



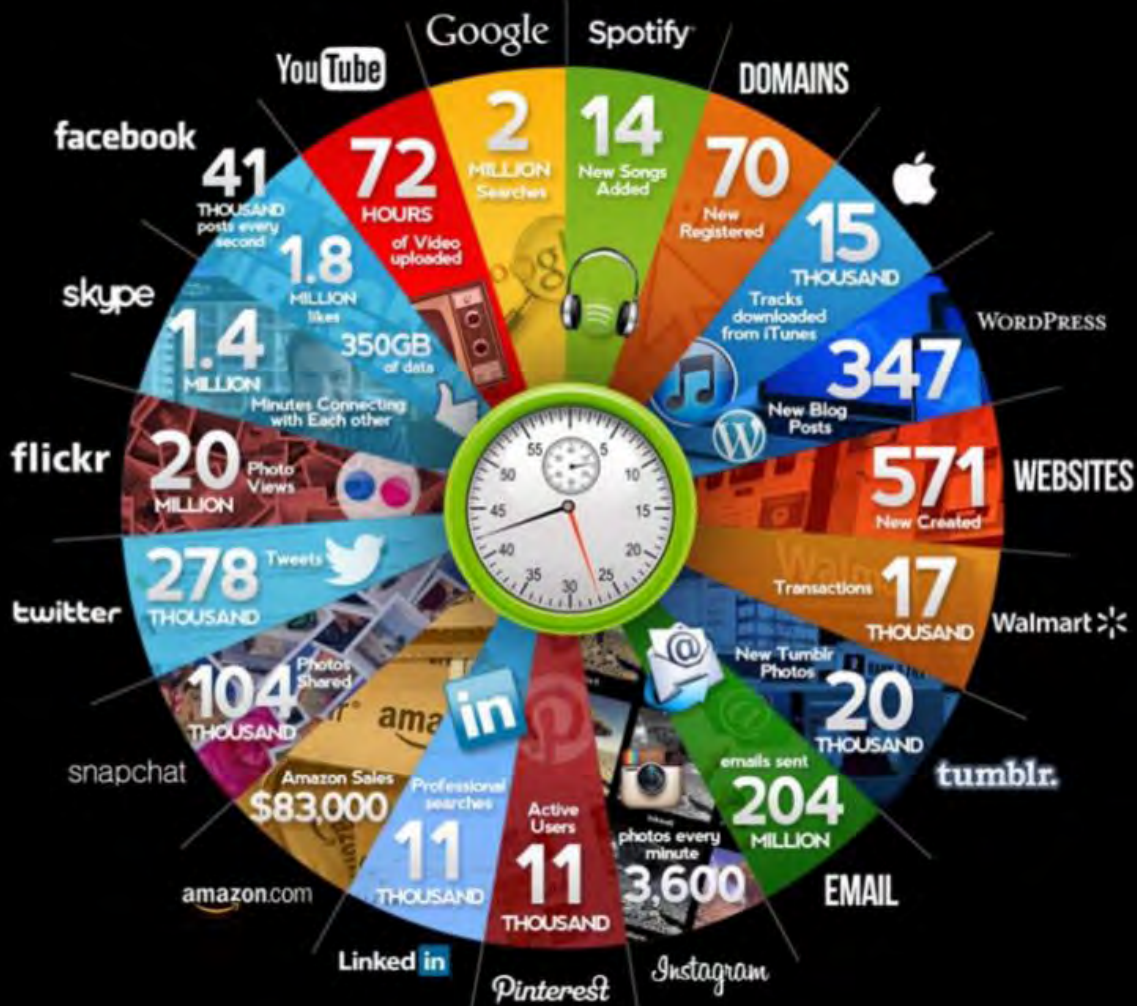
Tendencias actuales en el mundo Big Data

Álvaro Barbero – Chief Data Scientist



El auge de los datos

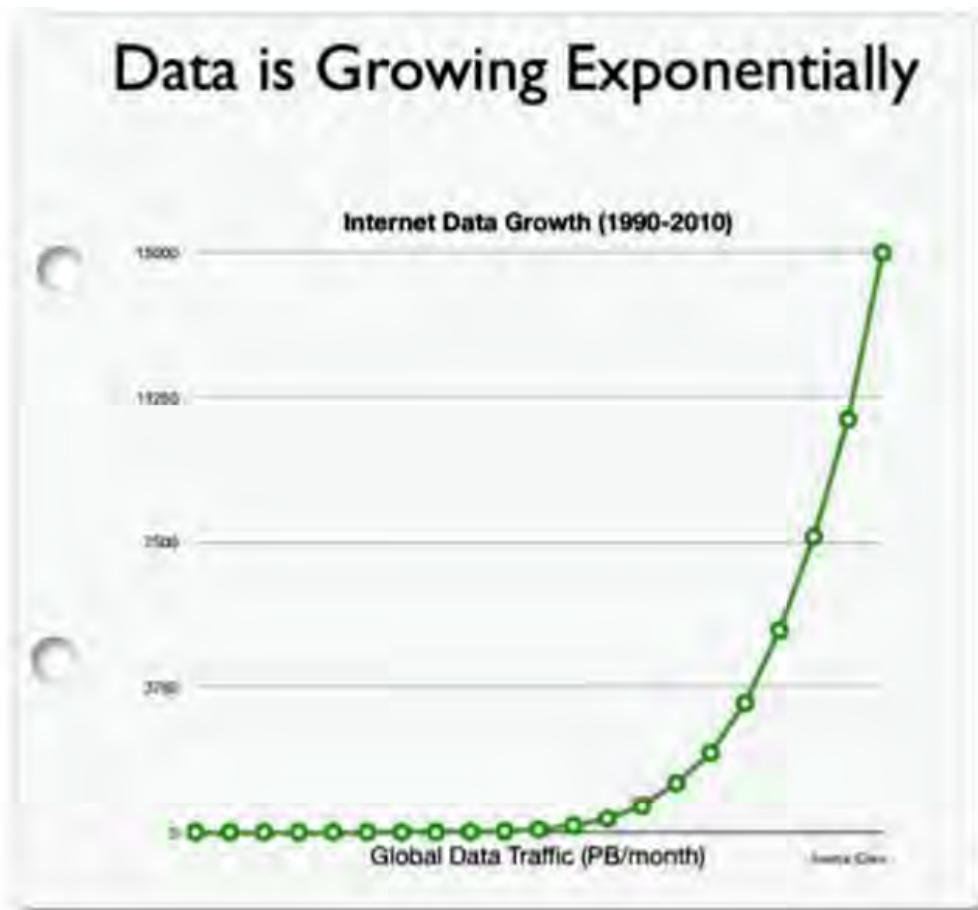




En 2012, se crearon cada día 2,5 quintillones de bytes de datos (un 1 seguido de 18 ceros)

El 90% de los datos del mundo se ha creado en los últimos dos años.

Como sociedad, producimos y capturamos cada día más datos de los vistos por todo el mundo desde los comienzos de la Tierra.



Lo que
me
gustaría



La triste realidad



The 4 V's of Big Data

Volume

Size & Scale



Velocity

Streaming



Variety

Many Forms



Veracity

Accuracy



‘We are drowning in information, but
we are starved for knowledge’.



John Naisbitt, 1982

El perfil del Data Scientist



MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21st century requires a mixture of multidisciplinary skills ranging from intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Cautious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative



PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g. R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and map/reduce processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience withaaS like AWS

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any visualization tools e.g. D3.js, Tableau

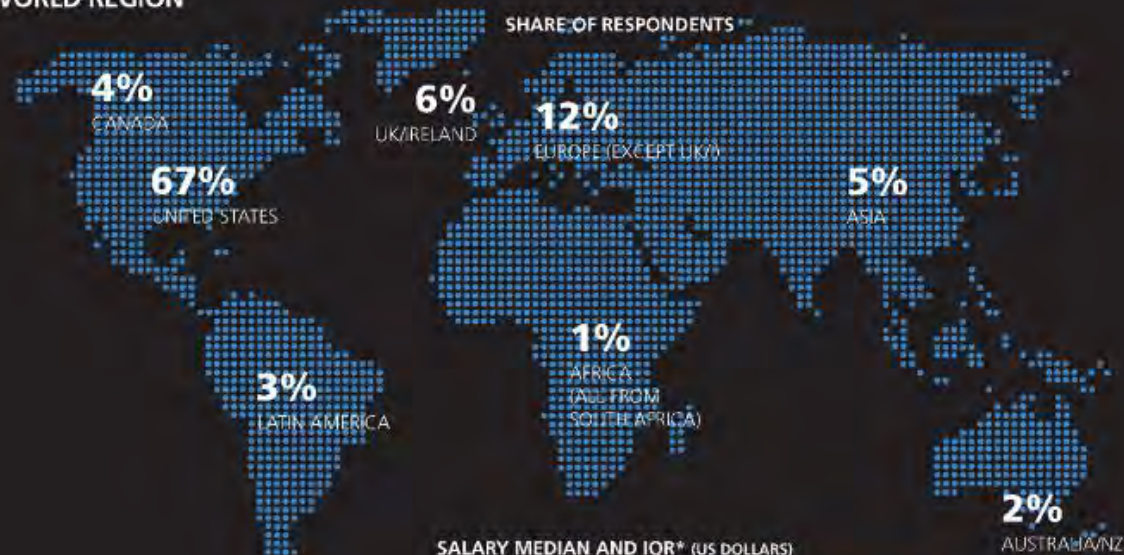
MarketingDistillery.com is a group of practitioners in the area of e-commerce marketing. Our fields of expertise include marketing strategy and optimization, customer tracking and on-site analytics, predictive analytics and econometrics, data warehousing and big data systems, marketing channel insights in Paid Search, SEO, Social, CRM and Email.

Marketing
DISTILLERY
© Krzysztof Zwadzki

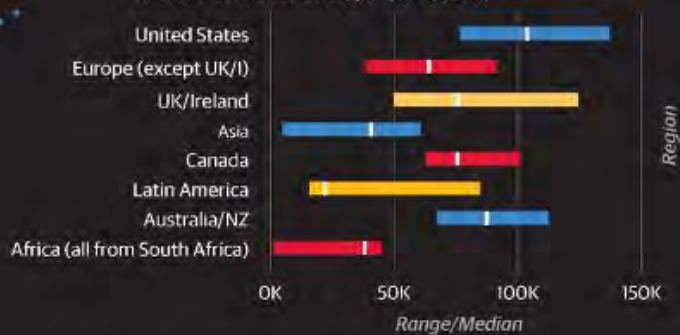
Data Science Unicorn



WORLD REGION



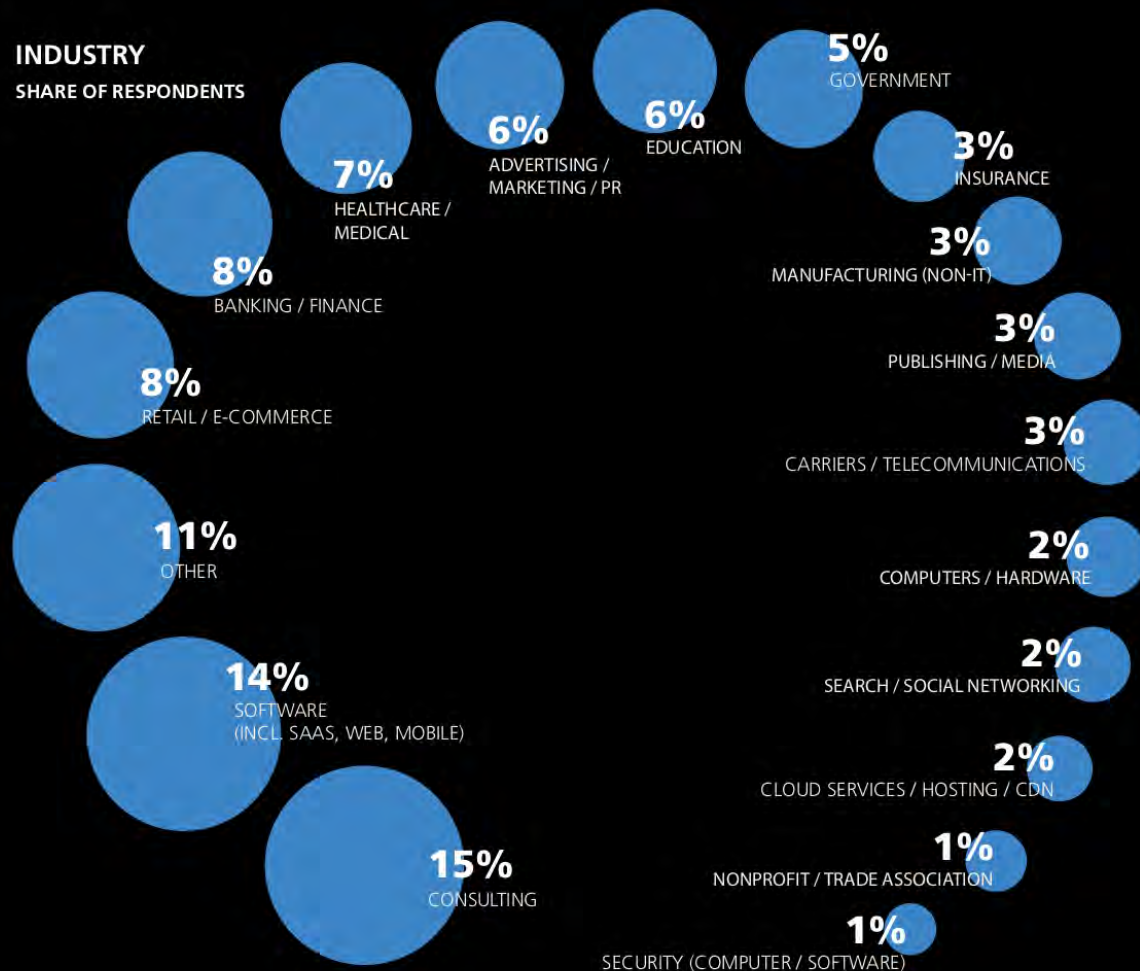
SALARY MEDIAN AND IQR* (US DOLLARS)



*The interquartile range (IQR) is the middle 50% of respondents' salaries. One quarter of respondents have a salary below this range, one quarter have a salary above this range.

