


ADVERTIMENT. L'accés als continguts d'aquesta tesi queda condicionat a l'acceptació de les condicions d'ús establertes per la següent llicència Creative Commons:  <https://creativecommons.org/licenses/?lang=ca>

ADVERTENCIA. El acceso a los contenidos de esta tesis queda condicionado a la aceptación de las condiciones de uso establecidas por la siguiente licencia Creative Commons:  <https://creativecommons.org/licenses/?lang=es>

WARNING. The access to the contents of this doctoral thesis it is limited to the acceptance of the use conditions set by the following Creative Commons license:  <https://creativecommons.org/licenses/?lang=en>

Tesi Doctoral

**Aplicació de models supervisats per a l'estudi de provenença
arqueològica i certificació d'origen de ceràmiques**

Memòria presentada per Anna Anglisano Roca per tal d'obtenir el títol de doctora en Geologia
per la Universitat Autònoma de Barcelona

Març 2023

Director de Tesi:

Dr. Lluís Casas Duocastella

Universitat Autònoma de Barcelona

Facultat de Ciències

Departament de Geologia

Resum (català):

Aquesta tesi doctoral se centra en l'estudi de ceràmiques i argiles. Ara bé, aquesta caracterització es planteja amb una finalitat molt concreta: la distinció d'argiles de jaciments amb contextos geològics similars i la transferència d'aquesta distinció a les corresponents produccions ceràmiques. Inicialment es van abordar diverses estratègies de caracterització que van conduir cap a l'aproximació geoquímica i el tractament estadístic de les dades; aquestes metodologies són les que varen resultar més efectives per classificar mostres (indistintament d'argiles o ceràmica) segons el seu origen. En particular, la metodologia estadística que s'ha demostrat més eficaç ha estat l'ús de models de classificació d'aprenentatge automàtic supervisat (*supervised machine learning*). L'adjectiu supervisat fa referència al fet que s'utilitzen dades geoquímiques de referència de les quals es coneix amb certesa l'origen. En una primera fase s'empren part de les mostres de referència (amb origen conegut) per entrenar els models de classificació. En una segona etapa es comprova si l'aprenentatge ha estat efectiu comparant les assignacions d'origen que fan els models a la resta de dades de referència amb el seu origen real. Cal destacar que les dades emprades en aquesta segona etapa no s'ha fet servir durant la fase d'entrenament. La tesi es presenta com a compendi de publicacions amb dos articles que exploten aquesta aproximació estadística de classificació supervisada.

En el primer article publicat (l'any 2020 a Minerals MDPI) es va treballar amb una base de dades de més de 200 anàlisis químiques d'argiles i ceràmiques de 6 centres de producció ceràmica local catalana (Breda, Quart, La Bisbal d'Empordà, Sant Julià de Vilatorrada, Esparreguera i Verdú). En aquest treball es demostra la capacitat dels models entrenats de classificar correctament les mostres (amb taxes d'èxit de l'ordre del 80%). Es proposa l'aplicació d'aquesta capacitat de discriminació per certificar l'ús d'argiles locals i autenticar produccions ceràmiques.

En el segon article (publicat l'any 2022 a Sustainability MDPI) es va ampliar la base de dades amb un nou centre productor de referència (Barcelona), arribant a reunir gairebé 300 anàlisis químiques d'argiles i ceràmiques. A l'article s'aplica la capacitat classificadora dels models a la determinació de la provenença (lloc de fabricació) de mostres arqueològiques de ceràmica. Es va treballar amb 38 mostres que conformen 5 grups tipològics: tres obtinguts al Castell de Montsoriu, un a Torre de la Mora (tots dos jaciments propers a Breda) i un altre grup recuperat a la Creueta (jaciment proper a Quart). Els models entrenats van coincidir a identificar un dels grups recuperats al Castell de Montsoriu com a producció barcelonina amb probabilitats de l'ordre del 90%, confirmant així la hipòtesi arqueològica. En canvi els altres 4 grups estudiats no procedirien de cap dels centres productors que conformen la base de dades de referència.

L'aproximació estadística supervisada s'ha demostrat efectiva per certificar l'origen de produccions ceràmiques modernes i per a la investigació de provenença de ceràmica arqueològica. La complexitat dels algorismes que contenen els models de classificació i la seva implementació pot ser un factor que juga en contra de la seva popularització. Per això en aquesta tesi també s'ha desenvolupat un codi de lliure distribució i fàcil de fer servir per afavorir l'adopció d'aquesta aproximació per part d'usuaris sense coneixements avançats de programació.

Es preveu continuar la recerca exportant aquesta metodologia a la investigació de provenença de mostres arqueològiques no ceràmiques, en particular marbres.

Resumen (castellano)

Esta tesis doctoral se centra en el estudio de cerámicas y arcillas. Sin embargo, esta caracterización se plantea con un propósito muy concreto: la distinción entre arcillas de yacimientos con contextos geológicos similares y la transferencia de esta distinción a las correspondientes producciones cerámicas. Inicialmente se abordaron diversas estrategias de caracterización que condujeron hacia la aproximación geoquímica y el tratamiento estadístico de los datos; estas metodologías son las que resultaron más efectivas para clasificar muestras (arcillas o cerámica indistintamente) según su origen. En particular, la aproximación estadística que se ha demostrado más eficaz ha sido el uso de modelos de clasificación de aprendizaje automático supervisado (*supervised machine learning*). El adjetivo supervisado hace referencia al hecho que se utilizan datos geoquímicos de referencia de los cuales se conoce con certeza el origen. En una primera etapa se usan parte de las muestras de referencia (con origen conocido) para entrenar los modelos de clasificación. En una segunda etapa se comprueba si el aprendizaje ha sido efectivo mediante la comparación de las asignaciones de origen que hacen los modelos sobre el resto de datos de referencia. Cabe destacar que los datos usados en esta segunda etapa no se han utilizado durante la fase de entrenamiento. La tesis se presenta como un compendio de publicaciones con dos artículos que explotan la aproximación estadística supervisada.

En el primer artículo publicado (en el año 2020 en Minerals MDPI) se trabajó con una base de datos de más de 200 análisis químicos de arcillas y cerámicas de 6 centros de producción cerámica local catalana (Breda, Quart, La Bisbal d'Empordà, Sant Julià de Vilatorrada y Verdú). En el artículo se demuestra la capacidad de los modelos entrenados de clasificar correctamente las muestras (con tasas de éxito del orden del 80%). Se propone la aplicación de esta capacidad de discriminación para certificar el uso de arcillas locales y autenticar producciones cerámicas.

En el segundo artículo publicado (en el año 2022 en Sustainability MDPI) se amplió la base de datos con un nuevo centro productor de referencia (Barcelona), llegando a reunir casi 300 análisis químicos de arcillas y cerámicas. En el artículo se aplica la capacidad clasificadora de los modelos para determinar la proveniencia (lugar de fabricación) de muestras arqueológicas de cerámica. Se trabajó con 38 muestras que conforman 5 grupos tipológicos (tres obtenidos en el Castell de Montsoriu, uno en Torre de la Mora, ambos yacimientos cercanos a Breda) y otro grupo recuperado en la Creueta (yacimiento cercano a Quart). Los modelos entrenados coincidieron en identificar uno de los grupos recuperados en el Castell de Montsoriu como una producción barcelonesa con probabilidades del orden del 90%, confirmando así la hipótesis arqueológica. En cambio, los otros 4 grupos estudiados no tendrían como origen ninguno de los centros productores que conforman la base de datos de referencia.

La aproximación estadística supervisada se ha demostrado efectiva para certificar el origen de producciones cerámicas modernas y para la investigación de proveniencia de cerámica arqueológica. La complejidad de los algoritmos que contienen los modelos de clasificación y su implementación puede ser un factor que juega en contra de su popularización. Por eso, en esta tesis también se ha desarrollado un código de distribución libre y fácil de usar para favorecer la adopción de esta aproximación por parte de usuarios sin conocimientos avanzados de programación.

Se prevé continuar la investigación iniciada para exportar esta metodología a la investigación de proveniencia de muestras arqueológicas no cerámicas, en particular mármoles.

Abstract (English)

This doctoral thesis focuses on the study of ceramics and clays. However, this characterization is approached with a very specific purpose: the distinction between clays bearing similar geological contexts and the transfer of this distinction to the corresponding ceramic productions. Initially, a number of different characterization strategies were attempted and progressively we have been directed towards a statistical approach based on geochemical data, as an effective methodology to classify the samples (either clays or ceramics) according to their origin. In particular, the statistical approach that has proven to be more successful is the use of supervised machine learning classification models. The term supervised refers to the fact that this approach uses geochemical samples from known origin as reference class-labelled samples. In a first step, the models are trained with part of the class-labelled samples. In a second step, the efficiency of the training process is verified using the rest of the reference samples by comparing the class assignments made by the models to their actual origin. It is worth to mention that the data that is used in the second step has not participated in the training process. The thesis is presented as a compendium of publications with two articles that exploit the supervised statistical approach.

The first published article (in 2020 within Minerals MDPI) manages a database of over 200 chemical analyses of clays and pottery from 6 Catalan ceramic production centers (Breda, Quart, La Bisbal d'Empordà, Sant Julià de Vilatorrada, and Verdú). The article demonstrates the ability of the trained models to classify samples correctly (with success rates of around 80%). The application of this discrimination capacity is proposed as a way to certify the use of local clays and authenticate ceramic productions.

In the second published article (in 2022 within Sustainability MDPI) the database was expanded by adding a new reference production center (Barcelona) and reaching around 300 entries of pottery and clay chemical analyses. In this paper the ability to classify samples according to their provenience (place of manufacture) is applied to archaeological ceramic samples. A total of 38 samples from 5 typological groups were studied, three were retrieved in the Castell de Montsoriu, one in Torre de la Mora, both sites near Breda and another group was recovered in La Creueta (a site near Quart). The different trained models agreed to assign one of the groups from the Castell de Montsoriu to the Barcelona class with probabilities of around 90%, thus confirming the archaeological hypothesis. On the other hand, the other 4 groups studied would not belong to any of the classes (production centers) contained within the reference database.

The supervised statistical approach has proven to be effective both to certify the provenance of modern ceramic productions and to explore the provenience of archaeological pottery. The complexity of the algorithms used by the classification models and their implementation can be a discouraging factor to their popularization. To counteract this, the thesis contains a freeware and easy-to-use code to help users without advanced programming knowledge to apply the supervised approach.

We plan to pursue the developed research by exporting the supervised approach to provenance research on non-ceramic archaeological samples, in particular marbles.

| | |
|--|----|
| Índex..... | 9 |
| Introducció | 11 |
| De les argiles a la producció de ceràmica local..... | 11 |
| La ceràmica, un producte amb denominació d'origen (DO) | 12 |
| La ceràmica, una font d'informació del passat | 13 |
| Justificació | 15 |
| Objectius | 17 |
| Materials i Metodologies | 19 |
| Materials | 19 |
| Metodologies analítiques..... | 22 |
| Publicacions..... | 29 |
| Anglisano et al 2020 | 29 |
| Anglisano et al 2022 | 31 |
| Resultats i Discussió | 73 |
| Exploració de diverses aproximacions analítiques per caracteritzar i distingir els diferents grups de mostres corresponents a diverses localitats..... | 73 |
| Definició i distinció dels grups de referència mitjançant mètodes supervisats..... | 75 |
| Utilització de mètodes supervisats per a la classificació de mostres de ceràmica arqueològica com a eina per a establir-ne la provenença..... | 76 |
| Supervised Provenance Analysis. Una nova eina per a estudis de provenença de restes arqueològiques..... | 77 |
| Conclusions finals..... | 79 |
| Línies de Futur | 81 |
| Agraïments | 83 |
| Índex de taules i figures | 85 |
| Bibliografia | 87 |

Introducció

De les argiles a la producció de ceràmica local

La paraula argila s'utilitza en diferents àmbits per referir-se a conceptes diferents per bé que relacionats, això pot comportar un cert grau de confusió. En primer lloc en un àmbit científico-geològic la paraula i des d'un punt de vista granulomètric, "argila" fa referència a qualsevol partícula inferior a 2 µm, també fa referència al grup de fil·losilicats microcristal·lins (aluminosilicats hidratats) d'hàbit laminar o fibrós. Aquests presenten una gran capacitat de bescanvi i de retenció d'aigua i de líquids orgànics [1,2]. En segon lloc, la terminologia argila popularment ha fet referència a qualsevol material d'una certa plasticitat, generalment d'una tonalitat vermella i que s'utilitza en el món industrial per a la fabricació de materials ceràmics [3,4]. En aquest sentit, les argiles es poden entendre com un sediment dotat d'una certa plasticitat format per la degradació natural de roques preexistents i és en aquest sentit que cal interpretar el terme argila al llarg d'aquest document a no ser que s'especifiqui el contrari. Aquest sediment està format principalment per fil·losilicats resultants de l'alteració de minerals de composició silicatada, generalment tenen una disposició orientada paral·lelament que sovint contribueixen a definir l'estratificació del dipòsit sedimentari. Altres components de les argiles poden ser quars, feldspats i miques que constitueixen la fracció més grollera del sediment. Els minerals de les argiles formen part de la fracció més fina d'un dipòsit d'argiles tot i que no necessàriament han de ser d'una mida inferior a 2 µm segons la classificació granulomètrica estàndard [5].

Les argiles són un recurs natural que s'ha utilitzat al llarg de la història en molts àmbits diferents. Com ja s'ha esmentat, una de les aplicacions que avui dia encara segueix molt vigent és la utilització d'argiles per a la producció de ceràmica. Cada centre productor, fins i tot cada artesà utilitza la seva pròpia tècnica per a la producció de la pasta per fer ceràmica a partir de les argiles. Aquest procés ha variat al llarg de la història: actualment s'extreuen les argiles de manera industrial i amb grans maquinàries, antigament l'extracció era manual i es feia amb pics per desterrassar els sediments. Per processar aquests sediments argilosos plens d'impureses (arrels, sorres, fragments de roca, etc.), els terrissers utilitzaven mètodes de decantació o bé trituradores per incloure la majoria d'elements com a desgreixants, és a dir components que fan disminuir la plasticitat excessiva de l'argila pura [6–8]. Aquest procés per aconseguir la pasta ceràmica pot alterar la composició química del producte final diferenciant-lo de les argiles extretes.

Un altre procés que altera visiblement l'argila és la cuita que la transforma en ceràmica. En una fase prèvia al procés de cocció és important deixar que l'argila s'assequi i perdi bona part de la humitat, aquest procés implica una reducció del volum del producte i és important dominar-lo per tal d'evitar tensions que puguin causar esquerdes durant aquest procés d'assecatge. Un cop comença la cocció als forns, s'observen dues fases molt ben diferenciades: i) en primer lloc una fase d'evaporació de l'aigua al forn, que escalfa les peces de ceràmica fins a temperatures inferiors als 100 °C durant un temps considerable (el temps dependrà del gruix de les peces); ii) a continuació inicia el procés pròpiament dit de cuita on a partir d'argiles (maleables i dispersables en aigua) s'obtenen ceràmiques (rígides i no dispersables). En el procés de cocció es produeixen reaccions químiques entre els components de l'argila per bé que globalment la

composició química de les argiles i la ceràmica resultant són pràcticament idèntiques. Els únics components que es poden perdre durant la cocció són alguns volàtils (formats durant la cocció o preexistents) com ara l'aigua, el diòxid de carboni o el diòxid de sofre. En contrapartida es pot incorporar oxigen addicional a partir de reaccions d'oxidació. A banda dels petits canvis de composició, es produeixen molts canvis a nivell mineralògic per reordenament de la matèria. Els canvis mineralògics s'emmarquen en un procés de recristal·lització molt complex, amb certes analogies amb els que es produeixen en els processos metamòrfics naturals. El resultat final del procés dependrà principalment de la temperatura de cocció i de la composició de les argiles cuites. Alguns dels canvis estan molt ben acotats tèrmicament, a tall d'exemple la caolinita $[Al_2Si_2O_5(OH)_4]$ es degrada a 550 °C [9] transformant-se en la seva versió anhidre (metacaolinita, $[Al_2Si_2O_5]$) mentre que la plagiòclasi $[NaAlSi_3O_8 - CaAl_2Si_2O_8]$ es degrada a 950 °C fonent-se. En canvi el quars $[SiO_2]$ i l'hematites $[Fe_2O_3]$ es mantenen estables al llarg de tot el procés de cuita (tot i que en general la seva proporció en el producte cuit haurà variat) que rarament supera els 1100°C. La mullita $[Al_6Si_2O_{13}]$, feldspats potàssics $[KAlSi_3O_8]$ i cristobalita $[SiO_2]$ son minerals habituals de nova formació que apareixen a temperatures superiors als 950 °C [10].

La ceràmica, un producte amb denominació d'origen (DO)

Els sediments, també les argiles, que s'acumulen en un dipòsit sedimentari presenten una composició que està en relació a la de les diverses àrees font erosionades. A més, la composició final del sediment també pot estar relacionada amb el quimisme de l'aigua circulant i de determinats processos geoquímics que es produeixen en la zona de dipòsit. En particular, la distribució de terres rares en sediments com ara argiles pot aportar informació sobre processos geològics. En estudis sedimentològics, les concentracions relatives de terres rares són útils com a traçadors dels processos geoquímics i canvis ambientals que s'esdevenen a l'aigua i els sediments de sistemes aquàtics [11,12]. Atenent al fet que la ceràmica (argila cuita) hereta de forma força intacta la composició de l'argila original, es dedueix que la ceràmica també conté aquesta informació geoquímica. A més, cal destacar que aquesta informació està lligada al territori. L'argila d'una determinada contrada té característiques distintives que es transmeten a la ceràmica que pot produir-se coent-la.

Actualment a Catalunya encara segueix vigent la fabricació de productes ceràmics. En la història recent, fins a una setantena de localitats han tingut una producció de ceràmica local molt important [13]. Tot i això, avui dia només una desena d'aquestes localitats mantenen una producció local activa de productes tradicionals tot i que de fet no sempre s'utilitzen argiles locals per a fer aquests productes. Aquestes produccions locals s'han convertit en una activitat artesana amb capacitat per promoure un turisme sostenible. El turisme sostenible és aquell que es desenvolupa minimitzant l'impacte en el medi ambient. L'Organització Mundial del Turisme anomena tres factors clau per a que una activitat turística pugui qualificar-se de sostenible: i) que optimitzi els recursos mediambientals; ii) que consideri l'autenticitat de la cultura local; iii) que provoqui distribució de riquesa. El turisme basat en la potenciació de la producció local de ceràmica pot actuar com a vertebrador del territori, mantenint una activitat pròpia del lloc i contribuint a la creació de nous llocs de treball relacionats amb la divulgació i promoció del patrimoni cultural propi de cada regió. A Catalunya actualment ja hi ha uns quants museus dedicats a l'activitat terrissera tradicional. A més, també s'organitzen fires i esdeveniments

relacionats amb aquesta activitat, com per exemple la Fira del Tupí de Sant Julià de Vilatorrada (Osona), la Fira de l'Olla de Breda (la Selva), la Festa del Fang de Quart (Gironès), etc. En aquestes fires sovint els artesans locals mostren els seus productes i els processos de fabricació.

No obstant això, és molt comú que alguns negocis aprofitin el reclam de la ceràmica local per vendre-hi productes d'importació sense indicar-ne l'origen. Aquesta pràctica va en contra de la potenciació del producte local, en pot malmetre la imatge i desvirtuar la qualitat del producte que es troba a les localitats amb tradició ceramista. La qüestió de l'origen de la ceràmica pot interpretar-se de moltes maneres, actualment s'entén per producció ceràmica local aquella que ha estat confeccionada i cuita localment, independentment de l'origen de l'argila. La realitat avui en dia és que els terrissers acostumen a importar argiles per a les seves produccions i els mateixos terrissers poden ser forans que s'han establert a una determinada zona aprofitant-ne el renom. L'única manera de lligar una producció local genuïna al territori és reproduir els procediments artesanals fent servir la matèria primera local. L'empremta geoquímica que les argiles transmeten a la ceràmica obre la porta a poder acreditar de forma científica la producció de ceràmica amb argiles locals. D'aquesta manera es podria certificar la ceràmica amb un segell de denominació d'origen realment vinculat al territori.

La ceràmica, una font d'informació del passat

Generalment les excavacions arqueològiques comporten la troballa d'una gran quantitat de restes ceràmiques que en bona part acaben numerades i emmagatzemades en caixes. Aquestes restes arqueològiques poden aportar molta informació sobre les societats que les van fabricar, els usos que en feien, la tecnologia emprada per fabricar-les [14,15], l'edat del jaciment [16], evidències de contactes entre cultures i canals de comerç, etc. És per aquest motiu que una de les principals branques de la arqueometria (que engloba totes les aproximacions científiques a l'estudi de materials arqueològics) se centra en l'estudi físic i geoquímic de restes ceràmiques [17,18], generalment per esbrinar-ne el lloc de manufactura [19,20] (de vegades anomenat provenença [21] per distingir-lo de la procedència entesa com a lloc de la troballa).

Un element crucial per a l'estudi la provenença de restes ceràmiques és la definició de grups de referència. Aquests grups comparteixen una sèrie de característiques determinants que els diferencien d'altres produccions ceràmiques. Davant un conjunt de ceràmica arqueològica, els grups de referència generalment es defineixen a partir de criteris petrogràfics [22–24], geoquímics [25–28], o de vegades combinant tots dos tipus de criteris [29–32].

A partir dels estudis petrogràfics es pot obtenir informació com ara la composició mineralògica i la textura de les ceràmiques. De tota manera, cal tenir present que la composició mineralògica d'una ceràmica és el resultat del reequilibri d'una determinada composició a unes condicions de cocció determinades. És a dir, una mateixa argila pot coure's a temperatures diferents generant dos conjunts de minerals diferents. De forma similar, la textura de la ceràmica (i sobretot la granulometria del desgreixant) pot ser el resultat d'un filtratge granulomètric de l'argila de partida o bé de l'addició controlada d'un determinat desgreixant. Per tant la informació petrogràfica ens informa més aviat sobre la tècnica de fabricació emprada (selecció de l'argila, additius, condicions de cocció, etc) i no tant sobre la provenença. En canvi, els estudis geoquímics són molt més informatius sobre la composició del material de partida (argila) i no tant sobre les

tècniques de fabricació, per bé que cal tenir en compte que alguns procediments de fabricació (filtració, addició de desgreixants) poden alterar la composició de la pasta ceràmica.

Actualment, en els estudis geoquímics de provenença ceràmica, la pràctica habitual és analitzar el material disponible i utilitzar mètodes estadístics per identificar les diferències i afinitats entre mostres. Fruit dels procediments estadístics, les mostres es classifiquen en grups, de vegades a aquests grups se'ls assigna hipotèticament una determinada provenença a partir de raonaments de tipus arqueològic o arqueomètric. Els raonaments de tipus arqueomètric passen per analitzar materials de referència amb provenença coneguda. Per exemple, a la ceràmica trobada en forns o bé ceràmica descartada pel fet de ser defectuosa (sobrecuïta, deformada, pobrament cuïta, etc) se li pot assignar directament la provenença amb seguretat i les seves característiques permetrien definir un grup i indirectament assignar la provenença també als altres fragments ceràmics que quedarien classificats en el mateix grup.

Justificació

El treball de final de grau (TFG) de l'autora d'aquesta tesi doctoral ja es va fer en relació a la temàtica de la present tesi i d'alguna manera en va significar la llavor. El TFG es va titular "Estudi i caracterització de les argiles explotades pels ollers a la zona de Breda" i fou dirigit pel mateix tutor i director d'aquesta tesi. En el TFG ja es va abordar la caracterització d'argiles mitjançant diverses tècniques analítiques, entre elles la fluorescència de raigs X, que han estat utilitzades sistemàticament en la present tesi doctoral.

L'objectiu inicial que es va plantejar en aquesta tesi fou la caracterització i distinció de les diverses produccions locals de ceràmica catalana de finals de segle XX. És a dir, d'entre l'amplíssim ventall de poblacions de tradició ceramista existent en època medieval es pretenia centrar l'atenció en les produccions artesanals que han perdurat gairebé fins als nostres dies. El focus s'ha situat en aquestes poblacions de forma premeditada i en part per dotar de discurs científic les corresponents iniciatives de promoció i manteniment d'aquesta activitat tradicional en perill de desaparició. Com a part del treball desenvolupat s'ha dut a terme una recerca inèdita sobre els jaciments locals d'argila d'aquest petit reducte de poblacions amb producció actual. Aquesta recerca ha comportat la interacció amb diversos agents privats i públics, una col·laboració que en alguns casos s'ha articulada oficialment a través de petits projectes que han proporcionat finançament a la tesi i han generat resultats útils per a les institucions. Així, per exemple s'ha creat la "Ruta de les Terreres de Breda", finançada per l'ajuntament de Breda i que ha fet que actualment es puguin visitar les terreres on s'hi poden trobar panells informatius. També s'han aconseguit dues beques que han donat suport a aquesta iniciativa de recerca. D'una banda, la beca de recerca de La Selva finançada per l'Institut d'Estudis Selvatans i que es va atorgar per investigar sobre les terreres de Breda i les tècniques de producció tradicional de terrissa. Com a resultat d'aquesta investigació es va publicar un llibre [6] que fou redactat en col·laboració amb l'historiador local Jordi Goñi i porta per títol: "La tradició terrissera de Breda (s.XV-s.XX), Terreres, obradors, forns, elaboració i vocabulari". D'altra banda, la beca de recerca Josep Romeu, finançada per l'ajuntament de Sant Julià de Vilatorrada, es va atorgar per tal de recuperar la memòria de les terreres del municipi, una informació que es vol utilitzar per crear una escola de terrissa local utilitzant argiles de la zona.

A l'hora d'afrontar l'objectiu inicial de la tesi, de seguida es va constatar que l'aproximació mineralògica a l'estudi de la ceràmica (anàlisi petrogràfica a partir de làmines primes o estudi per difracció de raigs X) no permet caracteritzar grups de referència segons procedència. D'una banda, la composició mineralògica de les argiles de partida sovint és molt similar independentment del sediment explotat a diferents terreres i, un cop cuita, la ceràmica presenta una composició mineralògica que s'allunya de l'argila no cuita. A més, centrant-nos en la composició de la ceràmica, les diferències en les diverses tècniques de producció són més determinants en les característiques texturals i mineralògiques que acaba tenint la ceràmica que no pas la composició inicial de les argiles. La correlació entre procedència i tècnica productiva no és gaire bona perquè una determinada població pot haver produït ceràmica mitjançant tècniques diverses i una mateixa tècnica pot haver-se utilitzat en diverses poblacions productores.

Atès el fracàs de l'anàlisi mineralògica com a eina per caracteritzar i distingir les produccions ceràmiques de les diverses localitats estudiades, es va voler abordar la problemàtica amb l'aproximació geoquímica, tal i com s'aplica sovint en ceràmica arqueològica. En aquest sentit,

va ser fonamental la formació en mètodes estadístics que es va realitzar a través d'una estada a la Universitat de Pàdua sota el guiatge de la Dra. Lara Maritan. A partir de l'estada a Pàdua, es va intentar l'aplicació dels dos mètodes estadístics més comunament utilitzats en l'anàlisi geoquímica de ceràmiques arqueològiques: l'anàlisi de components principals (PCA – Principal Component Analysis) i l'anàlisi jeràrquica de grups (HCA – Hierarchical Clustering Analysis). En tots dos casos, l'aplicació d'aquests mètodes estadístics per a la distinció de les diverses procedències de ceràmiques (i argiles) va demostrar-se ineficaç. Cal tenir present que, tot i que els mètodes PCA i HCA sovint són presentats com a mètodes de classificació de grups, en realitat no ho són; són simplement mètodes que exploren les similituds i diferències entre mostres individuals [33] i per tant no és sorprenent que siguin mètodes ineficaços a l'hora de distingir mostres de localitats diferents però amb composicions molt similars de les respectives argiles i ceràmiques.

A partir d'aquí, i sota la recomanació i ajut de l'expert especialitzat en anàlisi de dades Marc Anglisano, es va decidir orientar la recerca cap als mètodes estadístics de tipus supervisat. Els mètodes explorats fins aleshores (PCA i HCA) són mètodes de tipus no supervisat, és a dir els algorismes que s'apliquen per explorar les similituds i diferències entre mostres no fan servir informació sobre l'origen de les mostres, ni tan sols quan es disposa d'aquesta informació a partir de mostres de referència. En el cas que es pretenia abordar disposàvem de mostres de referència amb procedència segura i aquesta situació permet l'ús de mètodes supervisats. En aquest tipus de mètodes, que entren dins del que es coneix com a 'machine learning', s'optimitzen els paràmetres d'un model de classificació de forma dirigida a la diferenciació efectiva de grups a partir d'un conjunt de dades (mostres del grup d'entrenament) que serveixen per entrenar el model de classificació. Un cop entrenat, es pot avaluar si el model ha après a diferenciar els grups de forma més o menys reeixida, assignant procedència a un conjunt diferent de dades (mostres del grup de control) completament aliè a l'entrenament de les dades.

L'èxit d'aquesta aproximació de tipus supervisada ha donat cos a la present tesi fins al punt de transcendir l'objectiu inicial de la recerca, de manera que s'ha acabat explorant l'aplicació d'aquesta aproximació al món de l'arqueologia en l'àmbit dels estudis de procedència de restes ceràmiques.

Objectius

A continuació es presenten enumerats els principals objectius que conformen la tesi:

- 1.- Estudiar i caracteritzar les argiles i ceràmiques de les poblacions catalanes on actualment, i en la història recent, s'hi ha desenvolupat una indústria terrissera.
- 2.- Explorar diverses aproximacions de tipus analític per caracteritzar i distingir conjunts de mostres (ceràmiques subactuals i d'argiles) de diverses procedències geogràfiques però properes físicament i amb contextos geològics similars.
3. Comprovar l'efectivitat de l'aproximació supervisada en la discriminació geoquímica de conjunts de mostres (ceràmiques i d'argiles) de diverses procedències geogràfiques però properes físicament i amb contextos geològics similars.
4. Aplicar l'aproximació supervisada per a la classificació de mostres de ceràmica arqueològica com a eina per a establir-ne la seva provenença.
- 5.- Desenvolupar un paquet de codi amb instruccions que permeti aplicar de forma senzilla l'optimització de models de tipus supervisat per a la diferenciació de grups de mostres i la producció de probabilitats de pertinença als grups per a mostres de provenença desconeguda.

