

Módulo 1. Introducción a la ciencia de datos

Introducción a la Ciencia de Datos e Ingeniería de Datos

Jacinto Arias

Máster en Ciencia de Datos e Ingeniería de Datos en la Nube



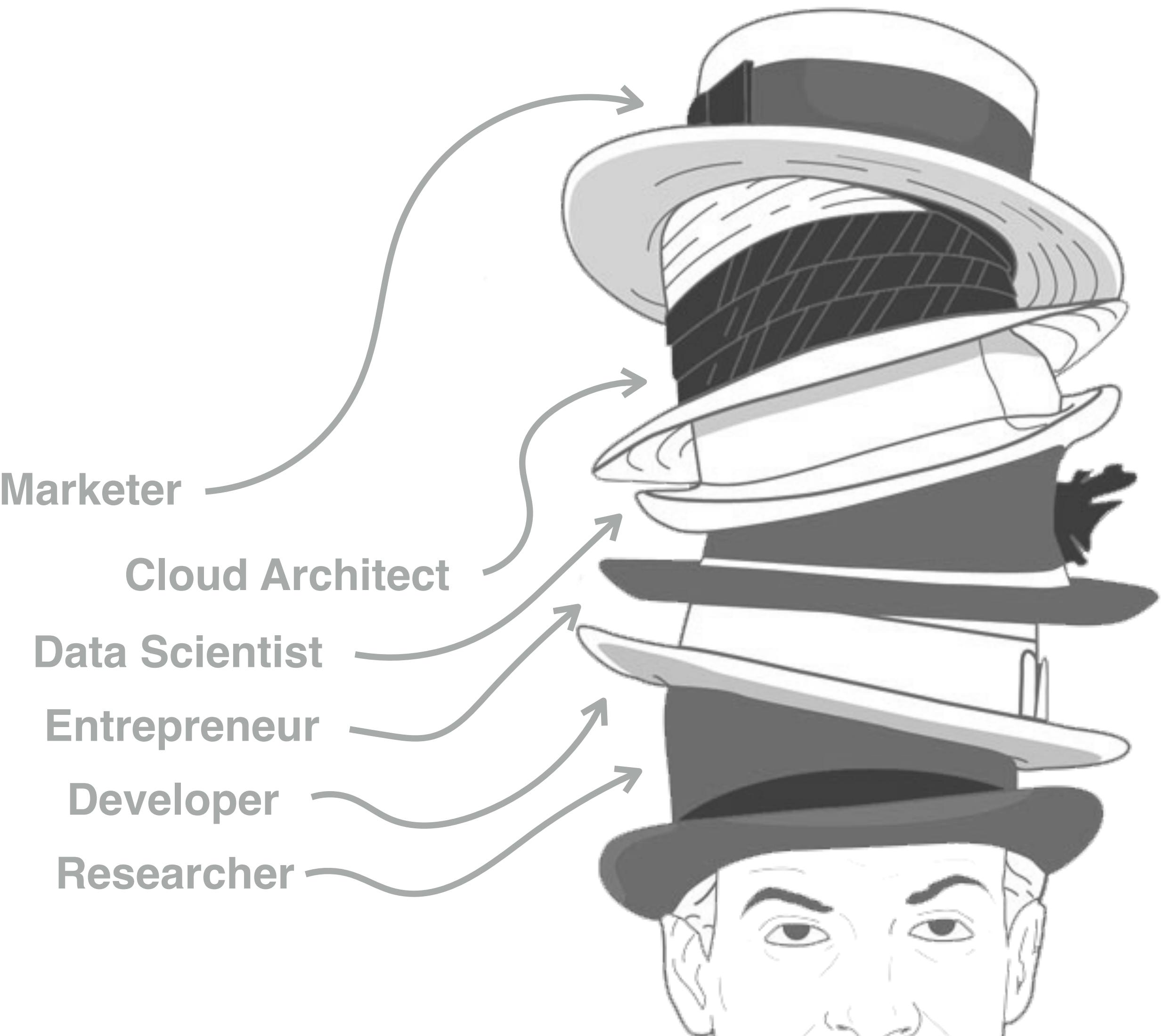


Sobre mi

- Informático
- Científico
- Desarrollador
- Emprendedor

 Visit taidy.cloud

 in/jacintoarias



A large central word "DATA SCIENCE" in bold blue letters, surrounded by a cloud of related terms in various sizes and colors, including:

- DETECTION
- SOCIAL MEDIA
- BIG DATA
- INFORMATION TECHNOLOGY
- PROMOTION
- COMPUTING
- CONTENT
- PROCESSING
- CONSUMER
- ORGANIZATION
- PLANNING
- E-MARKETING
- COMMUNICATION
- PROJECTS
- WWW
- SERVICES
- BRANDING
- COMPUTER
- MULTIMEDIA
- NETWORK
- CONSUMER DEMAND MARKETS
- PREDICTIVE
- PROGRAM
- ANALYTICS
- EVENTS
- PROGRAMMING
- SOFTWARE
- WEB MARKETING
- DATA MINING
- MACHINE LEARNING
- VISION
- ENGINEERING
- RESEARCH
- PROBABILITY
- COMPUTING
- WEB SERVICES
- WEB DEV
- KDD
- STRATEGY
- WORLDWIDE
- BIG
- DATA
- SERVICE
- VISUALIZATION
- PRICING
- BIG DATA
- SOLUTIONS
- MATHS
- PATTERN
- ENGINEERING
- PLANNING
- MEDIA
- STATISTICS
- MOBILE
- INFORMATION
- DIGITAL
- SEGMENTATION
- SOCIAL NETWORKS
- SOCIAL NETWORK

The image is a horizontal word cloud centered around the theme of "BIG DATA". The word "BIG" is the largest word in the center, flanked by "DATA" and "INFORMATION". Other prominent words include "TECHNOLOGIES", "ANALYSIS", "SEARCH", "PETABYTES", "VOLUME", "GOVERNMENT", "SOFTWARE", "LARGE", "STORAGE", "SETS", "FUTURE", "SCIENCE", "BILLION", "NEW", "INTERNET", and "RESULTS". The words are in various sizes and colors (blue, purple, red, yellow) and are arranged in a dense, overlapping pattern.



“I keep saying that the sexy job in the next ten years will be statisticians. And I’m not kidding.”

Hal Varian (2009)
Chief Economist @ Google



Buzzwords

DATA SCIENCE

DATA SCIENTIST

DATA ENGINEER

MACHINE LEARNING

DEEP LEARNING

BIG DATA



¿Qué es un Científico de datos?

