

# 5641 - Clasificación automática de textos y explotación BI

Javier Buill Vilches

# Estructura

- Motivación
  - Problemática
  - Metodología actual
  - Propuesta de solución
- Clasificación automática
  - Método de Bayes
- Business Intelligence (BI)
  - Qlikview
- Resultados
  - Clasificación
  - Visualización
- Conclusiones
  - Mejoras

# Motivación

- Extracción de datos y clasificación

Información de los medios

Clasificación manual



ID	Fecha	País	Independencia	Reduccion directa	Gobierno central	Gobierno autonómico	Elecciones	Actores1	Actores2	PSOE	PP	UPD	CS	ERC	BASCOP	ICV-EA M	PSC	PSC	Ciudadanos	Organizaciones	Interactiva
1	07/04/2012	Ans	0	1	1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0
2	04/04/2012	Ans	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	07/04/2012	Ans	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	07/04/2012	Ans	0	1	0	1	0	3	2	0	1	0	2	0	0	0	0	0	0	0	0
5	07/04/2012	Ans	0	1	0	1	0	3	0	0	0	0	2	0	0	0	0	0	0	0	0
6	07/04/2012	Ans	0	1	1	1	0	2	0	0	0	0	2	0	0	0	0	0	0	0	0
7	07/04/2012	Ans	0	0	1	0	0	2	0	0	2	0	1	0	0	0	0	0	0	0	0
8	07/04/2012	Ans	1	0	0	1	0	3	1	0	0	0	0	0	0	1	0	0	0	0	0
9	07/04/2012	Ans	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
10	07/04/2012	Ans	0	0	1	0	0	3	0	0	2	0	0	0	0	0	0	0	0	0	0
11	04/04/2012	Ans	0	1	1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0
12	04/04/2012	Ans	0	0	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
13	05/04/2012	Ans	0	1	0	1	0	3	0	0	0	2	0	0	0	0	0	0	0	0	0
14	05/04/2012	Ans	0	0	1	1	0	3	2	0	2	0	1	0	0	0	0	0	0	0	0
15	05/04/2012	Ans	0	1	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
16	05/04/2012	Ans	0	1	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
17	05/04/2012	Ans	0	0	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
18	05/04/2012	Ans	0	0	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
19	05/04/2012	Ans	0	0	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
20	05/04/2012	Ans	0	0	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
21	05/04/2012	Ans	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	05/04/2012	Ans	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	05/04/2012	Ans	0	1	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
24	07/04/2012	Ans	0	1	0	0	0	12	0	0	0	0	0	0	0	0	0	0	0	0	0
25	07/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	04/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
27	03/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
28	03/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
29	03/04/2012	Ans	0	1	0	0	0	0	0	0	2	0	2	0	0	0	0	0	0	0	0
30	03/04/2012	Ans	0	1	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0
31	04/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
32	03/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
33	04/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
34	03/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
35	03/04/2012	Ans	0	1	1	1	1	0	2	3	0	2	0	0	0	0	0	0	0	0	0
36	03/04/2012	Ans	0	1	1	1	1	0	2	2	0	0	2	0	0	0	0	0	0	0	0
37	03/04/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
38	28/03/2012	Ans	0	1	0	0	0	0	3	1	0	0	1	0	0	0	1	0	0	0	0
39	28/03/2012	Ans	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
40	27/03/2012	Ans	0	1	0	0	0	0	1	2	1	1	0	0	0	0	0	0	0	0	0
41	27/03/2012	Ans	0	1	1	1	1	0	2	2	0	2	0	0	0	0	0	0	0	0	0
42	03/04/2012	Ans	0	1	0	1	1	0	3	2	0	0	12	0	0	0	0	0	0	0	0
43	03/04/2012	Ans	0	1	1	1	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0
44	03/04/2012	Ans	0	1	1	1	1	0	2	0	0	0	0	0	0	0	0	0	0	0	0
45	03/04/2012	Ans	1	0	0	1	0	3	4	0	0	0	1	1	0	0	0	0	0	0	0
46	03/04/2012	Ans	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

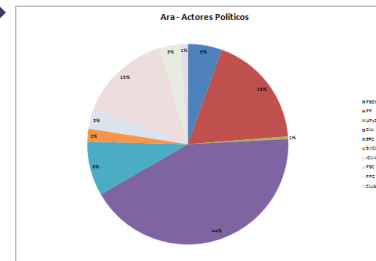
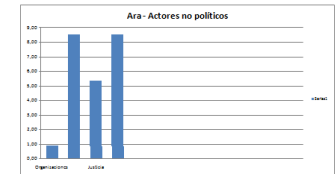
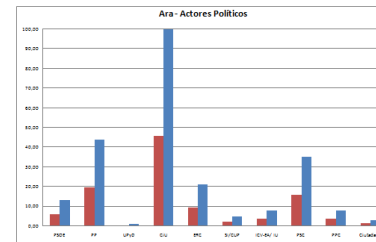
# Motivación