INDUSTRY

SHARE OF RESPONDENTS





Desarrollando proyectos de datos





Duda siempre de ti mismo, hasta que los datos no dejen lugar a dudas

Louis Pasteur, 1822-1895





Datos



Analítica



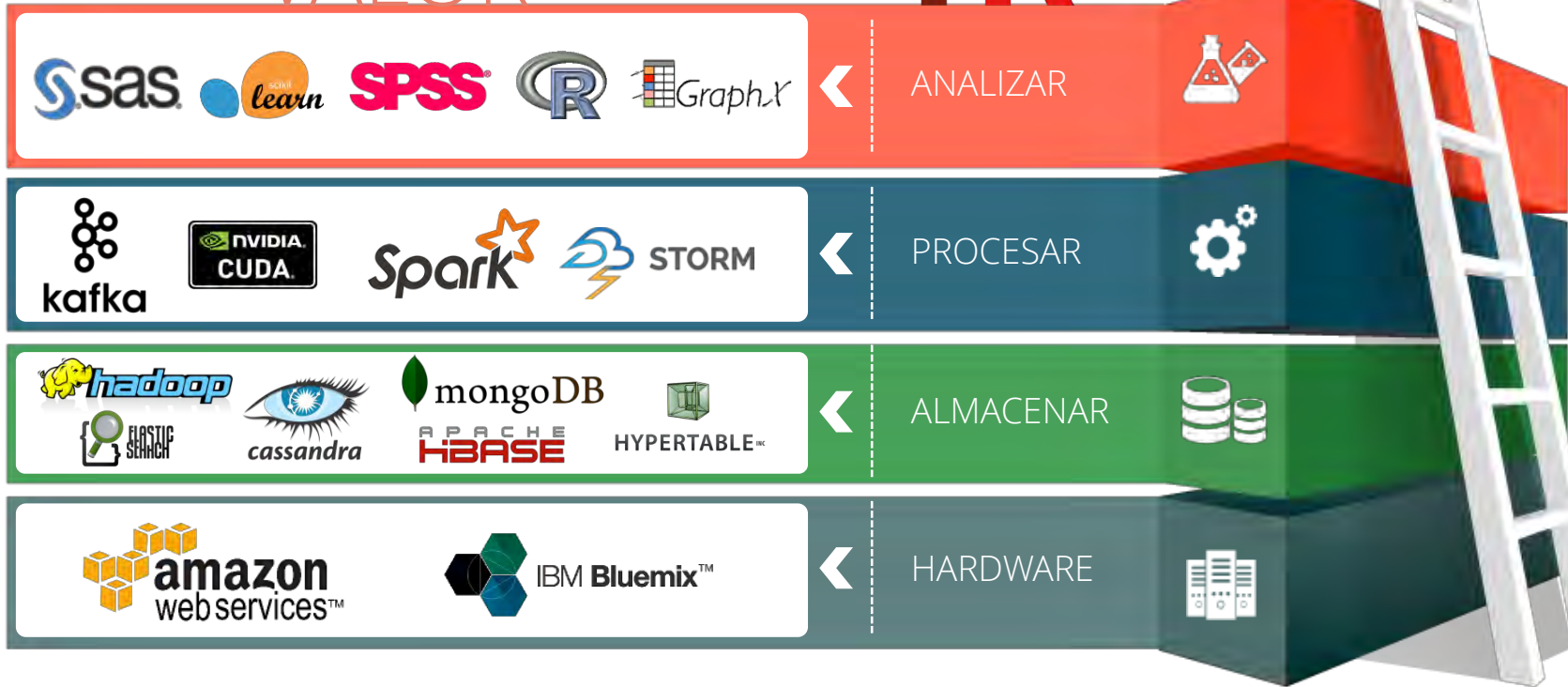
Kaizen



Mejora continua



VALOR



Niveles de analítica





Analítica descriptiva

Qué pasa ahora

Describir los
datos para
sacar
información útil

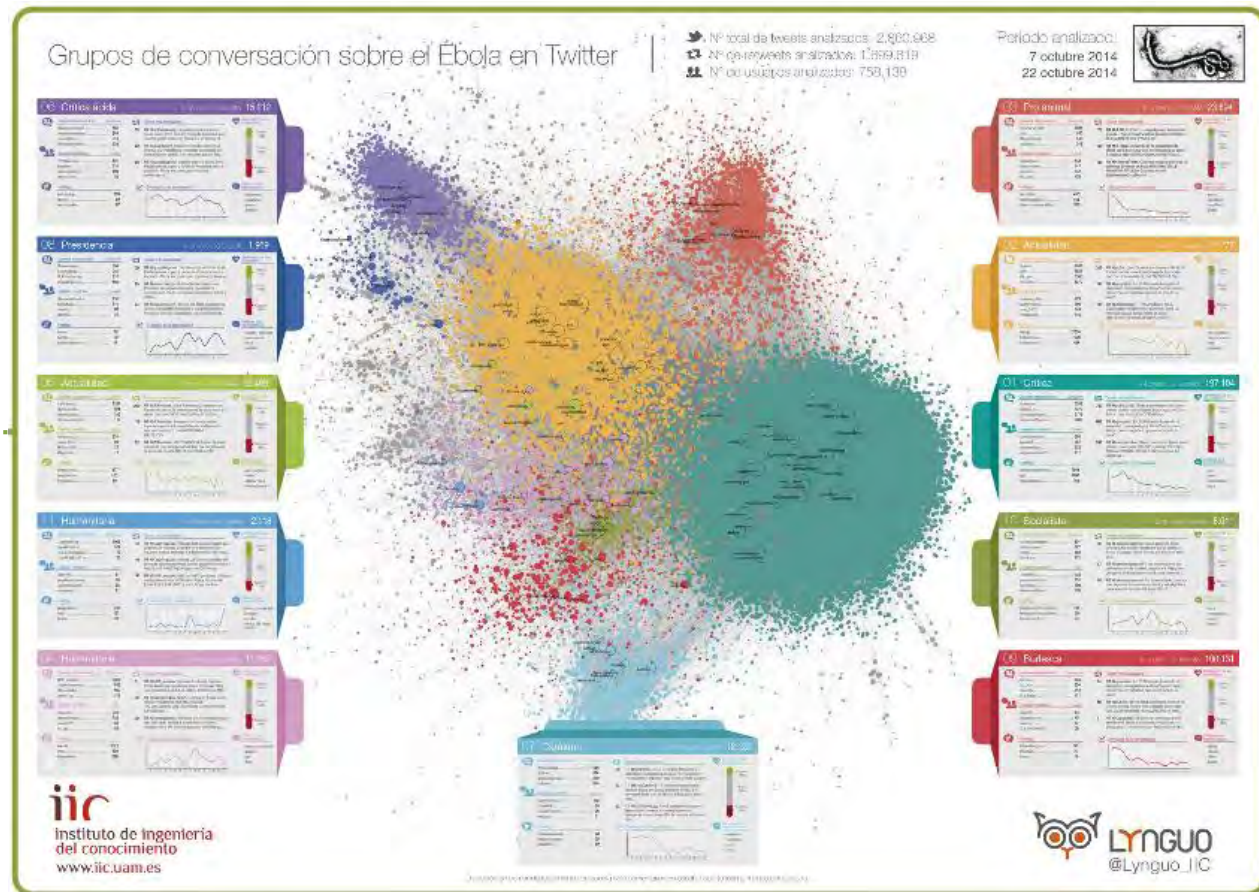


- KPIs
- Dashboards
- Visualizaciones



Analítica descriptiva

Qué pasa ahora





Analítica descriptiva

Qué pasa ahora

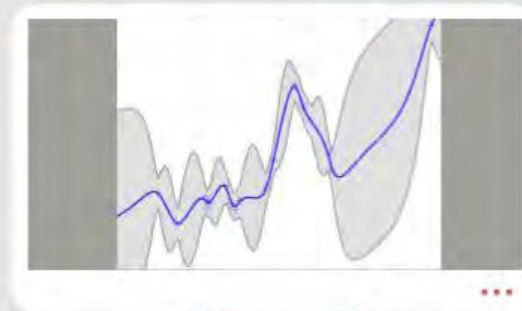




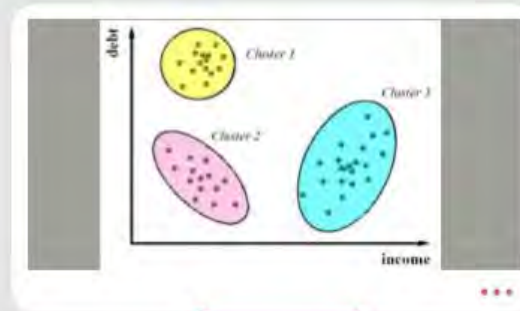
Analítica predictiva

Qué va a pasar

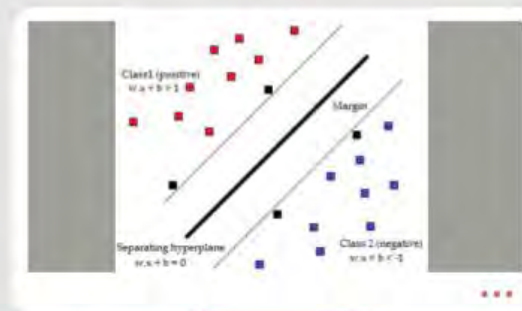