An Statistician who lives in San Francisco

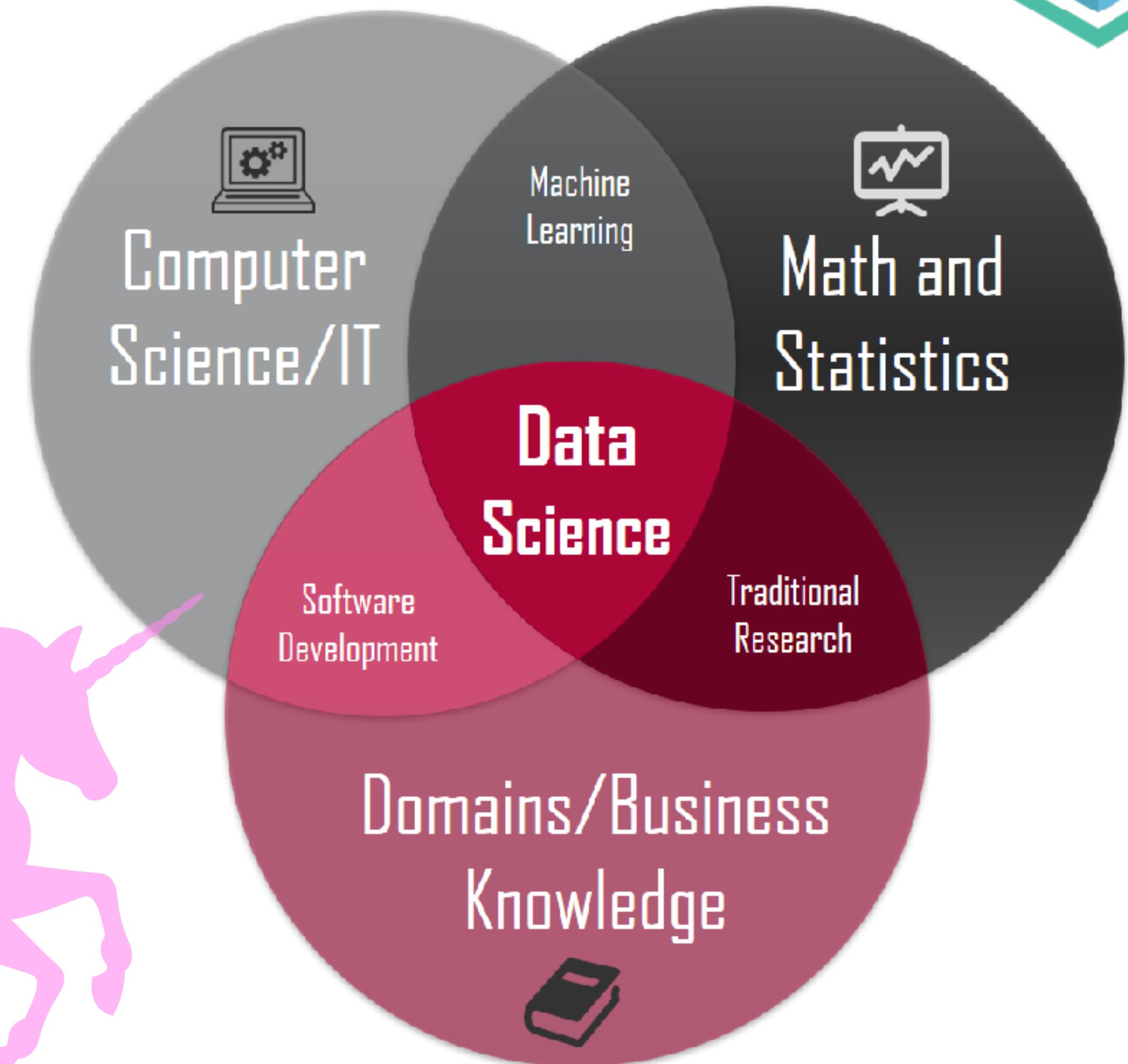
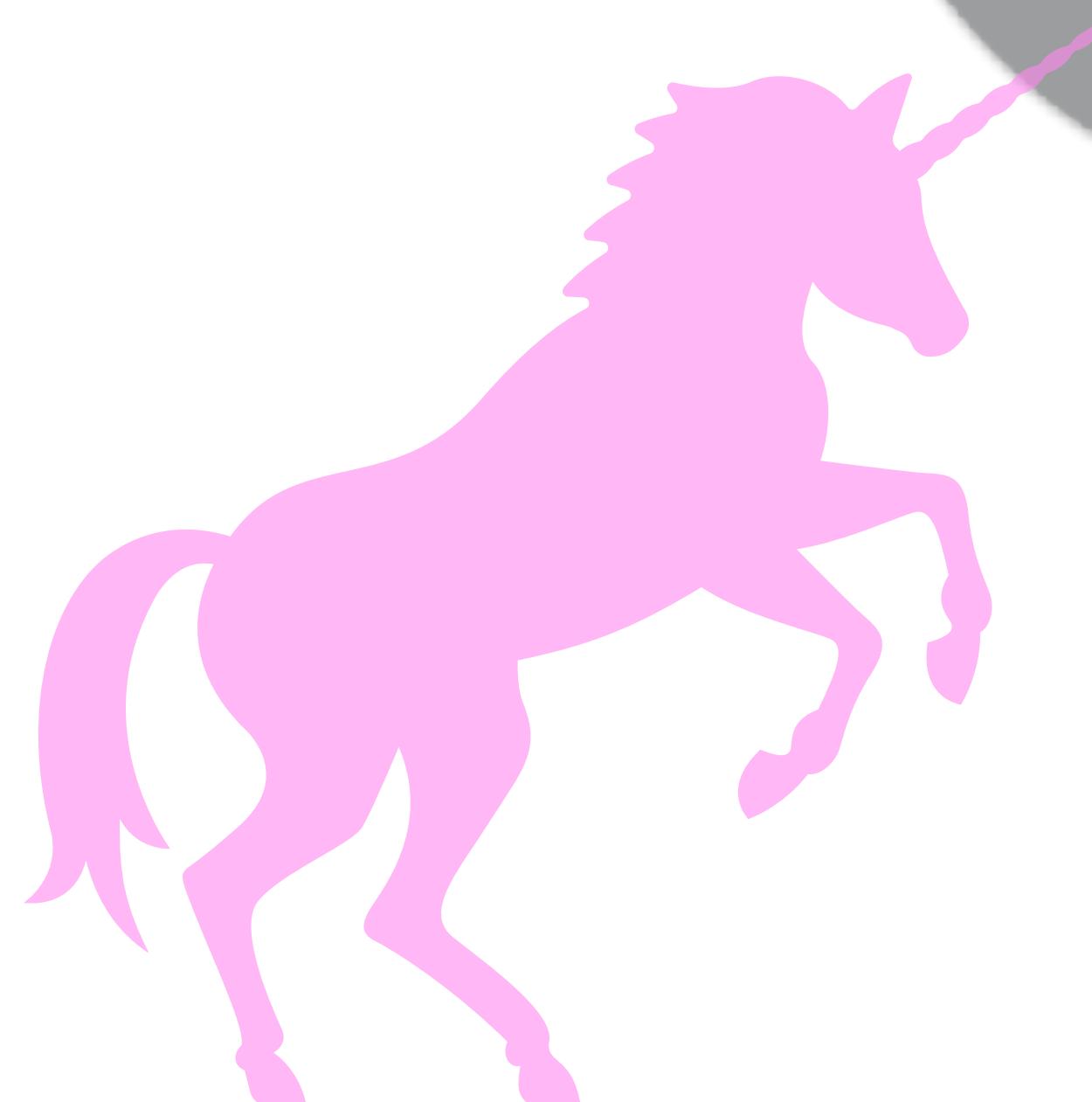
Statistics on a Mac

Someone who is better at statistics than a software engineer, and better at coding than any statistician



¿Qué es la Ciencia de Datos?