- Análisis de la información

## Clasificación manual

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
10	Excs	Foars	Indepndentis	Relucio disvts	Debitoris control	Debitoris atencioaris	Elccioaris	Actores	Actores		PSOE	PP	UPdI	CH	ERC	BNGO	ICV-EA	BI	PSC	ERC	Ciudadans	Organizaciones	Internavis			
1	0	0	0	1	1	0	0	2	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
2	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
4	0	0	1	0	1	0	3	2	0	0	1	0	2	0	0	0	0	0	0	0	0	0	0	0	0	
5	0	1	0	1	0	3	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	
6	0	1	1	1	1	0	2	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
7	0	0	1	0	0	0	2	0	0	0	2	0	1	0	0	0	0	0	0	0	0	0	0	0	0	
8	1	0	0	0	1	0	3	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	
9	0	0	0	0	0	0	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	
10	0	0	1	0	0	0	3	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
11	0	1	1	0	0	2	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
12	0	0	0	1	0	3	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
13	0	1	0	1	0	3	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
14	0	0	1	1	1	1	0	3	2	0	2	0	2	0	0	0	0	0	0	0	0	0	0	0	0	
15	0	1	1	1	1	0	3	2	0	0	2	0	2	0	0	0	0	0	0	0	0	0	0	0	0	
16	0	1	0	1	0	1	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
17	0	0	0	0	1	0	3	1	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
18	0	0	0	0	1	0	3	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
19	0	0	0	0	1	0	3	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
21	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
22	0	0	0	0	0	0	1	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
23	0	1	0	0	1	0	3	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
25	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
26	0	1	0	1	0	1	0	3	1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
27	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
28	0	1	0	1	0	1	0	3	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
29	0	1	0	0	0	0	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
30	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
31	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
32	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
33	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
34	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
35	0	1	1	1	1	0	2	3	0	0	2	0	2	0	0	0	0	0	0	0	0	0	0	0	0	
36	0	1	1	1	1	1	0	3	2	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
37	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
38	0	1	0	0	0	0	3	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	
39	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
40	0	1	0	0	0	0	1	2	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
41	0	1	1	1	1	1	2	3	0	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	
42	0	1	1	1	1	0	3	2	0	0	0	12	0	0	0	0	0	0	0	0	0	0	0	0	0	
43	0	1	1	1	1	1	0	3	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	
44	0	1	1	1	1	1	0	3	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
45	1	0	0	1	0	0	0	3	4	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	
46	1	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

## Análisis de los datos clasificados



# Problemática

## ¿Problema? TIEMPO



Para un texto de aproximadamente 4000 frases (~200 páginas) el tiempo estimado de codificación manual ronda los 2 meses en el trabajo de una persona a media jornada dedicada exclusivamente a esta tarea

# Metodología actual

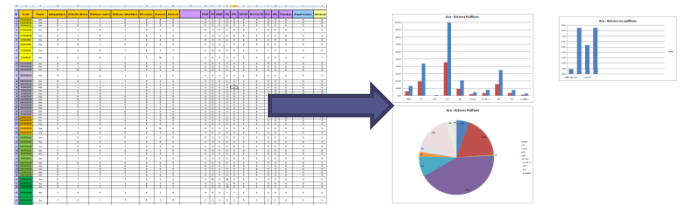
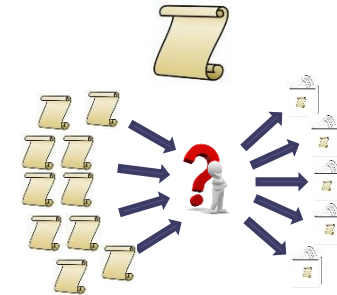
- Pasos de la actual metodología:
  - Formación de codificador
  - Lectura de texto (1)
  - Separación en cuasi-frases
  - Lectura de texto (2)
  - Clasificación manual
  - Recogida de datos
  - Estudio y análisis para los datos clasificados (sólo útil para éstos)



Hemos podido, en este contexto, desarrollar políticas de vertebración territorial para combatir los desequilibrios demográficos internos y luchar contra el despoblamiento rural,

Crear en una economía sostenible,

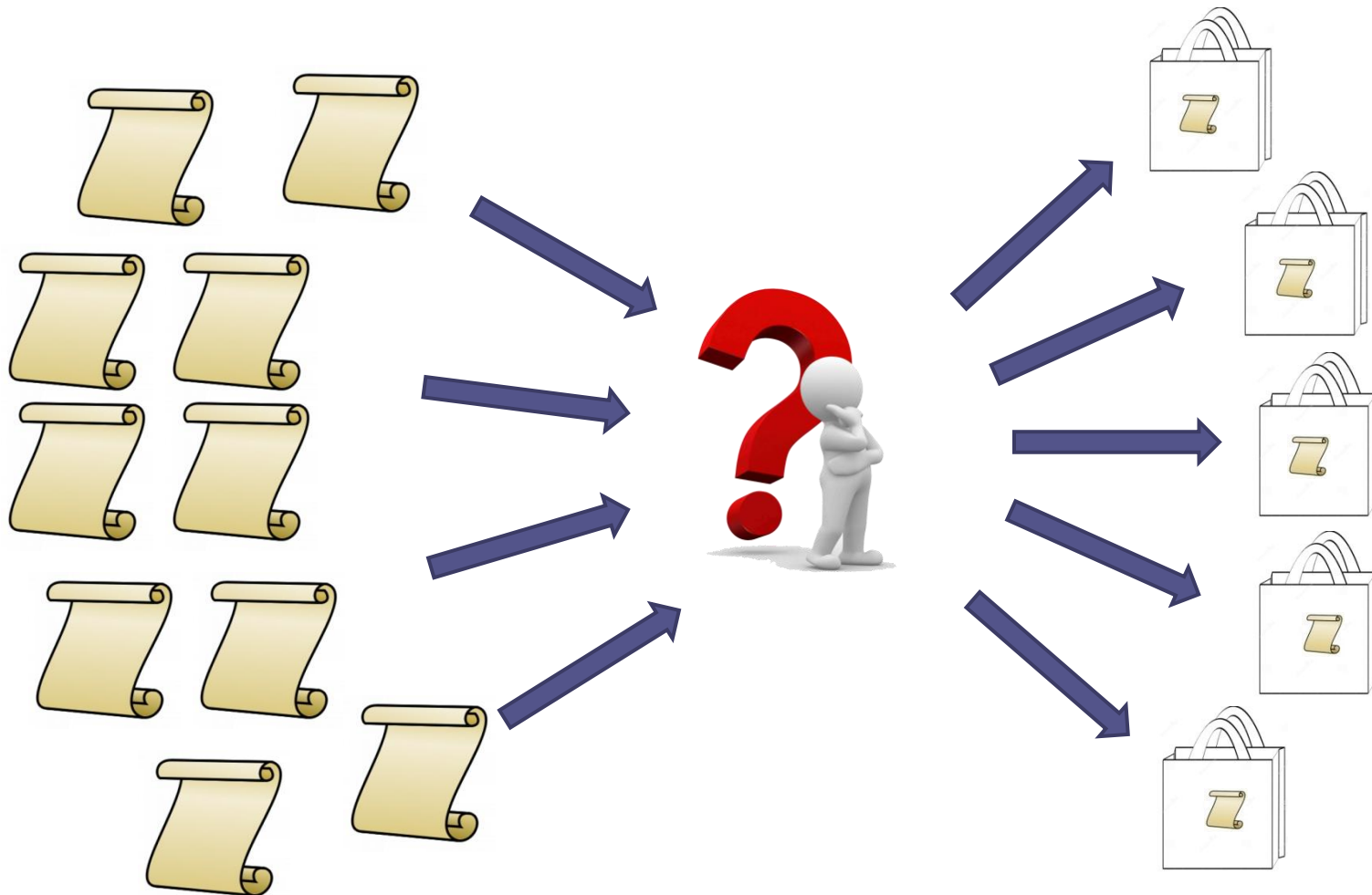
diversificada, equilibrada y tecnológicamente avanzada,