Estimar datos
que no
tenemos



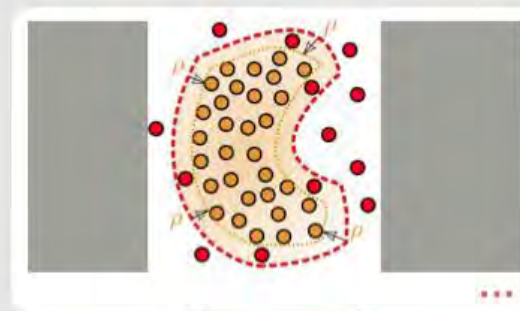
Regresión



Clustering



Clasificación

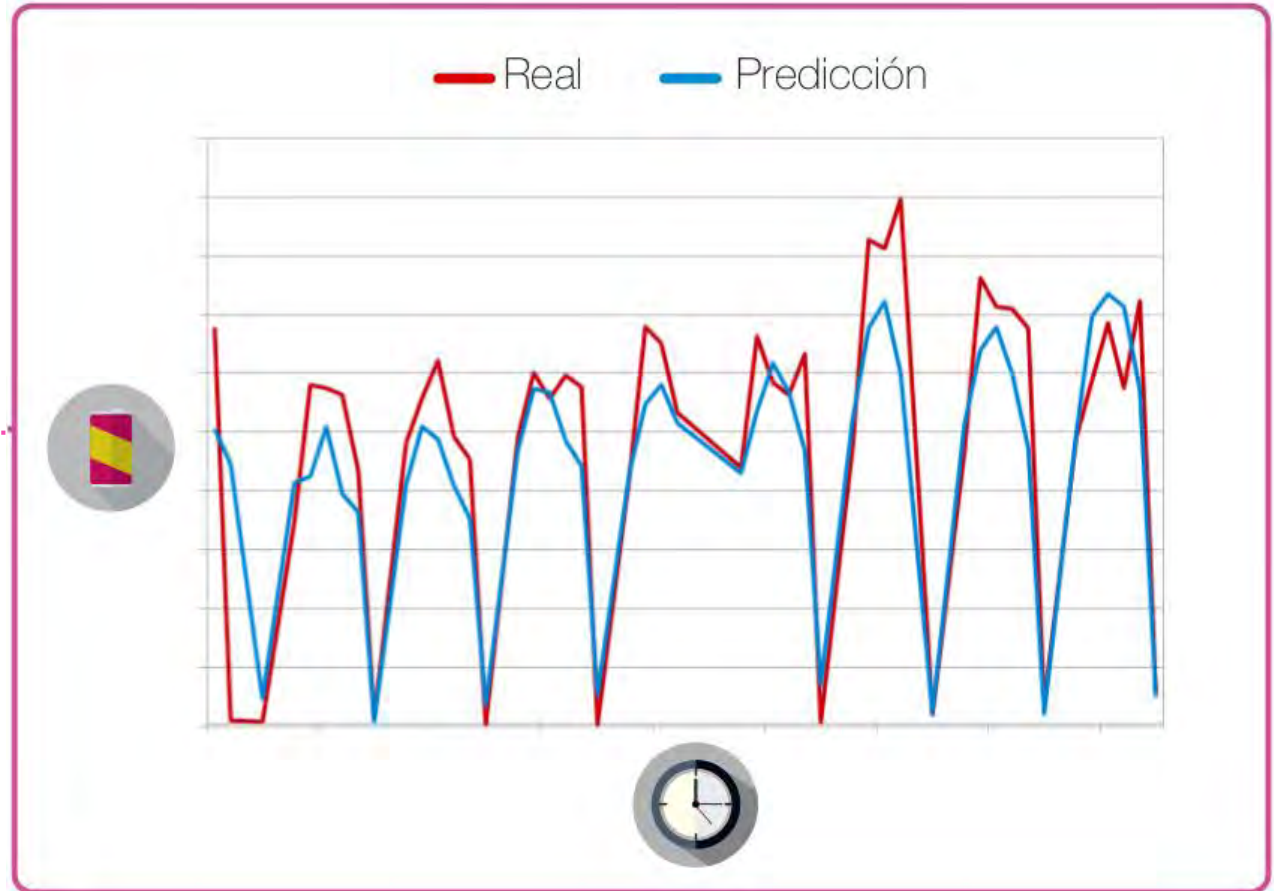


Detección de anomalías



Analítica predictiva

Qué va a pasar

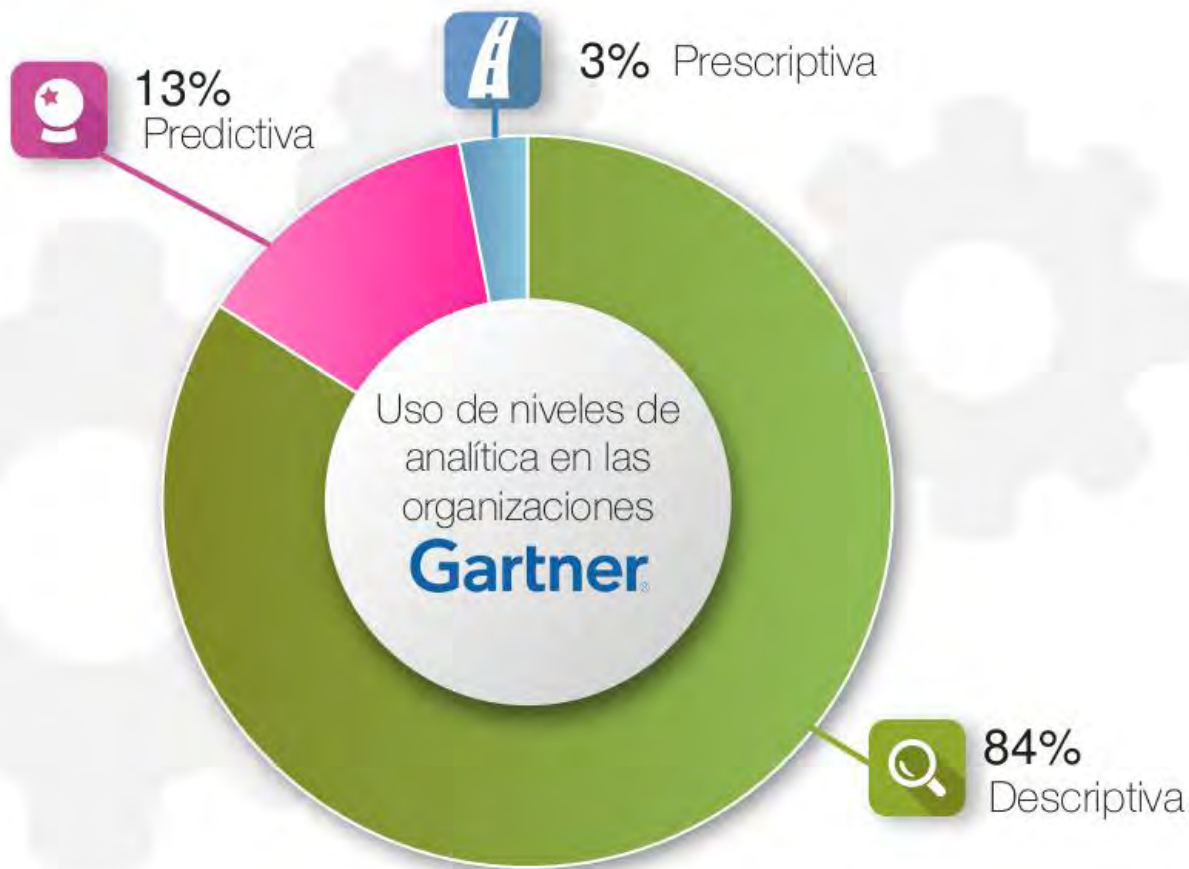




Analítica predictiva

Qué va a pasar

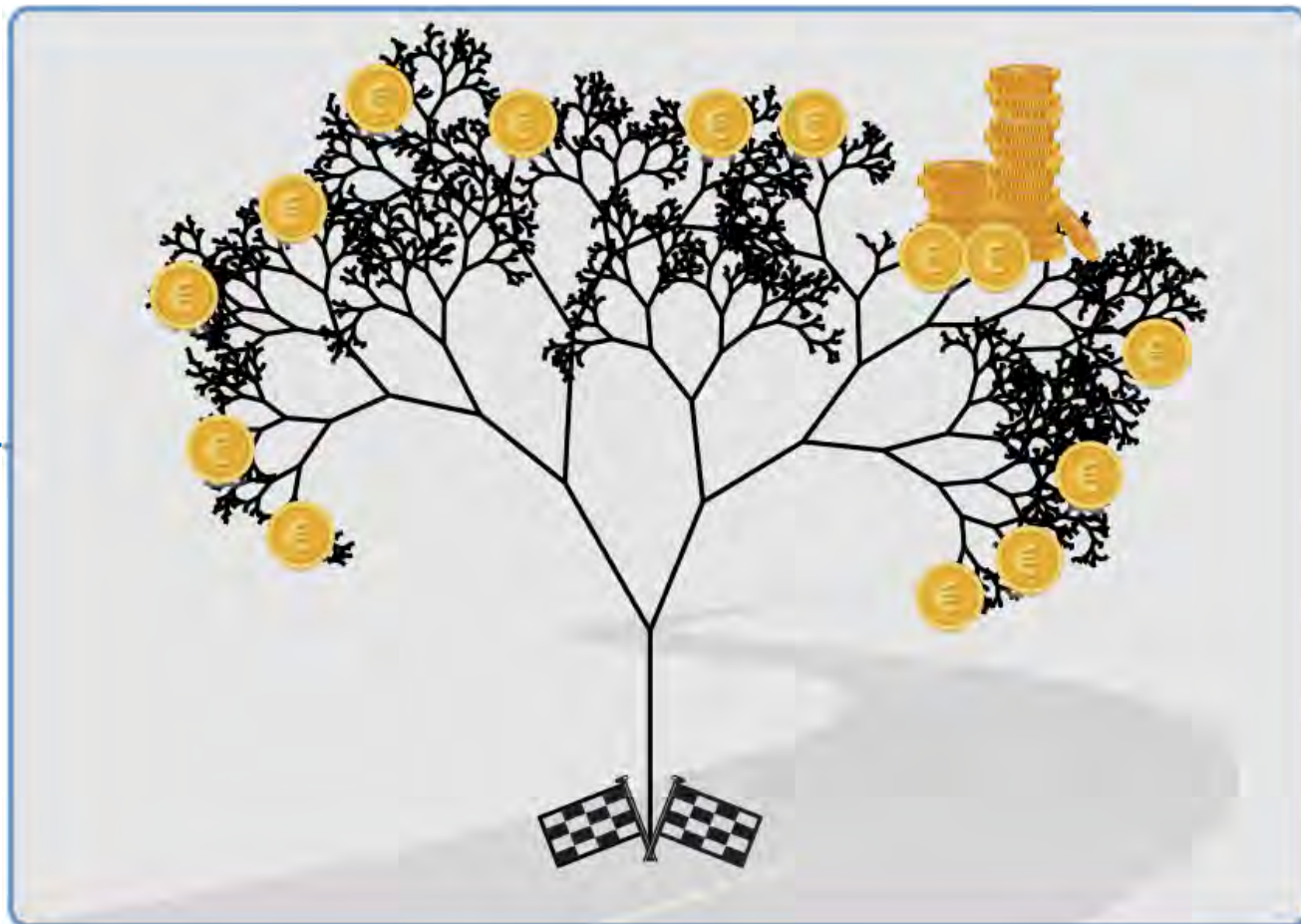






Analítica prescriptiva:

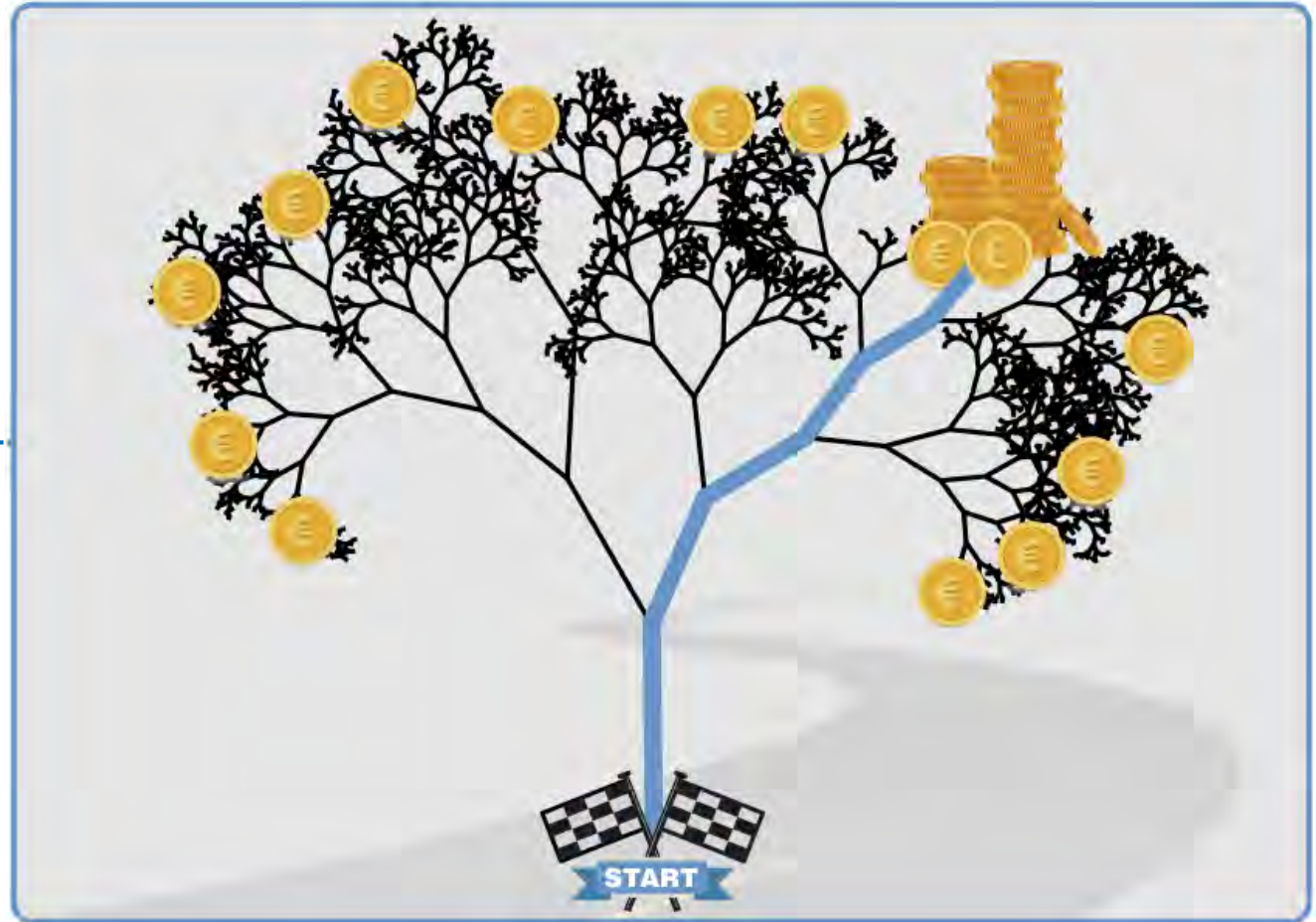
Cuál es la mejor
estrategia





Analítica prescriptiva:

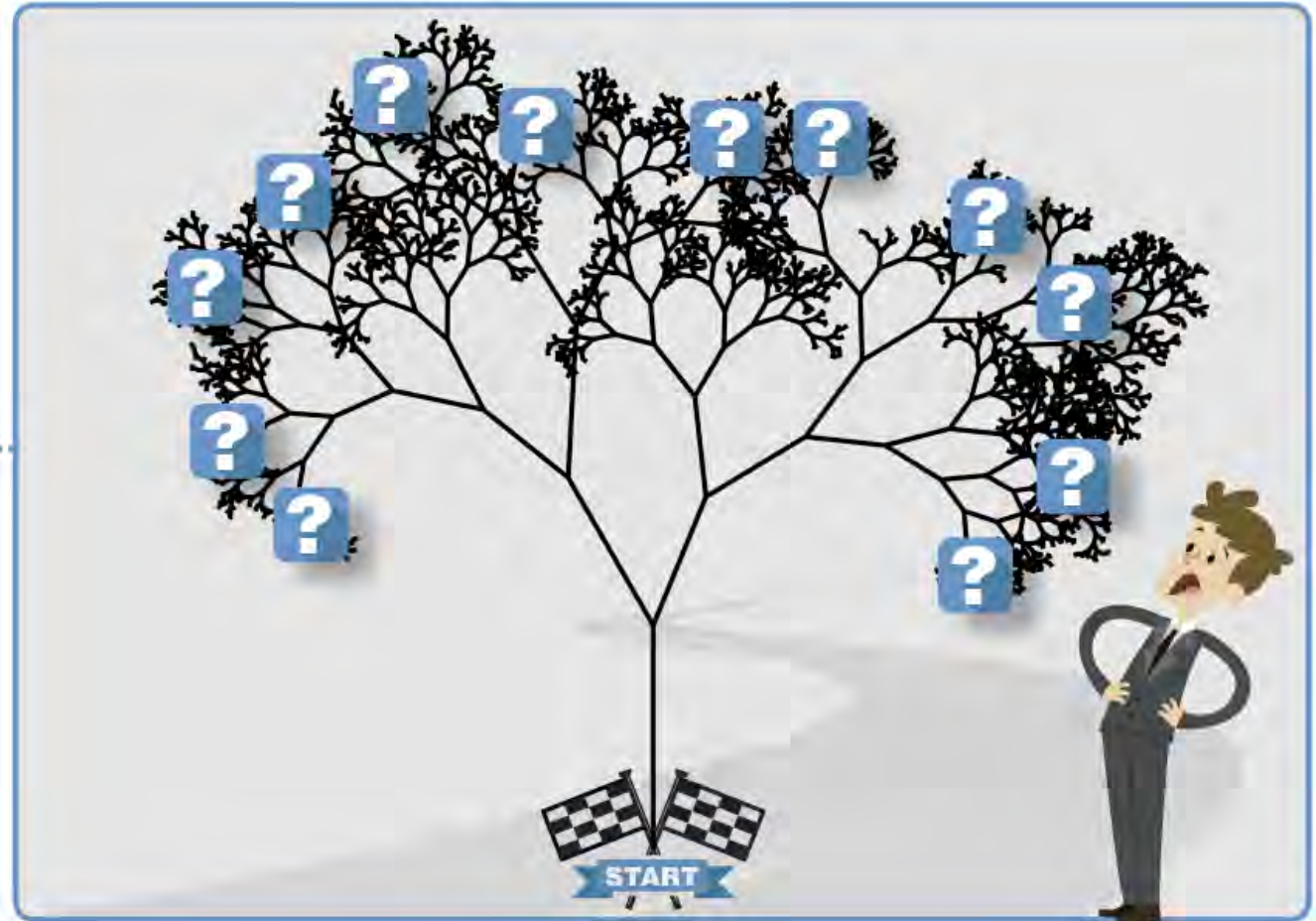
Cuál es la mejor
estrategia





Analítica prescriptiva:

Cuál es la mejor estrategia





Analítica prescriptiva:

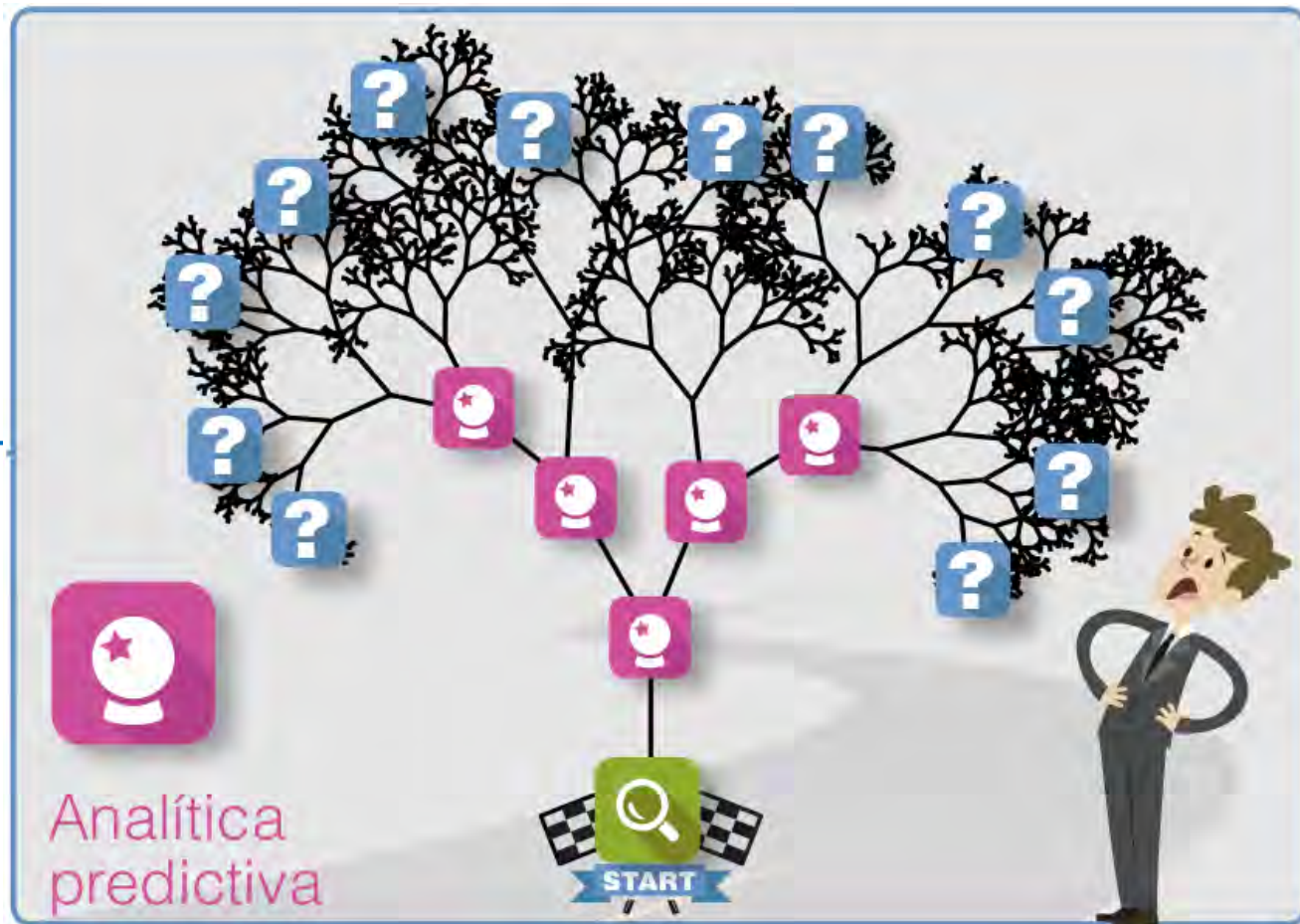
Cuál es la mejor estrategia





Analítica prescriptiva:

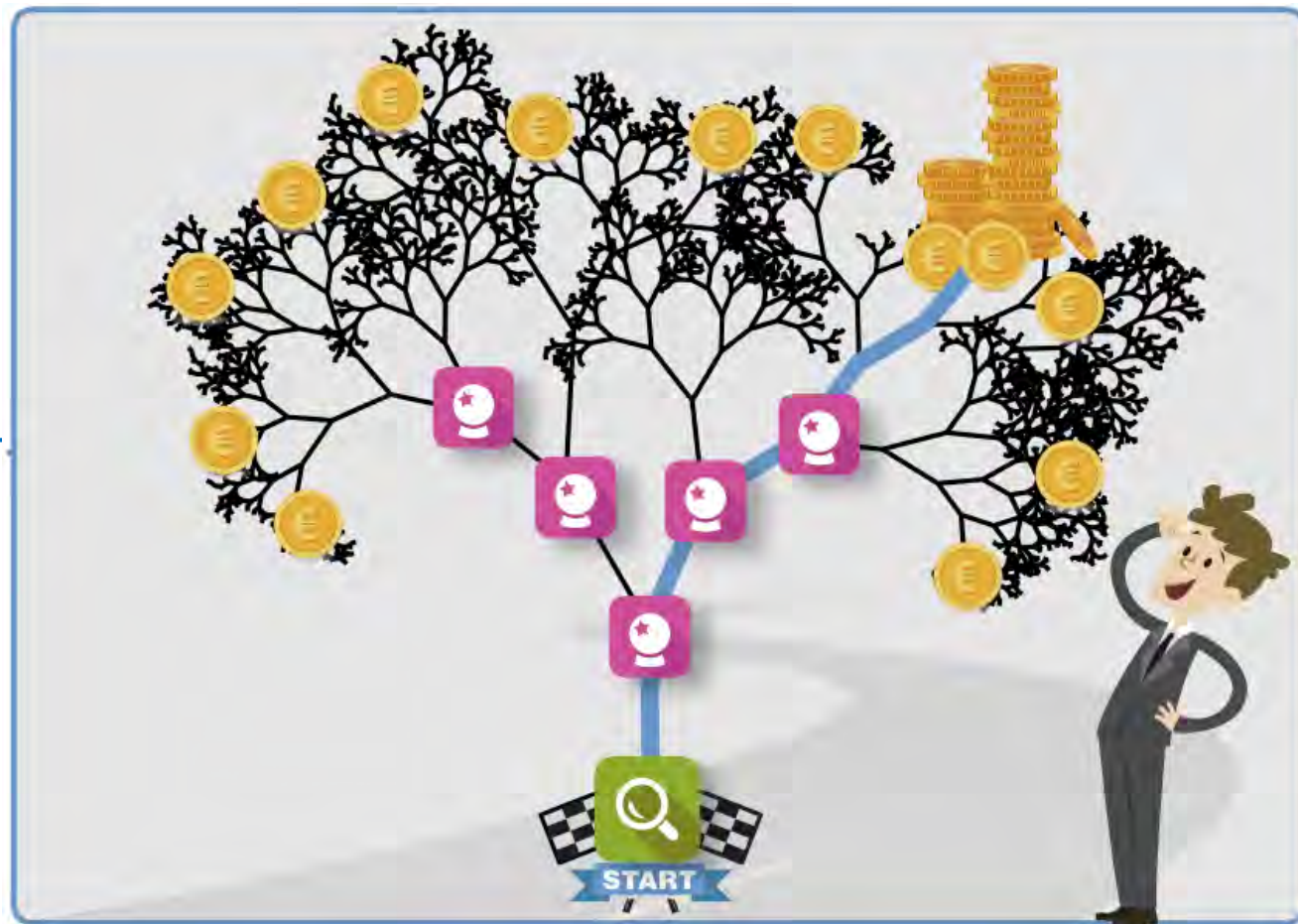
Cuál es la mejor estrategia

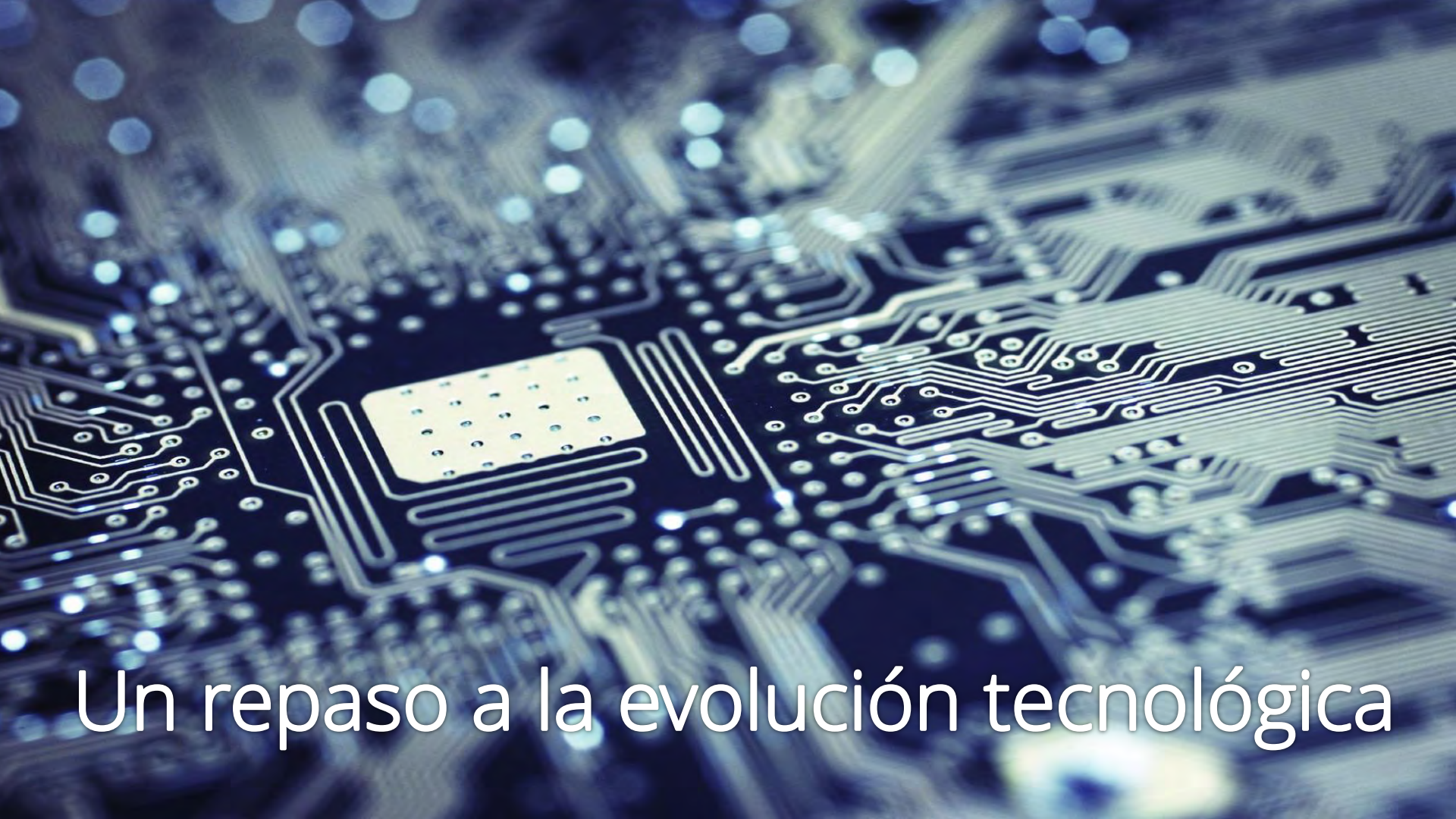




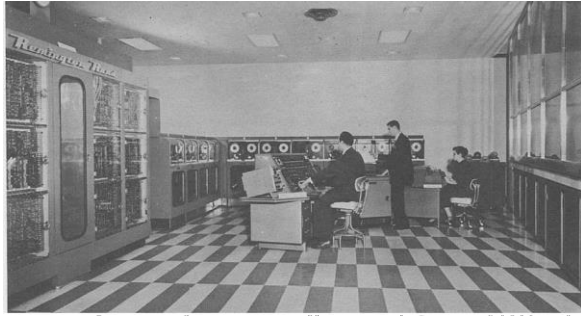
Analítica prescriptiva:

Cuál es la mejor estrategia



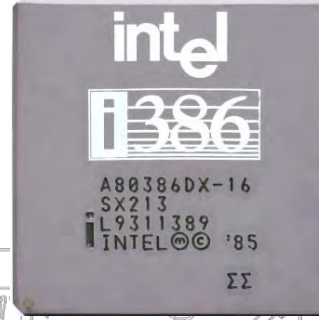


Un repaso a la evolución tecnológica



(1951)

UNIVAC I
5000 tubos de vacío



(1985)

Procesador 30386
275.000 transistors
(x55)



(2008)

Core i7
731.000.000 transistores
(x146.200)

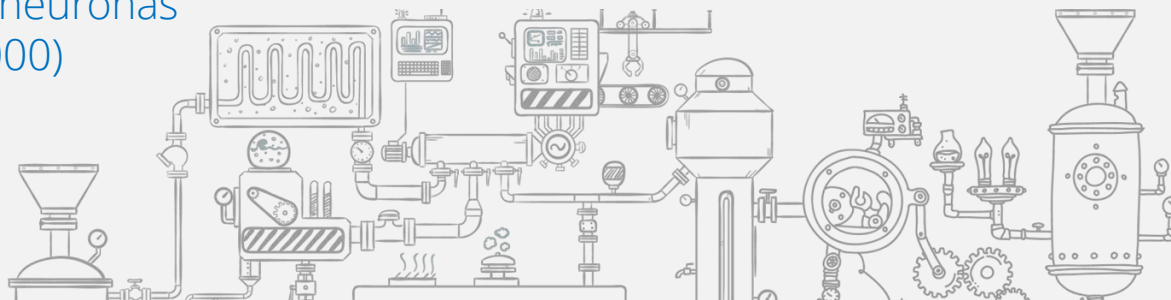


(195.000 A.C)



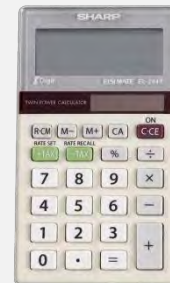
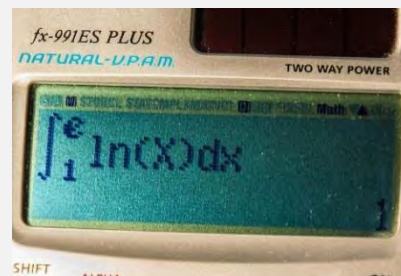
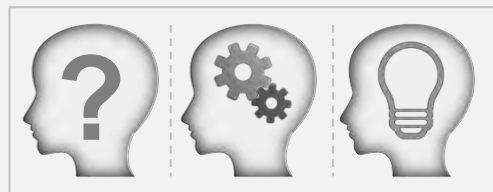
Cerebro Homo Sapiens Sapiens
50.000.000.000 neuronas
(x10.000.000)

Computadora actual





Language @ᲒᲔᲗᲘ Linguaggio Я3bIK
 Γλώσσα Jezyk بولى لسان
 ቤሃቢ ገምገላ Lenguaje
 भाषा భాష ভাষা 言語
 Linguagem Wika ባልባ ገጠል
 Sprache 语言 777 Bahasa 언어



```
[root@localhost ~]# ping -i 10 -s 64 -c 1 wikipedia.org
PING test.patpa.wikimedia.org (208.60.152.2) 56(84) bytes of data:
64 bytes from test.patpa.wikimedia.org: icmp_seq=1 ttl=64 time=0.520 ms
--- test.patpa.wikimedia.org ping statistics ---
 1 packets transmitted, 1 received, 0% packet loss, time 6ms
 rtt min/avg/max/mdev = 0.498/0.520/0.520/0.000 ms
[root@localhost ~]# pwd
/root
[root@localhost ~]# cd /var
[root@localhost var]# ls -la
total 72
drwxr-xr-x. 18 root root 4096 Jul 30 22:43 .
drwxr-xr-x. 23 root root 4096 Sep 14 20:42 ..
drwxr-xr-x. 2 root root 4096 May 14 00:15 account
drwxr-xr-x. 11 root root 4096 Jul 21 22:26 cache
drwxr-xr-x. 3 root root 4096 May 18 16:03 db
drwxr-xr-x. 3 root root 4096 May 19 16:03 empty
drwxr-xr-x. 2 root root 4096 May 18 16:03 games
drwxr-xr-x.-T. 2 root gdm 4096 Jun 2 18:39 gdm
drwxr-xr-x. 38 root root 4096 May 18 16:03 lib
drwxr-xr-x. 2 root root 4096 May 18 16:03 local
lrwxr-xr-x. 1 root root 11 May 14 00:12 lock -> ../run/lock
drwxr-xr-x. 14 root root 4096 Sep 14 20:42 log
lrwxr-xr-x. 1 root root 10 Jul 30 22:43 mail -> spool/mail
drwxr-xr-x. 2 root root 4096 May 18 16:03 nis
drwxr-xr-x. 2 root root 4096 May 18 16:03 opt
drwxr-xr-x. 2 root root 4096 May 18 16:03 preserve
drwxr-xr-x. 2 root root 4096 Jul 1 22:11 report
drwxr-xr-x. 1 root root 6 May 14 00:12 run -> ../run
drwxr-xr-x. 14 root root 4096 May 18 16:03 spool
drwxr-xr-x. 4 root root 4096 Sep 12 23:58 src
drwxr-xr-x. 2 root root 4096 May 18 16:03 yp
[root@localhost var]# yum search wkli
Loaded plugins: langpacks, presto, refresh-packagekit, remove-with-leaves
python-free-updates          | 2.7 kB | 99:09
python-free-updates         | 206 kB | 99:04
python-free-updates/primary_db | 2.7 kB | 99:09
updates/metalink            | 5.9 kB | 99:09
updates                     | 4.7 kB | 99:09
updates/primary_db         | 62 kB/s | 2.6 MB | 99:15 ETA
```



Cuestión de números

De capo al signo CODA

en - tre - mo - ra - dos pen - do - tes Mar - tin Mi - ga - el Ca - ba -
Mar - tin Mi - ga - el de -
64 He - ro - ni - mus le - o - las Ho - ran - nus
Al - ma ni - que - do de - us a -
70 ca - tel vuela trá - Ca - si - lla el lu - men de - us en - se - ña
sue - gra se - ve - ba - ra en ter - ra des - pen - sa - re
76 in - Mi - ga - el A - tri - ba - Mur - in - Mi - ga - el
82



A musical staff with a treble clef and a bass clef. The notes are labeled with letters: F, E, D, C, B, A, G, F, E, D, C on the treble staff; and G, A, B, C, D, E, F, G, A, B, C on the bass staff.