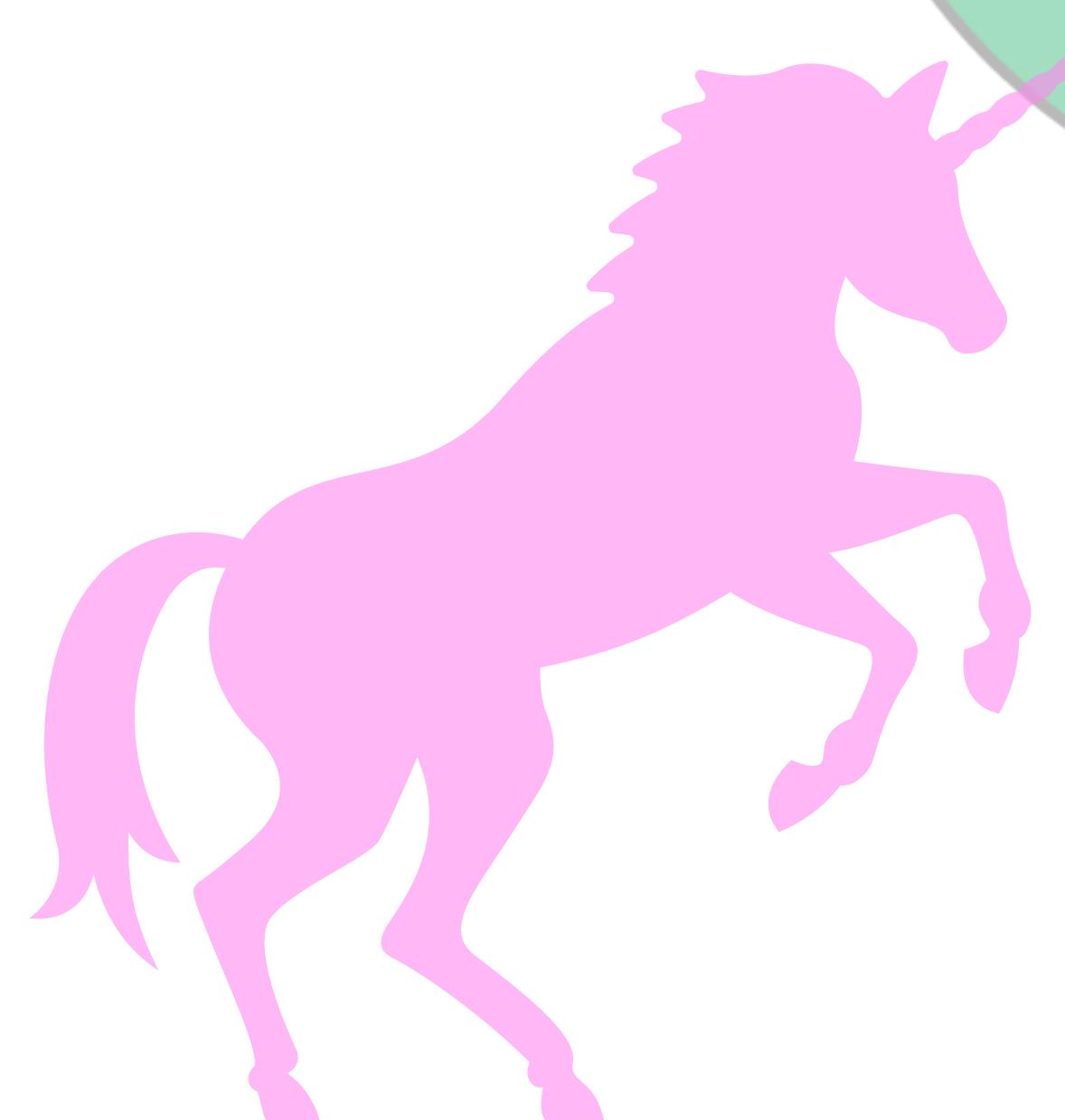
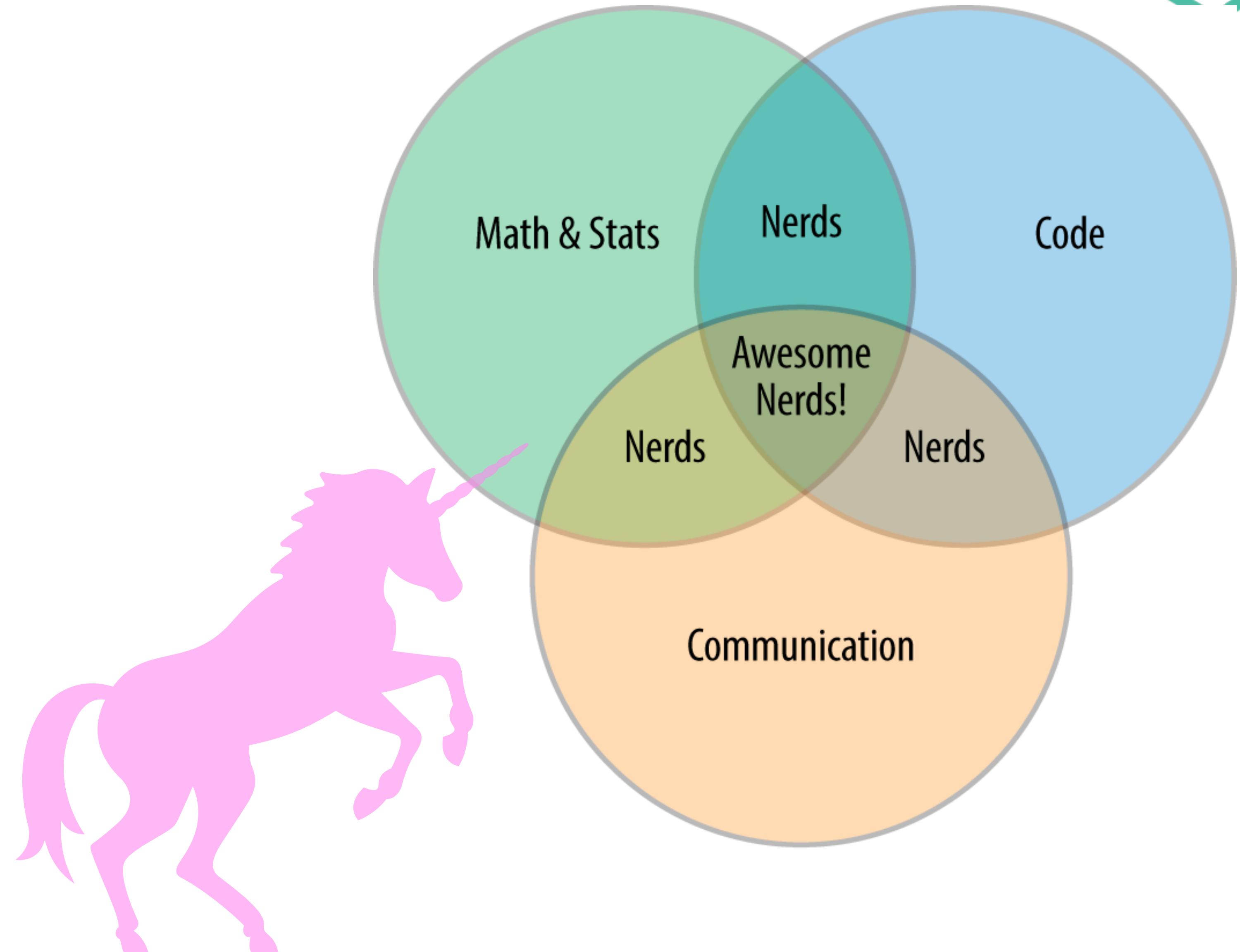
- Acuñado en 2008
(Patil@LinkedIn,
Hammerbacher@Facebook)
- “Data Science is what a Data Scientist does”
- “Data science is the discipline of making data useful.”