# Propuesta de solución - Objetivos

- Clasificador automático
  - Rápido
  - Alto índice de acierto
  - Fácil manejo
- Herramienta *Business Intelligence* (BI)
  - Cuadro de mando preparado para presentar esta información para análisis de datos

# Clasificación automática





# Algoritmo propuesto

Bayes

$$P(A_i | B) = \frac{P(A_i)P(B | A_i)}{P(B)}$$

Donde:

$P(A_i)$  → Probabilidad a priori

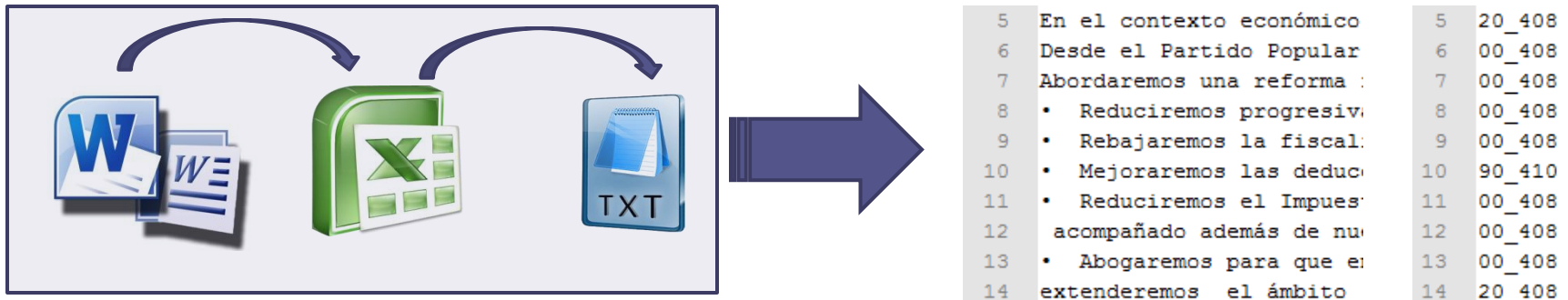
$P(B|A_i)$  → Probabilidad condicional

$P(B)$  → Probabilidad total

$P(A_i|B)$  → Probabilidad a posteriori

# Metodología desarrollada (1)

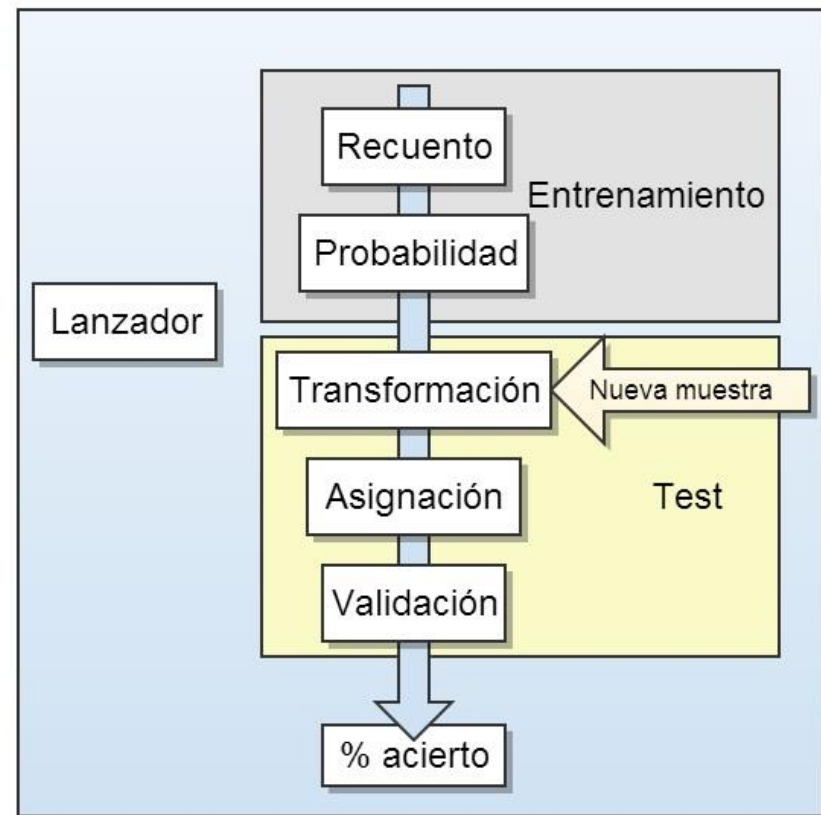
## Preparación y transformación de los datos