Materials i Metodologies

Materials

Mostres dels grups de referència

Un cop establert el llistat de poblacions amb tradició de producció ceràmica que serien objecte d'estudi, es van organitzar diferents campanyes de camp per tal de mostrejar les àrees d'exploració d'argiles i recopilar material ceràmic de diferents centres de producció de cada localitat estudiada. Part d'aquest estudi es va realitzar amb la col·laboració d'estudiants de la UAB i voluntaris. En total es van mostrejar un conjunt de 28 terreres de 7 poblacions catalanes diferents (Fig. 1).



Figura 1: Diferents terreres que s'han caracteritzat: a) La Bisbal d'Empordà, b) Breda, c) Esparreguera, d) Quart, e) Sant Julià de Vilatorrada i f) Verdú.

La recerca sobre la localització dels afloraments d'argiles ha resultat especialment necessària en poblacions com Breda, Quart, Sant Julià de Vilatorrada i Esparreguera, ja que les terreres no constaven descrites en cap document. Gràcies a la memòria oral de les persones del territori i a partir de fotografies aèries antigues s'ha pogut accedir a aquest coneixement que corria el risc de perdre's irremissiblement. Més enllà de la tesi i projectes accessoris ja duts a terme, aquesta tasca de recuperació de la memòria dels usos del territori es vol posar en valor amb nous projectes, com per exemple una escola de terrissa a Sant Julià de Vilatorrada utilitzant argiles locals o col·laborant amb diversos artesans ceramistes que investiguen amb argiles locals en desús.

En el cas de Barcelona, la població té tradició de producció ceràmica però els afloraments han desaparegut completament sota la trama urbana de la ciutat. Les argiles de Barcelona es van aconseguir a partir dels materials argilosos que apareixen en sondejos del subsol. Es van utilitzar onze sondejos diferents, els corresponents nivells d'argiles s'han agafat com a representatius de les argiles de Barcelona que podrien haver afluït en zones properes.

En total es van recollir múltiples mostres de cada aflorament i una seixantena de mostres ceràmiques (Taula 1) que es van analitzar com a part del treball de tesi. Addicionalment, per a diversos afloraments es va comprovar que les argiles eren adients per a la fabricació de ceràmica coent-ne una part en un forn industrial de ceràmica fins a formar mostres ceràmiques (Fig. 2).

Aquesta ceràmica fabricada en el marc de la tesi es va fer servir també per avaluar els canvis químics i mineralògics que es produeixen en el procés de cocció de l'argila. A més, les anàlisis químiques d'algunes d'aquestes mostres ceràmiques també es van incorporar a la base de dades geoquímiques de les mostres de referència.



Figura 2: Algunes de les mostres d'argila un cop cuites

Pel que fa a l'obtenció de mostres històriques de ceràmica fetes amb argiles locals també s'ha necessitat la col·laboració ciutadana per accedir a aquest material, ja que no sempre hi ha museus o entitats locals relacionades amb la indústria local amb qui poder establir col·laboracions. Atès que les anàlisis a realitzar en les ceràmiques sovint impliquen la destrucció de la ceràmica, ha calgut disposar de fragments o peces malmeses sense prou valor patrimonial per a ser exposades en museus. Es van obtenir més de 150 fragments de ceràmica provinents de les diferents produccions locals estudiades. (Taula 1).

Taula 1: Resum del total de mostres que formen part de l'estudi

| Localitat | Argiles | Argiles cuites | Ceràmiques locals | Total |
|-------------------------|------------|----------------|-------------------|------------|
| Barcelona | 20 | - | 64 | 84 |
| La Bisbal d'Empordà | 17 | 4 | 15 | 36 |
| Breda | 11 | 4 | 18 | 33 |
| Esparreguera | 10 | 5 | 18 | 33 |
| Quart | 11 | 3 | 19 | 33 |
| Sant Julià de Vilatorça | 8 | 6 | 22 | 36 |
| Verdú | 23 | 6 | 9 | 37 |
| TOTAL | 100 | 28 | 165 | 292 |

Mostres arqueològiques

Per complir amb un dels objectius del doctorat s'ha reunit un conjunt de mostres arqueològiques que hipotèticament poguessin tenir com a provinença algun dels grups de produccions locals catalanes prèviament caracteritzats.

Es van seleccionar 5 conjunts arqueològics de característiques i cronologies diferents. Tres d'aquests provinents del Castell de Montsoriu (CM) i un de Torre de la Mora (TM). Ambdós jaciments situats molt a prop del centre de producció de Breda. També es van analitzar un conjunt del jaciment de la Creueta (C), molt proper al centre de producció de Quart (Fig. 3). De cada conjunt es van seleccionar diferents fragments de ceràmica local (entre 6 i 10 en funció de la disponibilitat, Taula 2) per mirar d'assignar-los provinença a partir de les seves característiques geoquímiques.



Figura 3: Mostres representatives dels conjunts de mostres arqueològiques de provinença desconeguda i que formen part dels conjunts analitzats. D'esquerra a dreta: ceràmica gris de cuïta reductora, ceràmica de cuïna amb vidrat de plom, ceràmica de cuïna amb vidrat acolorit en verd gòtic, tots tres conjunts procedents del castell de Montsoriu (CM). Ceràmica de cuïna de Torre de la Mora (TM) i finalment ceràmica feta a mà del jaciment de La Creueta (C).

Taula 2: Conjunt de mostres arqueològiques estudiades

| Jaciment arqueològic | Tipologia | Cronologia | Ref. | Nº mostres |
|----------------------|---|-------------------------------------|------|------------|
| Castell de Montsoriu | Ceràmica gris de cuïta reductora | 1475-1560 DC | CM1 | 7 |
| | Ceràmica de cuïna amb vidrat de plom | | CM2 | 10 |
| | Ceràmica de cuïna amb vidrat acolorit en verd gòtic | | CM3 | 6 |
| Torre de la Mora | Ceràmica de cuïna | Finals s. IX DC – principis s. X DC | TM | 7 |
| La Creueta | Ceràmica feta a mà | S. IV AC | C | 8 |

Metodologies analítiques

Per a la realització dels primers objectius de la tesi doctoral es van dur a terme diverses aproximacions analítiques per l'estudi de ceràmiques i argiles. Finalment però, la major part del treball es va centrar en la producció de dades geoquímiques mitjançant fluorescència de raigs X i l'explotació d'aquestes dades.

Microscòpia Òptica (estudi petrogràfic)

Es van preparar làmines primes estàndard (gruix de 30 micròmetres) de mostres ceràmiques i de mostres d'argila, aquestes darreres van requerir un tractament previ d'inclusió en resina per a poder-les dotar de cohesió (Fig. 4). Totes les preparacions van ser elaborades pel servei de làmines primes de la UAB.

L'estudi de les làmines primes es va realitzar fonamentalment amb un microscopi petrogràfic Nikon Eclipse E600 POL, disponible al departament de Geologia de la UAB. Dels diversos modes de visualització es va utilitzar sobretot la modalitat de llum transmesa tant amb llum no analitzada com amb nícols encreuats.



Figura 4 : Porta làmines amb algunes de les làmines primes realitzades en el present estudi i una làmina prima d'una ceràmica subjectada amb la mà.

El principi de funcionament del microscopi petrogràfic es basa, com en el microscopi òptic convencional, en la capacitat d'augment de les lents convergents amb la particularitat que l'ús de polaritzadors permet analitzar diverses propietats de les substàncies anisòtropses, com són la majoria de minerals.

En el cas de les argiles i també en bona part de les ceràmiques, la mida de molts components és massa petita per a poder-la caracteritzar eficientment amb un microscopi òptic. No obstant això, sí que es poden obtenir força paràmetres com ara el color general, l'anisotropia de la matriu, la forma i distribució de la porositat o la mineralogia de les inclusions no plàstiques de mida més gran.

Difracció de raigs X (estudi mineralògic)

La mineralogia de mostres ceràmiques i d'argila es va investigar mitjançant la tècnica de difracció de raigs X de pols. En el cas de les argiles, per a algunes mostres, es va investigar específicament la mineralogia de les argiles (en el sentit mineralògic) presents a les mostres geològiques.

Per a l'estudi per difracció, les mostres es van preparar com a pols dipositada a sobre d'un portamostres. En el cas de les mostres ceràmiques vidrades, primerament es va eliminar el vidrat amb una mini-esmoladora elèctrica. Per a totes les mostres, prèviament a la molta, es va realitzar un

asseccament en un forn a 60°C. La molta es va realitzar fins assolir un particulat d'una mida de 125 µm amb un molí de boles (Pulverisette).

Per a l'estudi específic de la mineralogia de les argiles es van preparar mostres d'agregats orientats. Partint de la pols preparada a 125µm, es va portar la mostra a suspensió en aigua per agitació. Al cap d'unes hores bona part de la pols se sedimenta però romanen en la suspensió les partícules d'argila més fines. Amb una pipeta es va extreure part del líquid de la suspensió i es va dipositar sobre el porta-mostres. Un cop evaporat el líquid, la mostra conté una selecció orientada de les argiles. La posició dels

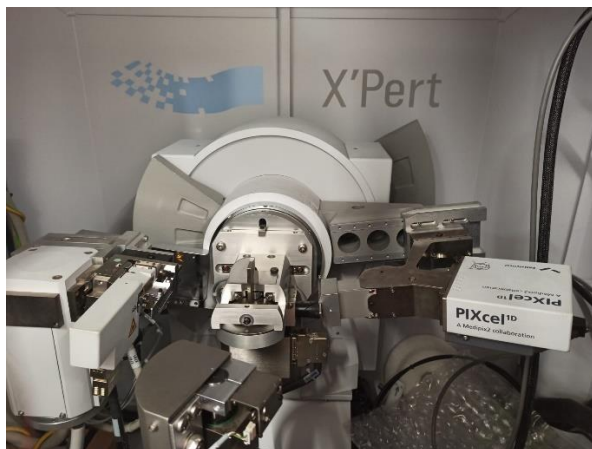


Figura 5: Equip de difracció de raigs X utilitzat.

pics de difracció que apareixen a baixos angles (entre 2° i 10° de 2θ) i la seva variació davant de tractaments de la mostra en una atmosfera d'etilenglicol, (C₂H₂OH)₂ permet discriminar el tipus d'argiles presents a la mostra.

El principi de funcionament de la difracció de raigs X de pols està relacionat amb la interferència constructiva de radiació X monocromàtica (λ) quan interacciona amb la matèria cristal·lina. És un fenomen amb certes analogies amb la reflexió de la llum però la radiació X no es reflecteix en un determinat pla físic de la matèria sinó en un pla reticular definit per l'estructura periòdica del mineral. La llei de Bragg ($\lambda=2d\sin\theta$) estableix en quines direccions es produeix la difracció, essent l'espaiat reticular i θ l'angle de difracció que forma el pla reticular que difracta el feix de radiació. A nivell experimental, l'equip conté bàsicament un generador de radiació, un porta-mostres i un detector de radiació. Aquests components es mouen de forma coordinada per anar recollint tots els feixos difractats en les diverses orientacions i finalment es genera un gràfic de la variació de la intensitat de radiació recollida en funció dels diversos valors angulars.

Les mesures de difracció es van realitzar amb un equip X'Pert Powder de Panalytical amb geometria θ - θ , ànode de coure per a la generació radiació X per impacte d'electrons accelerats, específicament es va utilitzar la radiació $K\alpha$ del coure ($\lambda=1.5406 \text{ \AA}$) i un detector PIXcel1D. L'equip es troba al Servei de Difracció de Raigs X de la Universitat Autònoma de Barcelona (Fig. 5). Per a la majoria de mostres es va realitzar un escanament entre els valors de 5° i 60° de 2θ. Per a la identificació de fases es va utilitzar el software X'Pert HighScore amb la base de dades de difractogrames PDF-2 (ICDD).

Fluorescència de raigs X (estudi químic)

Les primeres anàlisis geoquímiques de les mostres de ceràmiques i argiles es van realitzar a Pàdua (Itàlia) al laboratori d'espectrometria de la Università degli Studi di Padova en el marc d'una estada de doctorat auspiciada per la Dra. Lara Maritan. Posteriorment, quan el desenvolupament de la tesi va portar a la necessitat d'abordar sistemàticament l'anàlisi química de totes les ceràmiques i argiles mostrejades, es va arribar a un acord de col·laboració amb el

Dr. Ignasi Queralt (investigador del departament de geociències de IDAEA-CSIC), permetent la realització de les mesures al Parc Científic i Tecnològic de la Universitat de Girona.

Pel que fa a la preparació de les mostres, les mesures a Pàdua es van preparar com a perles per fusió amb tetraborat de liti ($\text{Li}_2\text{B}_4\text{O}_7$) portant les mostres a una temperatura de 1150°C . En canvi, per a la mesura sistemàtica de totes les mostres es van agafar 5 grams de pols de mostra barrejada amb 0.8 grams de resina Elvacite. La mescla es va comprimir sota un pistó a una pressió de 20 tones durant un minut generant-se pastilles coherents i homogènies de 4 cm de diàmetre.

La fluorescència de raigs X és una tècnica espectroscòpica en la que es produeix l'excitació electrònica de la mostra per irradiació amb raigs X. L'excitació electrònica consisteix en absorbir energia de la radiació per a fer saltar electrons dels àtoms de la mostra a orbitals més externs dels que ocupen i on romanen per un breu interval de temps abans guanyar estabilitat tornant a ocupar nivells més interns. Els salts electrònics a nivells més interns es produeixen emetent l'excés d'energia en forma de fotons (de radiació X). L'anàlisi de les característiques dels fotons emesos tant pel que fa a energia com a longitud d'ona permet identificar la naturalesa dels àtoms que han estat excitats i per tant aporta informació sobre la composició química de la mostra analitzada.



Figura 6: Espectròmetre utilitzat per a l'estudi.

L'equipament utilitzat per a l'anàlisi sistemàtica de totes les mostres objecte d'estudi ha estat un espectròscopi S2 Ranger de Bruker/AXS situat al laboratori de química analítica del Parc Científic i Tecnològic de la Universitat de Girona (Fig. 6). L'equip conté un tub de raigs X amb ànode de pal·ladi que pot treballar amb una potència màxima de 50W i un detector XFLASH LE Silicon Drift (SDD). Es tracta d'un equip que treballa per anàlisi de la dispersió d'energies (EDXRF) de manera que permet mesurar tots els elements químics simultàniament. Les mostres en forma de pastilla es van mesurar en

ambient de buit per a millorar la detecció d'elements lleugers. El temps de mesura de cada mostra ha estat de 400 segons. El software emprat per a transformar els espectres en valors de concentració del diversos elements químics ha estat el paquet SPECTRA.EDX (Bruker AXS). El software, de forma automatitzada, busca l'ajust a l'espectre experimental a partir de la combinació de subespectres teòrics i una funció polinòmica per al fons de radiació.

Tractament de les dades geoquímiques (mètodes estadístics)

S'han utilitzat dos tipus de mètodes per a tractar estadísticament els resultats de les anàlisis químiques obtingudes mitjançant fluorescència de raigs X. D'una banda, mètodes no supervisats pels quals s'investiguen les diferències i similituds entre les mostres individuals sense tenir en compte a quina classe (és a dir, a quina localitat) pertanyen. D'altra banda, els mètodes supervisats que busquen les diferències entre classes, és a dir entre els conjunts de mostres de cada classe, així doncs en els mètodes supervisats sí que s'utilitza com a informació la procedència de cada mostra.

Mètodes no supervisats

S'han utilitzats dos mètodes: l'anàlisi de components principals (PCA) i l'anàlisi de clústers jerarquitzada (HCA).

El mètode PCA consisteix en reduir el nombre de variables d'un determinat problema redefinint, a partir de les variables inicials, a un nou conjunt de variables compostes per combinació lineal de les inicials a través d'uns factors. Les noves variables (o components) es trien de manera que les primeres que es defineixen descriguin la major part possible de la variància mostral del conjunt de dades que es treballen. D'aquesta manera s'assumeix que considerant un conjunt petit de noves variables (els components principals) ja podem visualitzar les característiques del conjunt de dades analitzat sense haver de considerar un conjunt excessiu de variables. En el cas de les dades geoquímiques, les variables inicials són les concentracions dels diversos elements químics analitzats.

El mètode HCA consisteix a ordenar les diverses mostres d'acord amb la seva similitud. El càlcul de la similitud (o distància) entre les mostres es pot realitzar de diverses maneres però és habitual calcular-la com si es tractés d'una distància en un espai multidimensional on cada dimensió és cadascuna de les variables que defineixen les mostres. Partint de les posicions de les mostres en l'espai, es busquen les dues mostres més properes i s'agrupen formant una nova mostra composta a partir de la qual es torna a calcular quines són les dues més properes. Aquest procediment es va repetint fins que totes les mostres queden agrupades en una única mostra. Aquest procés d'agrupament de mostres genera un dendrograma que és un gràfic que va relacionant les diverses mostres i permet tenir un criteri a l'hora d'agrupar les mostres en el nombre de conjunts que es vulgui. En el cas de les dades geoquímiques les mostres són cadascuna de les mostres (de ceràmica o argila) i les variables són les concentracions dels diversos elements químics.

Mètodes supervisats

S'han utilitzat cinc mètodes de tipus supervisat: k-nearest neighbors analysis (kkNN), random forest (RF), artificial neural network (ANN), linear discriminant analysis (LDA), generalized linear models (GLM) i linear discriminant analysis (LDA). Per a tots aquests mètodes de classificació, els models s'entrenen fins assolir la configuració òptima dels seus paràmetres per tal de distingir de forma eficient les diverses classes a les quals pertanyen les mostres analitzades. La fase d'entrenament es fa amb un 80% de mostres i el 20% restant (grup de dades de control) es fa servir per a comprovar fins a quin punt el model parametritzat ha assolit la capacitat de classificar correctament les mostres. En el cas de les dades geoquímiques, les mostres són cadascuna de les mostres (de ceràmica o argila) i les classes són les localitats a les quals pertanyen aquestes mostres.

La capacitat de classificació dels models es calcula a partir dels encerts del model en les prediccions de classe que fa sobre el 20% de mostres que no han participat en la fase d'entrenament. Hi ha diversos paràmetres que redueixen a un únic valor l'anàlisi d'errors i encerts. Els que s'han utilitzat són:

-L'*accuracy* (de vegades traduïda com exactitud) que és la proporció d'encerts respecte al total de prediccions.

-La *balanced accuracy* que també parteix de la proporció d'encerts respecte a prediccions però es calcula de forma independent per a cada classe i a continuació se'n fa la mitjana aritmètica, és a dir se sumen els diversos valors obtinguts i es divideixen pel nombre de classes. És un paràmetre adequat quan les diverses classes contenen un nombre de mostres molt desigual.

-La puntuació F1 (o *F1 score*) que es calcula a partir de la mitjana harmònica entre la precisió (*precision*) i la sensibilitat (*recall*) calculada per a cada classe. Essent la precisió (P) el quocient entre el nombre d'encerts i el total de mostres assignades a la classe, i la sensibilitat (R) el quocient entre encerts i el total de mostres de la classe. Així per a una determinada classe la puntuació F1 seria:

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

Per obtenir un valor únic de F1 es fa la mitjana aritmètica entre tots els valors de F1 obtinguts per a cada classe.

A continuació es detallen molt breument les característiques dels cinc mètodes supervisats que s'ha utilitzat:

-K-nearest neighbors (kNN): Aquest model assigna la classe a una mostra examinant a quina classe pertanyen les mostres que té més a la vora (els seus veïns més propers). La manera de calcular la distància entre mostres i el nombre de veïns considerats pot canviar d'acord amb diversos criteris.

-Random forest (RF): Aquest model és de tipus *ensemble* (és a dir, que combina diversos submodels) on els submodels són arbres predictors o de decisió. Un arbre de decisió és bàsicament un conjunt de preguntes binàries (sí/no) aplicades a un subconjunt de dades i de les quals es dedueix una predicció de classe. La idea és que el model combina els resultats de diversos arbres de decisió que es van parametrizant (o entrenant) de forma independent.

Artificial neural network (ANN): Aquest model s'inspira en el funcionament del cervell. El model consisteix en múltiples unitats de processament d'informació (similars a neurones) organitzades en capes. Cada capa rep una sèrie de valors a partir de les variables inicials que caracteritzen les mostres, cada node transforma aquesta informació i la transfereix com a input a una nova capa de nodes fins que la darrera capa dona una predicció de pertinença a una determinada classe.

Linear discriminant analysis (LDA) : És molt similar a la PCA en el sentit que es redefeixen unes noves variables com a combinació lineal de les variables originals (és a dir les concentracions dels diversos elements químics) però les noves components principals en comptes de buscar la màxima dispersió de mostres individuals busquen la màxima separació entre els centroides corresponents a les diverses classes (és a dir a les diverses localitats estudiades).

Generalized linear models (Glmnet): És un grup de models que integra i amplia el concepte de regressió per relacionar dues variables a través d'una funció, anomenada funció d'enllaç, que s'obté per ajust a valors coneguts, de manera que la funció resultant permet predir valors desconeguts. En aquest cas les dues variables a relacionar són d'una banda el conjunt de concentracions dels diversos elements químics (l'input) i de l'altre la seva localitat o classe (l'output). La funció d'enllaç pot ser lineal o no lineal i sovint és més complexa que una simple funció.

Finalment, de forma similar als models que ja funcionen amb submodels, s'ha utilitzat un model de models que combina tots els models esmentats aquí en el que tècnicament s'anomena un *stack* (pila) de models. L'*stack* és un model que teòricament potencia les bondats de cada model i té l'avantatge que pot esquivar-ne les debilitats dels models sempre que no tots flauegin del mateix mal.

Publicacions

Anglisano, A., Casas, L., Anglisano, M., & Queralt, I. (2020). Application of Supervised Machine-Learning Methods for Attesting Provenance in Catalan Traditional Pottery Industry. *Minerals*, 10(1). <https://doi.org/10.3390/min10010008>

Anglisano, A., Casas, L., Queralt, I., & di Febo, R. (2022). Supervised Machine Learning Algorithms to Predict Provenance of Archaeological Pottery Fragments. *Sustainability*, 14(18). <https://doi.org/10.3390/su141811214>

Article

Application of Supervised Machine-Learning Methods for Attesting Provenance in Catalan Traditional Pottery Industry

Anna Anglisano ¹, Lluís Casas ^{1,*} , Marc Anglisano ² and Ignasi Queralt ³

¹ Department of Geology, Campus de la UAB, Autonomous University of Barcelona, 08193 Bellaterra, Catalonia, Spain; anna.ar.93@gmail.com

² Independent researcher, Professional Data Scientist, 17400 Breda, Catalonia, Spain; marcanglisano@gmail.com

³ Department of Geosciences, IDAEA-CSIC, Jordi Girona 18-26, 08034 Barcelona, Catalonia, Spain; ignasi.queralt@idaea.csic.es

* Correspondence: Lluís.Casas@uab.cat

Received: 18 November 2019; Accepted: 18 December 2019; Published: 20 December 2019



Abstract: The traditional pottery industry was an important activity in Catalonia (NE Spain) up to the 20th century. However, nowadays only few workshops persist in small villages where the activity is promoted as a touristic attraction. The preservation and promotion of traditional pottery in Catalonia is part of an ongoing strategy of tourism diversification that is revitalizing the sector. The production of authentic local pottery handicrafts aims at attracting cultivated and high-purchasing power tourists. The present paper inspects several approaches to set up a scientific protocol based on the chemical composition of both raw materials and pottery. These could be used to develop a seal of quality and provenance to regulate the sector. Six Catalan villages with a renowned tradition of local pottery production have been selected. The chemical composition of their clays and the corresponding fired products has been obtained by Energy dispersive X-ray fluorescence (EDXRF). Using the obtained geochemical dataset, a number of unsupervised and supervised machine learning methods have been applied to test their applicability to define geochemical fingerprints that could allow inter-site discrimination. The unsupervised approach fails to distinguish samples from different provenances. These methods are only roughly able to divide the different provenances in two large groups defined by their different SiO₂ and CaCO₃ concentrations. In contrast, almost all the tested supervised methods allow inter-site discrimination with accuracy levels above 80%, and accuracies above 85% were obtained using a meta-model combining all the predictive supervised methods. The obtained results can be taken as encouraging and demonstrative of the potential of the supervised approach as a way to define geochemical fingerprints to track or attest the provenance of samples.

Keywords: pottery industry; local products; clay; provenance; predictive modeling; supervised methods; geochemistry; XRF

1. Introduction

Spain is one of the leading countries in the EU within the ceramic sector, concentrating about half the total European production of wall and floor tiles (in particular within the Castelló region) and with a significant production of sanitary ware, bricks, roof tiles, and refractory materials. All these products are manufactured using rather uniformed extrusion processes. In contrast, handcrafted pottery produced using more or less traditional methods is becoming a receding activity. The traditional pottery industry was an important activity in many places within Catalonia (NE Spain), in the past and also in recent times. More than 70 different localities had a specific local production during the

20th century [1]. However, in most production centers the activity has stopped, or it has diminished drastically. The few locations where traditional pottery production persists are mainly small villages where the activity is promoted as a touristic attraction.

Tourism is a key economic sector for Spain (it accounts for 10–15% of the GDP). More than 70 million people are visiting the country every year, and Catalonia is among the main destinations. The preservation and promotion of traditional pottery in Catalonia is part of a strategy of tourism diversification undertaken by municipalities in cooperation with other higher governmental actors. The production of authentic local pottery handicrafts aims at attracting cultivated and high-purchasing power tourists that could counterbalance the expectable recession of mass tourism due to market saturation.

In this context, the production of traditional pottery should no longer be an uncontrollable activity. The authenticity of the production methods and, equally important, of the traditionally used raw materials (local clays) are crucial to provide an honest and genuine product offer. Similar approaches have demonstrated to be very effective to protect and certify products, particularly in the food and beverage sectors [2,3].

As far as we know, traditional pottery production in Catalonia is nowadays restricted to less than 100 workshops and potters in around 10 villages. Most of these villages produce pottery as a touristic attraction and events like ceramic festivals (including pottery workshops and demonstrations) are regularly organized (often annually or biannually). However, in these festivals, and also in local shops, the ceramic products that are presented for sale are not always produced locally (but imported from other places), sometimes they are not strictly produced following traditional methods (but using semi-industrial methods) or they are produced using imported clays instead of the local raw materials. Some local potters use stamps to characterize their productions, but besides certifying the ceramist, it would be very useful to develop a methodology that could certify the use of the local raw materials of every local production. In this way, quality production with certified local souvenirs (such as hand-crafted tableware and other clay-modelling art) would be possible.

Multi-element chemical analyses have a long tradition of use for the characterization of clays and pottery (e.g., [4–6]). However, the choice of the analytical method, the right set of elements to be analyzed or the multivariate method to process the obtained data are under discussion [7]. In this paper, we present a number of statistical approaches to characterize geochemically the local raw materials along with the corresponding ceramics from six Catalan villages. The goal is challenging as the six sites concentrate in a small area and some of them share essentially the same geological context.

The six selected villages are rural towns with local pottery being used as a touristic attraction. The traditional activity is particularly supported by their town councils but also by provincial and regional governments. The villages are Esparreguera, La Bisbal, Quart, Breda, Verdú and Sant Julià de Vilatorrada (Figure 1), and in all of them but one (Sant Julià) there are still active workshops and all of them but one (Esparreguera) belong to the Spanish 'Asociación de Ciudades de la Cerámica (AeCC)'. Extensive information on the past and present productions can be found in [1].



Figure 1. The selected traditional pottery centers (red dots) located around Barcelona (black dot).

- Esparreguera is the closest village to Barcelona (~35 km north-west of it) and the biggest one of the six selected, with a population of ~22,000 people. There is still one potter working in this area and from 2018 there is a museum in the village devoted to pottery. Many types of pieces were produced with the exception of cookware.
- La Bisbal d'Empordà (~100 km north-east of Barcelona, near the Costa Brava) has a population of about 11,000 people. The village is advertised as one of the leading pottery centers in Catalonia with many shopping areas focusing on pottery sales and a museum that promotes the local cultural heritage connected with pottery and ceramics [8]. There are 35 active pottery companies and small workshops. Besides the pottery market, clays from La Bisbal are exploited by four companies that sell it mainly within Spain. The town has been declared “craft area of interest” by the Catalan government and the EU registered trade mark “Ceràmica de la Bisbal” acts as a protected designation of origin [9]. However, the trademark only attests local producers but not the use of local clays.
- Quart is another of the selected villages (~80 km north-east of Barcelona), it is a smaller rural town (~3600 people) only ~15 km east of La Bisbal with the granitic Gavarres mountain ranges between both. A refurbished old brickyard hosts, from 2011, the local pottery museum. The ceramics tradition in Quart can be traced back to the 14th century [10]. The production is known by both red and black colored products. At present there are five active pottery producers in the village.
- Breda (~50 km north-east of Barcelona) is another small rural town (~3700 people) with active pottery workshops. However, the local clays have been currently replaced by imported clays. Comparatively, the pottery industry in Breda has been historically one of the most important in Catalonia. Pottery constituted practically half the total registered industries during late 18th century and early 19th century. However, nowadays there are only seven active workshops and small pottery industries in the village. From 2003 an old workshop hosts a cultural center devoted to the local history of pottery production.
- Verdú is even a smaller village (only about 1000 inhabitants) located ~90 km west of Barcelona, with six active pottery workshops, its pottery industry goes back to Roman times and it is well documented since the 13th century [1]. The typical pottery from Verdú is black colored (fired in reducing conditions) and the main produced item is the earthenware pitcher.
- Sant Julià de Vilatorça (~60 km north of Barcelona and 3000 inhabitants) in another of the selected villages. It has a long tradition on glazed pottery, in the early 20th century there were 32 active workshops but unfortunately nowadays there is no one and there is not available precise information on the extraction points of local clays. There are ongoing local initiatives to revive the pottery tradition of Sant Julià including ceramics festivals and a project to create a pottery school.

