

# Analizando las directivas de la EU – Tecnologías del Lenguaje y Big Data

25 de Abril  
DESARROLLO Y  
POTENCIAL DE LAS  
TECNOLOGÍAS DEL  
LENGUAJE

big data | analytics | cognitive

la interpretación de everis:



  
Cognitive Applications Development Platform · powered by ITAINNOVA

## HUMAN TECHNOLOGY

Understanding human knowledge to empower business

everisMoriarty is a platform for developing cognitive applications capable of understanding human knowledge and applying it to the decision making process.

○ ○ ○ ○

# La Comisión Europea

**Hay 32.966  
funcionarios  
generando  
documentos.**

**24 idiomas  
oficiales.  
Traducciones.  
Jerga legal.**

MALTÉS

PORTUGUÉS

ENGLISH

ITALIANO

ESPAÑOL

ελληνικά

FRANCÉS

LETÓN

BÚLGARO

DEUTSCH



# La complejidad...



Plan TL

Plan de Impulso de las  
Tecnologías del Lenguaje



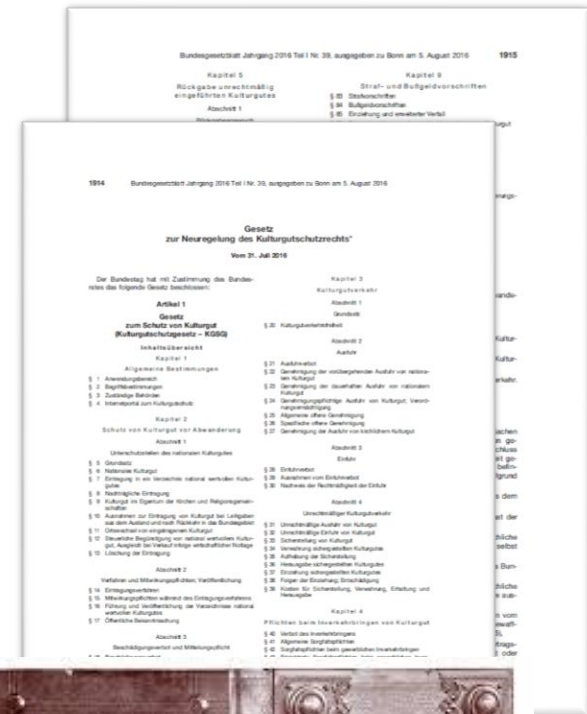
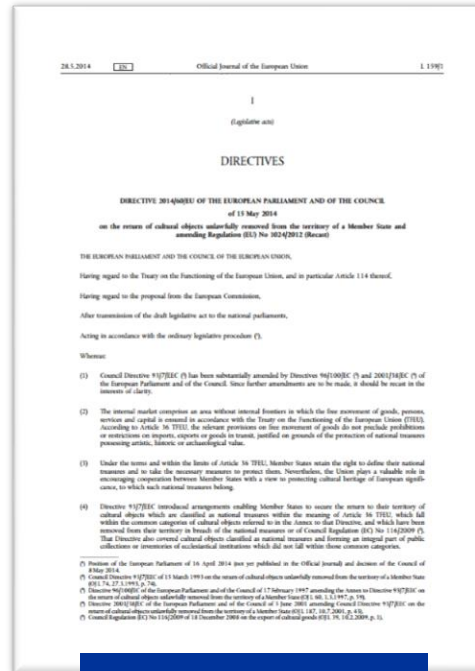
an NTT DATA Company

El proceso legislativo es largo, complejo.

La directivas nacen en Europa, se traducen a todos los idiomas oficiales.

Se trasladan al marco legal de cada estado miembro.