- Uso de formatos sencillos para usuarios ajenos a la informática
- Datos provistos en documentos de texto (.doc/.docx)
- Transformados con Excel
- Exportados a 2 archivos de texto .txt separados por:
  - Texto
  - Categoría

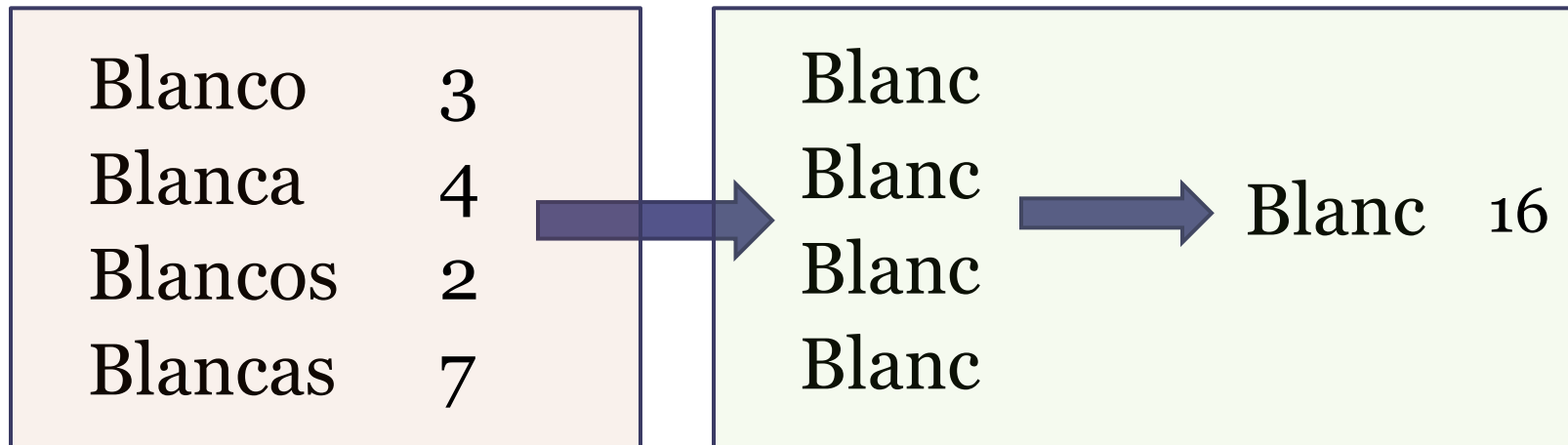
# Metodología desarrollada (2)

- Flujo de ejecución:
  - 2 Fases:
    - Entrenamiento
    - Test
- Entrenamiento con 32 configuraciones para cada texto combinando:
  - *Stemmed words*
  - *Stop words*
  - *Laplace Smoothing*
  - $\text{Prod}(\text{prob})/\text{sum}(\text{log})$
  - Probabilidad a priori



# Metodología desarrollada (3)

## Stemmed words



# Metodología desarrollada (3)

## Stop words

*Listado de palabras “sin significado” que no dan peso a ninguna categoría en particular.*

*Son las palabras que más aparecen*

*ahí, tal, de, aquí, allí, allá, la, que, el, en, y, a, los, del, se, las, por, un, para, con, no, una, su, al, lo, como, más, pero, sus, le, ya, o, este, sí, pues, decir, entonces, vez, porque, esta, entre, cuando, muy, sin, sobre, también, me, hasta, hay, donde, quien, desde, todo, nos, durante, todos, uno, les, ni, contra, otros, ese, eso, ante, ellos, e, esto, mí, antes, algunos, qué, unos, yo, otro, otras, otra, él, tanto, esa, estos, mucho, quienes, nada, muchos, cual, poco, ella, estar, estás, algunas, algo, nosotros, mi, mis, tú, te, ti, tu, tus, ellas, nosotras, vosotros, vosotras, os, mío, mía, míos, mías, tuyo, tuya, tuyos, tuyas, suyo, suya, suyos, suyas, nuestro, nuestra, nuestros, nuestras, vuestro, vuestra, vuestros, vuestras, esos, esas, ser, haber, tener, hacer, estar.*

# Metodología desarrollada (3)

## Laplace Smoothing

$$P(X_i = x_{ij} | Y = y_k) = \frac{\# D\{X_y = x_{ij} \wedge Y = y_k\} + 1}{\# D\{Y = y_k\} LM}$$

Corrección cuando el valor de recuento = 0

palabra	cat1	cat2	cat3	cat4	cat5	cat6
a	10	9	10	7	5	8
b	3	2	4	2	1	2
c	10	8	2	10	2	6
d	0	6	2	4	4	6
e	6	3	8	1	7	4
f	4	2	10	2	10	7
g	5	7	7	2	1	8
h	8	7	3	8	9	2
i	6	2	0	5	8	5
h	2	5	6	4	10	6
k	3	7	4	9	5	5
l	8	8	8	2	3	3
m	3	4	7	6	2	5
n	8	6	3	10	2	5
o	6	1	9	7	7	0
p	8	2	0	7	4	0
q	5	3	6	9	9	1
r	10	8	9	3	8	9
s	8	2	7	0	8	3
t	8	3	10	0	9	5
u	0	9	1	5	2	4
v	0	2	1	2	4	6
w	0	1	9	3	0	5
x	0	0	2	8	5	1
y	1	0	6	3	5	2
z	6	2	5	1	3	6

# Metodología desarrollada (3)

## Prod(prob) - Sum(log)

Al multiplicar muchos valores cercanos a 0, puede causar imprecisión numérica, generando ceros de manera similar al caso anterior.