¿Qué es la Ciencia de Datos?

- Acuñado en 2008
(Patil@LinkedIn,
Hammerbacher@Facebook)
- “Data Science is what a Data
Scientist does”
- “Data science is the discipline
of making data useful.”



¿Quién es un Científico de Datos?

- Suele utilizarse como título de un puesto de trabajo
- Acuñado en 2008 (Patil@LinkedIn, Hammerbacher@Facebook)

Es casi la descripción de un departamento entero

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21th century, requires a mixture of multidisciplinary skills ranging from an 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 little 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



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 parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience withaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

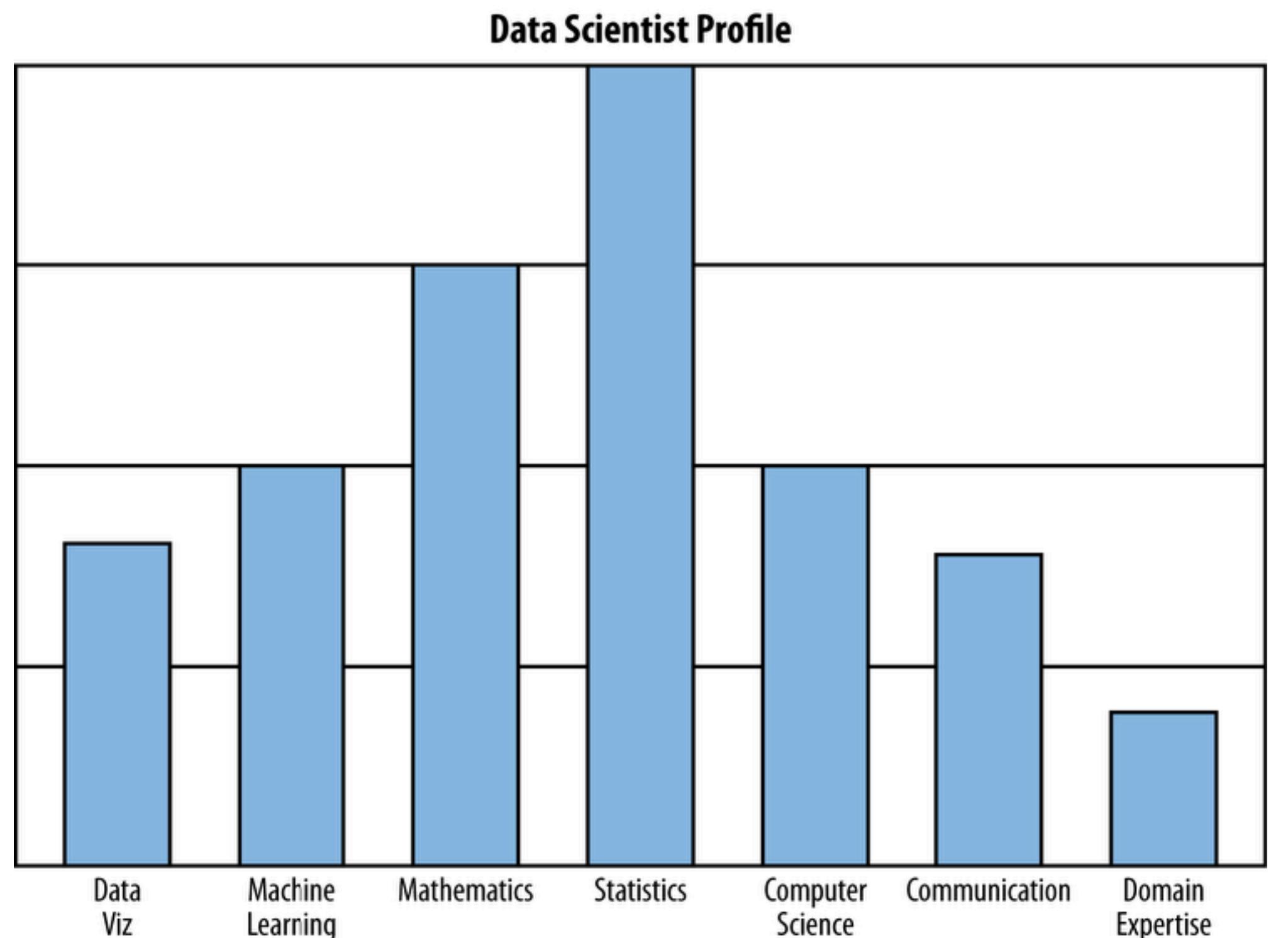
- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

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 of visualization tools e.g. Flare, D3.js, Tableau