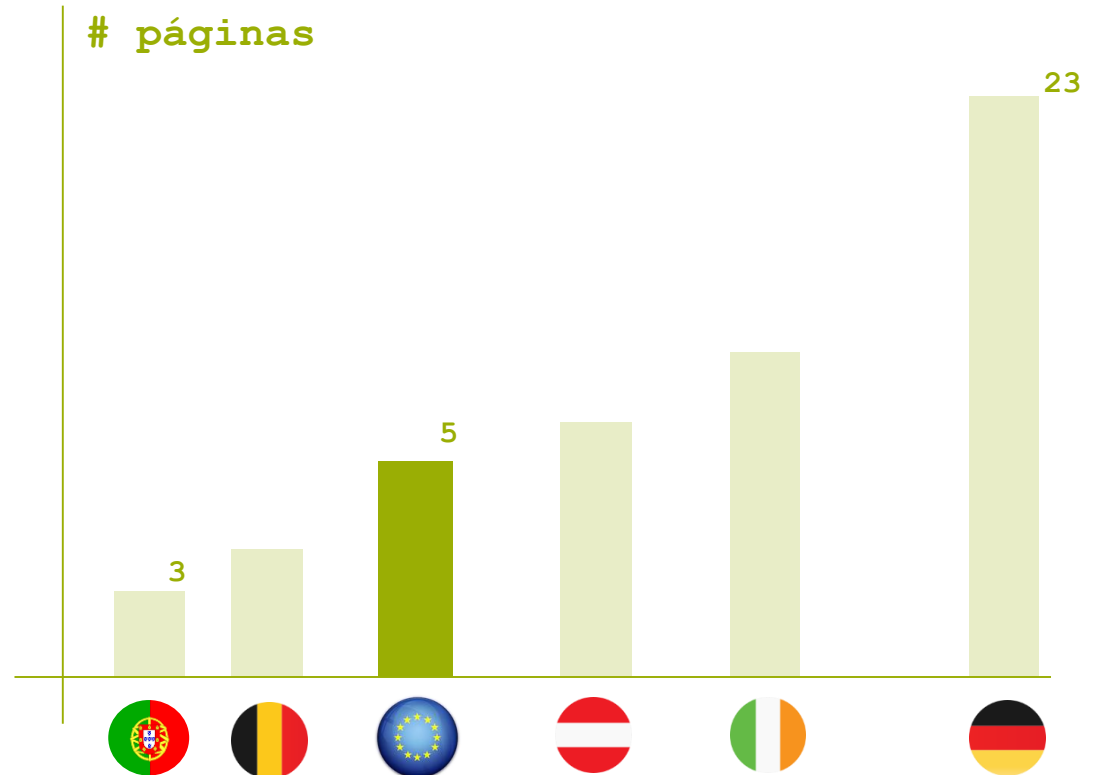
Cada estado informa (opcionalmente) a la Comisión de cada transposición.



# y la idiosincrasia...

## La directiva, sus transposiciones y el contenido.

Quizás es una medida del tamaño de las palabras, el amor por la literatura o la capacidad de síntesis en Europa.



QUE HABLE TODOS LOS IDIOMAS DE LA COMISIÓN

QUE INTERPRETE LAS TRANSPOSICIONES Y LAS COMPARE CON LA DIRECTIVA

“NECESITAMOS UNA MÁQUINA



QUE HAGA MAGIA CON TODOS ESTOS PAPELES.”



QUE MIDA LA CONFORMIDAD DE LA TRANSPOSICIÓN CUANTITATIVAMENTE

QUE DETECTE INCONFORMIDADES O GOLD-PLATING



# la primera: analítica clásica

## Natural Language Processing (NLP).

Implementamos un modelo con diccionarios, stop words, lematización ... para matrices de comparación “artículo vs párrafo”.