The aim of the work is finding a common geochemical fingerprint that could group together both clays and pottery from a given town and that could allow discrimination from similarly defined clusters for the other towns. The mineralogical changes that occur during the clay-to-pottery transformation prevent the use of other commonly used techniques for characterization of clays such as XRD or FTIR and Raman spectroscopies. Such techniques are sometimes indeed useful to characterize clays (e.g., [11,12]) or firing conditions in ceramic materials (e.g., [13]) but they cannot be used to group together clays and the corresponding ceramics. Therefore, chemical analysis is required, although some problems can also arise using elemental analysis. For instance, the concentration of some elements can decrease during the pottery production due to the formation of volatile compounds by thermal decomposition [14]. Other elements can have a high variance within a considered cluster (including both clays and ceramics) and some strong correlations between elements can provoke problems for statistics-based methods of analyses [15].

In the field of archaeological sciences, and specifically in provenance studies of pottery, data is often obtained by X-ray fluorescence analysis (XRF) or other elemental analysis techniques. The processing of the large chemical datasets that are obtained is usually addressed by application of statistical methods that predict groups or relations between samples [16,17]. The classical approach is to use unsupervised methods such as principal component analysis (PCA), hierarchical cluster analysis (HCA), k-means, or factor analysis (FA). Extensive literature can be found on the application of such techniques to address pottery provenance issues (e.g., [18–20]). In this paper, we show that these methods are not very useful to meet our objective and we explore the possibilities of a number of supervised methods (k-nearest neighbors analysis (kNN), random forest (RF), generalized linear models (Glmnet) and linear discriminant analysis (LDA)) that work much better for our purpose.

2. Materials and Methods

2.1. Sampled Areas and Materials

For this study, 80 samples of clays outcropping in the vicinity of the selected villages were taken into account (Table 1). All of them are geological Cenozoic formations, the exact extraction points of clay were determined through existing bibliography [21,22] and oral sources. When original extraction points were no longer available, sampling was performed on nearby equivalent geological outcrops (Figure 2). The sampled materials for every village were:

- **Esparreguera:** red clays (six samples) from decimetric layers inserted in proximal alluvial fan deposits belonging to the northern margin of the Vallès-Penedès basin (Miocene epoch, Vallesian age (11.2–8.9 Ma)). Additionally, four samples of the fault flour outcropping at the edge of the basin margin, in the main normal Vallès-Penedès fault, were also sampled as this material was traditionally added to the red clays [1].
- **La Bisbal:** ochre and red clays deposited in an alluvial plain environment (also Miocene epoch, Vallesian age) connected to the Gavarres massif. (17 samples)
- **Quart:** red clay levels within arkosic sandstones (possibly Pliocene) that belong to a system of alluvial fans linked to the Gavarres massif. (11 samples)
- **Breda:** red, black and white clays from relatively thin levels within sandstones (early Pliocene, ~5 Ma) deposited in an alluvial fan environment with predominance of igneous and metamorphic clasts from Variscan granitoids and the uplifted palaeozoic basement. (11 samples).
- **Verdú:** grey and red clayey and calcareous marls from an environment of distal alluvial fans and lacustrine limestones (Oligocene epoch, Chattian/Rupelian age (~28 Ma)) (23 samples).
- **Sant Julià:** red clays from relatively thin levels within sandstones and conglomerates (Eocene epoch, possibly Lutetian, ~45 Ma) deposited in an alluvial fan proximal environment. Initially 18 samples were obtained but 10 resulted to bear a high CaCO₃ content (>25 wt.%) inconsistent with the Ca content of the corresponding ceramics and therefore they were discarded. These 10 samples

possibly belong to a nearby formation of calcareous red marls. Therefore, eight samples were retained from Sant Julià.

Table 1. Summary of the retrieved and analyzed samples.

| Village | Clays | Pottery Shards | Ceramic Briquettes | Total Samples |
|--------------|-------|----------------|--------------------|---------------|
| Esparreguera | 10 | 18 | 5 * | 33 |
| Bisbal | 17 | 15 | 4 | 36 |
| Quart | 11 | 19 | 3 | 33 |
| Breda | 11 | 18 | 4 | 33 |
| Verdú | 23 | 9 | 5 | 37 |
| Sant Julià | 8 | 22 | 6 | 36 |
| | | | | 208 |

* 2 briquettes were produced using samples of fault flour.

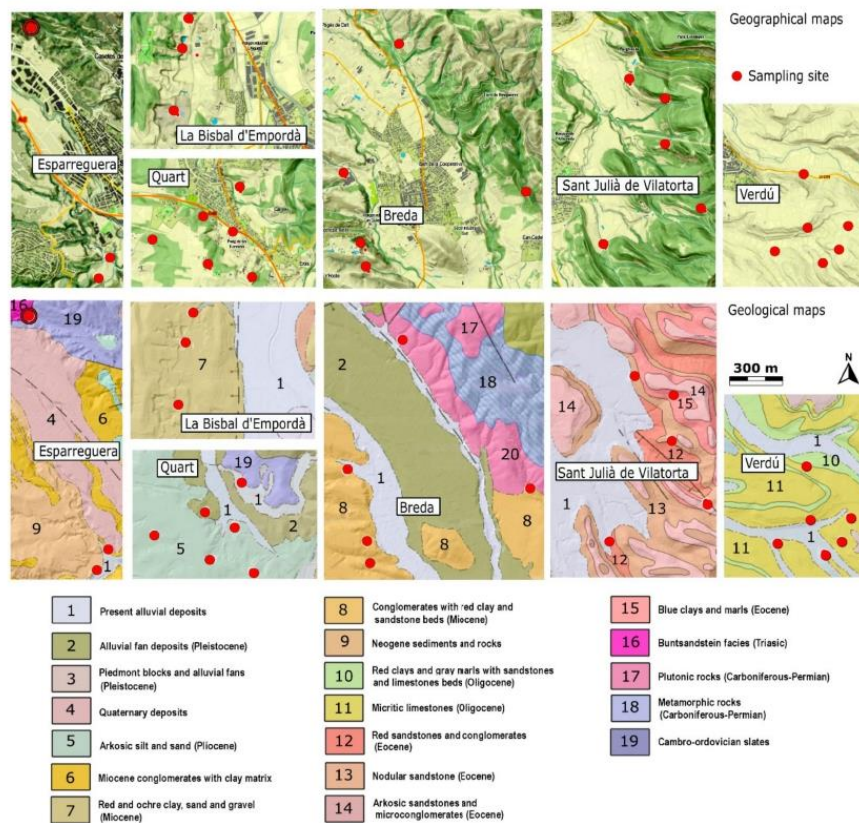


Figure 2. Geographical maps (top) and geological maps (bottom) of the sampling sites (red dots) around the studied pottery centers; scale is the same for all the maps. The concentric ring in the northern sampled site in Esparreguera indicates the sampled fault flour site.

Sampling of clays was performed by removing organic and soils layers and focusing on clays thin layers (avoiding coarser sized layers). However, most clay samples contain actually silt and very fine sand fractions. Nevertheless, the clay samples were neither sieved nor levigated because silt and sand fractions were also present on many pottery samples.

Besides the clay samples from the six selected areas, 101 samples of pottery likely produced by firing clays from the sampled formations were taken into account (Table 1). The pottery shards were obtained through local museums, active pottery workshops, local historians and retired potters. The selection criterion was the high probability that the pottery was really produced using local clays. For the Breda pottery samples, the presence of ceramic stamps even allowed the identification of the local potter and the production period (19–20th century) [23]. Pottery samples from Sant Julià are dated in the 17th century, and those from Esparreguera, Quart, and Bisbal are also from the 19–20th century like Breda. Finally, 27 extra newly produced pottery samples were added to samples from each locality (3–6 per studied site). These were prepared in form of small ceramic briquettes from local clays, firing them in a gasoil kiln during 10 h (including the heating ramp) reaching a maximum temperature of ~1000 °C for 2 h.

2.2. Mineralogical Analyses

A basic mineralogical and petrographic characterization was performed on several representative specimens of both clay and pottery samples from all the studied villages. Specimens were prepared as thin sections to be viewed using a petrographic microscope; other specimens were prepared from clay-suspension drops as well as ceramic powder that were deposited on aluminum discs to obtain X-ray diffraction (XRD) patterns. A Panalytical X'Pert powder diffractometer with θ - θ geometry, Cu anode X-ray tube and a PIXcel^{1D} detector has been used. Additionally, determination of the calcium carbonate content in clay specimens was determined by a volumetric method by adding hydrochloric acid and measuring the CO₂ released.

Petrographic results reveal that pottery productions from a given studied town can appear very different in terms of amount and size of the non-plastic inclusions (see Figure 3b). Perhaps the only clear trend is the predominance of fine-grained pastes for productions from Verdú. For pastes bearing large inclusions, quartz is dominant basically on all studied sites and therefore textural distinction between sites is not possible (Figure 3a). The XRD patterns of clays from the different sites share common features, exhibiting a mixture of characteristic peaks of mica-type and kaolinite-type minerals (Figure 3c). Only occasionally did calcite and/or quartz appear in the analyzed clay-fractions. In contrast XRD patterns of pottery shards can exhibit quite different signals, even for shards from a given town (Figure 3c), possibly this is the result of different firing conditions. Calcimetries revealed the presence of calcite in most clay specimens from Verdú, La Bisbal and Esparreguera.

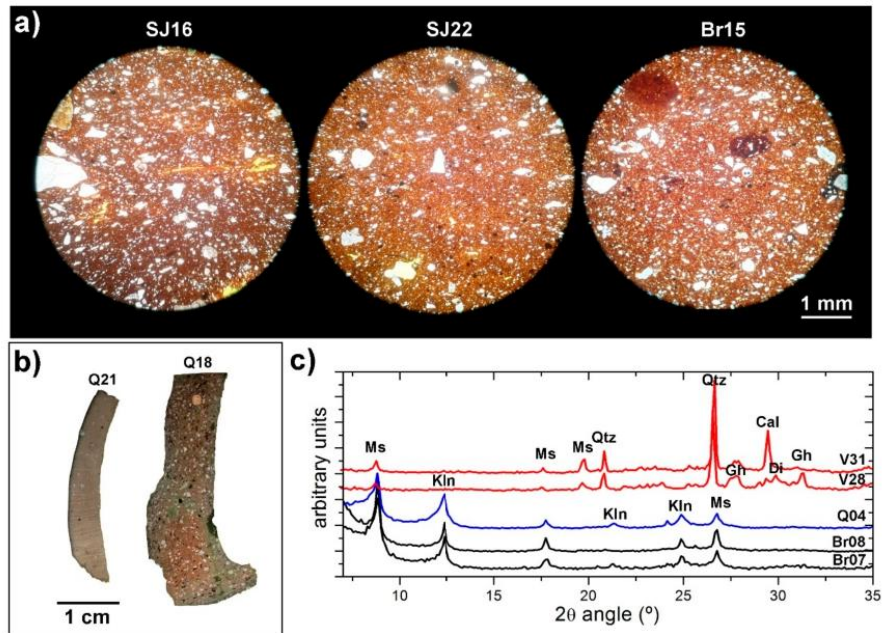


Figure 3. (a) Petrographic thin sections of three pottery samples, textural features of SJ22 (Sant Julià) and Br15 (Breda) appear quite similar. (b) Petrographic cross sections for two pottery samples from Quart exhibiting very different gran size. (c) XRD patterns from several clay (Q04, Br08 and Br07) and pottery (V28 and V31) samples with indication of the main phases: Ms: muscovite-type, Klin: kaolinite-type, Qtz: quartz, Cal: calcite, Gh: ghelenite, Di: diopside.

2.3. EDXRF Analysis

Besides basic mineralogical and petrographic characterization, the main analytical tool used has been Energy Dispersive X-ray fluorescence to obtain the required elemental data. That is the data that has been used to apply different statistical methods to identify a distinct geochemical fingerprint for each studied rural town (including their clays and pottery). Common sample preparation for geological materials and earthen objects has been undertaken consisting in the conversion of solid powdered samples into flat surface pellet specimens. Clays were first oven dried at 60 °C until constant weight. Glaze from glazed pottery shards was completely removed by scraping it using a drill. The shards were also oven dried as the clays and afterwards they were ground using a laboratory mill (Pulverisette™, Fritsch GmbH, Idar-Oberstein, Germany) to pass a 125 µm mesh.

Clay and pottery powders for analysis were prepared as pressed powder pellets following the methodology reported in [24] and using a methyl-methacrylate resin as a binding agent. 5 g of sample were mixed and homogenized with 0.8 g of binder (Elvacite™ commercial resin). The resulting powder was poured into a pressing die (40 mm in diameter) and pressed at a pressure of 20 T. The resulting pellet of tablet for analysis.

A commercially available benchtop EDXRF spectrometer (S2 Ranger, Bruker/AXS, GmbH, Karlsruhe, Germany) was used in the present study. This instrument is equipped with a Pd target X-ray tube (50 W power max.) and a XFLASH™ LE Silicon Drift Detector (SDD), ultra-thin beryllium window (0.3 µm thickness) with a resolution lower than 129 eV at Mn-K α line for a count rate of 100,000 counts-per-second (cps). In this LE configuration of SDD detectors, the intensities for Na K-alpha and Mg K-alpha are, respectively, close to around four times higher than the intensity recorded by conventional SDD detectors. The instrument is also equipped with nine primary filters that can be

used in front of the tube before X-ray beam reaches the sample surface in order to improve measuring conditions for the elements of interest and it can operate under vacuum conditions.

The software used to control the equipment, to build calibrations, and to perform spectral data treatment was SPECTRA.EDX package (Bruker AXS, GmbH, Karlsruhe, Germany). This software can perform the full line profile fitting, deconvolutions when lines overlap, intensity corrections for inter-elemental effects and qualitative, semiquantitative or full-quantitative routines.

All the samples were analyzed to obtain a spectrum for the identification of all the elements present in samples. Quantification was made by the assisted fundamental parameters approach included in the above-mentioned software using certified BAS BCS-315 and ECRM 776-1 firebricks as reference materials. Analysis was made in a vacuum atmosphere allowing better detection of low Z elements and using different conditions of voltage to properly excite low, medium and high atomic number elements existing at the samples. Current was automatically adjusted to obtain a fixed counting rate of 100,000 cps. Total measuring time was set at 400 s as a trade-off between an acceptable repeatability of measurements and total analysis time.

The net intensity of each analytical line was calculated by subtracting the theoretical background adjusted by a polynomial function to the obtained experimental spectra.

2.4. Data Processing

Data processing was performed using available scripts from RStudio (the integrated development environment for R software). In rough outlines, all the samples (both clays and baked clays) display similar assemblies of elements. From the obtained data, only those elements present (above their detection limit) have been taken into account. Besides that, Ca, S and Pb values have been disregarded, the first (Ca) because it shows a very strong inverse correlation with Si values, S because it shows a very high dispersion with many samples exhibiting values below the corresponding detection limit and Pb due to evidence of contamination from glazes in glazed pottery, even after glaze removal. The list of elements that have been taken into account to test several statistical machine learning approaches to classify geochemically the samples from each village have been: Al, Si, Fe, Na, Mg, Cl, K, Ti, Cr, Mn, Ni, Cu, Zn, Rb, Sr, Y, Zr, and Nb. Both unsupervised and supervised modeling of data has been applied to the datasets containing values for this list of elements.

Unsupervised learning is a way of organizing data that helps to find previously unknown patterns in datasets without pre-existing class labels. In contrast, supervised machine-learning relies on prior knowledge of the class labels. Some unsupervised methods, and particularly principal component analysis (PCA), are widely used in archaeometry to facilitate the identification of compositional groups and determining the chemical basis of group separation and extensive literature can be found on the subject. In contrast the use of supervised methods is quite scarce. However, some incursions within the supervised methods domain had also been done previously, from the pioneering works of [17,25] to much more recent papers [26–28] that focus particularly on the artificial neural network (ANN) approach.

In the processing of the obtained data some unsupervised models have been tested and, as it will be shown, all of them fail to produce data groups with good correlation with the real classes. The widely used hierarchical cluster analysis (HCA), k-means and PCA will be used to illustrate the low performance of such models. The HCA algorithm produces a tree diagram (dendrogram) according to a given metric and linkage criterion (e.g., [29,30]), the k-means algorithm identifies k clusters from a given dataset, every cluster is identified with a centroid and the corresponding data, the algorithm basically tries to keep inter-cluster data as similar as possible, while the centroids are as different as possible [31]. PCA logic is based on the concepts of linear correlation and variance. PCA is a dimensionality reduction technique, starting with the features (i.e., the chemical values) describing a set of objects (i.e., our samples), the target defines other variables that are linearly uncorrelated with each other. The output is a new set of variables defined as linear combinations of the initial features. The new variables are ranked on the basis of their relevance. The number of the new variables is less

than or equal to the initial number of features and it is possible to select the most relevant features. Then, it is possible to define a smaller set of features, reducing the problem dimension, see page 169 in [32].

Taking into account that the provenance of the analyzed samples is known (i.e., the class labels are known for each object) supervised models can also be used to process the data. Starting from the whole experimental dataset, it is possible to constitute a training dataset [26]. The labelled geochemical data is then used to build models that proxy for class characteristics. After optimization of the model parameters, the model is finally tested with new objects. The best performing predictive model can be selected looking at their ability to predict class memberships for these new objects. In this study, an 80% portion of the total dataset was used to train the models and the remaining data was used for model testing. The dataset was divided randomly using a specific seed in such a way that all the models tested in this study use the same train and test sets. The performance of the different tested models was evaluated using the widely employed confusion or error matrix [33], where each row represents the distribution of samples from an actual class among the predicted classes organized in columns (or vice versa). As the predicted and actual classes are presented in the same order, the successful predictions (usually called hits or true positives) concentrate along the main diagonal of the matrix, whilst unsuccessful predictions lie outside the diagonal. Additionally, an overall value of accuracy was computed as the ratio between hits and the total number of objects (i.e., samples). An overall accuracy value of 1 would indicate 100% of success and therefore no errors.

Predictive modeling of the training set was conducted using five machine learning algorithms and finally a combination of all them. These algorithms were called from packages within the freely available Caret R library:

1. Weighted k-nearest neighbors (kkNN) [34], its basic idea is that a new object will be classified according to the class that have their k-nearest neighbors. The R package used was `class`.
2. Random forest (RF) [35,36], this algorithm is based on the concept of decision tree (a series of yes/no questions asked to the data that in the end lead to a predicted class). The RF model deals with many decision trees (i.e., a forest) using random sampling to build the trees and random subsets of features when splitting nodes of the trees. The R package used was `randomForest`.
3. Artificial neural network (ANN) [37], a mathematical mimic of human learning where individual processing elements are organized in layers. The input layer receives the weighted values of the features of an object to produce new values through so called activation functions; these values will be also weighted and transferred to new layers until reaching the output which is made of as many elements as classes. The obtained values are used to assign a class to the object. The R package used was `nnet`.
4. Linear discriminant analysis (LDA) [38,39], similarly to the PCA logic, delineates a new set of variables defined as linear combinations of the initial features reducing the dimensionality of the problem, but instead of looking for the maximum variance, LDA maximizes the separability among classes (the distance between their means) and simultaneously minimizes the internal scatter within each class. The R package used was `lda`.
5. Generalized linear models (Glmnet) [40,41]. These are generalization models of a linear relationship between the output variable (class) and a set of input variables (features) where the distribution of the output variable can be non-normal and non-continuous and the function linking input and output variables can be more complex than a simple identity function. Specifically, the Glmnet algorithm incorporates regularization (i.e., reduction of variance) by the lasso and elastic-net methods to avoid overfitting (i.e., noise fitting). The R package used was `Glmnet`.
6. Stack of models. With the aim of improving the accuracy of the predictions, information from multiple models (i.e., a stack of models, see [42]) was used to generate a new model using a random forest approach to the predictions from different models. The R package used was `randomForest`.

All these models, including the stack of models, were optimized during the training step. This was done following a homogeneous approach for all the models and included a first phase of `traincontrol` using the `repeatedcv` method:

```
fit_control<-
trainControl(method="repeatedcv",number=10,repeats=2,savePredictions="final",classProbs=TRUE)
preproc=c("center","scale")
```

Then, in a second phase, the best optimizable parameters to improve the attained accuracy were automatically spotted with code lines as in the following example:

```
BASE.lrm<-
train(target~.,Dataset[train,],method="glmnet",metric="Accuracy",preProc=preproc,trControl=fit_control)
```

3. Results

3.1. EDXRF Analysis

The full set of chemical analyses can be downloaded from the Supplementary Material Section (Table S1). Table 2 shows a summary of the normalized chemical analyses with the major and minor element means (values are in weight percent of corresponding oxides) as well as those of the trace elements (values are in native element parts-per-million) for every studied locality. The associated standard deviation values give an indication of the internal variance within each group of data. Apart from Si and Ca values, for the remaining major elements, mean concentration values appear very similar in all the studied villages. Regarding the mean values for the minor elements, it is worth noting that almost all of them (Rb and Zr are the exceptions) show very high coefficients of variation (often above 40%) indicating a very high dispersion of values within each locality. Some major elements such as Na, S, Cl and occasionally also Ca also show very high variation coefficients. This highlights the difficulty to extract relevant data from the analyses, and therefore the use of statistical methods is indispensable.

Table 2. Arithmetic means (m) and standard deviations (σ) of the chemical composition for the samples (clay, pottery and briquettes samples merged) from every investigated location.

| Compound | Village | | | | | | | | | | | |
|------------------------------------|--------------|----------|-----------|----------|-------|----------|-------|----------|-------|----------|------------|----------|
| | Esparreguera | | La Bisbal | | Quart | | Breda | | Verdú | | Sant Julià | |
| | m | σ | m | σ | m | σ | m | σ | m | σ | m | σ |
| SiO ₂ (%) | 56.4 | 5.1 | 57.4 | 8.2 | 67.6 | 4.3 | 69.9 | 2.6 | 48.2 | 6.2 | 66.5 | 3.7 |
| Al ₂ O ₃ (%) | 11.5 | 1.9 | 11.5 | 1.9 | 13.1 | 1.2 | 12.5 | 0.9 | 9.9 | 1.4 | 12.3 | 1.1 |
| Fe ₂ O ₃ (%) | 8.6 | 1.7 | 8.1 | 1.5 | 8.6 | 1.9 | 8.3 | 1.5 | 9.2 | 1.7 | 8.5 | 1.7 |
| MgO (%) | 4.1 | 1.9 | 1.8 | 0.4 | 1.2 | 0.4 | 1.4 | 0.3 | 3.5 | 0.7 | 1.7 | 0.3 |
| CaO (%) | 11.9 | 7.2 | 14.6 | 10.9 | 2.0 | 3.2 | 0.8 | 0.2 | 22.6 | 9.8 | 3.1 | 4.2 |
| Na ₂ O (%) | 0.5 | 0.2 | 0.4 | 0.2 | 0.5 | 0.2 | 0.5 | 0.2 | 0.5 | 0.1 | 0.6 | 0.2 |
| K ₂ O (%) | 4.2 | 0.8 | 4.2 | 0.6 | 4.7 | 0.2 | 4.5 | 0.5 | 4.1 | 0.9 | 5.1 | 0.5 |
| TiO ₂ (%) | 1.3 | 0.2 | 1.4 | 0.2 | 1.5 | 0.3 | 1.5 | 0.2 | 1.2 | 0.1 | 1.2 | 0.2 |
| SO ₃ (%) | 0.5 | 0.7 | 0.1 | 0.2 | 0.1 | 0.2 | 0.0 | 0.0 | 0.1 | 0.1 | 0.1 | 0.1 |
| Cl (%) | 0.3 | 0.6 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.2 | 0.4 |
| Mn (ppm) | 386 | 55 | 372 | 171 | 335 | 71 | 259 | 58 | 352 | 39 | 314 | 54 |
| Cr (ppm) | 283 | 138 | 232 | 138 | 223 | 107 | 242 | 75 | 252 | 96 | 273 | 123 |
| Ni (ppm) | 53 | 53 | 46 | 51 | 44 | 29 | 36 | 24 | 33 | 37 | 52 | 45 |
| Cu (ppm) | 89 | 50 | 69 | 50 | 66 | 49 | 51 | 33 | 85 | 42 | 80 | 61 |
| Zn (ppm) | 189 | 87 | 164 | 67 | 157 | 74 | 168 | 108 | 200 | 64 | 161 | 95 |
| Rb (ppm) | 271 | 63 | 278 | 45 | 373 | 37 | 303 | 35 | 272 | 56 | 331 | 36 |
| Sr (ppm) | 233 | 84 | 320 | 188 | 133 | 74 | 122 | 26 | 479 | 125 | 183 | 53 |
| Y ₂ (ppm) | 30 | 14 | 39 | 19 | 54 | 18 | 45 | 12 | 29 | 13 | 37 | 17 |
| Zr (ppm) | 520 | 99 | 572 | 195 | 546 | 151 | 701 | 129 | 349 | 151 | 526 | 78 |
| Nb (ppm) | 20 | 10 | 22 | 19 | 34 | 13 | 24 | 11 | 17 | 10 | 23 | 9 |
| Pb (ppm) | 810 | 1431 | 369 | 566 | 145 | 136 | 719 | 678 | 67 | 38 | 877 | 1055 |

3.2. Unsupervised Modeling

HCA and k-means are both unsuccessful to group the different samples in their corresponding locations. Firstly, HCA could be an unsuitable classification method because the expectable clustering structure for the studied sites should not be particularly hierarchical but flat. In any case, the resulting HCA dendrogram can be cut at a certain level to produce a set of six clusters. These should contain a rather homogeneous number of samples (33 to 37) distributed according to the real distribution of samples per village (as appears in Table 1). However, the distribution of samples per cluster is clearly heterogeneous (Table 3). Two clusters (X1 and X3) contain more than 70% of the samples. Grosso modo X3 contains most of the samples from three real clusters (Breda, Sant Julià and Quart) and X1 contains most of the samples from Verdú and around half those from Esparreguera. The samples from the remaining real cluster (La Bisbal) appear scattered in the six predicted clusters. A slightly more balanced distribution of samples per class is obtained using the k-means model but still is far from being close to an acceptable result. The k-means model has been set to obtain six clusters (Table 4). This time three predicted clusters are big, amounting around 75% of the samples and the other three only contain the remaining 25%. Again, the predicted clusters do not contain samples from a preponderate location. For instance, the X3 cluster groups together most of the samples from Breda, half those from Quart and a third of those from Sant Julià, or X6 contains most of the samples from Verdú, half those from Esparreguera and a third of those from La Bisbal. Samples from La Bisbal, and this time also from Sant Julià, appear scattered within the six predicted clusters. It is worth of note that some similarities can be found between the predicted clusters using HCA and k-means, for instance many samples from Breda, Sant Julià and Quart tend to group in a single cluster and something similar occurs with those from Verdú and Esparreguera.

Table 3. HCA cluster prediction (cutting the dendrogram to produce six clusters).

| Predicted Class | Actual Class | | | | | | Samples within the Predicted Class |
|-----------------|--------------|--------|-------|-------|-------|------------|------------------------------------|
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià | |
| X1 | 18 | 11 | 1 | 0 | 26 | 0 | 56 |
| X2 | 7 | 9 | 1 | 0 | 1 | 3 | 21 |
| X3 | 7 | 9 | 25 | 25 | 0 | 29 | 95 |
| X4 | 0 | 2 | 6 | 8 | 0 | 4 | 20 |
| X5 | 0 | 2 | 0 | 0 | 6 | 0 | 8 |
| X6 | 1 | 3 | 0 | 0 | 4 | 0 | 8 |

Table 4. K-means cluster prediction (setting six clusters).

| Predicted Class | Actual Class | | | | | | Samples within the Predicted Class |
|-----------------|--------------|--------|-------|-------|-------|------------|------------------------------------|
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià | |
| X1 | 8 | 10 | 1 | 0 | 1 | 2 | 22 |
| X2 | 0 | 1 | 5 | 6 | 0 | 4 | 16 |
| X3 | 0 | 1 | 17 | 24 | 0 | 12 | 54 |
| X4 | 0 | 3 | 0 | 0 | 10 | 0 | 13 |
| X5 | 7 | 9 | 9 | 3 | 0 | 18 | 46 |
| X6 | 18 | 12 | 1 | 0 | 26 | 0 | 57 |

The results from PCA also fail to discriminate most of the clusters. Despite being a multivariate technique, the obtained result privileges mainly a single variable. From the new set of variables, PC1 alone can explain ~94% of the variance and PC2 nearly the 4%, therefore all the other new variables only hold the remaining 2% of variance. Besides that, looking at the definition of the two main principal components (Table 5 and inset in Figure 4) it is apparent that PC1 is basically the SiO₂ content and

PC2 a combination of the Fe_2O_3 and Al_2O_3 content. Figure 4 depicts the corresponding biplot where samples have been colored according to their known provenance. PC2 is actually a non-discriminant variable because samples from all the studied localities show a similar range of variability along the PC2 axis. Regarding PC1 (i.e., SiO_2), this single variable allows a net distinction between samples from Breda (high SiO_2 content) and Verdú (low SiO_2 content). Most of the samples from la Bisbal and Esparreguera are also relatively low in SiO_2 , and therefore they lie scattered basically in the same area along with the samples from Verdú. On the other hand, most of the samples from Sant Julià and Quart are comparatively richer in SiO_2 and they align with the samples from Breda, defining actually a straight line in the biplot (and therefore reveal a correlation between PC1 and PC2 for the samples from these three villages).