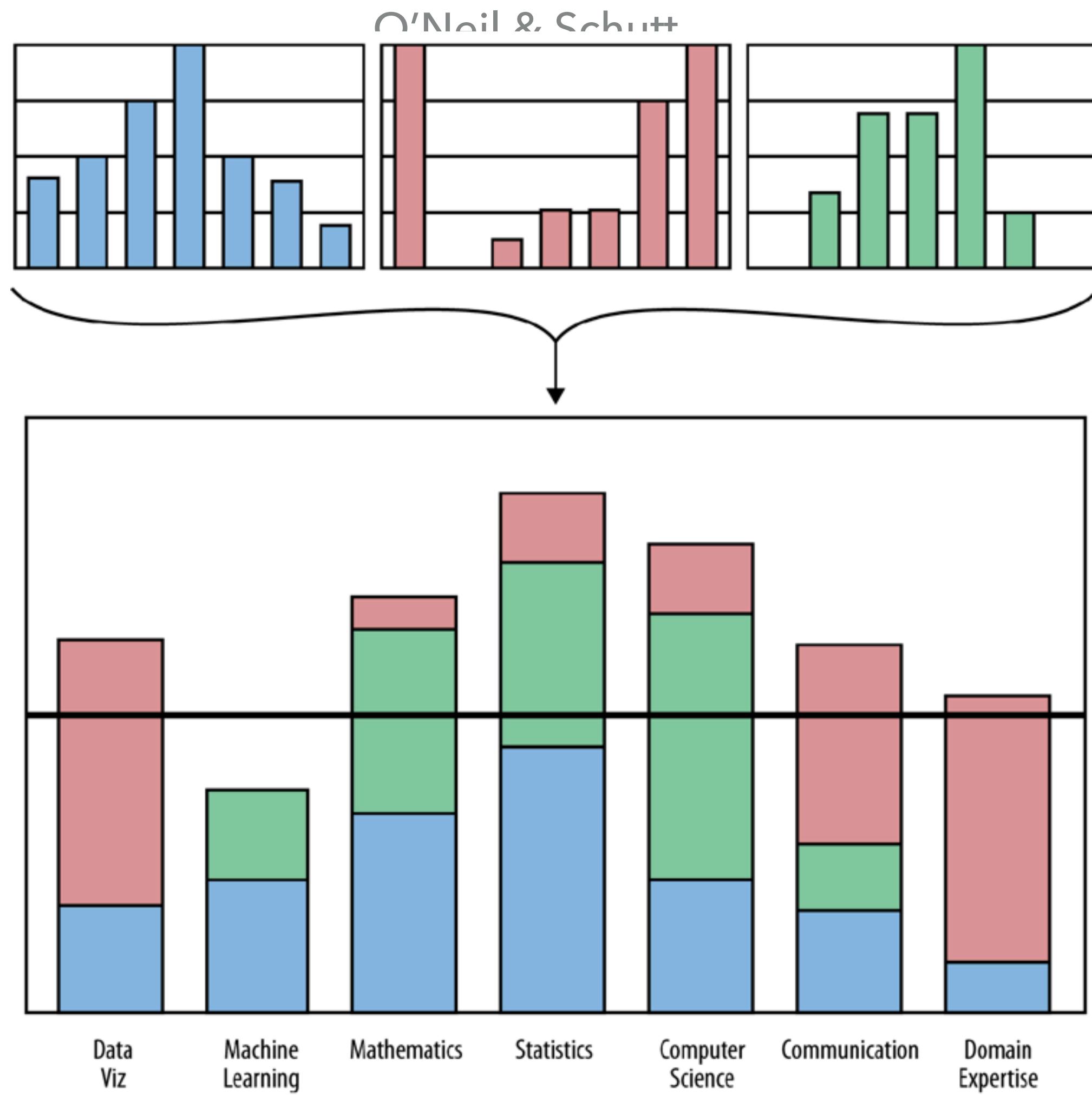
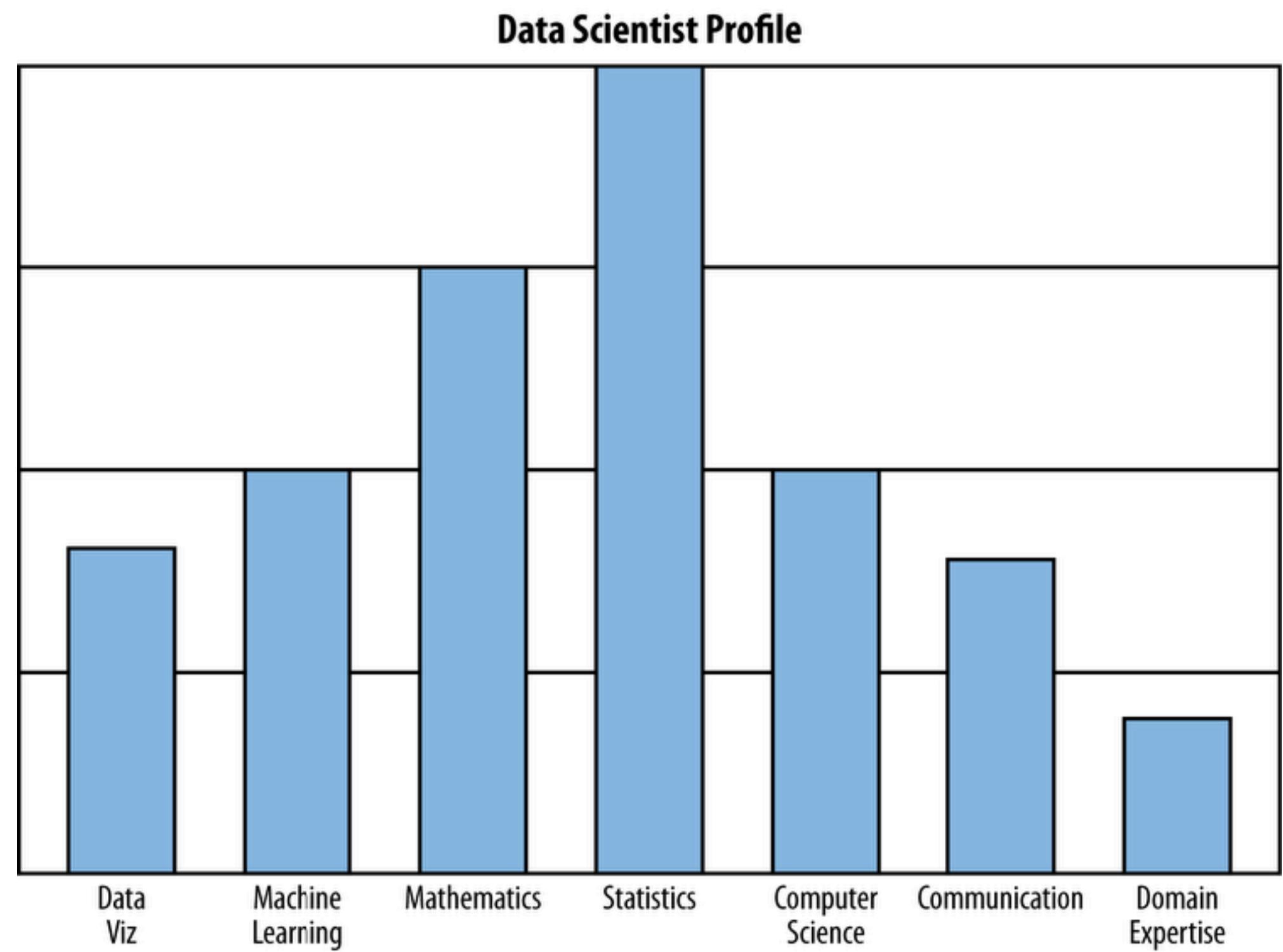


El Perfil de un Data Scientist





El Perfil de un Equipo de Data Science





Perdón, Data... ¿Qué?