Se puede cambiar el cálculo y en vez de hacer la multiplicación de probabilidades, calcular su logaritmo, y por tanto, queda como el sumatorio de los logaritmos

$$v = \operatorname{argmax}(v_j \in V (P(v_j) \prod_{a_i \in X} (P(a_i | v_j))))$$

$$v = \operatorname{argmax}(v_j \in V (\log P(v_j) + \sum_{a_i \in X} \log P(a_i | v_j)))$$

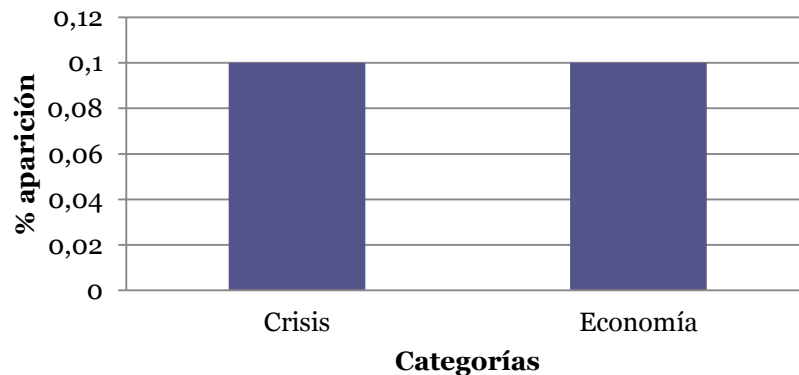
# Metodología desarrollada (3)

## Probabilidad a priori

$$v = \operatorname{argmax}_{v_j \in V} (P(v_j) \prod_{a_i \in X} (P(a_i | v_j)))$$

$$v = \operatorname{argmax}_{v_j \in V} (P(v_j) \prod_{a_i \in X} (P(a_i | v_j)))$$

### Palabra - euros



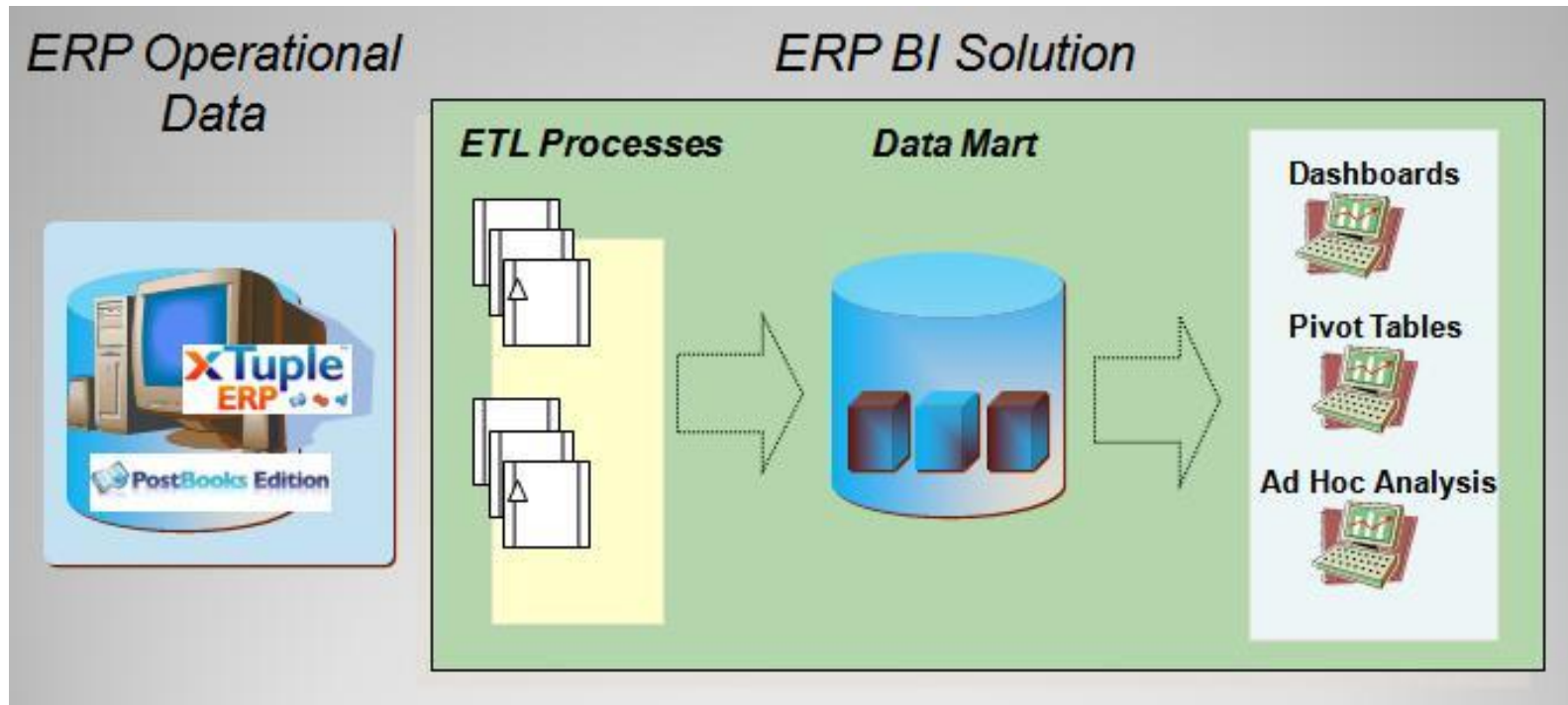
Categorías	# Palabras
Crisis	500
Economía	2500