Table 5. PCA coefficients for the two main principal components.

| | PC1 (93.69%) | PC2 (3.94%) |
|-----------|--------------|-------------|
| SiO_2 | -0.9863 | -0.0511 |
| Al_2O_3 | -0.1321 | 0.5063 |
| Fe_2O_3 | 0.0145 | 0.8187 |
| MgO | 0.0883 | 0.1522 |
| Na_2O | -0.0027 | -0.0083 |
| K_2O | -0.0398 | 0.2118 |
| TiO_2 | -0.0082 | 0.0512 |
| Cl | 0.0036 | -0.0017 |
| MnO | 0.0004 | 0.0014 |
| Cr_2O_3 | -0.0001 | 0.0037 |
| NiO | -0.0001 | 0.0006 |
| CuO | 0.0001 | 0.0002 |
| ZnO | 0.0001 | 0.0020 |
| Rb_2O | -0.0004 | 0.0018 |
| SrO | 0.0017 | -0.0001 |
| Y_2O_3 | -0.0001 | 0.0002 |
| ZrO_2 | -0.0014 | -0.0023 |
| Nb_2O_5 | -0.0001 | 0.0003 |

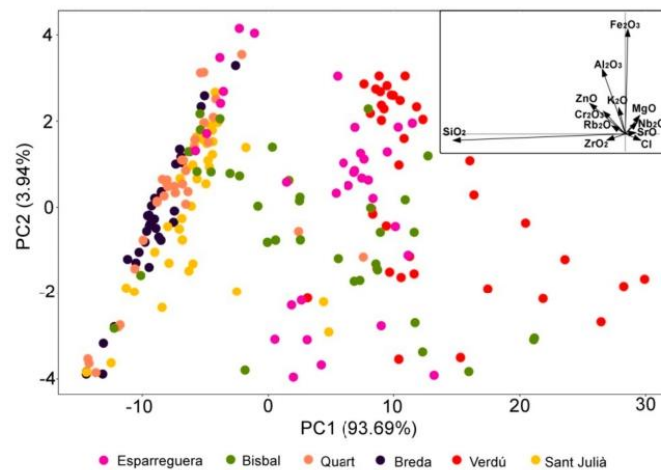


Figure 4. PCA biplot of factor scores for the first two principal components for all the processed samples. Inset: PCA biplot of the most relevant variables.

3.3. Supervised Modeling

Using the training dataset (a random 80% portion of the total dataset) every model was optimized to attain the maximum accuracy; this implied the automatic selection of different parameters for each model:

1. Weighted k-nearest neighbors (kkNN). Optimization was done by defining the values of three parameters of the model: Kmax = 5, distance = 2 and kernel = optimal; these parameters determine the way to define the neighbors and their distance to a given object.
2. Random forest (RF). The method rf gives the possibility to modulate three different parameters; mtry, splitrule and min.node.size. After optimization they were fixed as mtry = 19, splitrule = extratrees and min.node.size = 1. Mtry is the number of variables randomly sampled as candidates at each tree split.
3. Artificial neural network (ANN). After optimization the best architecture for the ANN classifier was found to be made of a single hidden layer of 5 units using a weight decay value of 0.1 (this is a multiplier factor for the weighted factors to avoid overfitting).
4. Linear discriminant analysis (LDA). No parameters were modulated for this straightforward model.
5. Generalized linear models (Glmnet). Only two parameters were optimized, $\alpha = 1$ and $\lambda = 0.00488$. The first implies that the lasso regularization was used and λ is the regularization penalty.

LDA is the only model without optimization of parameters during the training step. However, training of this model results in the definition of the new relevant variables as linear combination of the initial features (similarly to the unsupervised PCA model) that maximize the separation between clustered classes. Figure 5 is the corresponding 3D plot (for the three main discriminants) of the trained dataset and it becomes apparent that the clouds representing every class have little or no overlap. It is worth to note that, unlike PCA, the new variables are no longer essentially defined with a single compositional feature (see in Table 6, for instance that the main contributors to LD1 are SiO₂, MgO and K₂O whilst PC1 was essentially SiO₂ alone). Besides that, variance is not concentrated in a single variable. Therefore, the LDA model really uses a multidimensional space to separate the different clusters and to take advantage of all the available data. In contrast, PCA was so strongly biased to SiO₂ that basically only one dimension was used to separate the samples and yet unsuccessfully.

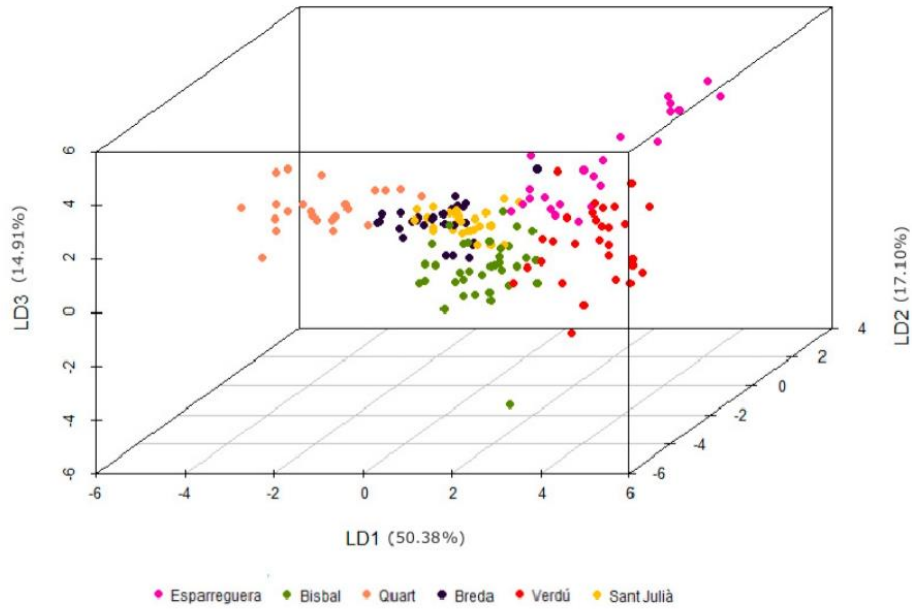


Figure 5. 3D scatter plot using the three main LDA linear discriminants.

Table 6. LDA coefficients for the four main linear discriminants.

| Compound | Linear Discriminants | | | |
|--------------------------------|----------------------|--------|--------|--------|
| | LD1 | LD2 | LD3 | LD4 |
| SiO ₂ | -1.190 | 1.210 | 0.467 | -0.635 |
| Al ₂ O ₃ | -0.363 | -0.595 | 0.352 | 0.940 |
| Fe ₂ O ₃ | 0.299 | 0.329 | -0.302 | -0.038 |
| MgO | 0.764 | 0.678 | 1.216 | -0.081 |
| Na ₂ O | 0.139 | 0.141 | 0.056 | -0.404 |
| K ₂ O | 0.648 | 0.776 | -1.216 | -0.391 |
| TiO ₂ | 0.300 | -0.202 | 0.262 | 0.202 |
| Cl | -0.295 | 0.113 | 0.127 | 0.079 |
| MnO | 0.071 | -0.160 | 0.338 | 0.315 |
| Cr ₂ O ₃ | 0.115 | 0.511 | -0.138 | -0.273 |
| NiO | -0.111 | -0.666 | -0.413 | 0.524 |
| CuO | 0.035 | 0.030 | 0.375 | -0.611 |
| ZnO | -0.073 | 0.283 | 0.006 | -0.109 |
| Rb ₂ O | -0.549 | -1.052 | 1.075 | -0.529 |
| SrO | 0.064 | -0.415 | -0.898 | -0.138 |
| Y ₂ O ₃ | -0.398 | -0.394 | -0.110 | -0.407 |
| ZrO ₂ | 0.295 | 0.412 | -0.811 | 0.718 |
| Nb ₂ O ₅ | -0.139 | 0.138 | 0.030 | 0.197 |

After the training step, each parametrized model was tested using the remaining 20% of the total dataset not used during the training step (precisely named test dataset). Table 7 contains the corresponding confusion matrices.

Table 7. Confusion matrix and accuracy for each supervised model (kkNN, RF, ANN, LDA and Glmnet).

| Actual class | kkNN-Based Provenance Classification | | | | | |
|-----------------------|--|--------|-------|-------|-------|------------|
| | Predicted class | | | | | |
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 5 | 0 | 0 | 0 | 1 | 1 |
| Bisbal | 1 | 7 | 0 | 0 | 1 | 0 |
| Quart | 0 | 0 | 2 | 2 | 0 | 0 |
| Breda | 0 | 0 | 1 | 5 | 0 | 0 |
| Verdú | 0 | 0 | 0 | 0 | 8 | 0 |
| Sant Julià | 0 | 0 | 0 | 0 | 0 | 8 |
| Accuracy 83.3% | | | | | | |
| Actual class | RF-based provenance classification | | | | | |
| | Predicted class | | | | | |
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 4 | 0 | 0 | 0 | 2 | 1 |
| Bisbal | 0 | 9 | 0 | 0 | 0 | 0 |
| Quart | 0 | 0 | 4 | 0 | 0 | 0 |
| Breda | 0 | 0 | 1 | 5 | 0 | 0 |
| Verdú | 0 | 0 | 0 | 0 | 8 | 0 |
| Sant Julià | 0 | 1 | 0 | 0 | 0 | 7 |
| Accuracy 88.1% | | | | | | |
| Actual class | ANN-based provenance classification | | | | | |
| | Predicted class | | | | | |
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 6 | 0 | 0 | 0 | 1 | 0 |
| Bisbal | 0 | 6 | 1 | 1 | 1 | 0 |
| Quart | 0 | 0 | 4 | 0 | 0 | 0 |
| Breda | 1 | 1 | 0 | 4 | 0 | 0 |
| Verdú | 1 | 1 | 0 | 0 | 6 | 0 |
| Sant Julià | 0 | 0 | 0 | 0 | 0 | 8 |
| Accuracy 81.0% | | | | | | |
| Actual class | LDA-based provenance classification | | | | | |
| | Predicted class | | | | | |
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 3 | 0 | 0 | 0 | 2 | 2 |
| Bisbal | 0 | 7 | 0 | 1 | 0 | 1 |
| Quart | 0 | 0 | 4 | 0 | 0 | 0 |
| Breda | 0 | 0 | 0 | 6 | 0 | 0 |
| Verdú | 0 | 1 | 0 | 0 | 7 | 0 |
| Sant Julià | 0 | 0 | 0 | 0 | 0 | 8 |
| Accuracy 83.3% | | | | | | |
| Actual class | Glmnet-based provenance classification | | | | | |
| | Predicted class | | | | | |
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 5 | 0 | 0 | 0 | 2 | 0 |
| Bisbal | 0 | 7 | 1 | 1 | 0 | 0 |
| Quart | 0 | 1 | 3 | 0 | 0 | 0 |
| Breda | 0 | 1 | 0 | 5 | 0 | 0 |
| Verdú | 0 | 1 | 0 | 0 | 7 | 0 |
| Sant Julià | 0 | 0 | 0 | 0 | 0 | 8 |
| Accuracy 83.3% | | | | | | |

All the models exhibit a high rate of successful class prediction for the 42 samples from the test dataset, with accuracies > 80% and RF attains the highest (88.1%). Esparreguera and La Bisbal are the classes that seem harder to predict successfully, in contrast Sant Julià is a class particularly well predicted by all the testes models. However, through the analysis of a correlation matrix, it can be seen that, besides these trends, there are not clear correlations between the results from each model. For instance, the LDA model is not very good at predicting the provenance of samples from Esparreguera (only three out of seven are well predicted) and in contrast both ANN is much more efficient at it (six out of seven). However, the LDA model is the best at predicting the provenance of samples from Breda. The lack of clear correlations between the models augurs well for the chances to increase the obtained accuracy using a combined approach (stack of models). Table 8 shows the corresponding confusion matrix for the stack of models.

Table 8. Confusion matrix and accuracy for the stack of models.

| Actual Class | Predicted Class | | | | | |
|-----------------------|-----------------|--------|-------|-------|-------|------------|
| | Esparreguera | Bisbal | Quart | Breda | Verdú | Sant Julià |
| Esparreguera | 4 | 0 | 0 | 0 | 0 | 0 |
| Bisbal | 0 | 9 | 0 | 0 | 0 | 1 |
| Quart | 0 | 0 | 4 | 1 | 0 | 0 |
| Breda | 0 | 0 | 0 | 5 | 0 | 0 |
| Verdú | 2 | 0 | 0 | 0 | 8 | 0 |
| Sant Julià | 1 | 0 | 0 | 0 | 0 | 7 |
| Accuracy 88.1% | | | | | | |

The stack is a meta-model or a combination of all models. Instead of using a single classification model based on experimental features (the chemical composition of samples), the stack approach uses the predictions of the different classifiers as features. However, using this combined approach the accuracy reaches 88.1%, a value already obtained using the RF approach.

4. Discussion

Unsupervised methods are widely used in provenance studies of pottery and particularly PCA is routinely applied [43] to define geochemical and/or petrographic groups within sampled materials from archaeological workshops and consumption centers. However, the presented results show that these methods would have failed to detect and distinguish a mixed ensemble containing pottery from the six studied production centers. The geographical proximity, inherent chemical variability and a similar geological context could explain the difficulty to distinguish geochemically the different sites.

The presented results for the three unsupervised methods (HCA, k-means and PCA) agree on detecting roughly two different classes: on the one hand samples with a relative high Si content (mainly those from Breda, Quart and Sant Julià) and on the other hand those with a comparatively lower Si content (predominantly samples from Verdú, Esparreguera and La Bisbal). The preponderant role of SiO₂ has been clearly illustrated by the PCA results, as the main composed variable (PC1) is basically defined as SiO₂ and it bears almost all the variance (94%). The higher SiO₂ content reflects a higher mineralogical abundance of quartz in the clays from Breda, Quart and Sant Julià, possibly correlated with a coarser grain size. The chemical analyses reveal an inverse correlation between Si and Ca, and indeed mineralogically there is a higher calcite content in the clays from Verdú, Esparreguera, and La Bisbal (and higher Ca-bearing minerals within the corresponding pottery). However, these are just general trends that cannot be used alone to identify productions from a given locality. Unsupervised methods do not produce distinguishable clusters that could be correlated with the actual provenance of the clay and pottery samples; the results only indicate roughly two different classes. When the number of classes is fixed to be six the different classes appear to be quantitatively

highly imbalanced and formed by a mixture of samples from different sites. Therefore, these methods cannot be used to define a geochemical fingerprint to track or certify the provenance of the samples.

In contrast, the tested supervised methods, through a machine learning approach, have been able to develop predictive models of provenance with accuracies above 80%, sometimes as high as 88.1%. Such level of accuracy can be considered very high, actually accuracies around 75% have been considered enough to rate a predictive model on soil prediction as successful [26]. This potentially opens the possibility for developing a tool that could predict the class of unknown samples, with an about 90% of accuracy, that could be used to certify the provenance of pottery productions with that level of probability.

Nevertheless, it should be noted that the accuracy predictions are performed using the test dataset, which has been defined as the 20% of the total dataset. As the full dataset contains the chemical analysis of 208 samples, there are only 42 samples within the test dataset. Taking into account that the samples come from six different localities it turns out that the capacity of the models to predict provenance is computed using only a very few test samples per site (for instance in the case of La Bisbal 9 test samples were used, and for Quart, only four). It is known that the success of machine learning methods depends on the amount and quality of available data [44] and a minimum total dataset size of 100 samples has been hypothesized as the lower limit to apply machine learning methods in materials research [45]. The presented results have been derived with a dataset of 208 (just above double the hypothesized minimum size), therefore an enlarged dataset would be required to increase the confidence on the obtained accuracy. However, a significant increase of the experimental dataset is time and cost consuming, the obtained results can be taken as encouraging and demonstrative of the potential of the supervised approach as a way to define geochemical fingerprints to track or certify the provenance of samples.

With the currently available dataset it is possible to further analyze the significance of the obtained results by repeatedly performing the full training and test process using different splits of training and test data. In the previous section it has been always used a particular split obtained using always the same random seed. Ten different seeds have been used to generate ten different splits that have been applied to every supervised model (and also to the stack of all the models), the obtained accuracies are not always exactly the same and therefore it is possible to analyze the distribution of the obtained accuracies for every tested method, and the results are shown in Figure 6.

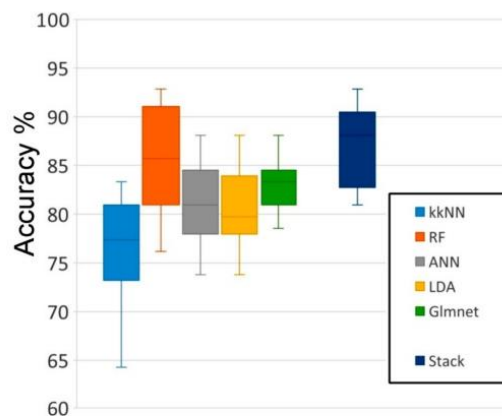


Figure 6. Boxplot with the accuracy variation for each model using different splits.

The obtained distributions of accuracy are relatively narrow. kkNN is the model that produces statistically lower accuracies, in some occasions even below 70%. Three models (ANN, LDA and Glmnet) yield accuracies with interquartile ranges between 78% and 85%. Finally, the best performing

models are RF and the stacked meta model, both exhibiting a distribution of accuracies that expand above 90%. However, the distribution of the stack of models is narrower and keeping its median value at 88.1%. Therefore, the results using different splits confirm that the ‘stack of models’ approach is the best classification approach.

5. Conclusions

Supervised machine learning methods have proven to be useful to extract geochemical fingerprints (for both clays and pottery from a given site) and these allow inter-site discrimination with accuracy levels of 80% and above.

Unsupervised methods are classically used in archaeometry to enable the identification of compositional pottery groups to distinguish between local products from a given workshop and different exports. Nevertheless, these methods have failed to distinguish the raw materials and pottery products from the six studied villages. The presented results should warn archaeometrists against the careless use of such methods, particularly if the distinction between closely related provenances is envisaged.

In the modern context of revitalization of the traditional pottery production in Catalonia, the presented approach based on supervised machine learning methods could be the useful to develop effectively a scientific protocol to control this industry. The protocol has the potential to make feasible the introduction of seals of quality and provenance to regulate the sector. Periodic chemical analyses (lead and cadmium) on ceramic products are already performed for articles intended to come into contact with foodstuffs to meet the European and Spanish regulations. A similar approach including an exhaustive compositional characterization could be implemented for those potters that would like to certify the use of local raw materials.

The geographical closeness and similar geological context for the six studied localities highlight the robustness of the presented approach that could easily be exported to other pottery centers and similar problematics.

Supplementary Materials: The following is available online at <http://www.mdpi.com/2075-163X/10/1/8/s1>, Table S1: Full set of chemical analyses (clays, pottery and briquettes) for the six studied localities.

Author Contributions: Conceptualization and fieldwork (clay and shard sampling) by A.A. and L.C., experimental work (clay firing) by A.A. and chemical analyses by A.A. and I.Q. Formal analyses of data by A.A. and M.A. writing—first draft preparation by A.A. and L.C.; writing—review and editing, all authors. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by the Spanish Ministerio de Economía y Competitividad, grant number CGL2013-42167-P.

Acknowledgments: We are grateful to Anna Pallàs, Eduard Recasens and Jenifer Obama for their contribution to fieldwork, sample preparation and experimental measurements. We want also to thank all the institutions that have contributed to the work with pottery samples: Ceràmiques Sedó, Terracotta museum, Terrissa de Quart museum, Terrissers de Quart association and Rocaguinarda museum. Finally, we would like to thank the editor as well as the anonymous reviewers for their valuable remarks and comments.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Romero, A.; Rosal, J. *La Terrissa a Catalunya*; Brau Edicions SL: Figueres, Spain, 2014; ISBN 841588513X.
2. Vandecandelaere, E.; Teyssier, C.; Barjolle, D.; Jeanneaux, P.; Fournier, S.; Beucherie, O. *Strengthening Sustainable Food Systems Through Geographical Indications*; Invest. Centre. Dir. Invest. eng no. 13; FAO: Rome, Italy, 2018.
3. González, A.; Llorens, A.; Cervera, M.L.; Armenta, S.; de la Guardia, M. Elemental fingerprint of wines from the protected designation of origin Valencia. *Food Chem.* **2009**, *112*, 26–34. [[CrossRef](#)]
4. Pollard, A.M.; Batt, C.M.; Stern, B.; Young, S.M.M. *Analytical Chemistry in Archaeology*; Cambridge University Press: Cambridge, UK, 2007; ISBN 9780511607431.
5. Pollard, A.M.; Heron, C. *Archaeological Chemistry*; The Royal Society of Chemistry: Cambridge, UK, 2008; ISBN 978-0-85404-262-3.

6. Mommsen, H. Short Note: Provenancing of Pottery—The Need for an Integrated Approach? *Archaeometry* **2004**, *46*, 267–271. [[CrossRef](#)]
7. Kuleff, I.; Djingova, R. Provenance study of pottery; choice of elements to be determined. *ArchéoSciences Rev. d'Archéométrie* **1996**, *20*, 57–67. [[CrossRef](#)]
8. Pagespetit, A.B. *La Ceràmica*; Diputació de Girona/Caixa de Girona: Olot, Spain, 1993; Volume 42, ISBN 84-8067-019-3.
9. Batista, D. La Bisbal crea una marca de denominación de origen para proteger sus productos cerámicos. *Expansión.com*, 4 April 2012.
10. Pairoli, M. *Quart: Natura, Història i Artesania*; Ajuntament de Quart: Girona, Spain, 1998; ISBN 84-923701-1-4.
11. Nayak, P.S.; Singh, B.K. Instrumental characterization of clay by XRF, XRD and FTIR. *Bull. Mater. Sci.* **2007**, *30*, 235–238. [[CrossRef](#)]
12. Zhou, X.; Liu, D.; Bu, H.; Deng, L.; Liu, H.; Yuan, P.; Du, P.; Song, H. XRD-based quantitative analysis of clay minerals using reference intensity ratios, mineral intensity factors, Rietveld, and full pattern summation methods: A critical review. *Solid Earth Sci.* **2018**, *3*, 16–29. [[CrossRef](#)]
13. Sanjurjo-Sánchez, J.; Montero Fenollós, J.L.; Barrientos, V.; Polymeris, G.S. Assessing the firing temperature of Uruk pottery in the Middle Euphrates Valley (Syria): Bevelled rim bowls. *Microchem. J.* **2018**, *142*, 43–53. [[CrossRef](#)]
14. Morgan, D.J. Thermal analysis—including evolved gas analysis—of clay raw materials. *Appl. Clay Sci.* **1993**, *8*, 81–89. [[CrossRef](#)]
15. Aitchison, J. *The Statistical Analysis of Compositional Data*; Blackburn Press: Caldwell, NJ, USA, 2003.
16. Panchuk, V.; Yaroshenko, I.; Legin, A.; Semenov, V.; Kirsanov, D. Application of chemometric methods to XRF-data—A tutorial review. *Anal. Chim. Acta* **2018**, *1040*, 19–32. [[CrossRef](#)]
17. Baxter, M.J. A review of supervised and unsupervised pattern recognition in archaeometry. *Archaeometry* **2006**, *48*, 671–694. [[CrossRef](#)]
18. Munita, C.S.; Paiva, R.P.; Alves, M.A.; de Oliveira, P.M.S.; Momose, E.F. Provenance Study of Archaeological Ceramic. *J. Trace Microprobe Technol.* **2003**, *21*, 697–706. [[CrossRef](#)]
19. Scarpelli, R.; Robustelli, G.; Clark, R.J.H.; De Francesco, A.M. Scientific investigations on the provenance of the black glazed pottery from Pompeii: A case study. *Mediterr. Archaeol. Archaeom.* **2017**, *17*, 1–10.
20. Buxeda, I.; Garrigós, J.; Cau Ontiveros, M.A.; Kilikoglou, V. Chemical Variability in Clays and Pottery from a Traditional Cooking Pot Production Village: Testing Assumptions in Pereruela*. *Archaeometry* **2003**, *45*, 1–17. [[CrossRef](#)]
21. Boleda Cases, R. *La Ceràmica Negra de Verdú. Cantirers i Terrissaires*; Grup de Recerques de les Terres de Ponent: Verdú, Spain, 2014; ISBN 9788461685875.
22. Rocas, X.; Roqué, C. Terres i terreres: la matèria primera de la indústria ceràmica bisbalenca. *Estud. del Baix Empordà* **2015**, *34*, 13–53.
23. Coll i Castanyer, J. *Breda Històrica i Actual*; Montblanc: Granollers, Spain, 1971.
24. Marguá, E.; Queralt, I.; Van Grieken, R. Sample Preparation for X-Ray Fluorescence Analysis. In *Encyclopedia of Analytical Chemistry*; Wiley: New York, NY, USA, 2016; pp. 1–25. ISBN 9780470027318.
25. Bell, S.; Croson, C. Artificial neural networks as a tool for archaeological data analysis. *Archaeometry* **1998**, *40*, 139–151. [[CrossRef](#)]
26. Oonk, S.; Spijker, J. A supervised machine-learning approach towards geochemical predictive modelling in archaeology. *J. Archaeol. Sci.* **2015**, *59*, 80–88. [[CrossRef](#)]
27. Barone, G.; Mazzoleni, P.; Spagnolo, G.V.; Raneri, S. Artificial neural network for the provenance study of archaeological ceramics using clay sediment database. *J. Cult. Herit.* **2019**, *38*, 147–157. [[CrossRef](#)]
28. Charalambous, E.; Dikomitou-Eliadou, M.; Milis, G.M.; Mitsis, G.; Eliades, D.G. An experimental design for the classification of archaeological ceramic data from Cyprus, and the tracing of inter-class relationships. *J. Archaeol. Sci. Rep.* **2016**, *7*, 465–471. [[CrossRef](#)]
29. Wilson, A.L. Elemental analysis of pottery in the study of its provenance: A review. *J. Archaeol. Sci.* **1978**, *5*, 219–236. [[CrossRef](#)]
30. Zhu, J.; Shan, J.; Qiu, P.; Qin, Y.; Wang, C.; He, D.; Sun, B.; Tong, P.; Wu, S. The multivariate statistical analysis and XRD analysis of pottery at Xigongqiao site. *J. Archaeol. Sci.* **2004**, *31*, 1685–1691. [[CrossRef](#)]





31. Kanungo, T.; Mount, D.M.; Netanyahu, N.S.; Piatko, C.D.; Silverman, R.; Wu, A.Y. An efficient k-means clustering algorithm: Analysis and implementation. *IEEE Trans. Pattern Anal. Mach. Intell.* **2002**, *24*, 881–892. [[CrossRef](#)]
32. Uselli, M. *R Machine Learning Essentials*; Packt Publishing: Birmingham, UK, 2014; ISBN 178398774X.
33. Stehman, S.V. Selecting and interpreting measures of thematic classification accuracy. *Remote Sens. Environ.* **1997**, *62*, 77–89. [[CrossRef](#)]
34. Hechenbichler, K.; Schliep, K. Weighted k-Nearest-Neighbor Techniques and Ordinal Classification. *Sonderforschungsbereich* **2004**, *386*, 1–16.
35. Barandiaran, I. The random subspace method for constructing decision forests. *IEEE Trans. Pattern Anal. Mach. Intell.* **1998**, *20*, 832–844.
36. Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45*, 5–32. [[CrossRef](#)]
37. Basheer, I.; Hajmeer, M.N. Artificial Neural Networks: Fundamentals, Computing, Design, and Application. *J. Microbiol. Methods* **2001**, *43*, 3–31. [[CrossRef](#)]
38. Tharwat, A.; Gaber, T.; Ibrahim, A.; Hassanien, A.E. Linear discriminant analysis: A detailed tutorial. *Ai Commun.* **2017**, *30*, 169–190. [[CrossRef](#)]
39. Kassambara, A. *Machine Learning Essentials: Practical Guide in R*; STHDA: Marseille, France, 2018.
40. McCullagh, P.; Nelder, J.A. *Generalized Linear Models*, 2nd ed.; Chapman & Hall/CRC: London, UK, 1989; ISBN 9780412317606.
41. Friedman, J.; Hastie, T.; Tibshirani, R. Regularization Paths for Generalized Linear Models via Coordinate Descent. *J. Stat. Softw.* **2010**, *33*, 1–22. [[CrossRef](#)]
42. Džeroski, S.; Ženko, B. Is Combining Classifiers with Stacking Better than Selecting the Best One? *Mach. Learn.* **2004**, *54*, 255–273. [[CrossRef](#)]
43. Angourakis, A.; Martínez Ferreras, V.; Torrano, A.; Gurt Esparraguera, J.M. Presenting multivariate statistical protocols in R using Roman wine amphorae productions in Catalonia, Spain. *J. Archaeol. Sci.* **2018**, *93*, 150–165. [[CrossRef](#)]
44. Schmidt, J.; Marques, M.R.G.; Botti, S.; Marques, M.A.L. Recent advances and applications of machine learning in solid-state materials science. *npj Comput. Mater.* **2019**, *5*, 83. [[CrossRef](#)]
45. Zhang, Y.; Ling, C. A strategy to apply machine learning to small datasets in materials science. *npj Comput. Mater.* **2018**, *4*, 25. [[CrossRef](#)]



© 2019 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).

Article

Supervised Machine Learning Algorithms to Predict Provenance of Archaeological Pottery Fragments

Anna Anglisano ¹, Lluís Casas ^{1,*}, Ignasi Queralt ² and Roberta Di Febo ¹¹ Department of Geology, Campus de la UAB, Autonomous University of Barcelona, 08193 Barcelona, Spain² Department of Geosciences, IDAEA-CSIC, Jordi Girona 18-26, 08034 Barcelona, Spain

* Correspondence: lluis.casas@uab.cat

Abstract: Code and data sharing are crucial practices to advance toward sustainable archaeology. This article explores the performance of supervised machine learning classification methods for provenancing archaeological pottery through the use of freeware R code in the form of R Markdown files. An illustrative example was used to show all the steps of the new methodology, starting from the requirements to its implementation, the verification of its classification capability and finally, the production of cluster predictions. The example confirms that supervised methods are able to distinguish classes with similar features, and provenancing is achievable. The provided code contains self-explanatory notes to guide the users through the classification algorithms. Archaeometrists without previous knowledge of R should be able to apply the novel methodology to similar well-constrained classification problems. Experienced users could fully exploit the code to set up different combinations of parameters, and they could further develop it by adding other classification algorithms to suit the requirements of diverse classification strategies.