Data Analyst

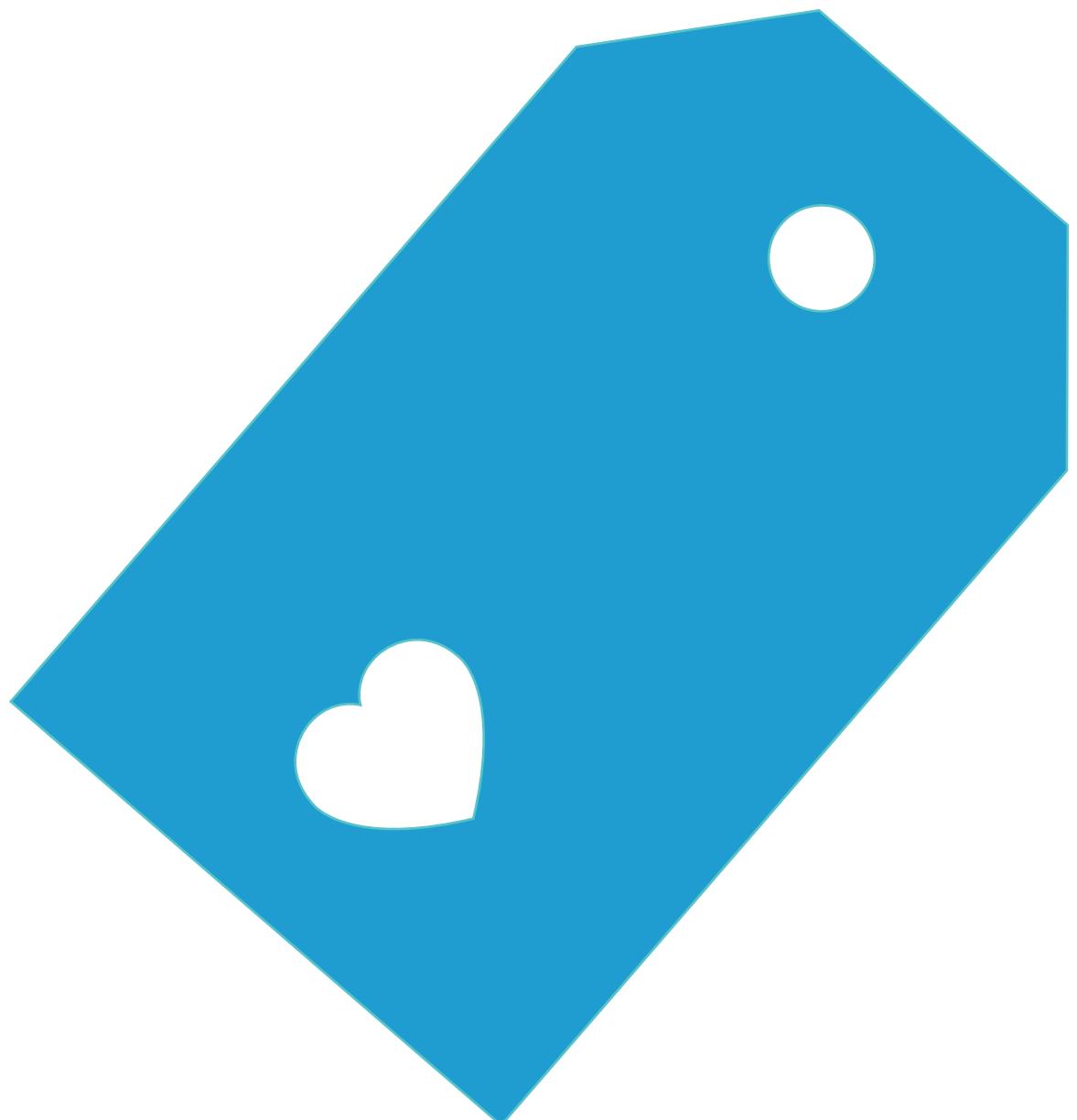
Data Scientist

Data Engineer

AI / ML Engineer

Data Translator

Citizen Data Scientist





“Making Data Useful”



Datos

Big Data. Experiencia técnica,
almacenamiento, modelado



Resultados

Métricas, nuevos datos,
iteración del ciclo





“Making Data Useful”





¿Para qué?

Hacer que los datos tengan utilidad

¿Cómo?



¿Para qué?

Para tomar decisiones...

“It’s through our decisions that we affect the world around us.”