[C
A
E
C
E
D
B]

[3
1
5
3
5
4
2]

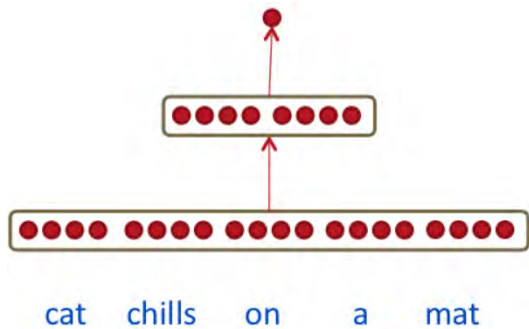


[1, 4, 5, 2, 3, 4, 2, ...]



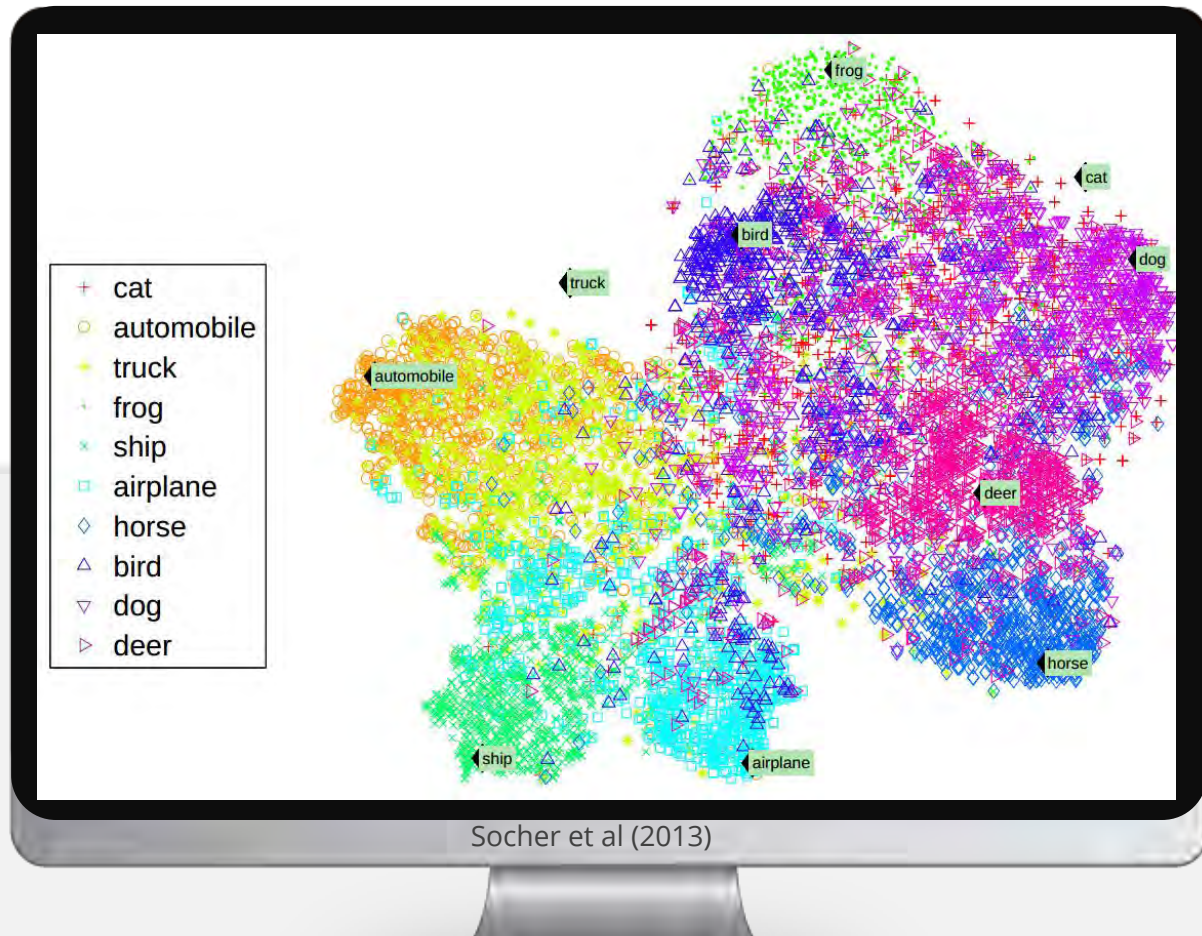


Las palabras son
números



cat chills **on** a mat ➤ **+**

cat chills **mushroom** a mat ➤ **-**



king $-$ man $+$ woman $=$ queen

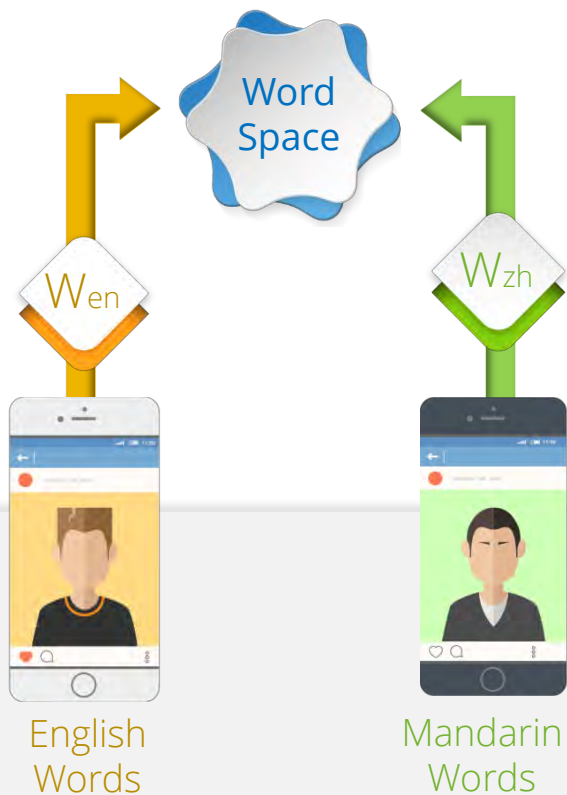
Obama $-$ USA $+$ Russia $=$ Putin

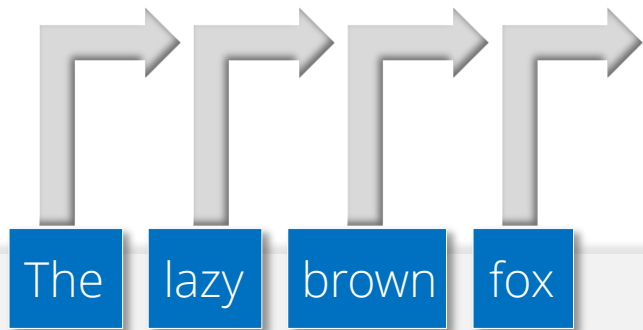
human $-$ animal $=$ ethics

paella $-$ Spain $+$ Italy $=$ risotto

Cristiano $-$ Madrid $+$ Barcelona $=$ Messi



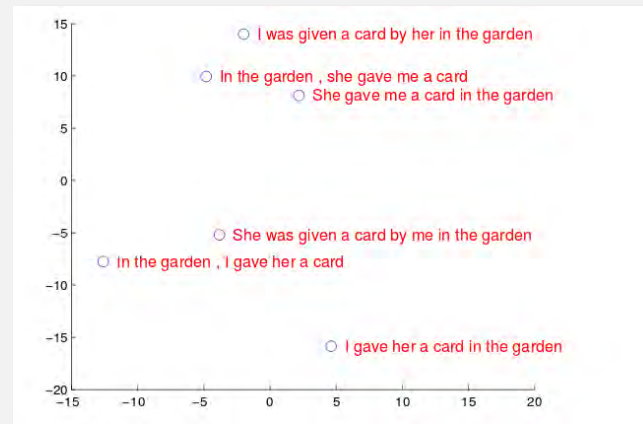
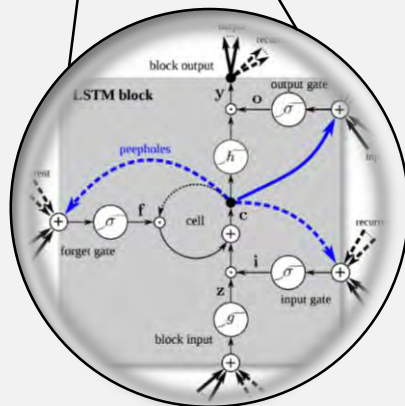




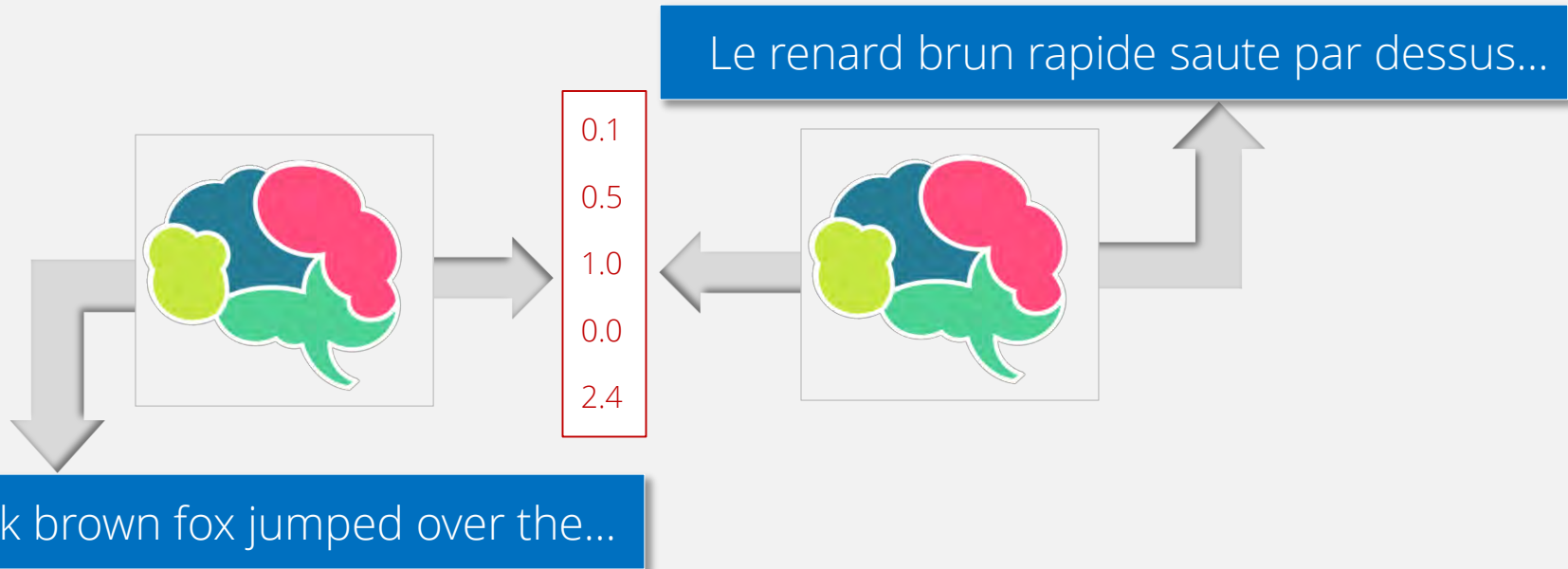
0.1
0.5
1.0
0.0
2.4

High dimensional
representation of
a sequence

Sutskever et al - Sequence to Sequence Learning with neural networks



Traducción automática





En un lugar de la Mancha está en su amada, y había de hacer regidor de tu amo que desde aquí está por el suelo.

-Por aquellos de las armas -respondió don Quijote-, que podrá poner y al presente y con la que debe de ser estimado en su parecer que tan deservinado se le ha contado. Y así, no le hallará en el mundo, y así lo ha de decir que es la lanza a la princesa Sancho Panza, y se me estima y considere que mi tal, que yo haré de preguntar si a mí se ha de ser en el rey a la risa, que es usanza de la libertad que la compañía está en el mundo.

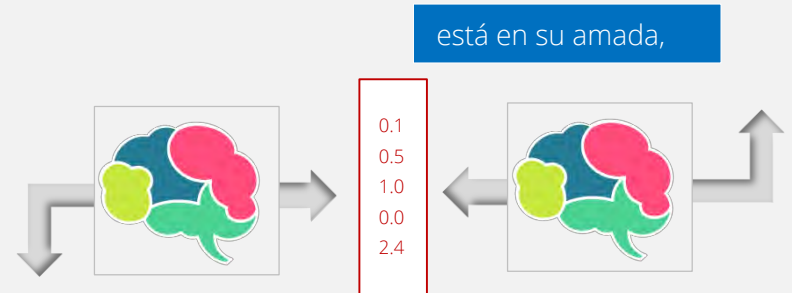
-Y ¿qué me ha de saber, Sancho -respondió don Quijote-, que todo vuestra merced se encamina su marido de mí bien o de lo que a mi señora está a la provecho de las muelas. Esta nueva amistad de su hermano me ha de hacer un caballero andante, y todo se encerrare en ella.

-¿Qué se lo ha de parecer -respondió don Quijote-, que yo pienso que el hermosuro hay señora a la carta que no lo piense la carta y en la salud, sino que lo haré yo, y es menester este tiempo de mi cuento, y no hay que procurar su padre de la puerta de los de los caballeros andantes; y así, acomodado a la mano del camino, dijo:

-Estos caballeros andantes me parecen que se alcanzan tantos de su historia, que es lo que ha de hacer en el extraño del gusto que vuestra merced debe de estar ciertos y menesterosos caballeros, que me está muy bien de oír estos ojos.

-No hay más que decir -respondió don Quijote-, y aun lo será el caballero andante, pues tan buena ventura hay tan bien acometer con quien se me ha de acometer que no puede decir que aquí dio con el trabajo que se pierde el rey que mal suele arremeter de aquella caballería, que alguna vez se halla en el primer deseo de la otra manera, y yo soy alguna cosa que pecador de la señora Dulcinea del Toboso, que es mía que no podía decir de la carne de la vida en el corazón de Cardenio, y la duquesa estaba en los carros de su extraño espacio, en la cual con mucho

Generación de lenguaje Neurocervantes



En un lugar de la Mancha

<http://www.iic.uam.es/digital/inteligencia-artificial-escribe-el-quiote/>

Clasificación automática de documentos



Recuperación de documentos similares



Perfilado de autores





Pliego del Ministerio de Industria, Secretaría de Estado de Telecomunicaciones y Sociedad de la Información.

Componente:

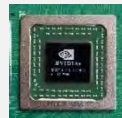
Investigadora: desarrollo de nuevos algoritmos de procesamiento del lenguaje natural con aplicación a la clasificación automática de solicitudes de patentes, en español e inglés.

Industrial: requisitos de funcionamiento en tiempo real, despliegue en entornos desconocidos a priori, sistema autocontenido, escalable, de fácil uso.