Keywords: pottery; provenance studies; supervised methods; machine learning; clustering; XRF; data sharing; open source software; heritage science



Citation: Anglisano, A.; Casas, L.; Queralt, I.; Di Febo, R. Supervised Machine Learning Algorithms to Predict Provenance of Archaeological Pottery Fragments. *Sustainability* **2022**, *14*, 11214. <https://doi.org/10.3390/su141811214>

Academic Editors: Andrea Zerboni, Francesco Carrer, Filippo Brandolini and Guido S. Mariani

Received: 31 July 2022

Accepted: 5 September 2022

Published: 7 September 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Pottery shards are possibly the most common artifacts among archaeological finds. These ceramic fragments can bring many types of information; among others, production technology [1,2], age [3] or evidence for cultural interactions through the identification of local and imported productions in the archaeological sites. In fact, one of the main branches of archaeometry regards the physical and geochemical analyses applied to provenance studies, mainly on pottery artifacts. The progressive increase in novel analytical techniques and approaches, along with digital methods and tools, has been argued to provoke a sustainability crisis due to the exponential growth of data. The generation of huge amounts of data is currently inherent to many other areas of present human activity, and its unsustainability would not be considered problematic per se. However, sustainability, as it applies to archaeology, can be understood in many different ways [4], and possibly the economic dimension of sustainability [5] is the one that should be particularly taken into account. In this sense, sustainable archaeology should promote the standardization of retrieved data, data sharing, open data and data recycling to minimize the number of required analyses to carry out an investigation.

The crucial step in provenance studies is the definition of reference groups, and this requires reference samples that only rarely are used by different authors. The typical approach is to use petrographical [6–8] and particularly chemical [9–12] analyses or a combination of both [13–16] to serve isolated research or, in some cases, multiple investigations directed by a given research group. Chemical data, commonly obtained using X-ray fluorescence (XRF) or neutron activation analysis (NAA), often consist of large datasets and imply lots of measurements. Their processing is usually addressed by the application

of statistical methods. The most common approach is the use of principal component analysis (PCA), hierarchical cluster analysis (HCA) and other unsupervised clustering methods [17–19]. These unsupervised methods are also occasionally used to process other kinds of analytical data, such as X-ray diffraction [20], infrared spectroscopic data [21], shard shapes or profiles [22] or even visible colors [23] also with the goal of differentiating ceramics having different provenance or production technologies.

By using unsupervised clustering methods, the data are not labeled before classification. Without labeling, it is basically assumed that the different groups or classes emerge naturally because the algorithms find the most useful variables to highlight differences between data. Usually, in every analyzed consumption center, different classes are identified based on the relative distances or similarities between data [24,25]. Analyses on different sites, including production centers, enable inferring inter-site connections and, ultimately, assigning provenance to most of the identified classes [14,26]. In addition to the characterization of pottery sampled at the production centers [27], reference groups are often also defined using kiln wasters [28,29], ceramic analogous materials [30,31] produced from clays fired under controlled conditions or even simply the presumably used clay deposits [32,33]. However, the most frequently used unsupervised method (PCA) is not strictly a classification method [34]. Unsupervised methods can easily fail to discriminate classes corresponding to provenance sites sharing similar features.

In contrast to unsupervised methods, supervised methods deal with data previously labeled with their corresponding class (i.e., the provenance of the reference samples is known). The models learn from a given training dataset, and after optimization of the model parameters, the model is tested with new labeled data. The best performing predictive model can be selected by looking at their ability to predict class memberships for these new data. This approach was recently tested on geochemical data from clays and modern baked clays from six local production centers of pottery [35]. Despite the geographical proximity and the common or similar geological contexts, the supervised approach proved to successfully classify the data with an accuracy above 80%. This high capability of inter-site discrimination has opened the door to applying supervised machine learning methods in archaeology and specifically within the field of pottery provenance studies. Predictive modeling using supervised methods is still largely an underexploited field within archaeology [36,37].

The use of machine learning applications in the field of archaeology is growing fast, in part due to the increasing accessibility and capability of the algorithms [38]. Its applications are expected to diversify and improve, gaining in usability and performance [37], covering vast areas of archaeology. Supervised machine and, in particular, deep learning (specifically deep convolutional neural network (CNN)) is commonly used to analyze images and recognize patterns. CNN was successfully applied in remote sensing applications within archaeological prospection [39] as well as artifact classification by site and period [40]. Classification of pottery fragments based on images was also explored [41], in particular for reconstruction purposes. The assembly criterion is often the matching between the shape of the fragments [42], and some sophisticated approaches even consider an imperfect matching due to erosion [43]. Alternatively, the classification criterion can be other than morphology; for instance, in [44], the pottery is classified according to its engravings. Non-ceramic archaeological items can also be classified using supervised methods, for instance, bone surface modifications [45]. However, archaeological machine learning contributions that do not deal with input images are not so commonly found, and in the particular field of provenance studies, these are rather scarce. Some pioneering works treated geochemical data using these methods to classify archaeological soils [46], clays [47], obsidians [48] and pottery [49]. Rare applications based on analytical data beyond elemental geochemistry can also be found, for example, gemstone provenance based on Raman data [50] or pottery provenance based on ultrasound data [51].

This paper illustrates the use of supervised modeling for provenancing archaeological pottery using chemical analyses. A chemical dataset from previous work [35] with data

from six nearby production centers was used and extended to an additional reference site representing the archaeological pottery produced in Barcelona (Catalonia, NE, Spain). After the training and optimization of a supervised discrimination models, the models were used to infer the provenance of several pottery samples retrieved from archaeological sites in Catalonia. Additionally, the models could potentially be applied to a large number of archaeological sites reported in this region. The presented approach is made available as an R-based code, including easy-to-read instructions to install and use.

The primary goal of this study was to evaluate the performance of supervised classification methods in those cases where the common unsupervised approach fails. Additionally, a secondary goal was to help the archaeological community easily implement the supervised clustering algorithms for provenancing pottery as a way to move a step forward from the common unsupervised practices.

2. Materials and Methods

2.1. Reference Sampled Materials and Geochemical Data

A previously published geochemical database (see Supplementary Materials within [35]) was used in the present study. The data were produced by Energy Dispersive X-ray fluorescence (EDXRF) analyses on 208 samples. These samples were clays (80 samples), pottery shards (101 samples) and ceramic briquettes produced in an oven (27 samples). They belong to six traditional pottery production centers relatively close to Barcelona (Figure 1):

- Esparreguera (~35 km northwest of Barcelona);
- Breda (~50 km northeast of Barcelona);
- Sant Julià de Vilatorrada (~60 km north from Barcelona);
- Quart (~80 km northeast of Barcelona);
- Verdú (~90 km northwest of Barcelona);
- La Bisbal d'Empordà (~100 km northeast of Barcelona).

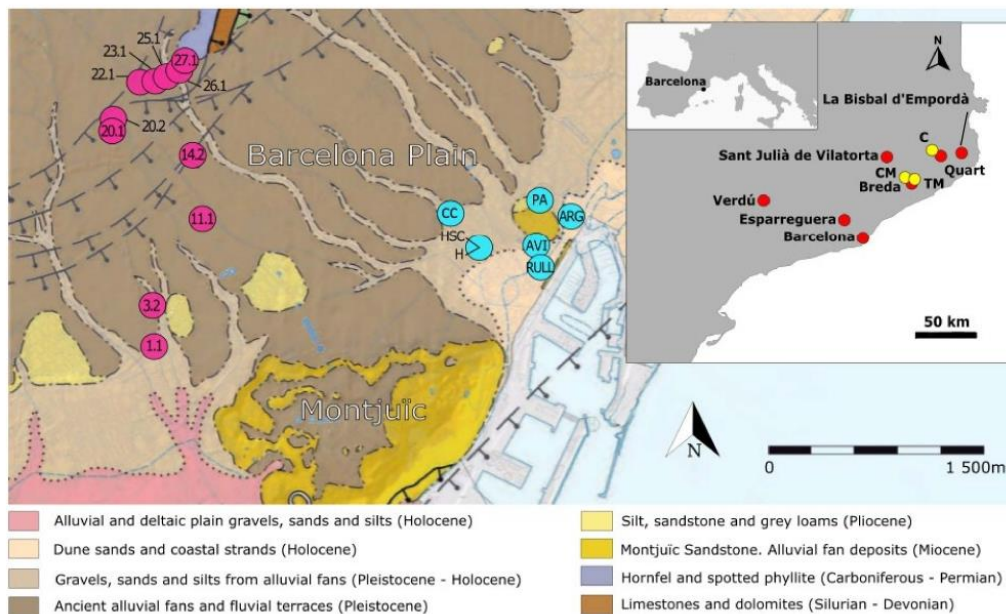


Figure 1. Geological map of the clay sampling sites (pink dots) and archaeological sites (blue dots) mainly around the *Sarrià-Sant Gervasi* and *Ciutat Vella* districts, respectively. The geological base map

was modified from ICGC. The top right corner shows a geographical map with the location of each characterized production center (red dots) and the location of the three archaeological sites that were studied (yellow dots).

These data can be labeled to their corresponding provenance class because the samples were directly extracted from their clay outcrops (clays and briquettes) or because they were known to have been produced from the corresponding local clays (pottery shards). To these six reference centers, an additional set of samples was also chemically characterized to add an extra class corresponding to the main center of the area, i.e., Barcelona itself. Since the foundation of the Crown of Aragon (12th century), Barcelona acted as the preeminent city in NE Iberia, and among many other activities, pottery production remains vastly attested from the 13th century onwards in many archaeological sites excavated in the city [52].

The extra reference dataset corresponding to Barcelona was produced by analyzing 84 samples comprising 64 archaeological samples and 20 clay samples (Figure 1). The archaeological samples were pottery shards from 7 archaeological sites (Table 1, for detailed data, see Table S1 within the Supplementary Materials). The sites cover a large time range (13th–19th centuries), and the shards were carefully selected from those constituting local productions and covering a wide range of typologies and techniques (non-glazed, glazed and glazed, including decoration). Concerning clays, as the Barcelona plain has been heavily urbanized, there is no trace of old clay pits. The 20 clay samples were extracted from eleven borehole cores drilled as part of geotechnical surveys on the Barcelona plain (Figure 1). The sampled clays were located at different depths (3 to 45 m, see Table S1 within the Supplementary Materials for detailed data), and they were selected as representative of the raw materials that could be used to produce pottery. Despite not corresponding to quarried clay pits, the sampled clays belong to levels that would have outcropped in places closer to the sea.

Table 1. Summary of the sampled pottery from Barcelona.

| Archaeological Site | N. of Samples | Age | References |
|---------------------------------|---------------|---------------------|------------|
| Avinyó (AVI) | 7 | 13th CE | [53] |
| Hospital (H) | 6 | 13th CE–14th CE | [54] |
| Rull (RULL) | 3 | 14th CE–18th CE | [55,56] |
| Hospital de la Santa Creu (HSC) | 27 | 17th CE | [54] |
| Casa de la Caritat (CC) | 9 | mid 17th CE–18th CE | [57] |
| Pia almoina (PA) | 6 | early 18th CE | [58] |
| Argenteria (ARG) | 6 | 19th CE | [59] |

Similar to the other reference sites [35], the pottery samples from Barcelona [53–59] appear to be petrographically heterogeneous, particularly in terms of grain size and inclusions/matrix ratio. However, as a general trend, the samples tend to be fine-grained, mainly including quartz, feldspar and metamorphic (phyllite and mica schist) inclusions. The petrographic heterogeneity reinforces the need to prioritize a geochemical approach to define the reference groups.

In order to obtain the geochemical composition, exactly as in other reference samples [35], both the clay and pottery samples from Barcelona were dried and ground using a laboratory mill (Pulverisette™, Fritsch GmbH, Idar-Oberstein, Germany) to pass a 125 µm mesh. The powders were then prepared in the form of pressed powder pellets using a methyl methacrylate resin as a binding agent (Elvacite™ commercial resin) under a pressure of 10 T. The composition of the pellets was measured by Energy Dispersive X-ray fluorescence (EDXRF) into an S2 Ranger system (Bruker/AXS, GmbH, Karlsruhe, Germany). The raw data were fitted using the SPECTRA.EDX package (Bruker AXS, GmbH, Karlsruhe, Germany). Quantification was made by the assisted fundamental parameters method. Analyses were made in a vacuum atmosphere for better detection of low Z elements and using different conditions of voltage to properly excite low, medium and high atomic number elements existing in the samples; the measuring time was set at 400 s.

With the extra samples from Barcelona, the full geochemical database comprises the elemental analyses of 292 samples. These are class-labeled samples for a total of seven classes or reference groups.

2.2. Archaeological Samples of Unknown Provenience

Five sets of pottery of unknown provenience were selected. They belong to three different archaeological sites and time periods. These unlabeled samples were chosen because they were retrieved from sites that lie in the vicinity of one of the reference groups, and the archaeologists that worked on them defined their production as local. Here, “local” means that the pottery was not imported from overseas or distant sites but produced near the archaeological site where it was found. Therefore, “local” here would include any of our seven reference groups and actually many other possible local workshops.

Three sets of samples were retrieved from the Montsoriu castle, an important gothic castle (10th–14th centuries) that is located on top of the homonymous hill, only ~4 km north of Breda, one of our labeled reference groups. By the end of the 20th century, the castle was an abandoned ruin, and it is currently in the process of restoration. From the several archaeological excavations on the site in 2007, the excavation of the filling of a cistern from the castle’s bailey produced many types of pottery, both imported and local [60]. Among the local pottery, there are common table and cooking ware, both glazed and non-glazed. Glazes are usually transparent (simple lead glazes) or green-colored. Pastes can be both reddish/ochre (oxidizing conditions) or grey (fired in reducing conditions). The ensemble of local pottery is dated between 1475 and 1560 [60] and despite some resemblance with later productions from Breda, it is often assumed that this pottery, in particular the green-colored glazed pottery, was produced in Barcelona’s workshops [61]. The three types of pottery from Montsoriu selected for the present investigation correspond (Table 2) to non-glazed grey pottery (7 shards), ochre-reddish pottery with a lead-glaze (10 shards), and a green-glaze (6 shards).

Table 2. Description of the five selected sets of samples of unknown provenience.

| Archaeological Site | Typology | Chronology | Tag | No. Samples | References |
|----------------------|---------------------------|---------------------------|-----|-------------|------------|
| Castell de Montsoriu | gray ware | 1475–1560 CE | CM1 | 7 | [60] |
| | lead-glazed cooking ware | | CM2 | 10 | |
| | green-glazed cooking ware | | CM3 | 6 | |
| Torre de la Mora | cooking ware | late 9th CE–early 10th CE | TM | 7 | [62] |
| La Creueta | hand-built cooking pots | 4th BCE | C | 8 | [63] |

Another set of studied samples was recovered from the Torre de la Mora site, a Roman watchtower that was reused during the Middle Ages. The tower, completely ruined at present, rose on top of a hill, also very close to Breda (~2 km northeast of it). Archaeological works undertaken in 1998–1999 produced Iberian and Early Medieval common ware. The Iberian pottery would be decontextualized, possibly integrated within rammed earth as part of the foundations. The Medieval pottery would be, according to morphological criteria, from the late 9th to the early 10th centuries [62]; 7 shards of this Medieval pottery were selected for the present study (Table 2).

Finally, another selected set of samples came from La Creueta. This is a prehistoric (Iberian) settlement that locates on top of a hill in the southern area of Girona, only 3km north of the village of Quart, i.e., one of our labeled reference groups. The site was discovered in the thirties of the 20th century, and it was repeatedly excavated. Different

types of pottery were retrieved from these excavation works, including Hellenistic and Punic imports and local wheel-made and hand-built pottery. From the typology of the identified pottery, the site is dated to the 4th century BCE [63]. The hand-built pottery would correspond to the most primitive type of pottery from the site, and it is more likely to truly represent local pottery. For this reason, 8 shards of this pottery typology were selected to take part in the present investigation (Table 2)

For detailed data of the analyzed samples of unknown provenience, see Table S2 within the Supplementary Materials. A basic petrographic and mineralogic characterization of the samples cut as thin sections was performed using a petrographic microscope (Eclipse E200, Nikon Instruments Inc., Tokyo, Japan). This was performed to check whether the samples share the same basic components as the reference samples or if any of them contain a particularly distinct mineralogical component. The chemical composition of all the shards from the five selected sets was obtained following exactly the same procedure described in Section 2.1.

2.3. Data Processing, Modelling and Class Prediction

In order to apply scripts to optimize the classification models and to perform class predictions, it is necessary to use a homogeneous set of variables. Therefore, the pastes of all the samples (both class-labeled and unlabeled) must be characterized exactly by the same ensemble of geochemical elements. All the samples were characterized using the following elements: Al, Si, Fe, Na, Mg, Cl, K, Ti, Cr, Mn, Ni, Cu, Zn, Rb, Sr, Y, Zr, and Nb. These are elements that appear in the samples above their detection limit; Ca values were disregarded because their values are strongly correlated with Si; on the other hand, Pb was also removed from the database to avoid problems related to contamination from different sources (in particular from glazes).

By using the same set of elements, the common principal component analysis (PCA) was used to illustrate the inability of unsupervised learning to classify the unlabeled samples. PCA is a method that reduces the dimensionality of the original dataset by creating a new set of variables avoiding correlated variables. By using only a few of the new variables (the main components), it is possible to plot the distribution of the dataset, maximizing the separation between samples.

The supervised classification models that were trained to learn accurate class prediction of unlabeled data were:

- Weighted k-nearest neighbors (kkNN); the samples are classified by taking into account the classes of their k-nearest neighbors;
- Random forest (RF); the algorithm is based on a combination of multiple and uncorrelated decision trees operating as an ensemble;
- Artificial neural network (ANN); this prognostic model uses an extensive network of nodes that exchange messages simulating the function of the human brain;
- Linear discriminant analysis (LDA); similar to PCA, is a linear transformation used for dimensionality reduction, but LDA maximizes the separation between classes and not between individual samples;
- Generalized linear models (GLM); this is a collection of regression models with the possibility to introduce a penalty term for the maximum likelihood (λ) to move from a pure ridge model to a pure lasso model.

These models were chosen among the most widely used supervised learning algorithms [64–66]. They are substantially different approaches that exhibit low correlation. This ensures the effectiveness of stacking models as an additional classification approach that takes advantage of the strong points from each model. The stacking technique was implemented through a random forest approach. The performance of all these models, including the stack of models, was improved using the repeated k-fold cross-validation technique during the training step. All the models are made available through the R-based code downloadable from a GitHub repository as part of this article. The algorithms corresponding to these models can be obtained freely from the following packages in the Caret

R library [67]: class (kkNN), randomForest (RF), nnet (ANN), lda (LDA), Glmnet (GLM). To make class predictions, the generic R function “predict” was used.

3. Results

A basic petrographic characterization reveals that all the archaeological samples from the five selected sets share the same main components: a ferric matrix characterized by silicate inclusions (quartz, feldspars, plagioclase and micas) and metamorphic rock fragments (mainly granitoid, phyllite and mica schist fragments). Some differences concern the state of oxidation that can change from well-oxidized (homogeneous red matrix for sets CM2 and CM3) to poorly oxidized (non-homogeneous dark-brown/grey matrix for the CM1 set), with all the other samples from other sets exhibiting matrixes with heterogeneous and intermediate oxidation. Regarding the aplastic inclusions, these are very abundant and coarser for the samples of the C set (reaching sizes of 4 mm and with a dominance of sericited feldspars and biotite) and those of the TM set (with slightly smaller grain sizes). Inclusions are moderately abundant for CM2 and CM3 samples and scarce for the more fine-grained CM1 set (with a pre-eminence of quartz). The petrographic characteristics cannot be used to connect any of these samples of unknown provenience with the geochemical reference groups. In addition to grain size and mineral abundance, the mineralogical components are basically the same for the five sets of archaeological samples and also for the samples within the reference groups, which bear clay and baked-clay reference samples with heterogeneous petrographic traits but invariably with similar minerals and rock fragments. Therefore, there are no mineralogical reasons that could prevent undertaking a geochemical classification method. The results of both unsupervised and supervised approaches are presented below.

3.1. Unsupervised Approach

Several configurations of PCA were attempted (centered alone, centered and scaled, log-transformed, etc.) using all the reference samples (208 from [35] and 84 to account for a new set representing Barcelona’s class). In all the tested configurations, the data appear distributed in a cloud where the samples from different reference sites appear completely intermixed. This overlapping is not surprising due to the similarity between the different reference sites. The use of non-scaled data produces a PCA with data from several groups distributed along a straight line (Figure 2a). The 95% confidence ellipses for every class were also drawn in the figure. Comparing these PCA results with a previously computed PCA (Figure 4 in [35]), it was apparent that the main components were almost unaffected by the addition of the new set of samples from Barcelona. PC1 (basically SiO_2) copes almost 93% of the variance, and PC2 (which correlates strongly with the Fe_2O_3 contents and Al_2O_3) around 4% of it.

The samples from three reference sites (Breda, Sant Julià de Vilatorca and Quart) are those that overlap and distribute along a straight line indicating a strong correlation between PC1 and PC2. However, the corresponding confidence ellipse is only very narrow for the Breda samples. A few outliers from Sant Julià and Quart cause wider confidence ellipses for these two reference sites. The rest of the samples from the other reference sites appear much more scattered, and the corresponding ellipses cover higher areas and appear with a very high degree of overlapping. In particular, the ellipse for the new class representing Barcelona appears completely enclosed within the ellipse representing La Bisbal, but there is actually a certain degree of overlapping with all the other ellipses. In fact, the only two ellipses with almost zero mutual overlapping are those representing Breda (high SiO_2 content) and Verdú (low SiO_2 content).

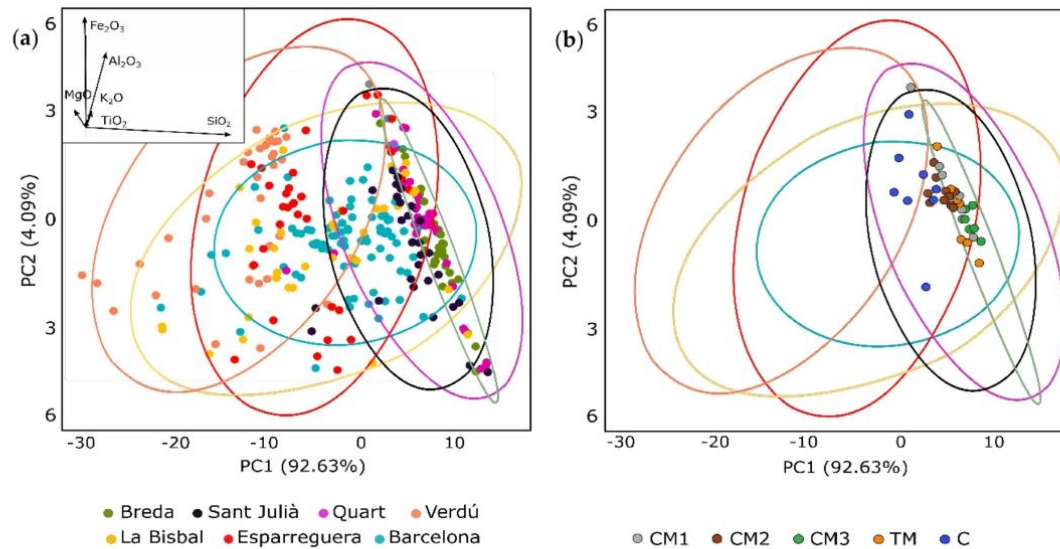


Figure 2. (a) PCA biplot of factor scores for the first two principal components for all the reference samples, 95% confidence ellipses were drawn for every class. Inset: PCA biplot of the most relevant variables. (b) The position of the samples of unknown provenience within the PCA biplot where the confidence ellipses were kept.

The position of the samples of unknown provenience within the PCA biplot is shown in Figure 2b, along with the confidence ellipses of the reference groups. It appears that the gray (CM1) and the green glazed (CM3) pottery samples from the Castell de Montsoriu, as well as those from Torre de la Mora (TM), appear in the position that defines the characteristic alignment of samples observed for Breda, Sant Julià and Quart. The lead-glazed samples from Castell de Montsoriu (CM2) appear in a similar position but do not clearly define an alignment; finally, the samples from La Creueta (C) appear a bit more scattered. In any case, the PCA biplot cannot be used to connect unambiguously any set of samples to a particular reference group. Reversely, only the CM1 and C sets appear rather disassociated from a reference site (Verdú and Breda, respectively).

3.2. Supervised Classification Models

The different tested supervised classification models were trained using 80% of the 292 reference samples. Dataset partition was performed, also keeping the 80% proportion for every labeled class; apart from this restriction, samples were randomly selected. After model optimization, the remaining 20% of the samples were used to test the performance of the classification models. Different random seeds can be used to obtain different splits of the database into the train (80%) and test (20%) subsets. Statistical values of accuracy (true positives divided by the total predictions) were computed as an indicator of the classification capability and also to check if this capability varies significantly with the number of seeds tested. The results (Figure 3) indicate that almost all models show little variation in the corresponding accuracy boxplots after computing it with the results from around ten seeds. In particular, the interquartile range remains rather stable, and only the ANN model exhibits more difficulty in stabilizing.

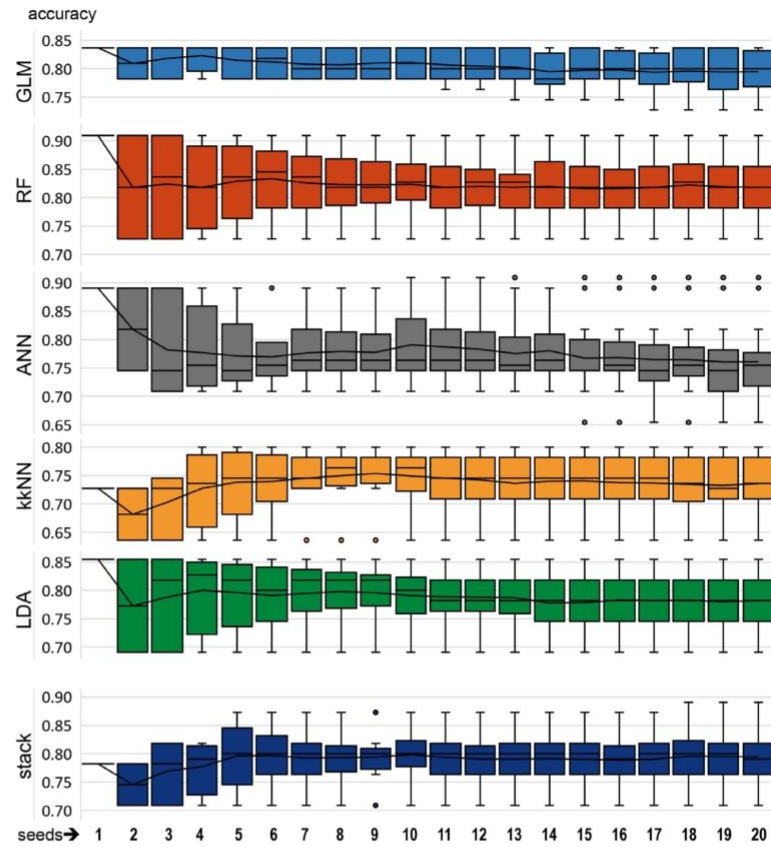


Figure 3. Boxplots of the accuracy variation as a function of the number of runs corresponding to different random seeds for all the tested classification models. The line connecting all the boxplots indicates the corresponding mean values.

By comparing the statistical distribution of accuracy for the different models as obtained using ten different random splits, it is apparent that the obtained accuracies (Figure 4a) range mostly between 0.75 and 0.85, with a moderate variation depending on the split. The low variability of the GLM model and the high accuracy values of the RF model is remarkable, and therefore these two models are less split-dependent. It is worth noting that with the addition of the new class (Barcelona), the targeted classes appear strongly imbalanced. The new class contains 84 reference samples whilst the others contain significantly fewer samples (from 33 to 37) [35]. Therefore, balanced accuracy [68] is a better indicator of the performance of the models. Additionally, an overall value of the F1-score was also computed as the simple arithmetic mean of the corresponding F1-scores per class. All three global performance indicators show very similar trends (Figure 4a), reinforcing the conclusion that ANN and, in particular, kkNN are the models with higher variability and lower mean performance.

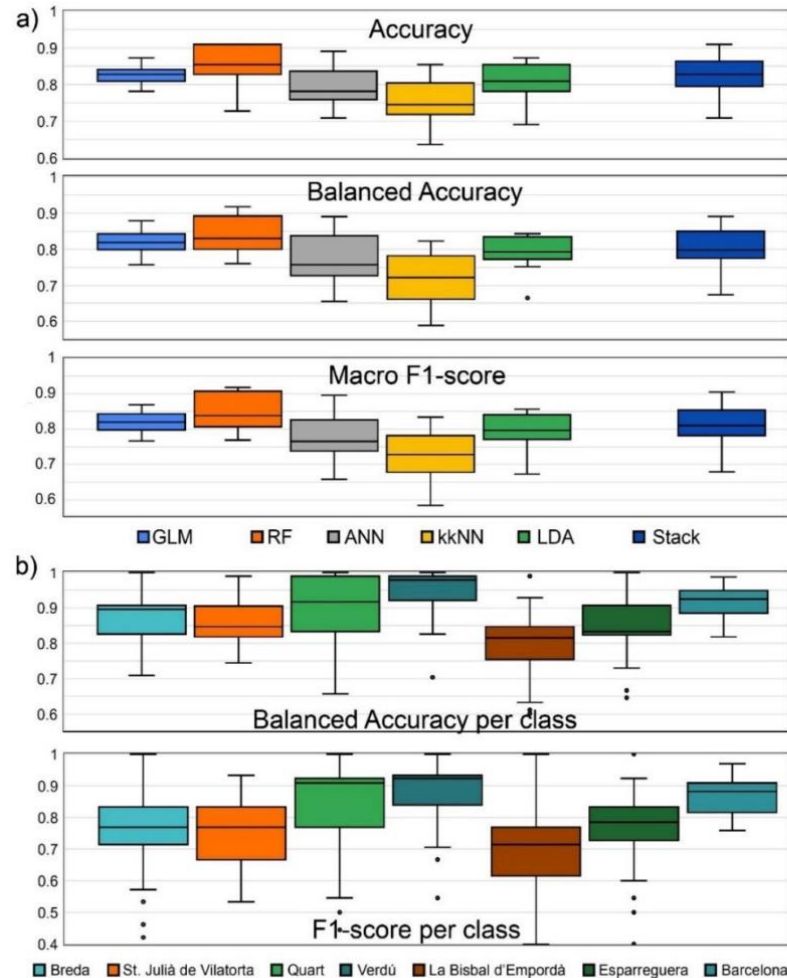


Figure 4. Different indicators of the performance of supervised models using ten different splits. (a) Boxplots corresponding to the accuracy, balanced accuracy and macro F1-score variations for each model. (b) Boxplots with indicators (balanced accuracy and F1-score) per class computed using the predictions from all the classification models.