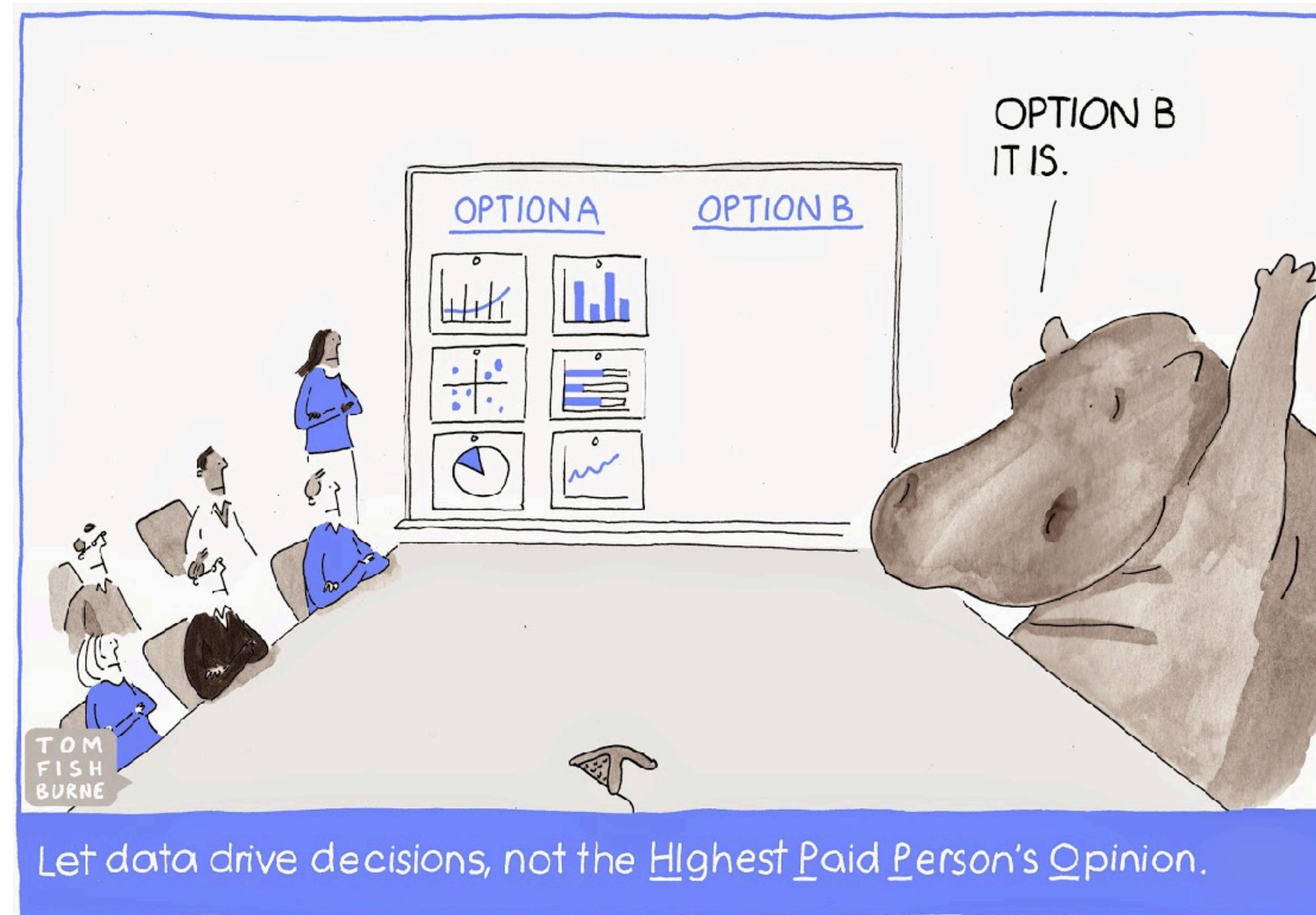
¿Para qué?

Para tomar decisiones...
...basadas en datos

¿Podemos hacerlo de otra manera?

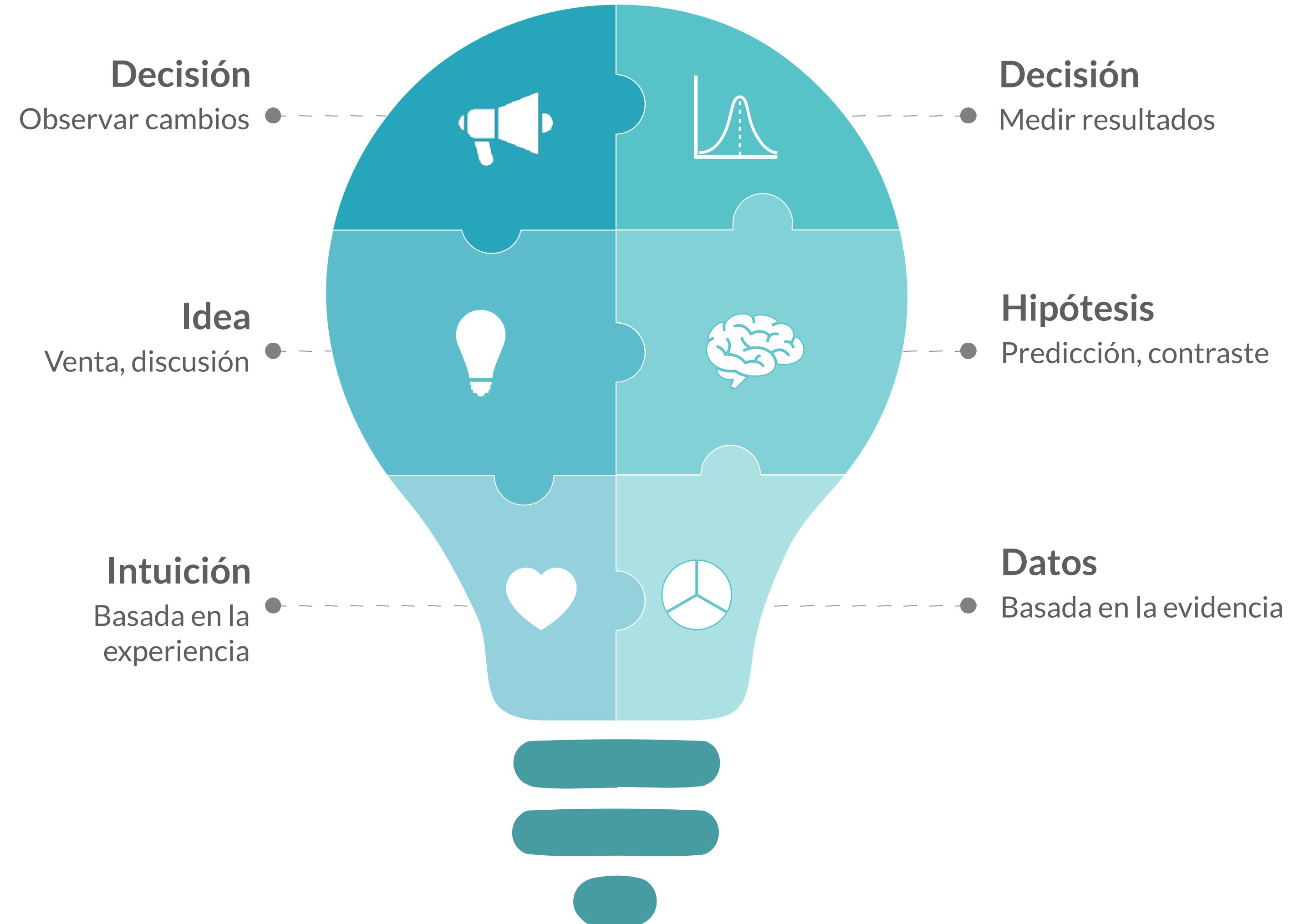


Antipatrón: The Hippo





Un cambio de paradigma





“In God we trust the rest bring data.”

W. EDWARDS DEMING (~1970)
Engineer and Scientist

**“If we have data, let's look at the data. If all we
have are opinions, let's go with mine.”**

JIM BARKSDALE (1998)
CEO @ Netscape



"Many people only use data to feel better about decisions they've already made."