	Art1	Art2	Art3	Art4	Art5	Art6	Art7	Art8	Art9	Art10
Paragraph1	14%	8%	8%	10%	12%	8%	8%	10%	12%	10%
Paragraph2	10%	8%	10%	10%	8%	8%	8%	10%	10%	18%
Paragraph3	8%	13%	7%	24%	13%	4%	8%	3%	6%	16%
Paragraph4	10%	8%	9%	8%	13%	10%	8%	13%	9%	12%
Paragraph5	10%	10%	10%	10%	10%	12%	10%	10%	10%	10%
Paragraph6	10%	10%	10%	10%	8%	15%	10%	7%	8%	13%
Paragraph7	12%	9%	9%	9%	9%	12%	9%	12%	9%	9%
Paragraph8	31%	9%	9%	6%	4%	6%	5%	9%	7%	11%
Paragraph9	10%	10%	10%	10%	10%	10%	10%	10%	10%	10%
Paragraph10	10%	10%	10%	8%	8%	10%	8%	12%	8%	12%
Paragraph11	20%	5%	8%	6%	9%	24%	8%	7%	7%	6%
Paragraph12	13%	9%	9%	11%	9%	11%	13%	9%	9%	11%
Paragraph13	14%	9%	9%	9%	12%	9%	9%	9%	9%	9%
Paragraph14	9%	9%	9%	9%	9%	9%	12%	9%	12%	12%
Paragraph15	10%	10%	10%	10%	10%	10%	10%	12%	10%	10%

Figure 5 - Example of compliance matrix

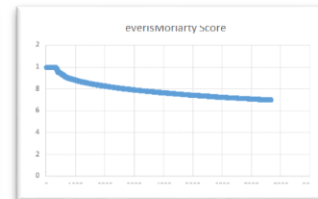
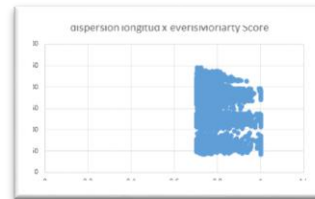
Esta ola nos pasó por encima...



# la segunda: big data



## Analítica y Big Data.



Implementamos un proceso con técnicas Big Data para procesar comparaciones de n-gramas de longitud variable de 1 a 250.

cadena_orig	everisMoriarty_score	longitud	cadena_origen	cadena_final
6301	0,768220769	176	that the period for payment fixed in the contract does not exceed the time limits provided for in paragraph 3, unless otherwise expressly agreed in the contract and provided it	period for payment fixed in a relevant contract shall not exceed the time limits provided for in section 6 unless otherwise expressly agreed contract and provided such
6344	0,768867325	204	means to increase awareness of the remedies for late payment among undertakings. 4. Member States may encourage the establishment of prompt payment codes which set out clearly defined payment time limits	Increase awareness of the remedies for late payment among undertakings. Establishment of codes. 17. The Minister may encourage the prompt payment codes which set out clearly defined payment time
6257	0,768786127	173	30 calendar days from the date of receipt of the goods or services, unless otherwise expressly agreed in the contract and any tender documents; and provided it is not grossly	exceed 30 calendar days from the date of receipt of the goods or services, unless otherwise expressly agreed in the relevant contract and not grossly unfair
2724	0,768707483	139	national central bank; (8) "amount due" means the principal sum which should have been paid within the contractual or statutory period;	substitute the following definitions: "amount due" means the principal sum which should have been paid within the contractual period
3862	0,768595041	120	date of receipt of the goods or services, unless otherwise expressly agreed in the contract and any tender documents and	from the date of receipt of the goods or services, unless otherwise expressly agreed in the relevant contract and provided
3906	0,768595041	119	fixed in the contract does not exceed the time limits provided for in paragraph 3, unless otherwise expressly agreed in	fixed in a relevant contract shall not exceed the time limits provided for in section 6 unless otherwise agreed
3932	0,768595041	113	objectively justified in the light of the particular nature or features of the contract, and that it in any event	agreed is objectively justified in the light of the particular nature or features of the relevant contract in question, and
3996	0,768595041	121	paid by the agreed date, interest and compensation provided for in this Directive shall be calculated solely on the basis	date, shall attract the interest and compensation provided for in this Act, and shall be calculated solely on the basis
5448	0,768361582	172	goods or services; (ii) where the date of receipt of the invoice or the equivalent request for payment is uncertain, 30 calendar days after the date of receipt of the	receipt of the invoice or an equivalent request for payment is uncertain, 30 calendar days immediately following the date of receipt of the supplied by the creditor
5868	0,768361582	172	goods or services; (ii) where the date of receipt of the invoice or the equivalent request for payment is uncertain, 30 calendar days after the date of receipt of the	supplied by the creditor
3267	0,767441886	120	days, unless otherwise expressly agreed in the contract and provided it is not grossly unfair to the creditor within the	services, unless otherwise expressly agreed in the relevant contract and provided that it is not grossly unfair to the creditor. (2)
			not exceed 30 calendar days from the date of receipt of the goods or services, unless otherwise expressly agreed in the contract and	of that procedure shall not exceed 30 calendar days from the date of receipt of the goods or services, unless otherwise expressly agreed