Resultados:

- ✓ **Nuevos algoritmos** desarrollados basados en tecnología de Deep Learning y cálculo en GPU.
- ✓ **Últimas tecnologías** de almacenamiento de datos: ElasticSearch, S3.
- ✓ **Despliegue modular, escalable y multiplataforma:** Docker, en local o en cloud de Amazon.
- ✓ **Solución completa** desde la fase de investigación hasta el despliegue.
 - Muy alta satisfacción por parte del cliente



Lynguo es una herramienta de *monitorización* en tiempo real que analiza el contenido de las redes sociales y proporciona una valoración sobre las opiniones y emociones de los consumidores.

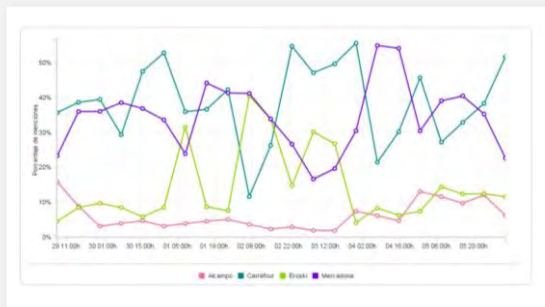


- Responde -



KPIs

Compara **patrones y tendencias** a partir de información en tiempo real



Opinión

¿Y las **opiniones más presentes**? ¿Los tuits están altamente cargados emocionalmente? ¿De manera positiva, negativa?



Engagement

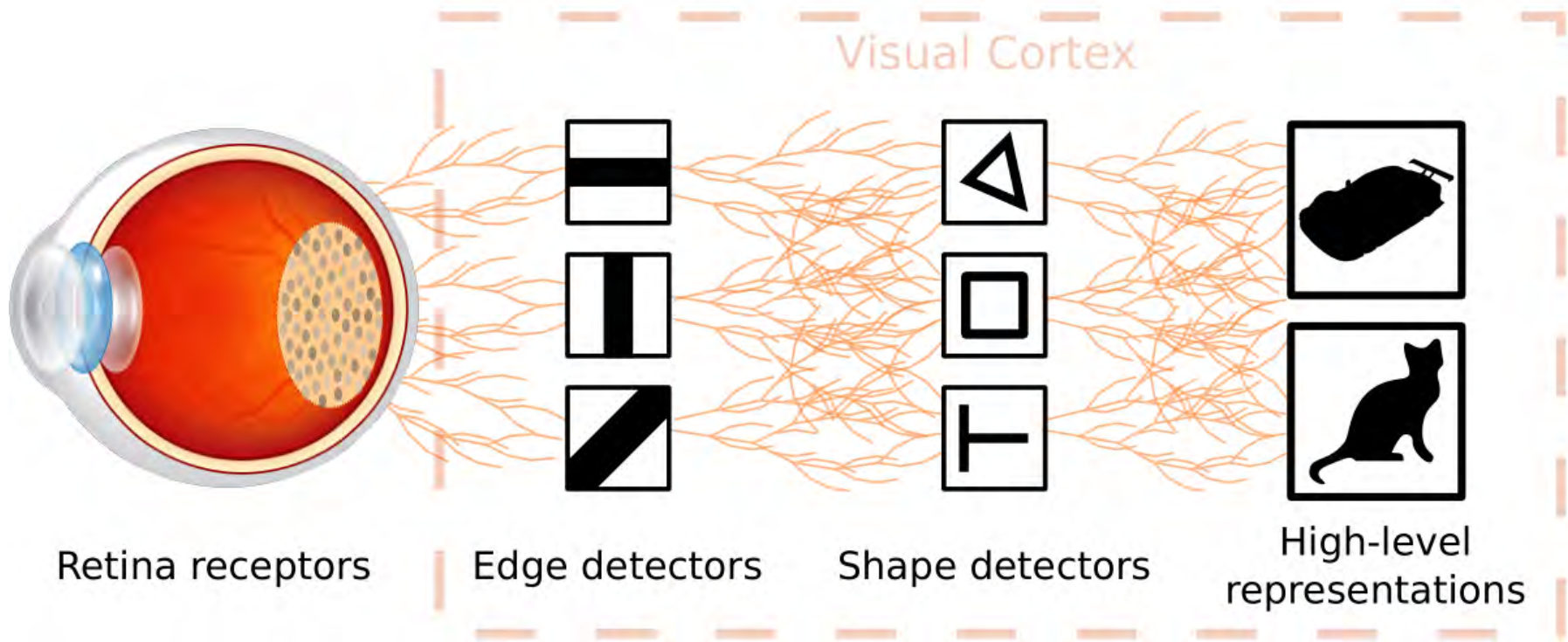
¿Cuáles son los **usuarios y temas estrella** sobre los que se está hablando en la red?

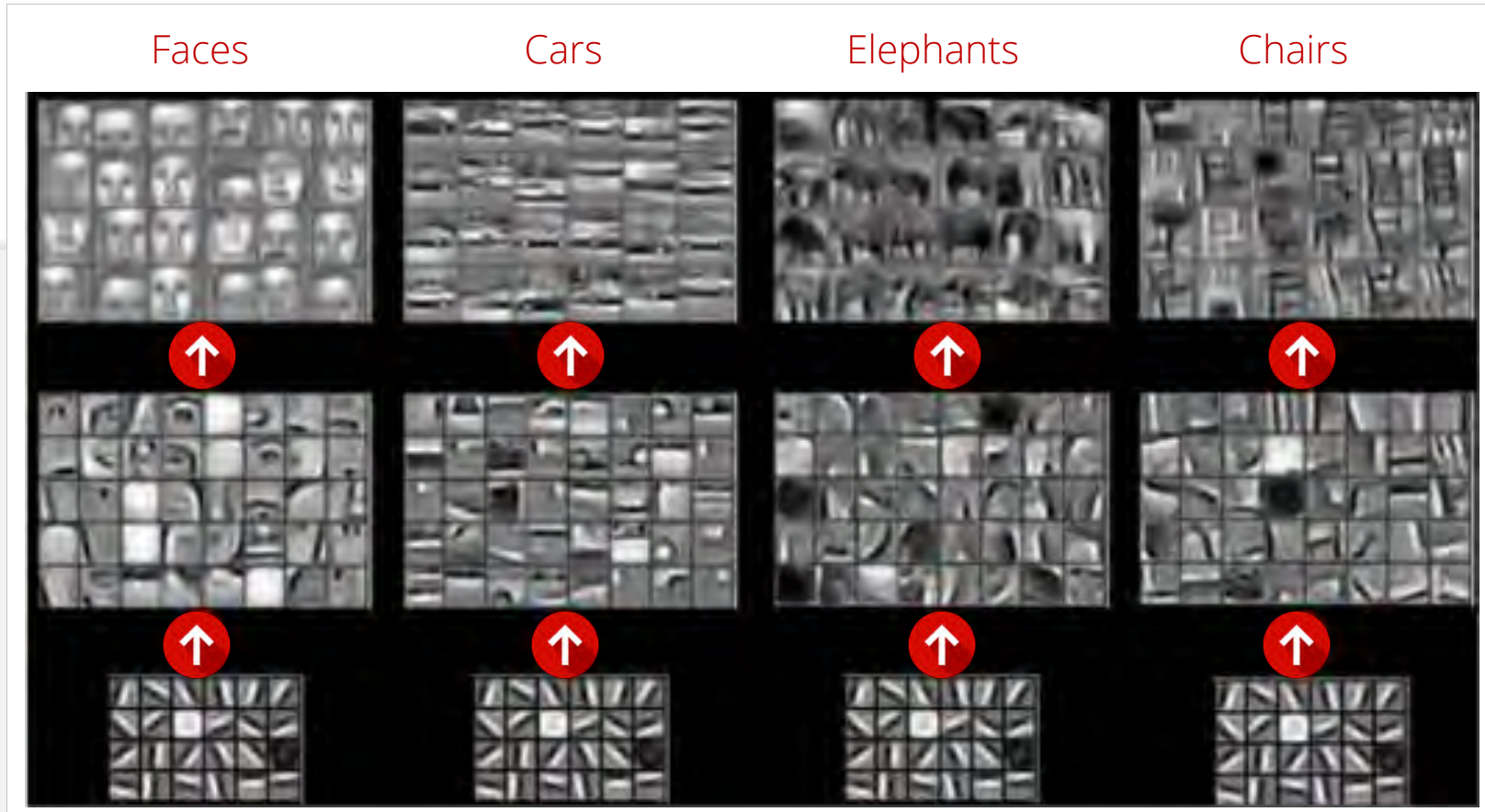


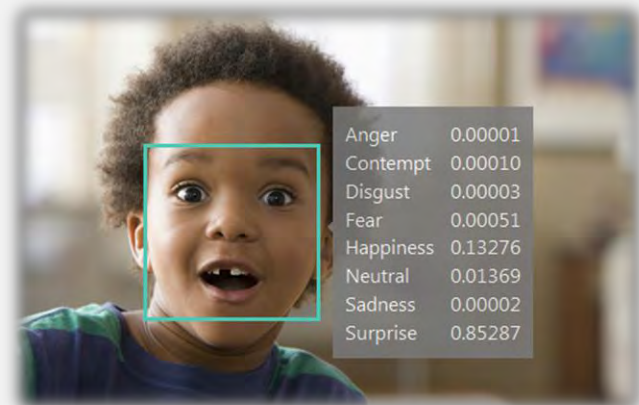
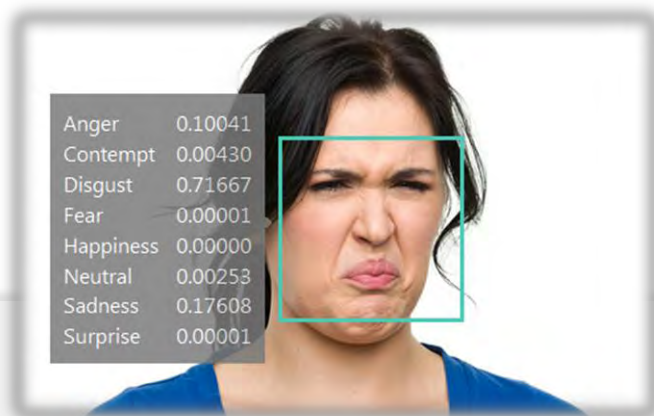


Las imágenes son
números





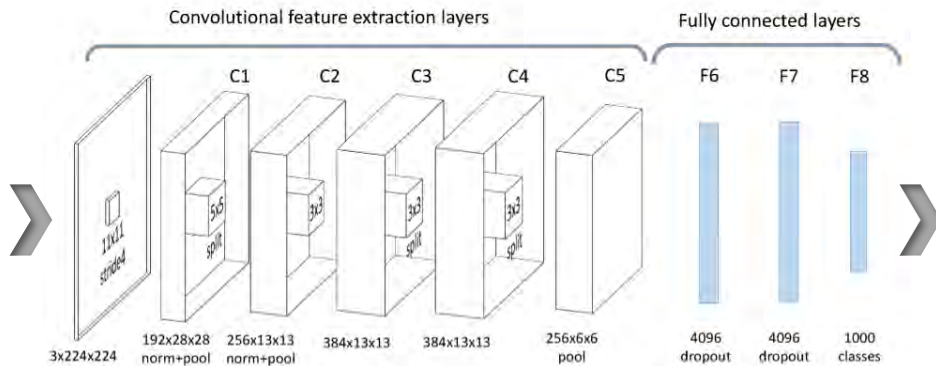




<https://www.microsoft.com/cognitive-services/en-us/emotion-api>



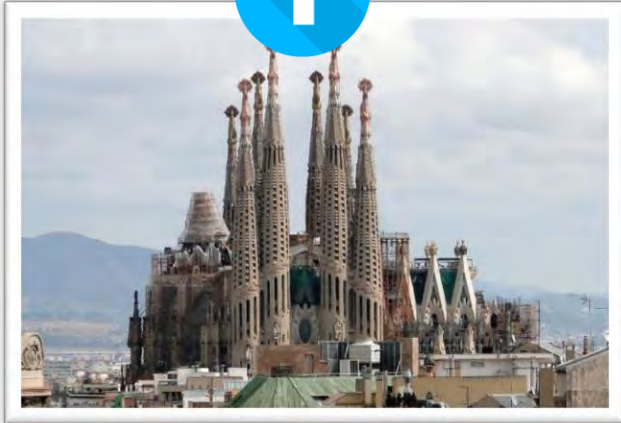
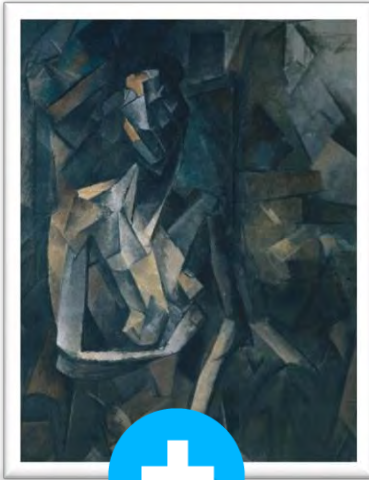
Low level observations

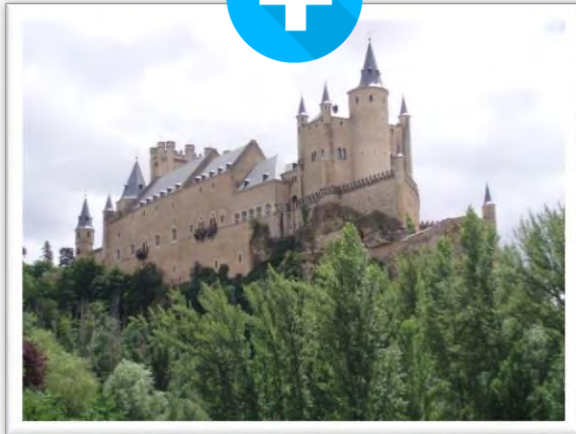
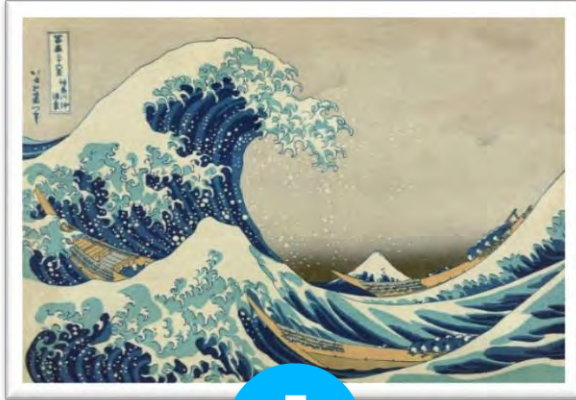