# Business Intelligence (BI)



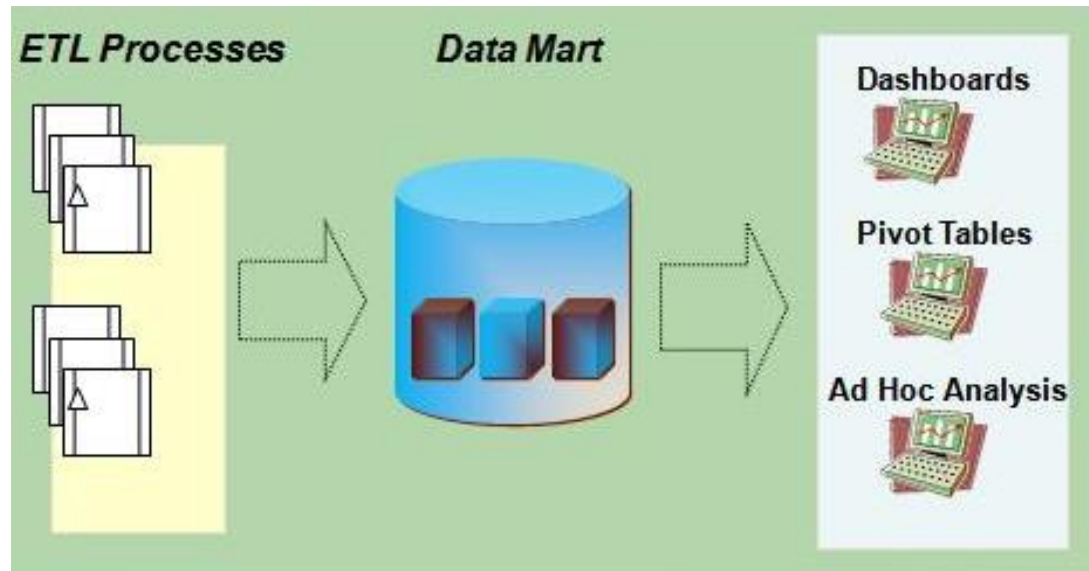
# Business Intelligence (BI)



Transformación de gran cantidad de datos provenientes del ERP a tablas y gráficos fácilmente interpretables

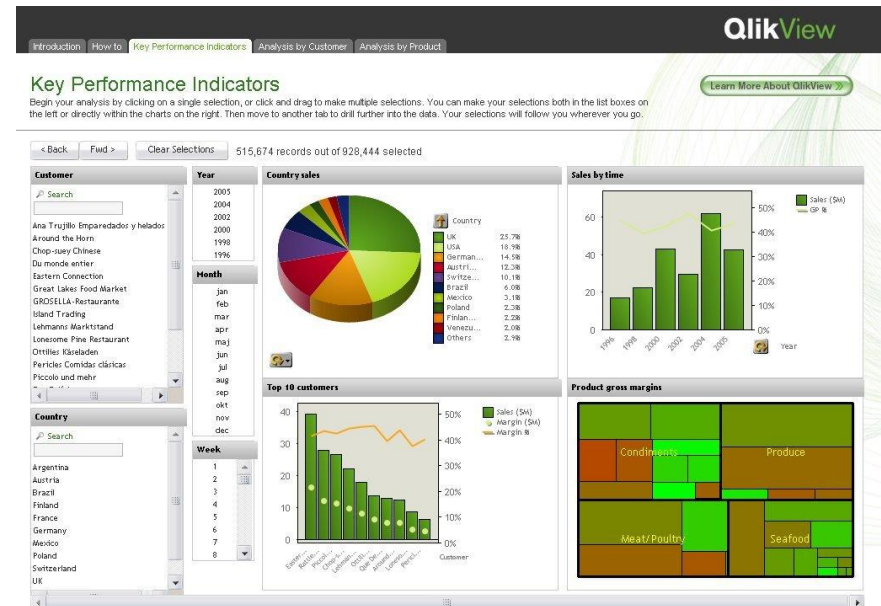
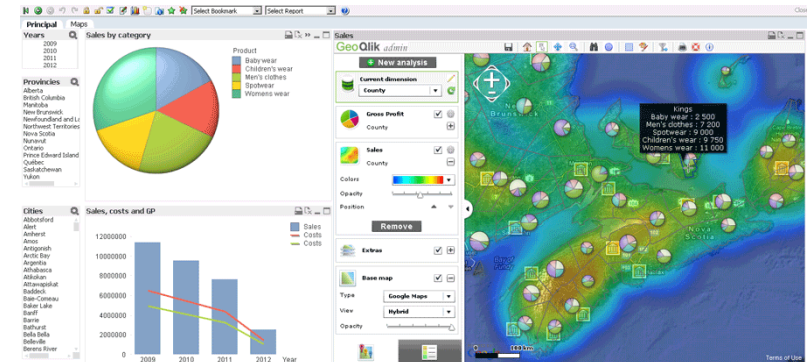
# Business Intelligence (BI)

- Proceso de analizar datos de una empresa y extraer conocimiento de ello
- Uso de almacén de información (*Datawarehouse*) como herramienta estratégica
- Habilidad de explorar y analizar datos para revelar la existencia de tendencias



# Qlikview

- Es una de las principales plataformas para el Business Discovery
- Aporta un lenguaje propio de modelado de datos
- Utiliza un modelo de datos asociativo que se carga en memoria