Finally, the statistics on balanced accuracy per class and F1-score per class (Figure 4b) allow identifying Quart, Verdú and Barcelona as the classes with a higher positive prediction score whilst La Bisbal appear to be the class that is more difficult to predict and yet, on average according to the accuracy, is correctly predicted in four out of every five cases and just a bit less looking at the F1-score. In summary, accuracies are almost always above 0.6 and often above 0.8, and other performance indicators show the same trends. This enables supervised methods to be used to predict classes for unlabeled samples.

3.3. Cluster Prediction

Not surprisingly, the unsupervised methods fail to classify the unlabeled samples (samples of unknown provenience) into a given reference cluster (Figure 2b), whilst trained,

supervised methods provide accuracies generally above 0.8. In this section, we show the classification results obtained by the application of several trained, supervised models onto the five types of unlabeled samples (CM1, CM2, CM3, TM and C).

Every training subset produces a trained model that, in a second stage, can be used to classify samples of unknown provenience. The classification results of an unlabeled sample are produced as a set of probabilities of this sample to belong to the different reference clusters. Despite the great performance of supervised classification models, it is worth mentioning the importance of applying different models, seeds and samples to obtain conclusions with statistical representativeness. By looking at the obtained results from a single run using only one sample from every unlabeled set and using a given trained model (e.g., LDA), we could preliminarily and misleadingly assume that all the proveniences can be unveiled. As shown in Table 3, a single tested sample for the CM3 and C sample types produces a 100% probability of belonging to the Breda cluster and 88% to the same cluster for a TM sample. A single sample of CM1 type seems to belong to the Quart cluster (85%), whilst the tested CM2 sample would be assigned to Barcelona with a moderate 69% probability. By using all the available samples from every provenanced set (and not just a single sample), most of the preliminary conclusions hold, although some appear less clearly supported. The C set persists unambiguously assigned to Breda (100%), whilst the CM2 set appears to be assigned to Barcelona with a 78% probability and the CM3 to Breda (75%). In contrast, the probability percentages corresponding to the CM1 and TM sets appear now much more distributed into different classes.

Table 3. Classification results in the form of percentages indicating probability of correspondence to every class. For every type of provenanced sample (CM1, CM2, CM3, TM and C), the first column indicates the probability obtained using a single sample in a single run, and the second column indicates the classification percentage mean (and uncertainty) also using single runs but here applied to all the available samples of every provenanced set. Only results of LDA model are shown.

| Model | Locality | CM1 | CM2 | CM3 | TM | C | | | | | |
|------------------------------|--------------|-----|---------|-----|---------|-----|---------|----|---------|-----|---------|
| Linear Discriminant Analysis | Breda | 1 | 23 ± 18 | 5 | 7 ± 9 | 100 | 75 ± 30 | 88 | 50 ± 25 | 100 | 100 ± 0 |
| | Sant Julià | 5 | 16 ± 12 | 2 | 1 ± 0 | 0 | 1 ± 0 | 1 | 27 ± 28 | 0 | 0 ± 0 |
| | Quart | 85 | 45 ± 25 | 1 | 1 ± 0 | 0 | 14 ± 26 | 0 | 2 ± 5 | 0 | 0 ± 0 |
| | Verdú | 1 | 0 ± 0 | 0 | 0 ± 0 | 0 | 0 ± 0 | 0 | 0 ± 0 | 0 | 0 ± 0 |
| | La Bisbal | 3 | 3 ± 0 | 10 | 7 ± 3 | 0 | 3 ± 3 | 0 | 15 ± 23 | 0 | 0 ± 0 |
| | Esparreguera | 3 | 6 ± 4 | 12 | 6 ± 3 | 0 | 0 ± 0 | 1 | 3 ± 7 | 0 | 0 ± 0 |
| | Barcelona | 2 | 5 ± 8 | 69 | 78 ± 13 | 0 | 7 ± 10 | 9 | 3 ± 5 | 0 | 0 ± 0 |

As seen before, different training sets (i.e., runs with different seeds) could result in slightly different classification results, and this also applies to samples of unknown provenience. Ten different trained configurations were used for every model to obtain probabilities with associated uncertainties. Computing the mean probabilities now clearly results in much more scattered probabilities (Table 4). By looking at the results from the LDA model, the probabilities associated with the CM2, CM3 and C sets still appear rather concentrated in a given class (Barcelona, Breda and again Breda, respectively). However, taking into account the results from different models, it becomes apparent that the provenance of the samples remains unknown for almost all the sample sets. For instance, the LDA model indicates that the C-type samples belong to the Breda cluster with a 97% probability, but the RF model distributes the probability into all the classes. The only case of systematic attribution of a set of samples to a given cluster is that of the CM2 set. For this set, the probability percentages always indicate provenance from Barcelona regardless of the classification model. For all the other sets, the probability percentages distribute into different classes. In the case of the CM3 set, two classes cope systematically with most of the probability, but in other cases, the concerned classes vary significantly from model to model.

Table 4. Classification results in the form of probability percentages (including uncertainties) obtained after running ten times the training and cluster classification code on all the samples from the 5 sets of samples of unknown provenience. Only results of LDA, RF and the stack of models are shown.

| Model | Locality | CM1 | CM2 | CM3 | TM | C |
|------------------------------|-------------|---------|---------|---------|---------|---------|
| Linear Discriminant Analysis | Breda | 29 ± 21 | 9 ± 13 | 76 ± 28 | 58 ± 25 | 97 ± 10 |
| | Sant Julià | 17 ± 12 | 1 ± 1 | 1 ± 1 | 22 ± 23 | 2 ± 7 |
| | Quart | 40 ± 24 | 0 ± 0 | 13 ± 23 | 1 ± 2 | 0 ± 0 |
| | Verdú | 2 ± 2 | 0 ± 0 | 0 ± 0 | 0 ± 0 | 0 ± 0 |
| | La Bisbal | 5 ± 3 | 7 ± 5 | 3 ± 5 | 11 ± 17 | 1 ± 2 |
| | Esparreguer | 4 ± 3 | 7 ± 5 | 0 ± 0 | 5 ± 12 | 0 ± 1 |
| | Barcelona | 4 ± 7 | 75 ± 16 | 7 ± 8 | 3 ± 5 | 0 ± 0 |
| Random Forest | Breda | 11 ± 10 | 8 ± 4 | 37 ± 17 | 25 ± 22 | 13 ± 16 |
| | Sant Julià | 20 ± 7 | 6 ± 2 | 5 ± 2 | 16 ± 8 | 21 ± 7 |
| | Quart | 33 ± 12 | 4 ± 2 | 46 ± 8 | 10 ± 4 | 11 ± 6 |
| | Verdú | 0 ± 0 | 0 ± 0 | 0 ± 0 | 0 ± 0 | 2 ± 1 |
| | La Bisbal | 16 ± 7 | 11 ± 5 | 5 ± 2 | 18 ± 7 | 17 ± 4 |
| | Esparreguer | 7 ± 4 | 3 ± 1 | 1 ± 1 | 3 ± 3 | 13 ± 6 |
| | Barcelona | 11 ± 5 | 68 ± 10 | 2 ± 1 | 21 ± 6 | 19 ± 7 |
| Stack of Models | Breda | 11 ± 16 | 0 ± 2 | 50 ± 37 | 57 ± 37 | 35 ± 31 |
| | Sant Julià | 17 ± 34 | 0 ± 0 | 0 ± 1 | 7 ± 20 | 21 ± 29 |
| | Quart | 62 ± 39 | 0 ± 0 | 48 ± 37 | 0 ± 2 | 0 ± 1 |
| | Verdú | 0 ± 0 | 0 ± 1 | 0 ± 2 | 0 ± 0 | 3 ± 7 |
| | La Bisbal | 6 ± 18 | 0 ± 2 | 0 ± 0 | 8 ± 21 | 3 ± 8 |
| | Esparreguer | 3 ± 8 | 0 ± 0 | 0 ± 1 | 1 ± 2 | 6 ± 15 |
| | Barcelona | 1 ± 5 | 99 ± 4 | 2 ± 6 | 28 ± 33 | 32 ± 33 |

4. Discussion

4.1. Cluster Prediction

The presented cluster-prediction results highlight the importance of increasing the data population to gain statistical significance for the obtained results. The predictions from a single run or using only a given classification model could be misleading.

In order to increase the size of the data, it is advisable to analyze different archaeological samples from a given typology instead of individual samples. Obviously, for a given archaeological provenience problem, perhaps only unique samples are available, and therefore, it would be impossible to work with a set of different samples. In these cases, a possibility could be to measure different specimens from the unique sample. Another way to increase the size of data, equally important, is to apply different classification runs of the presented programmed approach, starting from different splits and using different classification models.

In the case of the presented archaeological samples, the use of several samples for every set, as well as different classification runs and models, allows the production of statistically significant results. The only clear and systematic univocal prediction from the results displayed in Table 4 is the ascription of the CM2 sample set (lead-glazed cooking ware from the Montsoriu castle) to the Barcelona cluster. This origin was already hypothesized by archaeologists [61] as a possible provenience for the three analyzed sample sets retrieved from the Montsoriu castle. However, for the other two sets from Montsoriu (CM1 and CM3), the probability percentages corresponding to the Barcelona origin are very low, and there is no other reference cluster that systematically captures a portion of probability above 60%. Instead, for CM1 and CM3, the probability is greatly split between different groups, with Quart bearing the highest percentage for CM1 and both Breda and Quart for CM3. Regarding the other two unlabeled sets, the one retrieved in Torre de la Mora (TM) shows probabilities also split into different reference clusters, although, regardless of the

model, the higher percentage always corresponds to the nearby reference cluster of Breda. In contrast, the samples retrieved in La Creueta (C) do not show any particular affinity for the nearby site of Quart. In any case, the lack of systematic univocal prediction with high percentages (>60%) would indicate that all the unlabeled sets except CM2 do not really belong to any of the reference clusters. It is worth mentioning that the presented approach always produces probability percentages, even for samples that are known for sure not to belong to any of the labeled classes. We assumed that probability percentages greatly distributed into different clusters (either two or more) indicate that the true origin of the provenanced samples is an unknown site not included within the reference labeled samples.

For the successfully provenanced set, the probability percentages associated with the Barcelona cluster are always rather close to the measured accuracies using a test set. In particular, the percentage obtained using the stack of models is very high (99%). In fact, the stack of models uses, as input, the features from the different prediction models as a discrete income (1 or 0); therefore, the predictions from the stack of models also tend to be binary, and we should therefore expect a very high percentage for the labeled class that matches the true provenance. In the case of the unsuccessfully provenanced sets, the systematic concentration of probability into a given reference set or sets (see CM1 and Quart; C and Breda or CM3 and Quart/Breda) would only indicate that the true provenance site of these sets has a certain similarity with those reference sites, but the true provenance remains unknown. In order to unveil this provenance, the reference database of labeled samples should be expanded with new classes. New classes would require a given set of samples of known provenance that could be obtained and measured.

The presented supervised approach illustrates how, in the context of a very delimited archaeological classification problem, it is possible to assess the correctness of archaeological hypotheses statistically. As for the hypotheses that are not supported by the prediction results, it is apparent that archaeological assumptions, regardless of their plausibility, are not always accurate. Theoretically, the prediction capability of the trained models is not dependent upon the chronology dissimilarity between reference and archaeological samples. However, older samples could have been produced more easily from presently unknown claypits; therefore, the corresponding reference group would be impossible to produce unless kiln samples were available. Additionally, the correctness of the prediction results is highly dependent on the quality of the database of labeled samples. The models classify according to the categories they know about, and if the models do not incorporate the variability of the features being classified, they will be biased towards class subsets. Machine learning analyses are susceptible to missing the “forest for the trees” if the data used to train the models do not include sufficient information to distinguish between archaeologically relevant classes [38]. The database should therefore include representative samples covering all the internal variability of the reference site, including preferably also raw clay samples. Finally, another factor that should be considered is the experimental procedure. In order to avoid spurious correlations, all the geochemical values (for both labeled and unlabeled samples) should ideally be obtained using the same equipment and following the same experimental protocol.

4.2. Using and Exporting the Presented Approach to Other Contexts

Often, provenance studies on archaeological pottery deal with a large number of samples, if not all, of unknown origin, and PCA is only performed to divide the samples in different groups, but the provenance remains unknown. In order to attempt provenance determination, the groups are then compared with reference samples (from kilns or clay samples) using unsupervised methods (in particular PCA and HCA). In the absence of correlation between the samples of unknown origin and the reference samples, it is assumed that the samples are imports. The supervised approach presented here is intended to be used for specific and very delimited archaeological classification problems where the samples are very likely to belong to any of the labeled classes and with the particularity that the labeled classes appear to be undistinguishable using unsupervised

methods. In order to multiply the applicability of the approach, we envisage, in an ideal future, the building of a collective and standardized geochemical database corresponding to pottery/clay of known origin. This should constitute an increasing number of labeled classes that could be selected by any archaeologist interested in the application of the presented code. In the meantime, any researcher would have to obtain their own reference database. Keeping this in mind, the approach here is already available to any researcher that would like to apply it to a constrained classification problem, similarly to the one presented in Section 2.1. In order to facilitate this task, we provided the so-called “Supervised Provenance Analysis” (SPA) code that can be freely downloaded from a GitHub repository: https://github.com/AnnaAnglisano/SPA_Supervised_Provenance_Analysis.git (accessed on 1 September 2022).

This is an R language program in R Markdown format that can be used through the freely downloadable RStudio interface. In addition to the R Markdown (rmd) files, other files (Figure 5) are available with it, along with instructions to reproduce the illustrative example presented in this paper, even without previous knowledge of R projects. The files include both the faculty of training a model using a reference database and also that of predicting classifications from a database of unlabeled samples.

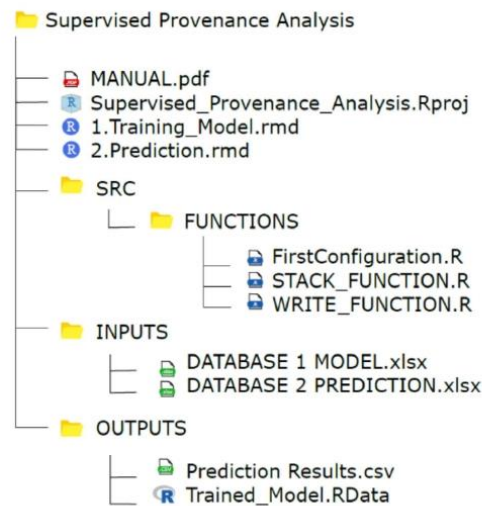


Figure 5. Complete tree of folders, subfolders and files required to perform the “Supervised Provenance Analysis”. Freely downloadable from a GitHub repository.

The downloadable materials contain a pdf file (manual.pdf) with detailed instructions on how to install and use the SPA code in all of its different options. The main document is the Supervised_Provenance_Analysis.Rproj file, which is an RStudio project file that is executed within an RStudio session and enables relative paths to read and save data as the programs within the project are executed. This allows the application to work properly on any computer regardless of the location of the files. The folder contains two R Markdown files (1.Training_Model.Rmd and 2.Prediction.Rmd) that should be executed one after the other (Figure 6). The first step tunes the models (through training and testing them), and once tuned, in the second step, the models are used to make class predictions. The use of rmd files within the RStudio project facilitates the use of the programming codes for non-specialist users through color codes. This file format combines code (actions and functions) that appear on a gray background with information or instructions on a white background. The results appear directly under the corresponding actions upon execution on a white background.

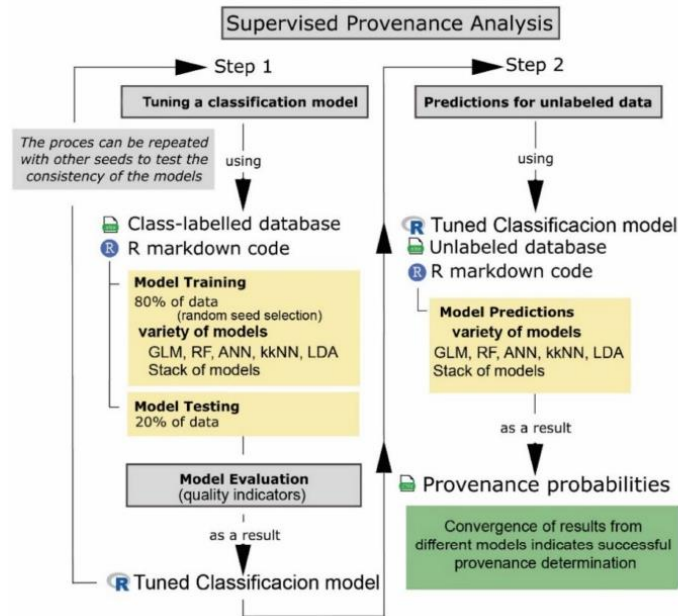


Figure 6. Schematic diagram of the two-step process (model tuning and predictions) to produce provenance probabilities for samples of unknown provenience using the R code to perform the “Supervised Provenance Analysis”.

Apart from the main files (Rproj and Rmd types), there are three folders: SRC, INPUTS and OUTPUTS. Within SRC, there are several codes that are called from the Rmd main programs. The INPUTS folder contains two Excel spreadsheet files. One of them, called here MODEL for short, is a database containing the features (geochemical values) of the class-labeled data that is used to train and test the models. The first two columns contain the numbers identifying the classes and the individual samples, respectively. The other spreadsheet, called here PREDICTION for short, contains the features of the non-class-labeled data that are intended to be classified. The data from this database can be altered by replacing them with new geochemical data to classify them within the clusters described in Section 2.1. Alternatively, and more likely, the user could change both the contents of MODEL and PREDICTION databases to fully export the approach to other contexts with completely different reference groups. Both databases should contain the same ensemble of features (i.e., the same set of analyzed chemical elements), and in any case, it is important to avoid empty fields within the geochemical values. Empty fields can be replaced by zeros, or alternatively, the whole column of the feature bearing empty fields could be deleted within the MODEL database. Regarding the OUTPUTS folder, a file named Trained_Model.Rdata appears in it after executing all the code within the 1.Training_Model.Rmd file. This newly created Rdata file overwrites any existing previous version of it within the folder, and it is required to produce cluster predictions using the 2.Prediction.Rmd file. In turn, the execution of this second rmd file (step 2 in Figure 6) creates another file within the OUTPUTS folder. This is a comma-separated file (Prediction_results.csv) containing a row for every unlabeled sample with the probabilities of belonging to any of the reference groups (numbered as within the MODEL spreadsheet) expressed on a per unit basis, and this is for every classification model (GLM = Generalised Linear Model; RF = Random Forest; NNET = Neuronal Networks; KNN = K-Nearest Neighbour; LDA = Linear Discriminant Analysis; STK = stack of models).

4.3. Contribution to Sustainable Archaeology

The presented study contributes to sustainable archaeology in various aspects. In this section, these are detailed, along with some reflections to promote sustainability in archaeology and particularly in research on the geochemical classification of pottery artifacts.

4.3.1. Free and Open-Source Software

In the current context of the continuous and unstoppable growth of software-based research, there are unsustainable practices that remain quite rooted. One of these is arguably the use of proprietary software. The use of such software often implies issues such as the impossibility of accessing the source code and the use of inaccessible file formats that make it difficult to export the obtained results; the results are often also difficult to reproduce, and the use of the software usually requires a license payment. The archaeological community is becoming increasingly aware that the use of free and open-source software (FOSS) is one of the steps toward the sustainable development of archaeology [69]. The presented classification and prediction approach was developed in R, a FOSS that is being widely used for statistical analysis, data mining and data visualization. After the release, in 2010, of RStudio [70], the usability of R greatly increased. The emergence of this integrated development environment enables archaeologists without any programming background to use R effectively. The complexity of the machine learning algorithms that were used is significant. The result of this complexity could be branded as a “black box” because their use relies on previously created classification models [38]. Archaeologists without a solid mathematic background accept their applicability to data without becoming too concerned over the mathematics and its possible limitations. However, the whole process is reproducible, and all its steps can be tracked.

In particular, the R Markdown file format that was used here is the antithesis of the black box approach of many proprietary software. In the classic black box software, after uploading the data, the parameters are set by a series of unrecordable clicks, and then, as if by magic, the results are produced. In contrast, R Markdown files are readable from RStudio or any common text editor; they contain editable R code blocks and text with instructions that guide the user through execution. Additionally, the functions are written one after the other in a chain of instructions (a so-called pipeline), which is the intuitive flow for most people new to programming. Therefore, it becomes clear the instructive implications of using R Markdown files and, in general, any kind of FOSS.

In addition to the pedagogical aspects, FOSS contributes to the economic dimension of sustainable archaeology and not only because of the lack of license payment. Reproducibility is essential to scientific progress and sustainable research. If research is non-reproducible, all the resources used to produce that research can be considered useless and wasted. The computer source code is critical for understanding and evaluating computer programs [71]. Full publication of the source code is a common demand in any scientific research involving computer codes. Reproducibility is one of the obvious criteria to assess the quality of archaeological research objectively [69]. The publication of the source code allows feedback in the form of collaborative peer review; the exchange of ideas can bring about improved, extended or customized versions of the code. A loyal open-source community results in insightful, careful and sustainable research development. The active online community of users of R is a paradigmatic example of collective development. Advanced users create contributed packages that help new users to be productive using R. Due to the large number of such packages, they were organized into lists relevant to specific areas of analysis and one of these areas is explicitly archaeology [72].

4.3.2. Open Access and Data Sharing

Publishing in open access is the best way to spread knowledge and allow that knowledge to be built upon. The move to open access is gradually changing working practices and helping the development of a sustainable future for data. The growth of data papers promotes new working practices enabling reuse and critical reassessment of primary data.

Archaeologists are often reluctant to share data, arguing that in doing so, they are giving away their research. However, their research is built upon data often gathered at the public's expense [37]. In our study, we not only present the potential of supervised classification methods but we also advocate geochemical data sharing. Our reference geochemical database with nearly 300 class-labeled analyses can be freely downloaded and used, and we prompt others to extend it to more samples and classes.

The quality of research results is highly dependent on the nature of the available data, issues of sustainability of digital data repositories, accessibility and reliability of data, standardization of data formats and management of property rights are currently widely debated [73]. In our approach, the data format is not really a big issue. Geochemical data consists of numbers commonly organized in columns that can easily be presented as a plain text table (CSV or Excel are suitable formats). However, it would be very important to monitor the quality of the data. In addition to the actual chemical analyses, a centralized repository should contain all details on the measuring conditions for every sample (technique and equipment used, measuring time, sample weight, sample format, etc.). Similar to other existing initiatives [74], a dedicated research project on a specific geochemical archaeological database for model training and cluster prediction would be helpful. However, the best way to guarantee the long-term sustainable archiving of the data would be to involve administrative authorities at either the national or transnational level. The current shared data initiatives are practices that should be followed by other data owners around the world in academia [37]. Publishing data papers and documenting good practices seem the most effective way to persuade reluctant archaeologists to share their data and the authorities to become involved in data curation.

5. Conclusions

The possibility of using supervised machine learning modeling for provenancing archaeological pottery was positively checked. A very delimited archaeological classification problem was used to illustrate this. From the five sets of pottery of unknown provenience but suspected to be produced locally in Catalonia, only the provenience of one (CM2) appears clearly defined. Therefore, the set of lead-glazed cooking ware retrieved from the Castell de Montsoriu site, dated between 1475 and 1560 CE, can be attributed to Barcelona according to the probability results of all the tested supervised models. In particular, a stack of models used repeatedly using different configurations produces an average estimate of the probability of 99% belonging to Barcelona for the CM2 set. For other provenanced sets, there is a lack of systematic univocal prediction, and we should conclude that these sets do not belong to any of the seven reference production centers.

In order to implement this supervised approach, it is essential to obtain a relatively large and reliable ensemble of reference samples for the several possible provenances of the samples that are intended to be provenanced. It is difficult to determine the minimum number of required labeled samples to define a reference group because it will depend on the homogeneity of its features. Clearly, the higher the number of reference samples, the better, although this would be an unsustainable solution. In the presented case, the reference group with fewer samples contained 33, and it is advised to use a balanced number of samples for all the reference classes.

After training and testing the classification models using the reference samples, it can be checked if the distinction between groups is feasible with an acceptable level of accuracy. In the case illustrated in this paper, all the models succeed with accuracies generally above 0.75 and relatively split-independent. Other performance indicators show the same trend. The approach could be exported to other similar classification problems following ideally the same train and test protocol to identify the useful (high accuracy) models to classify unlabeled samples (which should ideally consist of several samples for each sought classification). Classification conclusions should only be considered reliable in case of convergence of the predictions from all the high-accuracy models.

Archaeologists non-versed in statistical and machine learning techniques under R programming can use the files from the “Supervised Provenance Analysis” file folder that is made freely available with this paper and includes detailed installation and operation instructions. Interested archaeometrists without previous knowledge of R will be able to set up their own preloaded classification model using default parameters just by introducing the required databases. Additionally, experienced users could freely modify the preset parameters and adapt the codes to incorporate other classification models to suit the requirements of various classification strategies (not only necessarily based on geochemical data).

As a matter of principle, the presented approach contributes to sustainable archaeological practices as it is based on an open source R free software environment, and the user-friendly Rproject files are made freely available to any interested archaeologist. In the long term, generalized use of the presented approach and massive geochemical data sharing would result in a reduced number of analyses. Taken to the extreme, once established an exhaustive reference record for a given region, archaeologists should only analyze their samples of unknown origin without the need to provide reference samples. For the moment, our aim is to divulge the supervised approach and to help others to experiment with it.

Supplementary Materials: The following supporting information can be downloaded at: <https://www.mdpi.com/article/10.3390/su141811214/s1>: Table S1: detailed list of geological and archaeological samples that were used as reference samples to define the Barcelona reference cluster. Table S2: detailed list of the five sets of archaeological samples that were used to perform cluster predictions.

Author Contributions: Conceptualization and fieldwork—clay and shard sampling, A.A. and L.C.; experimental work—sample preparation, A.A. and R.D.F.; experimental work—petrographic analyses, A.A., L.C. and R.D.F.; experimental work—geochemical analyses, A.A. and I.Q.; code writing, A.A.; formal analyses of data, A.A. and L.C.; writing—first draft preparation, A.A. and L.C.; writing—review and editing, all authors. All authors have read and agreed to the published version of the manuscript.

Funding: IDAEA-CSIC is a Centre of Excellence Severo Ochoa (Spanish Ministry of Science and Innovation, Project CEX2018-000794-S). Article processing charges were partially supported by funds from the Grup de Recerca Aplicada al Patrimoni Cultural (GRAPAC).

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: The data presented in this study are openly available in the GitHub repository: https://github.com/AnnaAnglisano/SPA_Supervised_Provenance_Analysis.git (accessed on 1 September 2022).