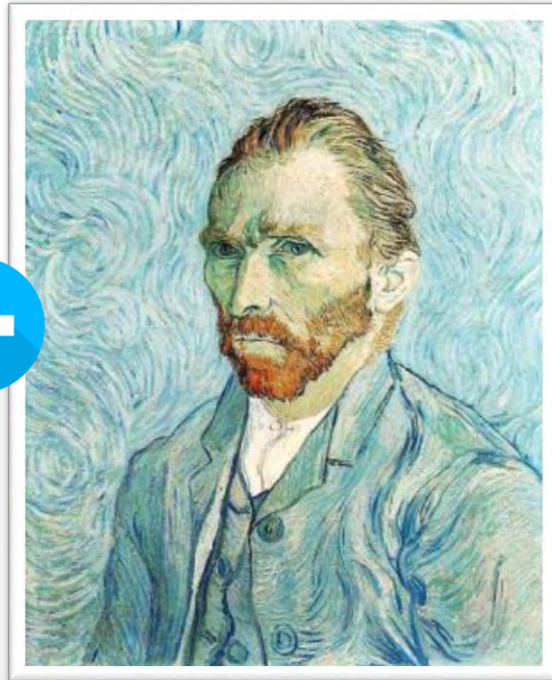
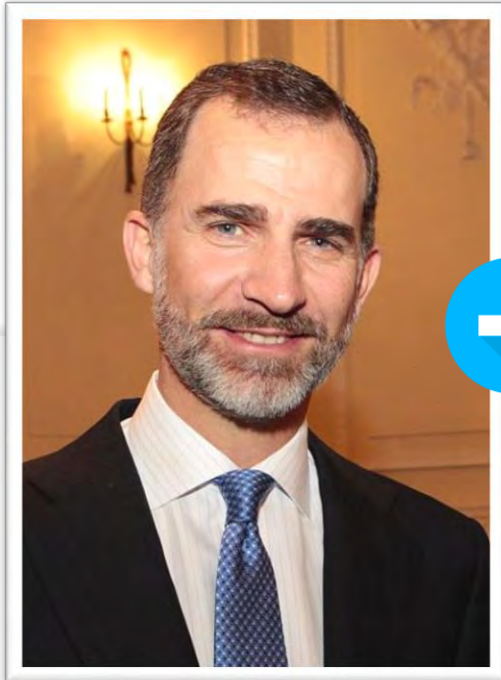
Style embedding

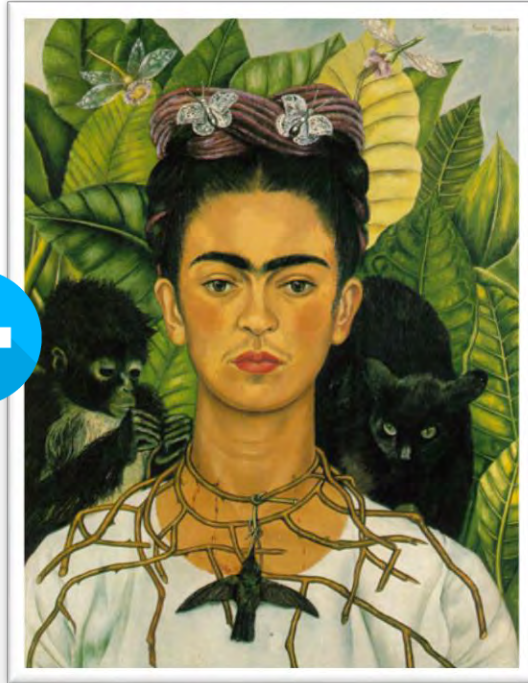
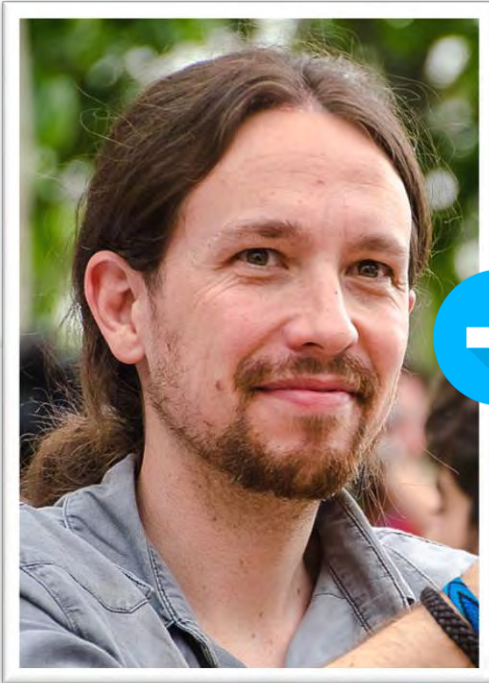
Gatys et al – A Neural Algorithm of Artistic Style
Bottou et al - Optimization Methods for Large-Scale Machine Learning

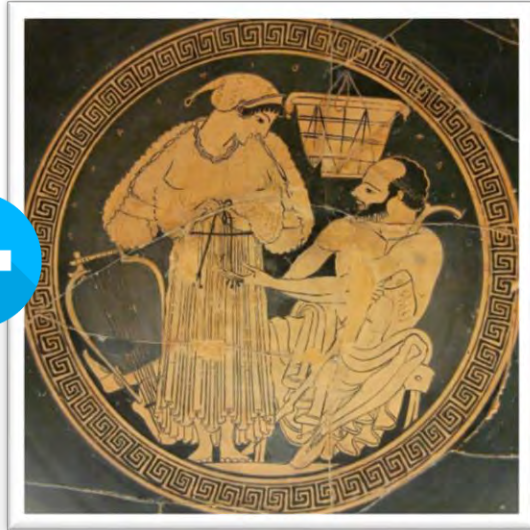


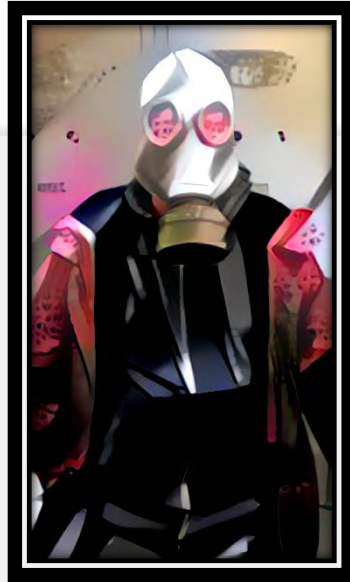
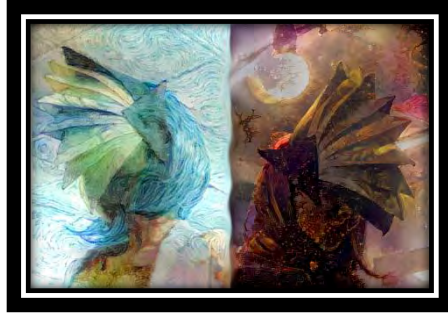
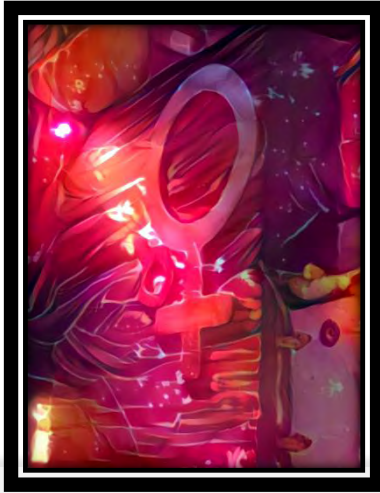


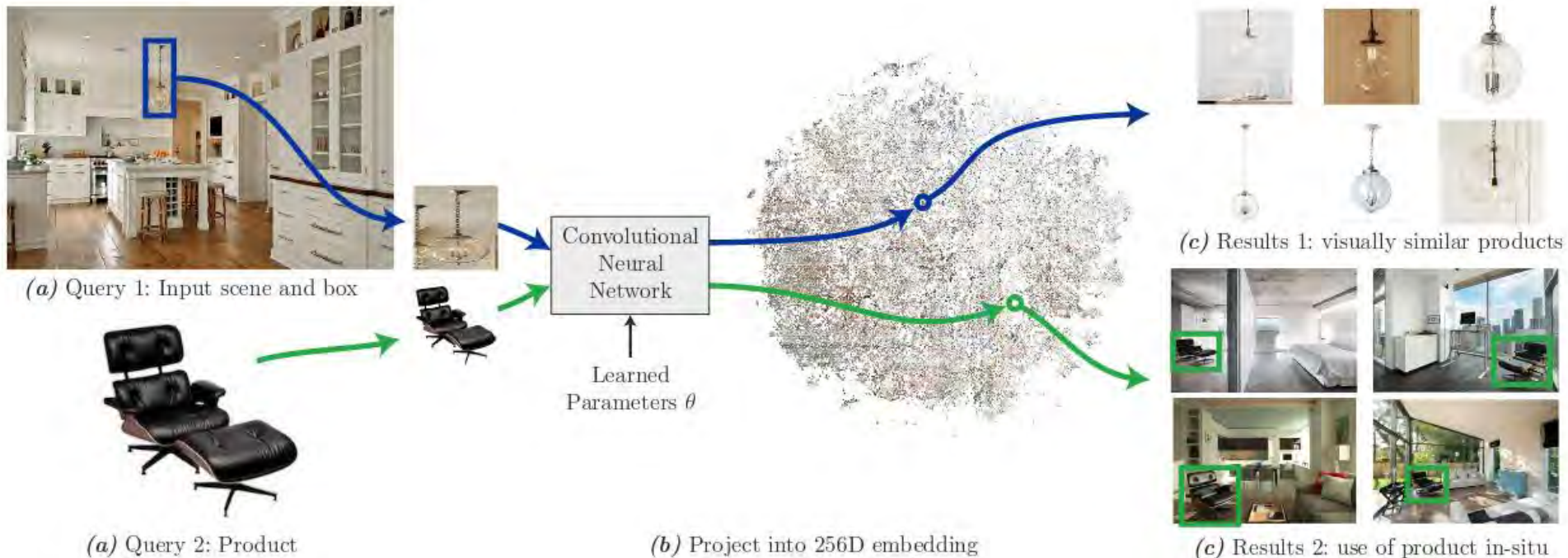










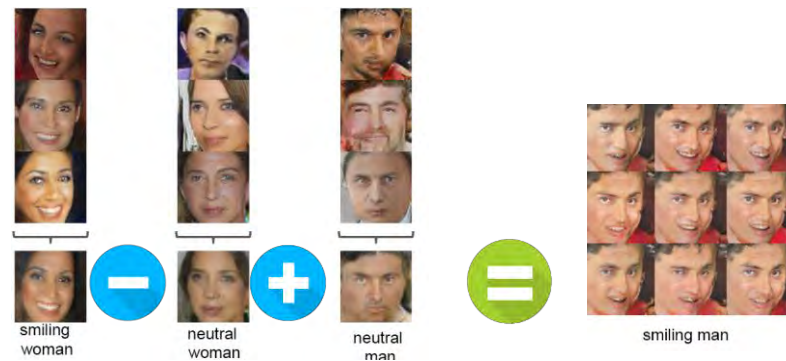


Bell and Bala - Learning visual similarity for product design with convolutional neural networks

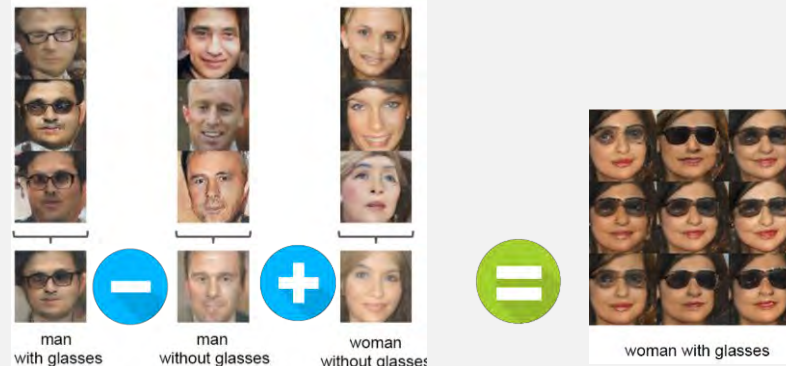
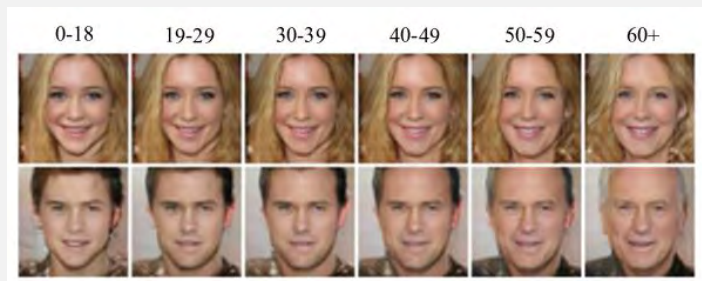
Restando ventanas



Modificando caras



Sumando años



Alec Radford, Luke Metz, Soumith Chintala - Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
Antipov et al - Face Aging with Conditional Generative Adversarial Networks

Coloreando fotos



Nguyen et al - Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space, <http://demos.algorithmia.com/colorize-photos/>

Volcano

Creando
imágenes
sintéticas



Nguyen et al - Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space, <http://demos.algorithmia.com/colorize-photos/>

Redshank

|

Ant

|

Monastery



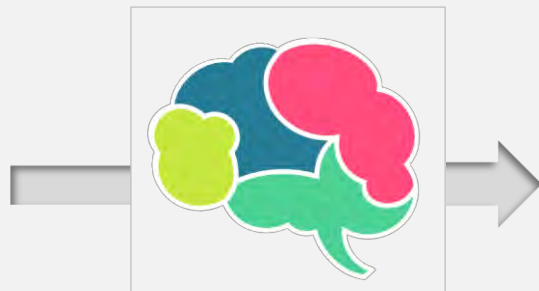
Creando
imágenes
sintéticas

Nguyen et al - Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space, <http://demos.algorithmia.com/colorize-photos/>



Mezclando números

Descripción de imágenes



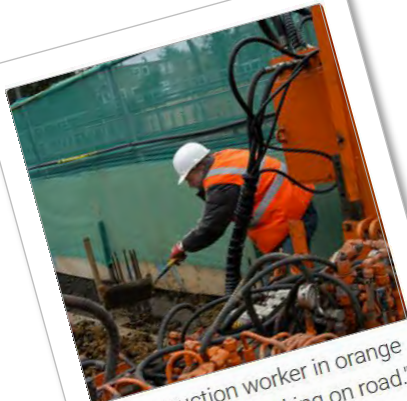
"A close up of a child holding a stuffed animal"

0.1
0.5
1.0
0.0
2.4





"man in black shirt is playing guitar."



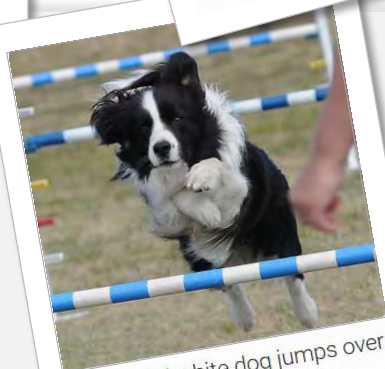
"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"girl in pink dress is jumping in air."