# Qlikview (2)

- ETL

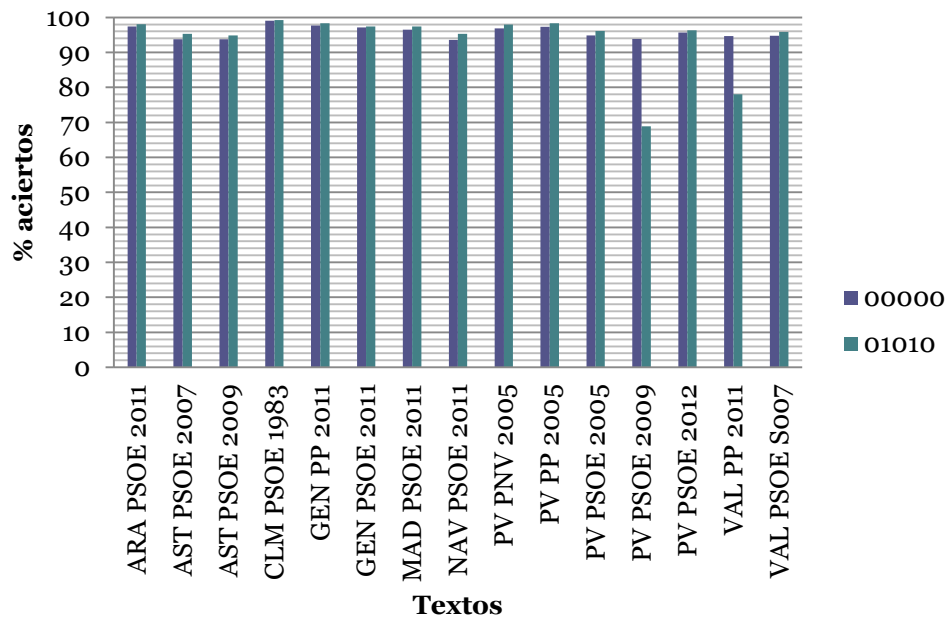


- Selección

The screenshot shows a selection pane in QlikView with three filters: 'Año', 'Partido', and 'Region'. The 'Año' filter has three buttons: 2012, 2013, and 2014. The 'Partido' filter has two buttons: PP and PSOE. The 'Region' filter has seven buttons: Aragon, Castilla la mancha, Cataluña, Madrid, Navarra, Pais vasco, and Valencia. The 'Aragon' button is highlighted in green.

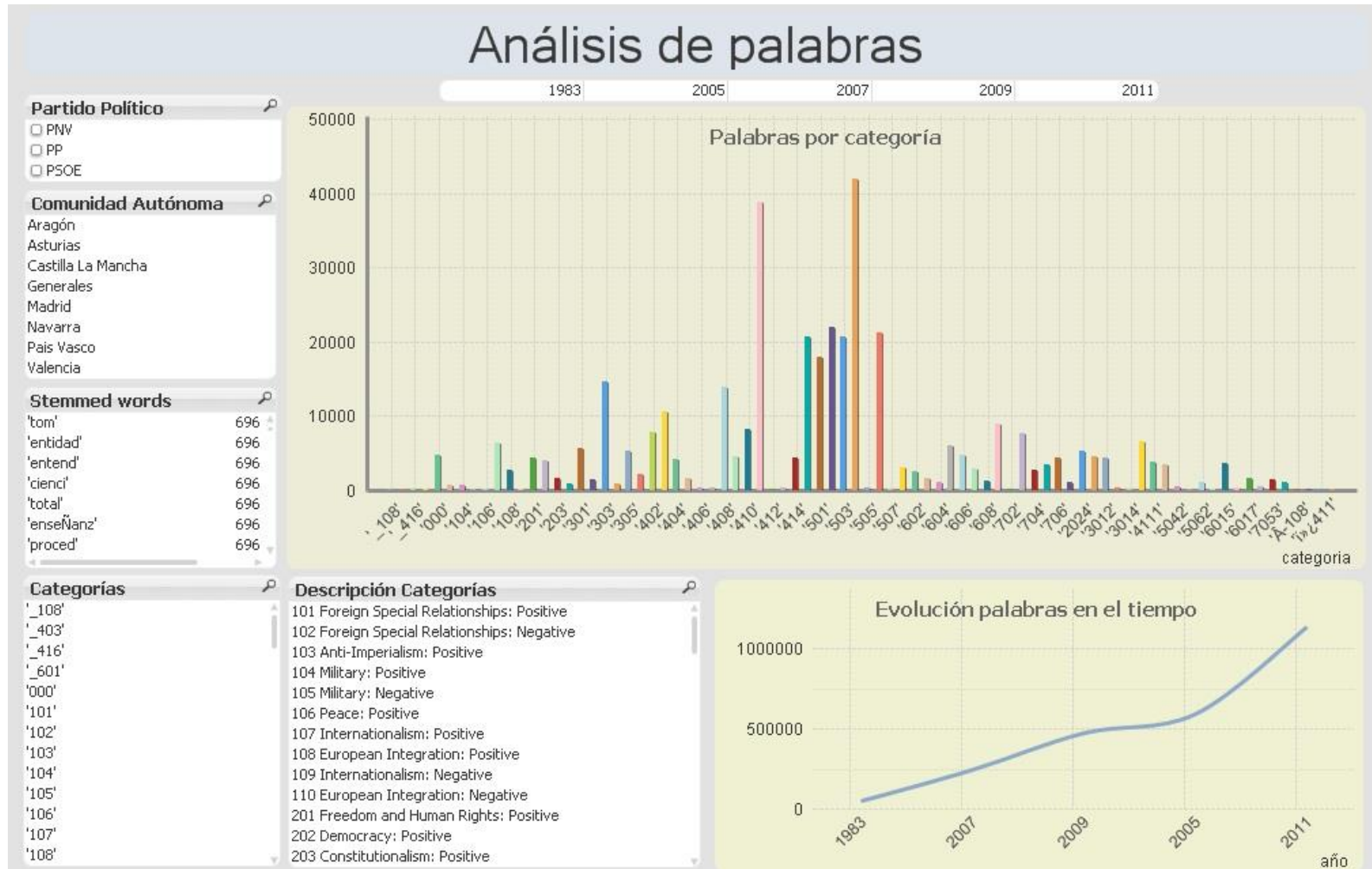
# Resultados clasificación

Mejores resultados para la configuración  
01010  
stop words + sum(log)