Acknowledgments: We are grateful to Núria Miró and Emili Revilla (Servei d’Arqueologia de Barcelona) for providing the archaeological samples from Barcelona, to Gemma Font and Jordi Tura (Museu Etnològic del Montseny) for providing the archaeological samples from Montsoriu castle, Torre de la Mora and la Creueta. Albert Ventayol (Bac and Ventayol Geoserveis) is acknowledged for providing the geological samples from borehole cores. We are also thankful to Marc Gabasa for helping with sample preparation and experimental measurements. We are also very grateful to Marc Anglisano for helping to develop the SPA code. Finally, we would like to thank the editor as well as the reviewers for their valuable remarks and comments.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Heimann, R.; Franklin, U. Archaeo-thermometry: The assessment of firing temperatures of ancient ceramics. *J. Int. Inst. Conserv.-Can. Group* **1979**, *4*, 23–45.
2. Holakoei, P.; Tessari, U.; Verde, M.; Vaccaro, C. A new look at XRD patterns of archaeological ceramic bodies. *J. Therm. Anal. Calorim* **2014**, *118*, 165–176. [[CrossRef](#)]

3. Aitken, M.J. Dating by archaeomagnetic and thermoluminescent methods. *Philos. Trans. R. Soc. London. Ser. A Math. Phys. Sci.* **1970**, *269*, 77–88.
4. Howard, S. Understanding the concept of sustainability as applied to archaeological heritage. *Rosetta* **2013**, *14*, 1–19.
5. Carman, J. Educating for sustainability in archaeology. *Archaeologies* **2016**, *12*, 133–152. [[CrossRef](#)]
6. Reedy, C.L. *Thin-Section Petrography of Stone and Ceramic Cultural Materials*; Archetype Publications Ltd.: London, UK, 2008. ISBN 9781904982333.
7. Quinn, P.S. (Ed.) *Interpreting Silent Artefacts: Petrographic Approaches to Archaeological Ceramics*; Archaeopress Publishing Ltd.: Oxford, UK, 2009. ISBN 9781905739295.
8. Quinn, P.S. *Ceramic Petrography: The Interpretation of Archaeological Pottery & Related Artefacts in Thin Section*; Archaeopress Publishing Ltd.: Oxford, UK, 2013. ISBN 978-1-905-73959-2.
9. Neff, H. (Ed.) *Chemical Characterization of Ceramic Pastes in Archaeology*; Prehistory Press: Madison, WI, USA, 1992. ISBN 0962911062.
10. Hein, A.; Tsolakidou, A.; Iliopoulos, I.; Mommsen, H.; Garrigós, J.; Montana, G.; Kilikoglou, V. Standardisation of elemental analytical techniques applied to provenance studies of archaeological ceramics: An inter laboratory calibration study. *Analyst* **2002**, *127*, 542–553. [[CrossRef](#)]
11. Baxter, M.J. *Exploratory Multivariate Analysis in Archaeology*; Eliot Werner Publications-Inc.: Clinton Corners, NY, USA, 2015.
12. Ricca, M.; Paladini, G.; Rovella, N.; Ruffolo, S.A.; Randazzo, L.; Crupi, V.; Fazio, B.; Majolino, D.; Venuti, V.; Galli, G. Archaeometric characterisation of decorated pottery from the archaeological site of villa dei quintili (Rome, Italy): Preliminary study. *Geosciences* **2019**, *9*, 172. [[CrossRef](#)]
13. Buxeda, I.; Garrigós, J.; Ontiveros, M.A.C.; Kilikoglou, V. Chemical variability in clays and pottery from a traditional cooking pot production village: Testing assumptions in pereruela. *Archaeometry* **2003**, *45*, 1–17. [[CrossRef](#)]
14. Brorsson, T.; Blank, M.; Fridén, I.B. Mobility and exchange in the middle neolithic: Provenance studies of pitted ware and funnel beaker pottery from Jutland, Denmark and the West Coast of Sweden. *J. Archaeol. Sci. Rep.* **2018**, *20*, 662–674. [[CrossRef](#)]
15. Papachristodoulou, C.; Oikonomou, A.; Ioannides, K.; Gravani, K. A study of ancient pottery by means of X-ray fluorescence spectroscopy, multivariate statistics and mineralogical analysis. *Anal. Chim. Acta* **2006**, *573–574*, 347–353. [[CrossRef](#)]
16. Aquilia, E.; Barone, G.; Mazzoleni, P.; Ingoglia, C. Petrographic and chemical characterisation of fine ware from three archaic and hellenistic kilns in gela, sicily. *J. Cult. Herit.* **2012**, *13*, 442–447. [[CrossRef](#)]
17. Munita, C.S.; Paiva, R.P.; Alves, M.A.; de Oliveira, P.M.S.; Momose, E.F. Provenance study of archaeological ceramic. *J. Trace Microprobe Tech.* **2003**, *21*, 697–706. [[CrossRef](#)]
18. Scarpelli, R.; Robustelli, G.; Clark, R.J.H.; Francesco, A.M.D. Scientific investigations on the provenance of the black glazed pottery from Pompeii: A case study. *Mediterr. Archaeol. Archaeom.* **2017**, *17*, 1–10.
19. Buxeda i Garrigós, J.; Kilikoglou, V.; Day, P.M. Chemical and mineralogical alteration of ceramics from a late bronze age kiln at Kommos, Crete: The effect on the formation of a reference group. *Archaeometry* **2001**, *43*, 349–371. [[CrossRef](#)]
20. Maritan, L.; Holakooei, P.; Mazzoli, C. Cluster analysis of XRPD data in ancient ceramics: What for? *Appl. Clay Sci.* **2015**, *114*, 540–549. [[CrossRef](#)]
21. Medeghini, L.; Mignardi, S.; Vito, C.D.; Conte, A.M. Evaluation of a FTIR data pretreatment method for principal component analysis applied to archaeological ceramics. *Microchem. J.* **2016**, *125*, 224–229. [[CrossRef](#)]
22. Parisotto, S.; Leone, N.; Schönlieb, C.-B.; Launaro, A. Unsupervised clustering of Roman potsherds via variational autoencoders. *J. Archaeol. Sci.* **2022**, *142*, 105598. [[CrossRef](#)]
23. Bratitsi, M.; Liritzis, I.; Vafiadou, A.; Xanthopoulou, V.; Palamara, E.; Iliopoulos, I.; Zacharias, N. Critical assessment of chromatic index in archaeological ceramics by Munsell and RGB: Novel contribution to characterization and provenance studies. *Mediterr. Archaeol. Archaeom.* **2018**, *18*, 175–212.
24. Visiedo, J.P.; Madrid i Fernández, M.; Buxeda i Garrigós, J. The case of black and green tin glazed pottery from Barcelona between 13th and 14th century: Analysing its production and its decorations. *J. Archaeol. Sci. Rep.* **2021**, *38*, 103100. [[CrossRef](#)]
25. Calparsoro, E.; Arana, G.; Iñáñez, J.G. Pottery from orduña village in the 17th–19th centuries: An archaeometrical approach. *J. Archaeol. Sci. Rep.* **2019**, *23*, 304–323. [[CrossRef](#)]
26. Baklouti, S.; Maritan, L.; Ouazaa, N.L.; Casas, L.; Joron, J.-L.; Kassaa, S.L.; Moutte, J. Provenance and reference groups of African Red Slip ware based on statistical analysis of chemical data and REE. *J. Archaeol. Sci.* **2014**, *50*, 524–538. [[CrossRef](#)]
27. Mackensen, M.; Schneider, G. Production centres of African red slip ware (2nd–3rd c.) in northern and central Tunisia: Archaeological provenance and reference groups based on chemical analysis. *J. Rom. Archaeol.* **2006**, *19*, 163–190. [[CrossRef](#)]
28. Monette, Y.; Richer-LaFlèche, M.; Moussette, M.; Dufournier, D. Compositional analysis of local redwares: Characterizing the pottery productions of 16 workshops located in southern québec dating from late 17th to late 19th-century. *J. Archaeol. Sci.* **2007**, *34*, 123–140. [[CrossRef](#)]
29. Montana, G.; Randazzo, L.; Tsantini, E.; Fourmont, M. Ceramic production at Selinunte (Sicily) during the 4th and 3rd century BCE: New archaeometric data through the analysis of kiln wastes. *J. Archaeol. Sci. Rep.* **2018**, *22*, 154–167. [[CrossRef](#)]
30. Maritan, L.; Gravagna, E.; Cavazzini, G.; Zerboni, A.; Mazzoli, C.; Grifa, C.; Mercurio, M.; Mohamed, A.A.; Usai, D.; Salvatori, S. Nile river clayey materials in Sudan: Chemical and isotope analysis as reference data for ancient pottery provenance studies. *Quat. Int.* **2021**, *in press*. [[CrossRef](#)]

31. Baklouti, S.; Maritan, L.; Casas, L.; Ouazaa, N.L.; Järrega, R.; Prevosti, M.; Mazzoli, C.; Fouzai, B.; Kassaa, S.L.; Fantar, M. Establishing a new reference group of keay 25.2 amphorae from Sidi Zahrani (Nabeul, Tunisia). *Appl. Clay Sci.* **2016**, *132*–133, 140–154. [CrossRef]
32. Montana, G.; Ontiveros, M.Á.C.; Polito, A.M.; Azzaro, E. Characterisation of clayey raw materials for ceramic manufacture in ancient sicily. *Appl. Clay Sci.* **2011**, *53*, 476–488. [CrossRef]
33. Gutsuz, P.; Kibaroglu, M.; Sunal, G.; Hacrosmanoğlu, S. Geochemical characterization of clay deposits in the Amuq Valley (Southern Turkey) and the implications for archaeometric study of ancient ceramics. *Appl. Clay Sci.* **2017**, *141*, 316–333. [CrossRef]
34. Efenberger-Szmechtyk, M.; Nowak, A.; Kregiel, D. Implementation of chemometrics in quality evaluation of food and beverages. *Crit. Rev. Food Sci. Nutr.* **2018**, *58*, 1747–1766. [CrossRef]
35. Anglisano, A.; Casas, L.; Anglisano, M.; Queral, I. Application of supervised machine-learning methods for attesting provenance in Catalan traditional pottery industry. *Minerals* **2020**, *10*, 8. [CrossRef]
36. Fiorucci, M.; Khoroshiltseva, M.; Pontil, M.; Traviglia, A.; Bue, A.D.; James, S. Machine learning for cultural heritage: A survey. *Pattern Recognit. Lett.* **2020**, *133*, 102–108. [CrossRef]
37. McKeague, P.; van't Veer, R.; Huvila, I.; Moreau, A.; Verhagen, P.; Bernard, L.; Cooper, A.; Green, C.; van Manen, N. Mapping our heritage: Towards a sustainable future for digital spatial information and technologies in European archaeological heritage management. *J. Comput. Appl. Archaeol.* **2019**, *2*, 89–104. [CrossRef]
38. Bickler, S.H. Machine learning arrives in archaeology. *Adv. Archaeol. Pract.* **2021**, *9*, 186–191. [CrossRef]
39. Cheng, G.; Han, J. A survey on object detection in optical remote sensing images. *ISPRS J. Photogramm. Remote Sens.* **2016**, *117*, 11–28. [CrossRef]
40. Resler, A.; Yeshurun, R.; Natalio, F.; Giryas, R. A deep-learning model for predictive archaeology and archaeological community detection. *Humanit. Soc. Sci. Commun.* **2021**, *8*, 295. [CrossRef]
41. Navarro, P.; Cintas, C.; Lucena, M.; Fuertes, J.M.; Delrieux, C.; Molinos, M. Learning feature representation of iberian ceramics with automatic classification models. *J. Cult. Herit.* **2021**, *48*, 65–73. [CrossRef]
42. Wilczek, J.; Monna, F.; Navarro, N.; Chateau-Smith, C. A computer tool to identify best matches for pottery fragments. *J. Archaeol. Sci. Rep.* **2021**, *37*, 102891. [CrossRef]
43. Derech, N.; Tal, A.; Shimshoni, I. Solving archaeological puzzles. *Pattern Recognit.* **2021**, *119*, 108065. [CrossRef]
44. Chetouani, A.; Treuillet, S.; Exbrayat, M.; Jesset, S. Classification of engraved pottery sherds mixing deep-learning features by compact bilinear pooling. *Pattern Recognit. Lett.* **2020**, *131*, 1–7. [CrossRef]
45. Domínguez-Rodrigo, M.; Cifuentes-Alcobendas, G.; Jiménez-García, B.; Abellán, N.; Pizarro-Monzo, M.; Organista, E.; Baquedano, E. Artificial intelligence provides greater accuracy in the classification of modern and ancient bone surface modifications. *Sci. Rep.* **2020**, *10*, 18862. [CrossRef]
46. Oonk, S.; Spijker, J. A supervised machine-learning approach towards geochemical predictive modelling in archaeology. *J. Archaeol. Sci.* **2015**, *59*, 80–88. [CrossRef]
47. Barone, G.; Mazzoleni, P.; Spagnolo, G.V.; Raneri, S. Artificial neural network for the provenance study of archaeological ceramics using clay sediment database. *J. Cult. Herit.* **2019**, *38*, 147–157. [CrossRef]
48. Lopez-García, P.A.; Argote, D.L.; Thrun, M.C. Projection-based classification of chemical groups for provenance analysis of archaeological materials. *IEEE Access* **2020**, *8*, 152439–152451. [CrossRef]
49. Ma, Q.; Yan, A.; Hu, Z.; Li, Z.; Fan, B. Principal component analysis and artificial neural networks applied to the classification of Chinese pottery of neolithic age. *Anal. Chim. Acta* **2000**, *406*, 247–256. [CrossRef]
50. Díez-Pastor, J.F.; Jorge-Villar, S.E.; Arnaiz-González, Á.; García-Osorio, C.I.; Díaz-Acha, Y.; Campeny, M.; Bosch, J.; Melgarejo, J.C. Machine learning algorithms applied to Raman spectra for the identification of variscite originating from the mining complex of Gavà. *J. Raman Spectrosc.* **2020**, *51*, 1563–1574. [CrossRef]
51. Salazar, A.; Safont, G.; Vergara, L.; Vidal, E. Pattern recognition techniques for provenance classification of archaeological ceramics using ultrasounds. *Pattern Recognit. Lett.* **2020**, *135*, 441–450. [CrossRef]
52. Buxeda, J.; Iñáñez, J.; Madrid, M.; Beltrán, J. La ceràmica de Barcelona. Organització i producció entre els segles XIII i XVIII a través de la seva caracterització arqueomètrica. *Quarhis* **2011**, *7*, 192–207.
53. Serra, J. Ceràmica de rebuig al carrer d'Avinyó. Un possible nou taller barceloní en el primer quart del segle XIII. *Quad. D'arqueologia Història Ciutat Barcelona. Quarhis* **2016**, *12*, 194–209.
54. Miró, N. Excavació de les voltes de la sala de reserva de la biblioteca de Catalunya, antic hospital de la Santa Creu, Barcelona (el Barcelonès). In *15 Anys D'Intervencions Arqueològiques: Mancanes i Resultats, Proceedings of 1r Congrés d'Arqueologia Medieval i Moderna a Catalunya, Igualada, Spain, 13–15 November 1998*; Associació Catalana per a la Recerca en Arqueologia Medieval: Barcelona, Spain, 2000; pp. 168–176. Available online: <https://dialnet.unirioja.es/servlet/libro?codigo=782515> (accessed on 30 July 2022).
55. Nebot, N. La botiga de Josep Barba: Un terrisser a la Barcelona del segle XVIII. *Quad. D'arqueologia Història Ciutat Barcelona. Quarhis* **2015**, *11*, 184–199.
56. Madrid, M.; Marcos, C.F.D.; Barrachina, C.P.; Heredia, J.B.D.; Escribano-Ruiz, S.; Ibáñez, J.G.; Ferrer, S.G.; Febo, R.D.; Amores, F.D.; Buxeda, J. Ceràmica, tecnologia i transferències. Els centres productius del projecte tecnològic. *Quad. D'arqueologia Història Ciutat Barcelona. Quarhis* **2017**, *13*, 16–67.

57. Caixal, A.; Fierro, X.; López, A. Resultats de l'excavació arqueològica en la galeria alta del pati Manning de l'antiga Casa de Caritat. In *Actuacions en el Patrimoni Edificat Medieval i Modern (Segles X al XVIII) = Actuaciones en el Patrimonio Edificado Medieval y Moderno (Siglos X al XVIII)*; Servei del Oatrimoni Arquitectònic: Barcelona, Spain, 1991; pp. 13–15.
58. Oriol, J. *Memòria de la Intervenció Arqueològica a Pia Almoïna, Barcelona*; Generalitat de Catalunya: Barcelona, Spain, 1993.
59. Miró, N. *Memòria de la Intervenció Realitzada als Carrers de l'Argenteria i Manresa de Barcelona (Barcelonès)*; Ajuntament de Barcelona: Barcelona, Spain, 1997.
60. Font, G.; Mateu, J.; Pujadas, S.; Tura, J.; Llorens, J.M. Montsoriu al Segle XVI. Testimonis Arqueològics de L'abandonament d'un Gran Castell. *Tribuna D'arqueologia* 2011–2012. 2014, pp. 244–263. Available online: <http://calaix.gencat.cat/handle/10687/91795#page=1> (accessed on 30 July 2022).
61. Tura, J.; Font, G.; Pujadas, S.; Mateu, J.; Llorens, J.M. El conjunt arqueològic del segle XVI localitzat a la cisterna est del castell de Montsoriu. *Rodis J. Mediev. Post-Mediev. Archaeol.* **2022**, 25–46.
62. Tura, J.; Mateu, J. Torre de la Mora o del Far (Sant Feliu de Buixalleu, la Selva): Una ocupació alt-medieval al Montseny. In *Fars de L'islam: Antiques Alimares d'al-Andalus, Proceedings of the Jornades Científiques Ocorde*; Barcelona, Spain, 9–10 November 2006; Martí, R., Ed.; Edar Press: Barcelona, Spain, 2008; pp. 139–154. Available online: <https://cataleg.parc.s.diba.cat/cgi-bin/koha/opac-detail.pl?biblionumber=10093> (accessed on 30 July 2022).
63. Pericot y García, L.; Corominas Planellas, J.M.; Oliva Prat, M.; Riuró Ilapat, F.; Padrol Salellas, P. *La Labor de La Comisaria Provincial de Excavaciones Arqueológicas de Gerona. Informes y Memorias*; Ministerio de educación nacional. Comisaria general de excavaciones arqueológicas: Madrid, Spain, 1952; Volume 27.
64. Zhao, Y. R and Data Mining. In *R and Data Minig*; Zhao, Y., Ed.; Academic Press: Cambridge, MA, USA, 2013; Chapter 5; pp. 41–50. ISBN 978-0-12-396963-7.
65. Praveena, M.; Jaiganesh, V. A literature review on supervised machine learning algorithms and boosting process. *Int. J. Comput. Appl.* **2017**, 169, 32–35. [[CrossRef](#)]
66. Kotsiantis, S.B. Supervised machine learning: A review of classification techniques. In *Proceedings of the 2007 Conference on Emerging Artificial Intelligence Applications in Computer Engineering: Real Word AI Systems with Applications in EHealth HCl, Information Retrieval and Pervasive Technologies*, Amsterdam, The Netherlands, 10 June 2007; IOS Press: Amsterdam, The Netherlands, 2007; pp. 3–24.
67. Kuhn, M. Building predictive models in r using the caret package. *J. Stat. Softw. Artic.* **2008**, 28, 1–26. [[CrossRef](#)]
68. Grandini, M.; Bagli, E.; Visani, G. Metrics for multi-class classification: An overview. *arXiv* **2020**, arXiv:2008.05756.
69. Bibby, D.; Ducke, B. Free and open source software development in archaeology. Two interrelated case studies: GvSIG CE and survey2GIS. *Internet Archaeol.* **2017**, 43. [[CrossRef](#)]
70. Van der Loo, M.P.J.; de Jonge, E. *Learning RStudio for R Statistical Computing*; Packt publishing: Birmingham, UK, 2012. ISBN 1782160604.
71. Morin, A.; Urban, J.; Adams, P.D.; Foster, I.; Sali, A.; Baker, D.; Sliz, P. Research priorities. shining light into black boxes. *Science* **2012**, 336, 159–160. [[CrossRef](#)]
72. Marwick, B. CRAN Task View: Archaeological Science. Available online: <https://github.com/benmarwick/ctv-archaeology> (accessed on 29 August 2022).
73. Kintigh, K. The promise and challenge of archaeological data integration. *Am. Antiq.* **2006**, 71, 567–578. [[CrossRef](#)]
74. Derudas, P.; Dell'Unto, N.; Callieri, M.; Apel, J. Sharing archaeological knowledge: The interactive reporting system. *J. Field Archaeol.* **2021**, 46, 303–315. [[CrossRef](#)]

Resultats i Discussió

Fins al moment actual la recerca endegada en aquesta tesi doctoral ha permès publicar els resultats obtinguts més rellevants en dos articles en revistes amb factor d'impacte important en els àmbits de la mineralogia, la mineria i el processament de minerals, i els estudis i ciències ambientals. També s'han anat presentant resultats a congressos internacionals com European Meeting on Ancient Ceramics (2017, Bordeus i 2019, Barcelona). Malgrat tot, hi ha molta feina que, ja sigui perquè els resultats no foren els esperats o perquè encara està en fase de redacció per donar-hi format d'article, encara no s'han publicat. En aquest apartat es presenten els resultats més destacats que s'han anat obtenint ordenats de manera cronològica per tal de mostrar el procés a través del qual s'ha arribat a poder assolir els objectius plantejats en aquest projecte de recerca.

Exploració de diverses aproximacions analítiques per caracteritzar i distingir els diferents grups de mostres corresponents a diverses localitats

La primera aproximació analítica que es va provar va ser la més pròpiament geològica: l'estudi amb microscopi petrogràfic sobre mostres preparades en forma de làmina prima complementat amb l'estudi mineralògic per difracció de raigs X de pols. Això no obstant, durant el Treball de Final de Grau que va precedir aquesta tesi, es van explorar també altres tècniques complementàries a l'estudi mineralògic, com ara calcimetries per poder quantificar el CaCO_3 present a les mostres i l'espectrometria Mössbauer per a caracteritzar l'estat d'oxidació del ferro. Globalment, aquesta aproximació va portar a la conclusió que un estudi petrogràfic o mineralògic no és eficient a l'hora de distingir les diverses procedències dels materials que han estat objecte d'estudi. Les característiques texturals de la ceràmica depenen fortament de la granoselecció i la mescla d'argiles que hagi efectuat el mestre terrisser en cada producció concreta i pel que fa a la mineralogia sovint trobem els mateixos components en argiles i ceràmiques d'indrets diferents.

Per il·lustrar la dificultat de l'anàlisi petrogràfica per a la identificació de procedència, a la figura 7a es pot veure com dues mostres de ceràmica (Q21 i Q18) produïdes a Quart presenten granulometries completament diferents que les fan aparèixer diferents fins i tot a ull nu. Observades en làmina prima (Fig. 7b) mostres de ceràmica de localitats diferents (Quart i Sant Julià de Vilatorça) però amb una granulometria similar acostumen a presentar el mateix conjunt de minerals principals (molt sovint quars i feldespats) de manera que les descripcions petrogràfiques acaben essent molt similars malgrat que les mostres vinguin de localitats diferents.

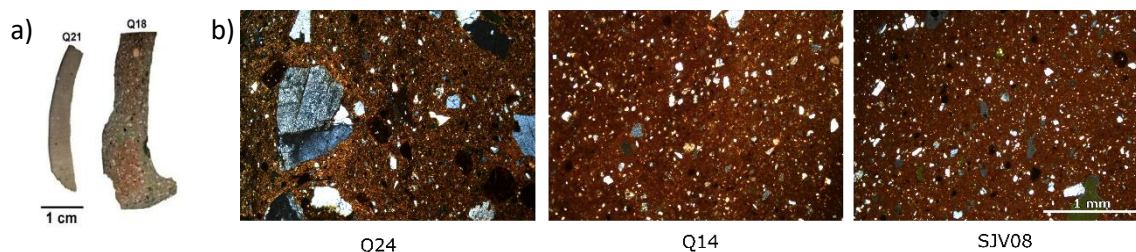


Figura 7: a) dues làmines primes de ceràmica local de Quart escanejades. b) fotografies de tres làmines primes de ceràmica local de Quart (Q24 i Q14) i Sant Julià de Vilatorça (SJV08).

Els estudis de la mineralogia mitjançant difracció de raigs X també poden mostrar resultats molt similars per a mostres de procedències diferents tal com es posa de manifest a la figura 8.

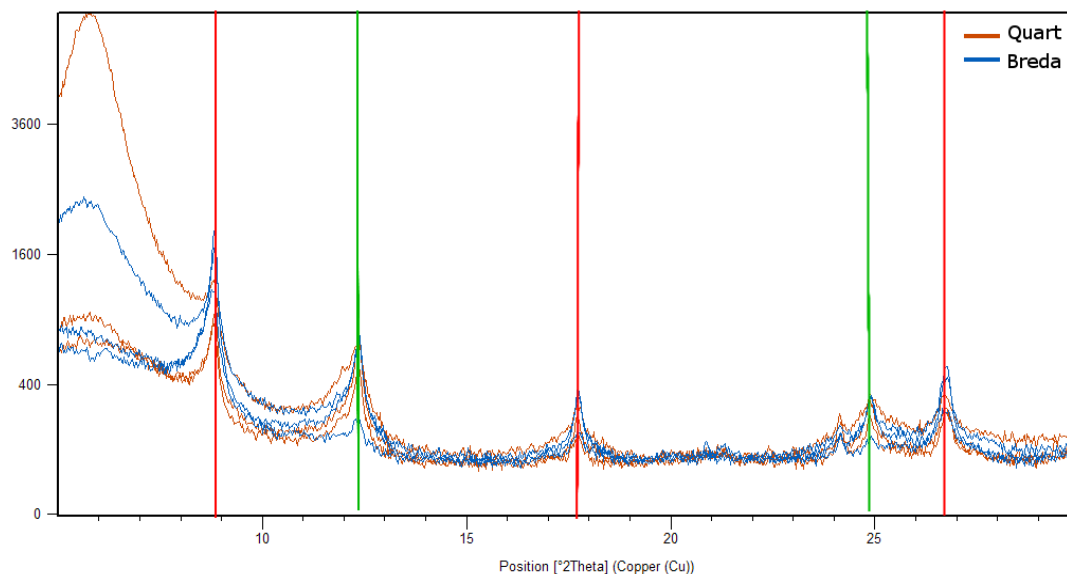


Figura 8: Quatre difractogrames gairebé idèntics de mostres d'argila, dues procedents de Breda (blau) i dues de Quart (carabassa). S'han ressaltat els pics corresponents a fil·losilicats del grup de la mica – il·lita (vermell); i del grup de la caolinita – serpentina (en verd).

L'ús de l'aproximació petrogràfica i mineralògica pot ser perfectament vàlida per a distingir produccions en determinats contextos. En particular, la presència de determinats minerals o litologies poc habituals en les inclusions no plàstiques de ceràmiques pot ser diagnòstica de procedència. Fins i tot la morfologia d'inclusions de mineralogies molt freqüents pot arribar a ser un indicador de procedència. Un exemple clar d'això és la presència de quars eòlic com a indicador de produccions tunisianes [34,35]. Ara bé, en el context d'aquesta tesi es pot concloure que l'aproximació petrogràfica i mineralògica no és útil i això en bona part és degut al context geològic similar de a totes les produccions estudiades.

Descartada l'aproximació petrogràfica i mineralògica, la geoquímica va anar guanyant pes en el projecte. Ara bé, la distinció de produccions en base a l'anàlisi química tampoc és d'entrada evident, el nombre d'elements analitzats és elevat i genera grans volums d'informació difícil de processar. Cal tenir present que s'han analitzat 21 elements químics per a gairebé 300 mostres de referència, generant més de 6000 valors numèrics a tenir en compte.

A més, l'anàlisi química té la dificultat afegida que conjunts mineralògics diferents poden arribar a donar lloc, almenys teòricament, a resultats idèntics d'anàlisi química. Els mètodes habituals de definició de grups (o si més no d'establiment de diferències entre mostres) tampoc van permetre la distinció entre les diverses procedències. Des dels gràfics en què es comparen parelles d'elements fins als més elaborats mètodes d'anàlisi no supervisat (PCA i HCA), tots ells fracassen (Fig. 9). Això és degut a que la variabilitat de concentracions dins d'una determinada producció és major o similar a les diferències que hi ha entre mostres de diferents procedències. Per tant, els intents de mostrar gràficament diferències sempre donen com a resultat gràfics en què els núvols de punts formats per mostres de diferents procedències se solapen completament o en bona mesura.

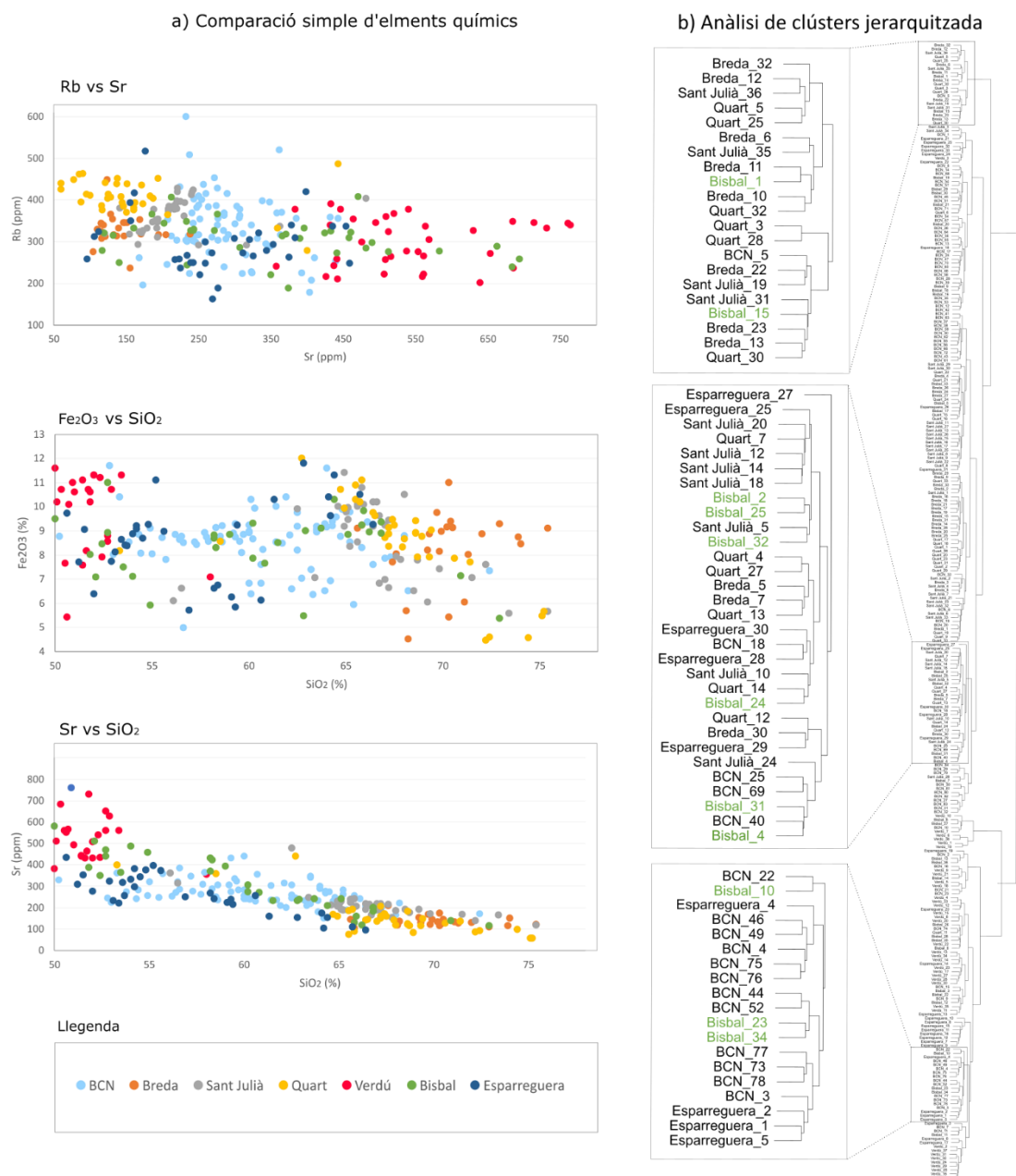


Figura 9: a) Comparació a través de parelles d'elements químics de les mostres dels grups de referència on es posa de manifest l'elevat grau de solapament. b) Gràfic HCA on s'observa que les mostres d'un dels grups (la Bisbal) no queden agrupades sinó repartides en zones allunyades del dendrograma resultant.