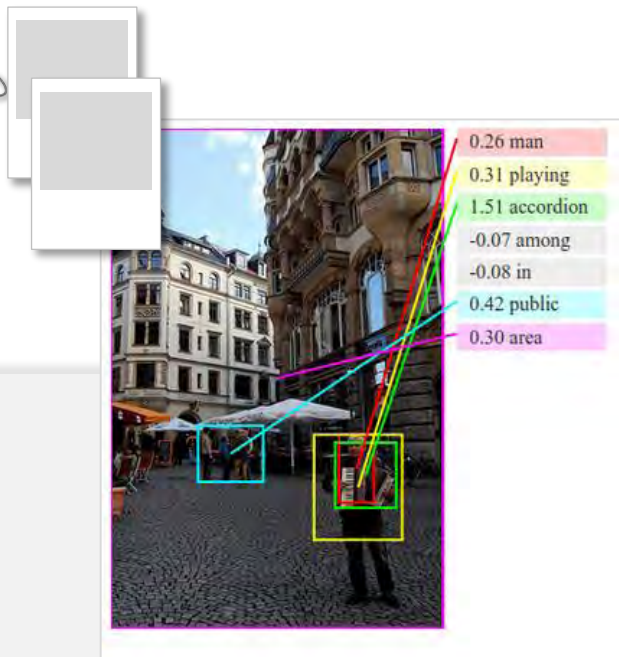
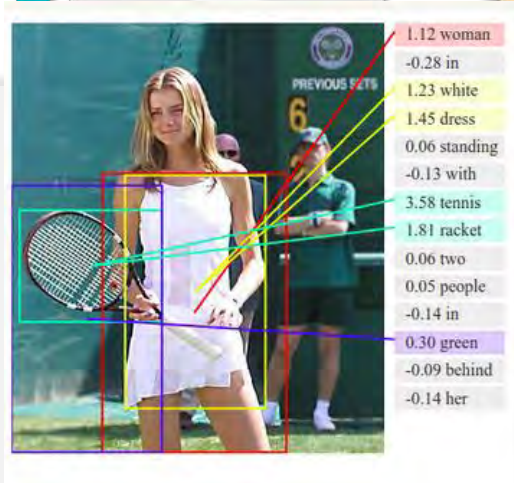
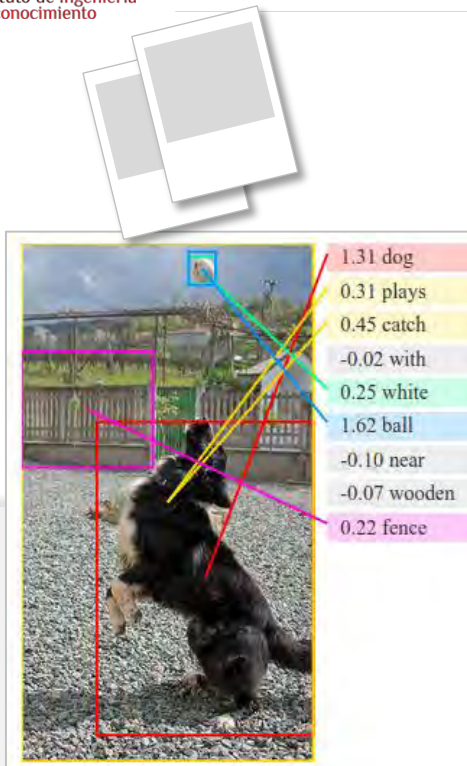
"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

Deep Visual-Semantic
Alignments for Generating
Image Descriptions
[http://cs.stanford.edu/people/
karpathy/deepimagesent/](http://cs.stanford.edu/people/karpathy/deepimagesent/)





Deep Visual-Semantic Alignments for
 Generating Image Descriptions
<http://cs.stanford.edu/people/karpathy/deepimagesent/>



a red car parked on the side of a road



a blue car parked on the side of a road



a pizza on a plate at a restaurant



someone is just about to cut the pizza



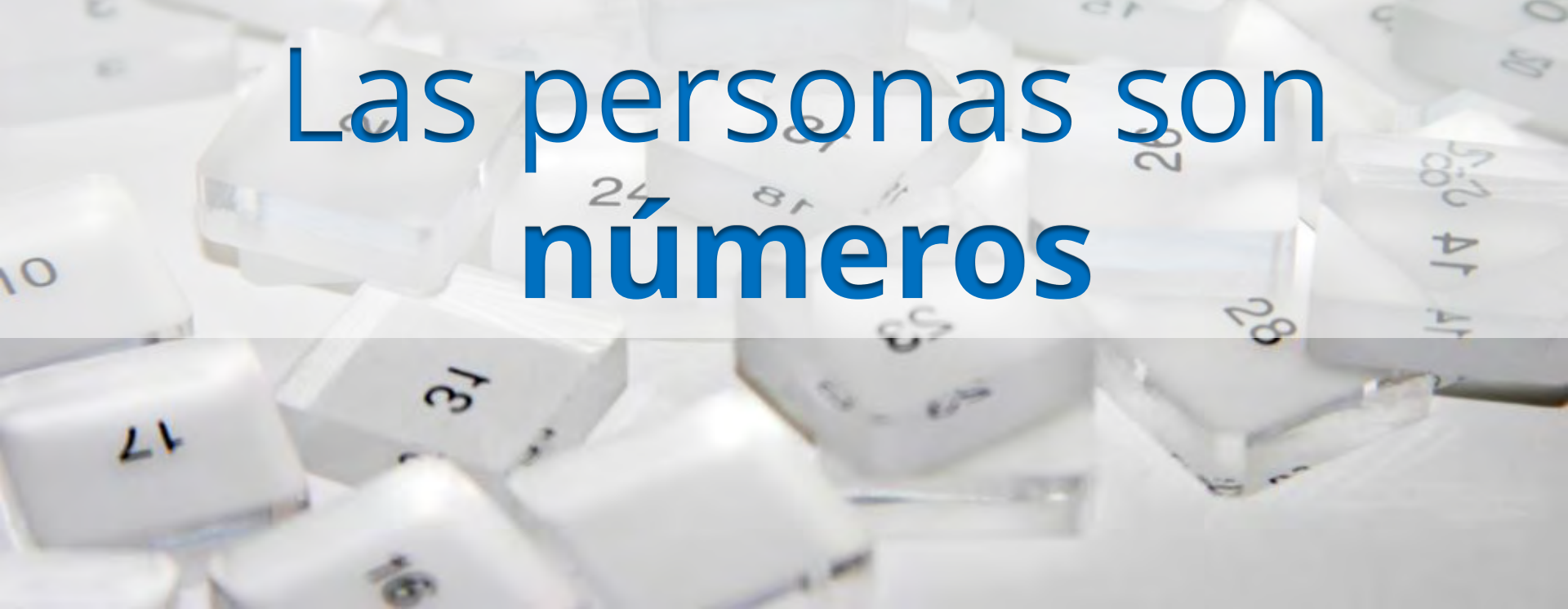
oranges on a table next to a liquor bottle



a pile of oranges sitting in a wooden crate



Las personas son
números



Una persona **no** es un número...



¡Son **muchos** números!





Objetivo: perfilado completo del cliente y detección de comportamientos fraudulentos

Comercios



Banca por internet



Oficinas



Banca telefónica



Cajeros

Operaciones financieras:

- Compras
- Retirada de efectivo
- Transferencias
- Recibos
- Recarga de móviles
- Préstamos, etc.

Operaciones **no** financieras:

- Consulta de saldo
- Cambio de PIN
- Consulta de movimientos, etc.



Criba cv



Llamada
Tfónica.



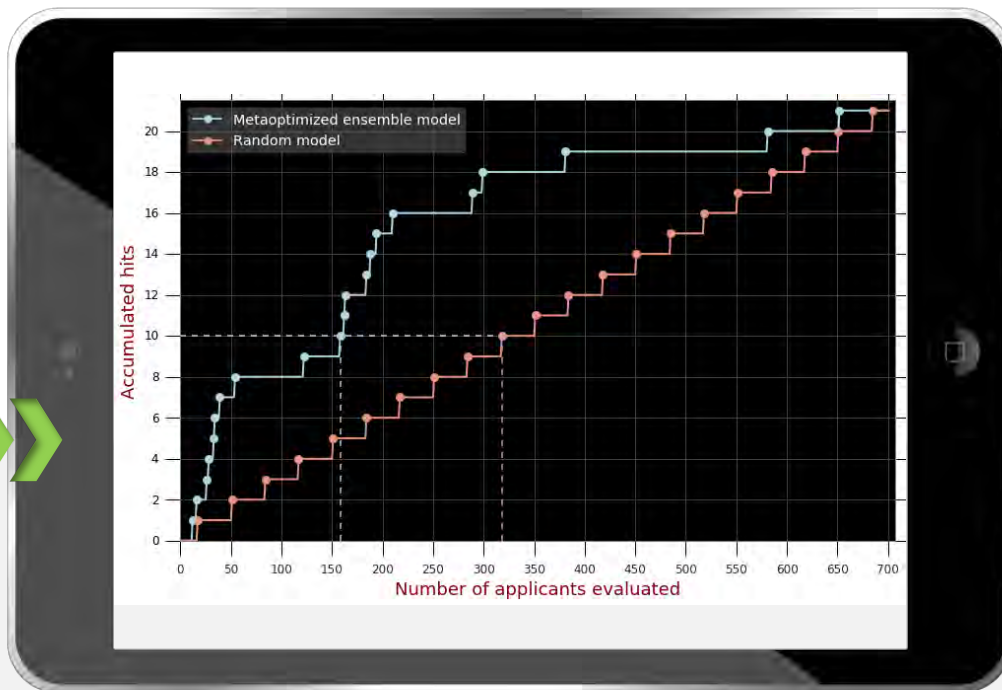
Pruebas
Psicométricas



Breve
entrevista



Dinámica



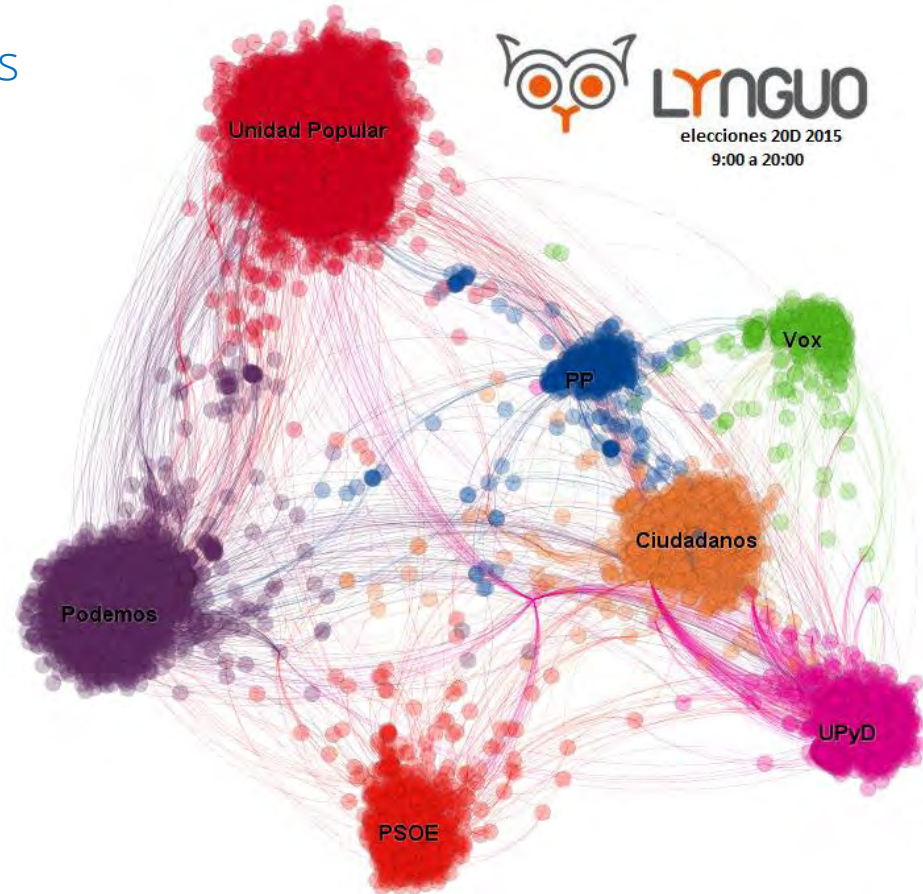
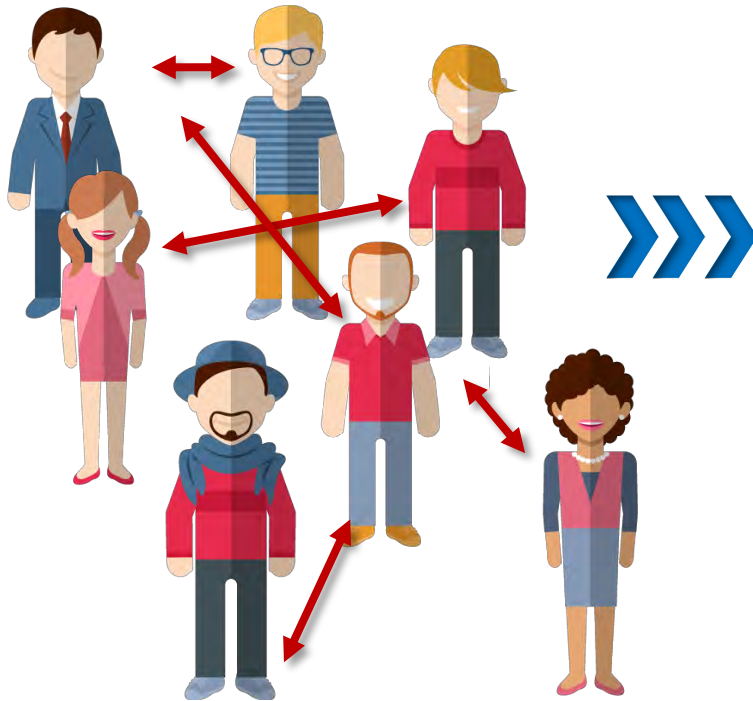
Foco

Modelo Azar

Modelo Propuesto

Sólo puedo revisar 160 CVs	5 Aptos	10 Aptos
Seleccionar 10 Aptos	Reviso 320 CVs	Reviso 160 CVs

Grupos de personas y relaciones entre ellas



AROS y eAROS es un servicio de consultoría para una gestión eficaz de la red social corporativa de las empresas a través de su monitorización y análisis.





En resumen



Integración de toda clase de fuentes de datos, incluyendo imágenes, textos, video, datos de personas...



Soluciones de datos de mayor nivel de análisis:
predictivas y prescriptivas



Nuevas tecnologías de análisis multimedia:
prácticas y efectivas



Solamente necesitas los expertos adecuados...

¡Cuenta con nosotros!





Álvaro Barbero Jiménez

PhD, PMP, Chief Data Scientist at
Instituto de Ingeniería del Conocimiento (IIC)



[Alvaro Barbero](#)

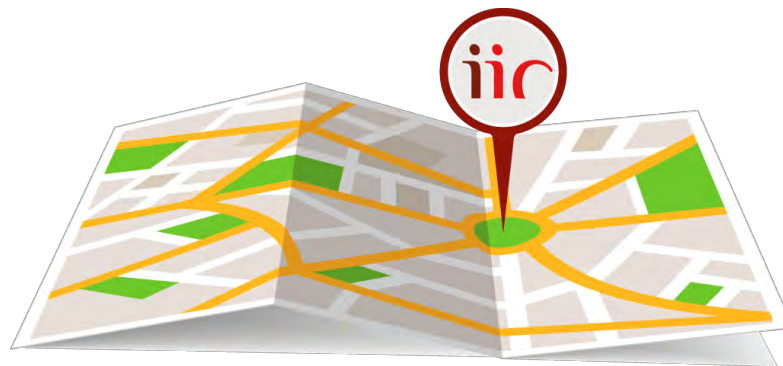


@albarjip



[albarji.deviantart.com](#)

www.iic.uam.es



Elementos gráficos de apoyo obtenidos en:

designed by  freepik.com

pixabay 