Texto	% acierto 01010
ARA PSOE 2011	98,06
AST PSOE 2007	95,35
AST PSOE 2009	94,87
CLM PSOE 1983	99,28
GEN PP 2011	98,28
GEN PSOE 2011	97,38
MAD PSOE 2011	97,39
NAV PSOE 2011	95,31
PV PNV 2005	97,98
PV PP 2005	98,3
PV PSOE 2005	96,16
PV PSOE 2009	68,86
PV PSOE 2012	96,35
VAL PP 2011	78,03
VAL PSOE S007	95,86

# Dashboard

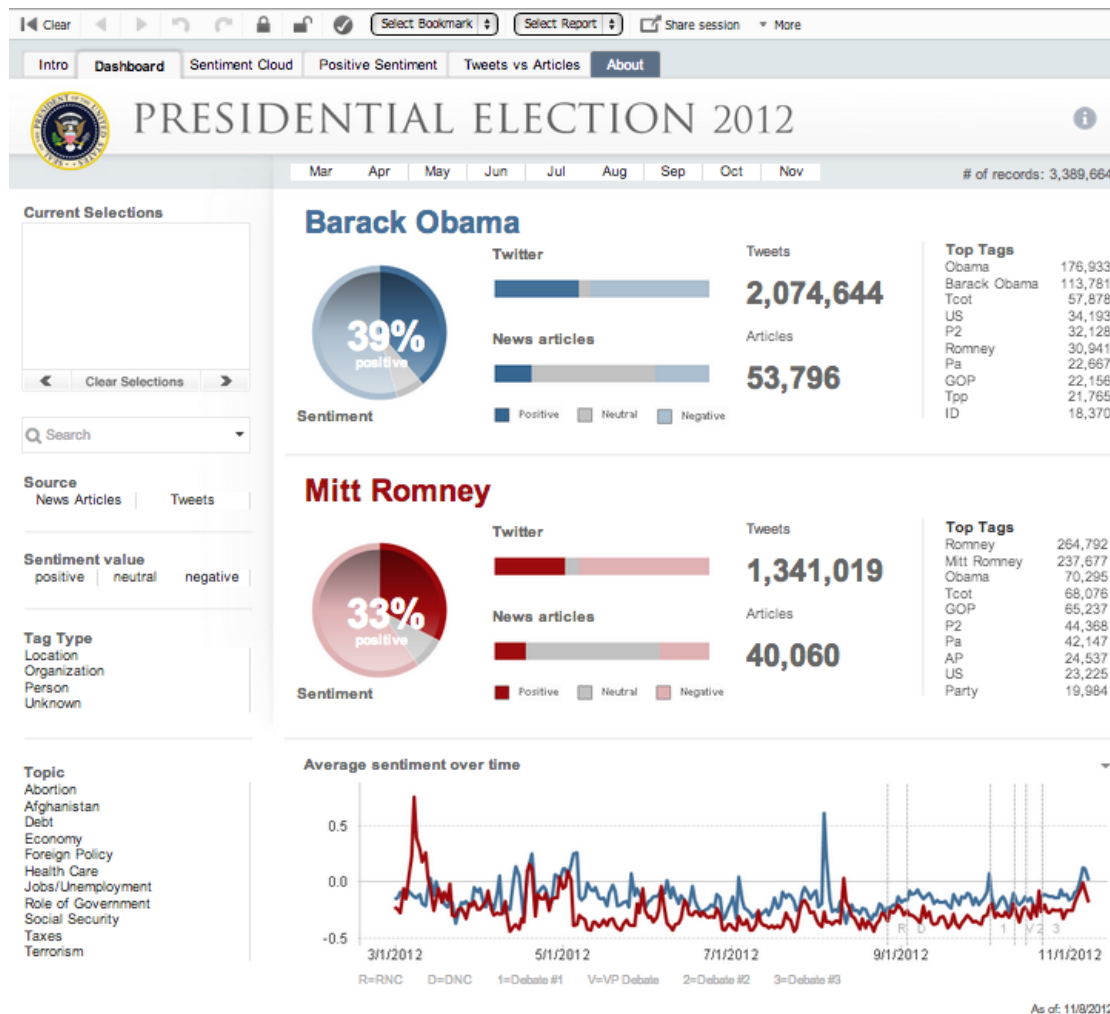


# Planificación





# Propuesta continuación



# Conclusiones

- El tiempo invertido en la clasificación automática es prácticamente nulo y se debe considerar como la principal ventaja.
- Altos índices de acierto en la clasificación, llegando hasta 99.2%
- Cuadro de mando preparado para datos formateados, permite estandarizar. Simplemente al cargar nuevos datos se dispone de información ya procesada para analizar con eficacia.

# Mejoras

- Mejorar el porcentaje de acierto (a través de muchos más datos de entrenamiento)
- Multi-idioma (actualmente sólo castellano)
- Mejorar el método de stemming
- Añadir palabras a la lista de stop words
- Mejorar el tiempo de clasificación
- Realizar un estudio a medida de las necesidades para mostrar en Qlikview para explotar al máximo la información proveniente de los datos

Gracias