## Analítica, Big Data y Ciencia.

Apliquemos nuevos avances en investigación analítica e introduzcamos las mejoras que nos lleven a una solución óptima.

### Citation Pattern Matching Algorithms for Citation-based Plagiarism Detection: Greedy Citation Tiling, Citation Chunking and Longest Common Citation Sequence

Bela Gipp

Plagiarism Detection Reduced to String Matching

Norman Meuschke

U, Germany & UC Berkeley, USA  
meuschke@berkeley.edu

### Winnowing: Local Algorithms for Document Fingerprinting

Saul Schleimer  
MSCS  
University of Illinois, Chicago  
saul@math.uic.edu

Daniel S. Wilkerson  
Computer Science Division  
UC Berkeley  
dsw@cs.berkeley.edu

Alex Aiken  
Computer Science Division  
UC Berkeley  
aiken@cs.berkeley.edu

**Abstract.** The number of students steadily increasing at the same time is readily available. There is a significant variety of reasons – take advantage illicitly copy other students' projects, disguise their deception. Software plagiarism is therefore of immense assistance prevent – such abuse. We design to such software. Our algorithm very efficient means to detect plagiarism over sparse suffix trees that allow program file against all the files in

#### ABSTRACT

Digital content is for copying: quotation, revision, plagiarism, and file-sharing all create copies. Document fingerprinting is concerned with accurately identifying copying, including small partial copies, within large sets of documents.

We introduce the class of *local* document fingerprinting algorithms, which seems to capture an essential property of any fingerprinting technique guaranteed to detect copies. We prove a novel lower bound on the performance of any local algorithm. We also develop *winnowing*, an efficient local fingerprinting algorithm, and show that winnowing's performance is within 33% of the lower bound. Finally, we also give experimental results on Web data, and report experience with MOSS, a widely-used plagiarism detection service.

A do run run run, a do run run  
(a) Some text from [7].

adorunrunrunadorunrun  
(b) The text with irrelevant features removed.

adoru dorun orunr runru unrun nrunr runru  
unrun runru runad unadador adoru dorun  
orunr runru unrun  
(c) The sequence of 5-grams derived from the text.

77 72 42 17 98 50 17 98 8 88 67 39 77 72 42  
17 98  
(d) A hypothetical sequence of hashes of the 5-grams.

- Scoring global.
- Comparación real-time para el usuario: funcionalidad de subrayado.
- Filtrado de ruido textual.
- Optimización del rendimiento.
- Optimización del balance de los filtros eM Score - longitud
- Optimización adaptada a cada lenguaje.



**Plan TL**

Plan de Impulso de las  
Tecnologías del Lenguaje

everis

an NTT DATA Company

1. Cada caso de uso tiene su camino. Y a veces no es evidente.
2. Se requiere imaginación y la combinación de ciencia, tecnología y analítica.



# GRACIAS