Definició i distinció dels grups de referència mitjançant mètodes supervisats

No hi havia cap garantia que la resolució del problema plantejat en la tesi tingués solució. Afortunadament però, els mètodes supervisats d'agrupament de dades geoquímiques sí que han permès assolir la pretesa distinció.

La superioritat dels mètodes supervisats està relacionada amb el fet que la cerca d'una solució està orientada a través del procés d'entrenament de dades. Així, l'algorisme subjacent al model

va provant diverses configuracions i pot anar-se perfeccionant automàticament per aconseguir assignar a totes les mostres, la seva procedència. En aquest procés de perfeccionament es diu que el model aprèn a fer allò que se li demana i si el procés és efectiu, després és capaç de determinar correctament la procedència de mostres alienes al procés d'entrenament. La manera com el model determina la procedència depèn de com està dissenyat el seu algoritme però sempre té l'avantatge que aprofita totes les dades disponibles i que pot establir condicions molt elaborades difícilment imaginables sense el suport dels algorismes.

Cal destacar que tots els models utilitzats, després de passar pel procés d'entrenament, han tingut taxes d'èxit notables valorades ja a partir de la taxa d'encerts en la predicció de la procedència de les mostres del grup de control (el 20% de les dades de referència). Tots els models presenten valors mitjans molt elevats d'exactitud i de tots els altres paràmetres d'anàlisi d'errors i encerts. Així, els percentatges d'encert se situen estadísticament al voltant del 80% per a tots els models excepte per al model kNN, amb valors mitjans lleugerament inferiors (~75%). Els valors d'exactitud obtinguts al primer article (publicat a Minerals MDPI) són lleugerament superiors als obtinguts en el segon article (publicat a Sustainability MDPI). Aquesta pèrdua d'exactitud (d'aproximadament un 3%) pot estar relacionada amb la incorporació d'un nou grup de referència (Barcelona) i el consegüent increment de complexitat que això comporta. Així doncs, una de les limitacions que podria tenir la metodologia explorada és el nombre de grups de referència que es poden arribar a gestionar conjuntament.

Utilització de mètodes supervisats per a la classificació de mostres de ceràmica arqueològica com a eina per a establir-ne la provenença.

Els bons indicadors de la capacitat dels models degudament entrenats d'assignar correctament la provenença de les mostres van donar confiança per afrontar el següent objectiu. El següent pas ha estat aplicar aquesta capacitat a mostres d'origen realment desconegut. Malgrat tot, per poder demostrar aquesta capacitat no n'hi ha prou d'analitzar qualsevol tipus de mostres arqueològiques i aplicar els models de classificació. La classificació només té possibilitats d'èxit si es compleix l'assumpció que les mostres que es volen classificar pertanyen realment a alguna de les classes (és a dir, localitats) definides a partir de les mostres de referència. Cal fer una selecció molt acurada de les mostres per tal que això passi. En aquest sentit, la metodologia presentada no pretén substituir cap de les tècniques utilitzades per caracteritzar, arqueològicament i arqueomètricament, les restes de ceràmica. De fet és l'anàlisi del conjunt d'informacions prèvies, com ara el context arqueològic, la cronologia, la tipologia o l'aspecte de les mostres el que permet tenir majors possibilitats de satisfer l'assumpció.

Els resultats presentats a l'article publicat a Sustainability MDPI mostren que dels cinc conjunts de mostres arqueològiques, només un procedeix clarament d'una de les classes (concretament Barcelona). Amb aquest estudi de cas, també s'ha posat de manifest que per tal d'aplicar els models de classificació a ceràmiques arqueològiques cal prendre diverses precaucions. Cal que les mostres arqueològiques s'analitzin exactament seguint el mateix protocol analític que les mostres de referència (el mateix tipus de preparació de mostra, el mateix conjunt d'elements químics analitzats i si és possible, el mateix aparell per fer les mesures). També convé que el conjunt de mostres arqueològiques siguin representatives de la tipologia ceràmica sobre la que es pretén determinar la provenença. Com que l'aproximació que s'utilitza és de tipus estadístic, cal cercar la provenença d'un conjunt de mostres, si és possible amb un grup d'almenys 7 mostres

i interpretar la validesa dels resultats també amb criteri estadístic. L'estadística en els resultats no només prové del fet que s'analitza un conjunt de mostres sinó també de l'ús de diversos models de classificació i de l'ús de diversos conjunts d'entrenament (és a dir diverses particions 80/20 generades aleatòriament). Finalment cal tenir en compte que els models sempre assignaran una classe de referència (és a dir localitats) a qualsevol mostra que introduïm, no tenen la capacitat d'identificar mostres que no pertanyen a cap de les classes de referència.

En el cas d'estudi presentat es conclou que s'ha pogut determinar la provenença d'un dels conjunts arqueològics (la ceràmica de cuina amb vidrat de plom procedent del jaciment arqueològic del Castell de Montsoriu) perquè tots els models apunten estadísticament (amb percentatges elevats, de l'ordre del 70% o superiors) a assignar la classe 'Barcelona' a aquest conjunt de mostres. En els altres conjunts analitzats no s'observa una tendència tan sistemàtica a assignar estadísticament una determinada classe a les mostres. Els corresponents percentatges d'assignació de classe són molt més repartits entre les diverses classes i varien força en funció de quin model de classificació s'apliqui. Així doncs, quan no s'observa una tendència clara s'ha de concloure que les mostres no provenen de cap de les localitats de referència. És destacable que malgrat haver seleccionat 5 conjunts de mostres arqueològiques procedents de jaciments propers a alguna de les localitats de referència i amb hipòtesis arqueològiques que apuntaven a procedències en alguna de les localitats de referència, la metodologia només ha pogut confirmar una de les hipòtesis arqueològiques. Així, per exemple, la provenença del conjunt de mostres de ceràmica vidrada acolorida en color verd gòtic (també trobada al Castell de Montsoriu) no s'ha pogut acreditar malgrat que era el conjunt per al qual arqueològicament s'apuntava amb més seguretat a un origen de Barcelona. Això indicaria que l'origen d'aquestes mostres és una localitat no recollida entre les classes de referència, o bé que les mostres que conformen la classe barcelonina no són completament representatives d'aquesta localitat.

Supervised Provenance Analysis. Una nova eina per a estudis de provenença de restes arqueològiques.

Finalment, un altre resultat de la tesi ha estat posar a disposició de la comunitat científica un codi preparat per a reproduir de forma relativament senzilla el tipus d'aproximació supervisada presentat en els dos articles que s'han publicat. Es pretén que els usuaris puguin aprofitar la base de dades creada i en puguin generar de noves amb altres localitats de referència. És a dir, que puguin utilitzar els models supervisats que s'han presentat però que també puguin afegir-ne de nous. Tot plegat amb l'objectiu final que en els estudis de determinació de procedència de ceràmica arqueològica s'hi afegixi la metodologia presentada com a eina addicional, una eina que permet obtenir resultats allà on altres aproximacions fracassen.

En programació existeixen diferents tipus de llenguatge, en aquest cas s'ha utilitzat el llenguatge R. Aquest, igual que la resta, no és fàcil ni intuïtiu per a un profà en programació i pot arribar a resultar complex de fer servir. Es requereix experiència i dedicació per poder generar codi, ni que sigui només per a enllaçar codis d'altri o instal·lar i explotar les funcionalitat de les diverses llibreries d'R. En aquest sentit, s'ha generat un document de projecte R (.Rproj) que s'utilitza a partir d'un codi escrit en un format més senzill, RMarkdown (.Rmd), aquest combina fragments de text amb instruccions i informació per a l'usuari entre els fragments de codi R. Tots els canvis que s'hagin de fer en el codi per fer-lo servir en les seves diverses opcions estan explícitament

indicats a les instruccions, tot i que els canvis sempre són opcionals. Les úniques modificacions que són veritablement necessàries són la modificació de les bases de dades geoquímiques (del conjunt de mostres de referència i del conjunt de mostres arqueològiques sobre el qual es vulgui determinar la provenença). Aquestes dades es poden modificar simplement canviant els corresponents fulls Excel que serveixen d'input. La lectura dels resultats també es pot fer a través d'un fitxer Excel que es genera automàticament amb l'execució del codi. També s'ha posat a disposició de qui ho necessiti un manual per a facilitar la instal·lació del programari lliure R i tot allò que es necessita per a la correcta execució del projecte R.

El potencial d'aplicació d'una eina d'aquests tipus és molt ampli però requereix acotar bé el problema de determinació de provenença, és a dir, cal tenir clar quines són les localitats d'origen plausibles i que sigui viable fabricar-ne una base sòlida de dades de referència. Amb aquestes condicions, la metodologia es pot aplicar a qualsevol lloc del món, a mostres de cronologies diverses. Per exemple, a ceràmica subactual amb finalitats d'autenticació i a ceràmica antiga amb finalitats de caire arqueològic. A més, l'objecte d'estudi no s'ha de limitar estrictament a mostres de ceràmica i podria obrir-se a altres materials.

En definitiva, el recorregut que s'ha fet al llarg d'aquest doctorat ha donat prou bons resultats. Els objectius plantejats s'han assolit i amb ells es pretén que altres segueixin i puguin aprofitar la feina feta. L'ampli ventall de possibilitats que s'obre seria inabastable en un únic projecte de tesi, però retrospectivament es pot pensar que han quedat pendents algunes investigacions complementàries a allò que s'ha presentat que sí que s'hi podrien haver inclòs, o més aviat que hauria estat desitjable haver-les abordat. Així, en el marc de les produccions ceràmiques subactuals, objecte d'estudi en l'article publicat a Minerals MDPI, la caracterització hauria d'ampliar-se també, com a mínim, a Miravet (Riviera d'Ebre) amb tallers encara actius avui dia. Si no s'ha inclòs en aquest estudi ha estat per motius de proximitat i manca de recursos. Es va prioritzar abordar el repte de la distinció entre centres productors molt propers i amb argiles de contextos geològics similars. D'altra banda, en el marc de l'estudi de conjunts de mostres arqueològiques i al constatar que dels cinc estudiats només la provenença d'un ha quedat ben determinat, hauria calgut cercar altres conjunts i sotmetre'ls al mateix estudi. Finalment, també hauria estat interessant demostrar que la metodologia és exportable a altres tipus de materials, confirmació que arribarà ben aviat en forma d'article.

Conclusions finals

- Les anàlisis petrogràfica i mineralògica no són eficients a l'hora de relacionar produccions ceràmiques i sediments argilosos amb les seves corresponents zones d'origen, si aquestes zones d'origen comparteixen contextos geològics similars.
- Els mètodes estadístics de tipus no supervisat (PCA i HCA) sobre conjunts de dades geoquímiques de ceràmiques i sediments argilosos no són eficients a l'hora distingir diversos grups de mostres corresponents a les diverses zones d'origen. Aquesta conclusió hauria de posar en alerta els arqueòlegs que utilitzen aquests mètodes per a definir grups de referència.
- Els mètodes estadístics de tipus supervisat, a través d'un procés d'entrenament automàtic de dades, permeten assignar correctament l'origen a mostres, les dades de les quals no han format part del procés d'entrenament (amb percentatges d'encert de l'ordre del 80%).
- La metodologia presentada, amb la definició de grups de referència mitjançant l'ús dels models de tipus supervisat, té el potencial d'aplicar-se per a l'acreditació de l'ús d'argiles locals en la indústria de la ceràmica tradicional. Atès que les empreses de ceràmica ja han de fer anàlisis químiques periòdiques per temes de control de qualitat i salut, l'aplicació d'aquesta metodologia no hauria de representar un sobrecost per a la indústria.
- S'ha pogut aplicar i demostrar l'efectivitat dels models supervisats per a estudis de provenença a través d'un cas d'estudi amb cinc conjunts de ceràmica arqueològica.
- S'ha corroborat amb una rigorosa aproximació arqueomètrica la hipòtesi dels arqueòlegs sobre la provenença d'un conjunt de ceràmica de cuina amb vidrat de plom procedent del Castell de Montsoriu i datada entre els anys 1475-1560 DC. Es pot afirmar amb molta fiabilitat que es tracta d'una producció local de Barcelona. En canvi les hipòtesis de provenença dels altres quatre conjunts queden en entredit.
- Per tal d'utilitzar la metodologia presentada, amb l'ús dels models de tipus supervisat, és imprescindible disposar d'un nombre relativament ampli de mostres, tant per a la definició de les classes de referència com per a la investigació de la provenença de mostres arqueològiques. És difícil d'establir quin és el mínim de mostres ja que dependrà del grau de similitud entre les classes i de la seva variabilitat interna. En el cas que s'ha presentat, les classes de referència amb menor nombre mostres en contenen 33. I els grups de mostres arqueològiques com a mínim en contenen 6.
- Les funcionalitats de la metodologia presentada, amb l'ús dels models de tipus supervisat, tant per a la definició de les classes de referència com per a la investigació de la provenença de mostres arqueològiques es poden exportar fàcilment a altres contextos i es posen lliurement a disposició de la comunitat científica mitjançant un projecte R amb instruccions autoexplicatives, descarregable des del portal de codi Github:

https://github.com/AnnaAnglisano/SPA_Supervised_Provenance_Analysis/

Línies de Futur

L'optimització de models de classificació de tipus supervisat s'ha aplicat a una extensa base de dades geoquímica per a: i) la distinció de produccions ceràmiques subactuals; ii) identificació de la provenença de restes ceràmiques d'origen desconegut. La base de dades que s'ha produït cobreix un extens nombre de mostres d'un total de 7 centres productors (i les seves corresponents argiles). Com a línies de futur es planteja l'ampliació de la base de dades a nous centres de producció de la mateixa zona estudiada (Catalunya) però també a zones del món completament diferents. És una línia de treball que no necessàriament ha de ser desenvolupada pels promotors de l'aplicació d'aquesta metodologia supervisada, sinó que es pretén donar-la a conèixer i proporcionar les eines per tal que qui vulgui pugui contribuir a exportar la metodologia a altres contextos geogràfics.

A banda de continuar la seva aplicació en els àmbits d'estudi explorats en la present tesi també es planteja com a línia de futur, l'aplicació de la metodologia a l'anàlisi de mostres d'altres tipologies, més enllà de les ceràmiques (i argiles). En aquest sentit, actualment s'està treballant en l'aplicació de l'aproximació supervisada per a la identificació de la provenença de marbres arqueològics. El cas d'estudi que s'està abordant pretén determinar si diversos marbres arqueològics tenen la seva provenença geològica en afloraments de Ceret o Gualba. Els marbres d'ambdues localitats presenten unes característiques molt similars. Les tècniques arqueomètriques que generalment s'utilitzen per als estudis de provenença de marbres (petrografia, isotopia estable de C i O, catodoluminescència) no permeten la distinció clara entre les dues localitats. Partint d'una base de dades geoquímica formada amb anàlisis de mostres de marbre de les pedreres d'ambdues localitats es podran entrenar els models per fer que puguin distingir els dos orígens. Si s'assoleix èxit en l'entrenament, s'espera poder aconseguir la determinació de provenença dels marbres arqueològics. La feina es troba molt avançada i es preveu culminar-la en forma de publicació.

Agraïments

En primer lloc vull agrair l'ajut de totes les persones que han col·laborat en l'estudi de la caracterització de les diferents localitats estudiades: Albert Egea, Jordi Goñi, Anna Pallàs, Eduard Recasens i Jenifer Obama. També a les diferents institucions que han col·laborat en la recerca dels materials (ceràmiques i argiles) que conformen els diferents grups de referència i els conjunts arqueològics estudiats: Glòria Sedó (ceràmiques Sedó, Esparreguera), Xavier Mir (Museu Terracotta, La Bisbal d'Empordà), David Compte (Museu Rocaguinarda de Terrissa dels Països Catalans i el Bisbat de Vic), Núria Miró i Emili Revilla (Servei d'Arqueologia de Barcelona), Gemma Font i Jordi Tura (Museu Etnològic del Montseny), Albert Ventayol (Bac i Ventayol Geoserveis), Joan Vicens (Museu de Terrissa de Quart) i a (Josep Mestres) l'Associació de Terrissers de Quart.

A banda de les institucions i el seu personal, vull ressaltar el paper crucial d'alguns mestres terrissers: Josep Mestres (1957-2021) per la seva col·laboració en la caracterització de la zona Quart, ja que fou ell qui dirigí la campanya de camp a aquesta localitat i es va posar en contacte amb l'associació dels terrissers locals per tal d'aconseguir les mostres necessàries de ceràmica produïda amb seguretat amb argiles locals. També, els bredencs Josep Samón i Saurí (1932-2022), Antoni Majó Manresa (1954-2021) i Carme Serra Majó (1938-2021). Sense la seva saviesa, memòria i sobretot amor a l'ofici s'hauria perdut tot el coneixement que finalment s'ha pogut recopilar al llibre monogràfic sobre la Tradició Terrissera de Breda.

Més persones que han resultat imprescindibles per a l'èxit d'aquest projecte han estat el Dr. Ignasi Queralt, que des del departament de Geociències de l'IDAEA-CSIC ha prestat suport humà i l'equipament per a tot el procés d'anàlisi química de les mostres; Marc Anglisano que ha participat en la selecció dels models d'aprenentatge supervisat i en l'elaboració del codi; Dra. Roberta Di Febo que ha ajudat en els processos de preparació de les mostres i l'indispensable Dr. Lluís Casas pel suport tècnic, per posar els mitjans necessaris i per tota la supervisió, revisió i suport en totes les fases d'aquesta recerca. Finalment agrair també el suport, consells i revisions dels companys de l'AEPECT Dr. David Brusi i Xavier Juan.

Tampoc s'hagués pogut dur a bon terme el projecte sense el suport incondicional de la família (Pere Anglisano, Margarita Roca i Marc Gabasa) que en els moments més difícils sempre han fet el possible per aixecar els ànims i m'han fet creure i estimar el projecte que va arrencar el ja llunyà 2016.

Finalment agrair també a totes les persones que conformen el tribunal d'avaluació d'aquest treball, tant titulars com suplents per la seva total disposició a formar-ne part.

Les despeses que s'han derivat de les analítiques que ha comportat aquesta recerca han estat sufragades, en part, pel projecte de recerca amb referència CGL2013-42167-P del Ministerio de Economía y Competitividad, pel Grup de Recerca Aplicada al Patrimoni Cultural (GRAPAC) i pel projecte amb referència CEX2018-000794-S de l'IDAEA-CSIC (Ministeri de Ciència i Innovació).

Índex de taules i figures

Figures

Figura 1: Diferents terreres que s'han caracteritzat: a) La Bisbal d'Empordà, b) Breda, c) Esparreguera, d) Quart, e) Sant Julià de Vilatorça i f) Verdú.....17

Figura 2: Algunes de les mostres d'argila un cop cuites.....18

Figura 3: Mostres representatives dels conjunts de mostres arqueològiques de provenença desconeguda i que formen part dels conjunts analitzats. D'esquerra a dreta: ceràmica gris de cuïta reductora, ceràmica de cuïna amb vidrat de plom, ceràmica de cuïna amb vidrat acolorit en verd gòtic, tots tres conjunts procedents del castell de Montsoriu (CM). Ceràmica de cuïna de Torre de la Mora (TM) i finalment ceràmica feta a mà del jaciment de La Creueta (C)19

Figura 4: Porta làmines amb algunes de les làmines primes realitzades en el present estudi i una làmina prima d'una ceràmica subjectada amb la mà.....20

Figura 5: Equip de difracció de raigs X utilitzat.....21

Figura 6: Espectròmetre utilitzat per a l'estudi22

Figura 7: a) dues làmines primes de ceràmica local de Quart escanejades. b) fotografies de tres làmines primes de ceràmica local de Quart (Q24 i Q14) i Sant Julià de Vilatorça (SJV08).....73

Figura 8: Quatre difractogrames gairebé idèntics de mostres d'argila, dues procedents de Breda (blau) i dues de Quart (carabassa). S'han ressaltat els pics corresponents a fil·losilicats del grup de la mica – il·lita (vermell); i del grup de la caolinita – serpentina (en verd).....74

Figura 9: a) Comparació a través de parelles d'elements químics de les mostres dels grups de referència on es posa de manifest l'elevat grau de solapament. b) Gràfic HCA on s'observa que les mostres d'un dels grups (la Bisbal) no queden agrupades sinó repartides en zones allunyades del dendrograma resultant.75

Taules

Taula 1: Resum del total de mostres que formen part de l'estudi.....20

Taula 2: Conjunt de mostres arqueològiques estudiades.....21

Bibliografia

1. Stanley E., M. *Chimica Dell'ambiente*; Piccin-Nuova Libreria, 2000;
2. Sen, T. *Clay Minerals: Properties, Occurrence and Uses*; Nova Science, 2017;
3. Searle, A.B. *The Natural History of Clay*; Cambridge: University Press, 1912;
4. Ion, R.-M.; Fierascu, R.-C.; Teodorescu, S.; Fierascu, I.; Bunghez, I.-R.; Turcanu-Carutiu, D.; Ion, M.-L. Ceramic Materials Based on Clay Minerals in Cultural Heritage Study. In *Clays, Clay Minerals and Ceramic Materials Based on Clay Minerals*; InTech, 2016.
5. Sabatino, G.; Franzone, M.; Martinelli, M.C.; Rondinella, M.T.; Italiano, F.; Caccamo, M.T.; Mezzatesta, F.; Magazù, S.; Tripodo, A.; Di Bella, M. From Clays to Pottery: Role of Geomaterials in the Social-Technological Development of the Messina Territory (Sicily, Italy) and Archaeological-Historical Information on the Main Kilns. . *Classe di Scienze Fisiche, Matematiche e Naturali* **2021**, 99.
6. Anglisano, A.; Goñi, J. *La Tradició Terrissera de Breda (s.XV-s.XX), Terreres, Obradors, Forns, Elaboaració i Vocabulari*; Estudis i Textos, 2021;
7. Santanach, J.; Rosal, J.; Suñol, M. *La Ceràmica de Quart En La Memòria Viva: Els Obradors*; Ajuntament de Quart, Associació de Terrissaires Artesans de Quart, 1998;
8. Bover, A. *La Ceràmica*; Diputació de Girona.; Quaderns de la Revista de Girona, 1993; Vol. 42;.
9. Yan, K.; Guo, Y.; Fang, L.; Cui, L.; Cheng, F.; Li, T. Decomposition and Phase Transformation Mechanism of Kaolinite Calcined with Sodium Carbonate. *Appl Clay Sci* **2017**, 147, doi:10.1016/j.clay.2017.07.010.
10. Ouahabi, M. El; Daoudi, L.; Hatert, F.; Fagel, N. Modified Mineral Phases During Clay Ceramic Firing. *Clays Clay Miner* **2015**, 63, doi:10.1346/CCMN.2015.0630506.
11. Sholkovitz, E.R. Chemical Evolution of Rare Earth Elements: Fractionation between Colloidal and Solution Phases of Filtered River Water. *Earth Planet Sci Lett* **1992**, 114, doi:10.1016/0012-821X(92)90152-L.
12. Borrego, J.; López-González, N.; Carro, B.; Lozano-Soria, O. Origin of the Anomalies in Light and Middle REE in Sediments of an Estuary Affected by Phosphogypsum Wastes (South-Western Spain). *Mar Pollut Bull* **2004**, 49, doi:10.1016/j.marpolbul.2004.07.009.
13. Romero Vidal, A.; Rosal Sagalés, J. *La Terrissa a Catalunya*; Brau Edicions SL, 2014;
14. Heimann, R.; Franklin, U.M. Archaeo-Thermometry: The Assessment of Firing Temperatures of Ancient Ceramics. *Journal of the Institute of Conservation - Canadian Group* **1979**, 4, 23–45.
15. Holakooei, P.; Tessari, U.; Verde, M.; Vaccaro, C. A New Look at XRD Patterns of Archaeological Ceramic Bodies. *J Therm Anal Calorim* **2014**, 118, doi:10.1007/s10973-014-4012-z.

16. Aitken, M.J. Dating by Archaeomagnetic and Thermoluminescent Methods. *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences* **1970**, 269, doi:10.1098/rsta.1970.0087.
17. Pollard, A.M.; Heron, C. *Archaeological Chemistry*; 2nd ed.; Royal Society of Chemistry,: Cambridge, UK, 2008;
18. Pollard, A.M.; Batt, C.M.; Stern, B.; Young, S.M.M. *Analytical Chemistry in Archaeology*; Cambridge University Press, 2007; ISBN 9780521652094.
19. Mommsen, H. Short Note: Provenancing of Pottery- The Need for an Integrated Approach? *Archaeometry* **2004**, 46, doi:10.1111/j.1475-4754.2004.00156.x.
20. Kuleff, I.; Djingova, R. Provenance Study of Pottery; Choice of Elements to Be Determined. . *ArchéoSciences, revue d'Archéométrie* **1996**, 20, 57–67.
21. Alzate Gallego, L.A. *Arqueología Histórica y Arqueometría Para El Estudio de La Cerámica Colonial En Fundaciones de Terra Firme - Siglo XVI*. Tesi doctoral, Universitat de Barcelona: Barcelona, 2016.
22. Reedy, C.L. *Thin-Section Petrography of Stone and Ceramic Cultural Materials*; Archetype Publications : London, 2008;
23. Queen, P.S. *Ceramic Petrography: The Interpretation of Archaeological Pottery & Related Artefacts in Thin Section*; Archaeopress: Oxford, UK, 2013;
24. Queen, P. *Interpreting Silent Artefacts: Petrographic Approaches to Archaeological Ceramics*; Queen, P.S., Ed.; Archaeopress Publishing: Oxford, UK, 2009;
25. Hein, A.; Tsolakidou, A.; Iliopoulos, I.; Mommsen, H.; Buxeda i Garrigós, J.; Montana, G.; Kilikoglou, V. Standardisation of Elemental Analytical Techniques Applied to Provenance Studies of Archaeological Ceramics: An Inter Laboratory Calibration Study Electronic Supplementary Information (ESI) Available: Five Tabular Appendices Giving Element Concentrations Measured in Reference Materials. See [Http://Www.Rsc.Org/Suppdata/an/B1/B109603f/](http://www.rsc.org/Suppdata/an/B1/B109603f/). *Analyst* **2002**, 127, doi:10.1039/b109603f.
26. Baxter, M.J. *Exploratory Multivariate Analysis in Archaeology*; EliotWerner Publications-Inc.: Clinton Corners, NY, USA, 2015;
27. Ricca, M.; Paladini, G.; Rovella, N.; Ruffolo, S.A.; Randazzo, L.; Crupi, V.; Fazio, B.; Majolino, D.; Venuti, V.; Galli, G.; et al. Archaeometric Characterisation of Decorated Pottery from the Archaeological Site of Villa Dei Quintili (Rome, Italy): Preliminary Study. *Geosciences (Basel)* **2019**, 9, doi:10.3390/geosciences9040172.
28. Neff, H. *Chemical Characterization of Ceramic Pastes in Archaeology*; Neff, H., Ed.; Prehistory Press: Madison, WI, USA, 1992;
29. Buxeda i Garrigós, J.; Cau Ontiveros, M.; Kilikoglou, V. Garrigós, Jaume Buxeda i, Miguel Ángel Cau Ontiveros and Vassilis Kilikoglou. "Chemical Variability in Clays and Pottery from a Traditional Cooking Pot Production Village: Testing Assumptions in Pereruela*." *Archaeometry* 45 (2003): 1-17. *Archaeometry* **2003**, 45, 1–17.

30. Brorsson, T.; Blank, M.; Fridén, I.B. Mobility and Exchange in the Middle Neolithic: Provenance Studies of Pitted Ware and Funnel Beaker Pottery from Jutland, Denmark and the West Coast of Sweden. *J. Archaeol. Sci. Rep.* **2018**, *20*, 662–674.
31. Papachristodoulou, C.; Oikonomou, A.; Ioannides, K.; Gravani, K. A Study of Ancient Pottery by Means of X-Ray Fluorescence Spectroscopy, Multivariate Statistics and Mineralogical Analysis. *Anal Chim Acta* **2006**, *573–574*, doi:10.1016/j.aca.2006.02.012.
32. Aquilia, E.; Barone, G.; Mazzoleni, P.; Ingoglia, C. Petrographic and Chemical Characterisation of Fine Ware from Three Archaic and Hellenistic Kilns in Gela, Sicily. *J Cult Herit* **2012**, *13*, doi:10.1016/j.culher.2012.02.005.
33. Granato, D.; Santos, J.S.; Escher, G.B.; Ferreira, B.L.; Maggio, R.M. Use of Principal Component Analysis (PCA) and Hierarchical Cluster Analysis (HCA) for Multivariate Association between Bioactive Compounds and Functional Properties in Foods: A Critical Perspective. *Trends Food Sci Technol* **2018**, *72*, doi:10.1016/j.tifs.2017.12.006.
34. Ben Tahar, S.; Capelli, C. L’atelier Céramique d’Henchir Chouggaf (Ouedhref, Tunisie). *Antiquités africaines* **2018**, doi:10.4000/antafr.995.
35. Waksman, Y.; Gragueb, S.; Trégliia, J.-C.; Capelli, C. Jarres et Amphores de Sabra Al-Mansūriya (Kairouan, Tunisie). *École Française de Rome* **2011**, 197–220.