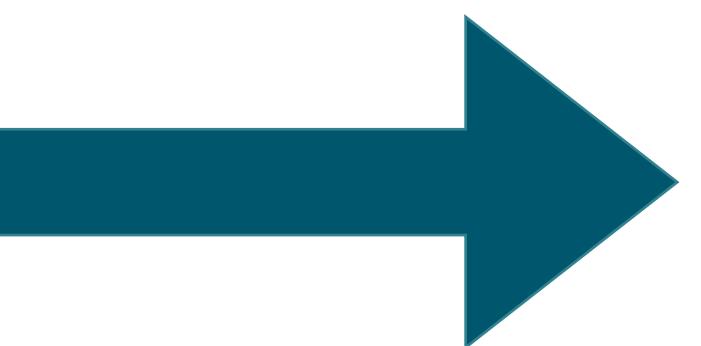
CASSIE KOZYRKOV (2018)
Chief Decision Maker @ Google





¿Para qué?

Hacer que los datos tengan utilidad





El Ecosistema



Personas



Procesos



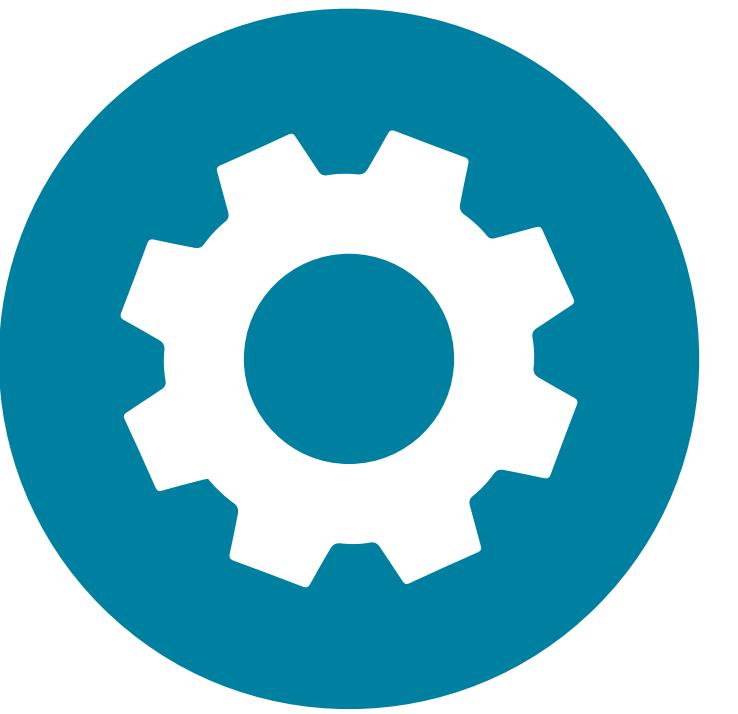
Tecnología



El Ecosistema



Personas



Procesos



Tecnología



¿Cuál es la diferencia con...?

- Estadísticos
- Analistas de negocio
- Desarrolladores de software
- Administradores de Bases de Datos
- Administradores de sistemas
- ...





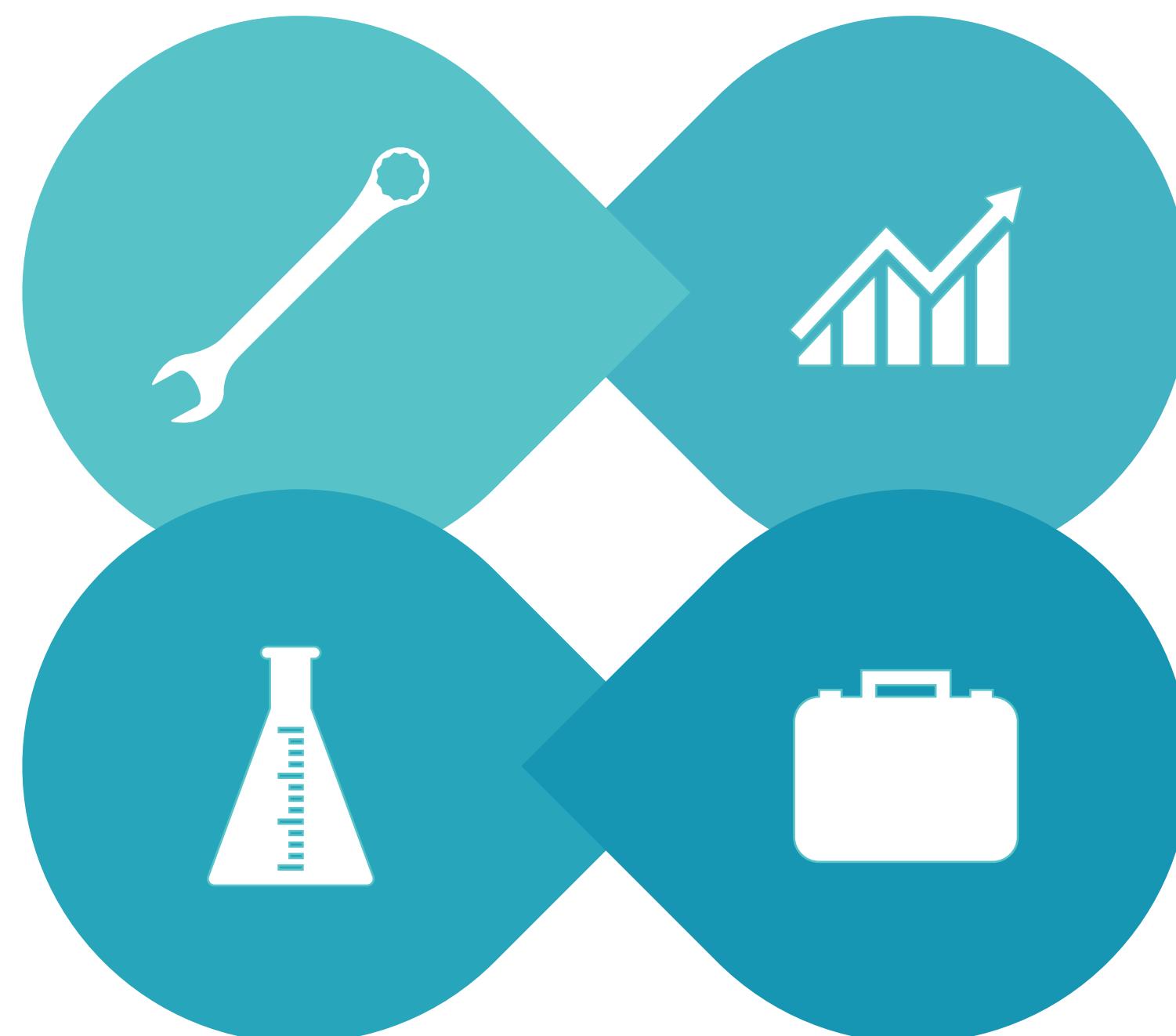
Combinar y adaptar

Data Engineers

Expertos en tecnología y en datos

Data Scientists

Machine Learning y curiosidad



Business Analyst

Traducir e interpretar los resultados

Decisión Makers

Liderazgo basado en datos



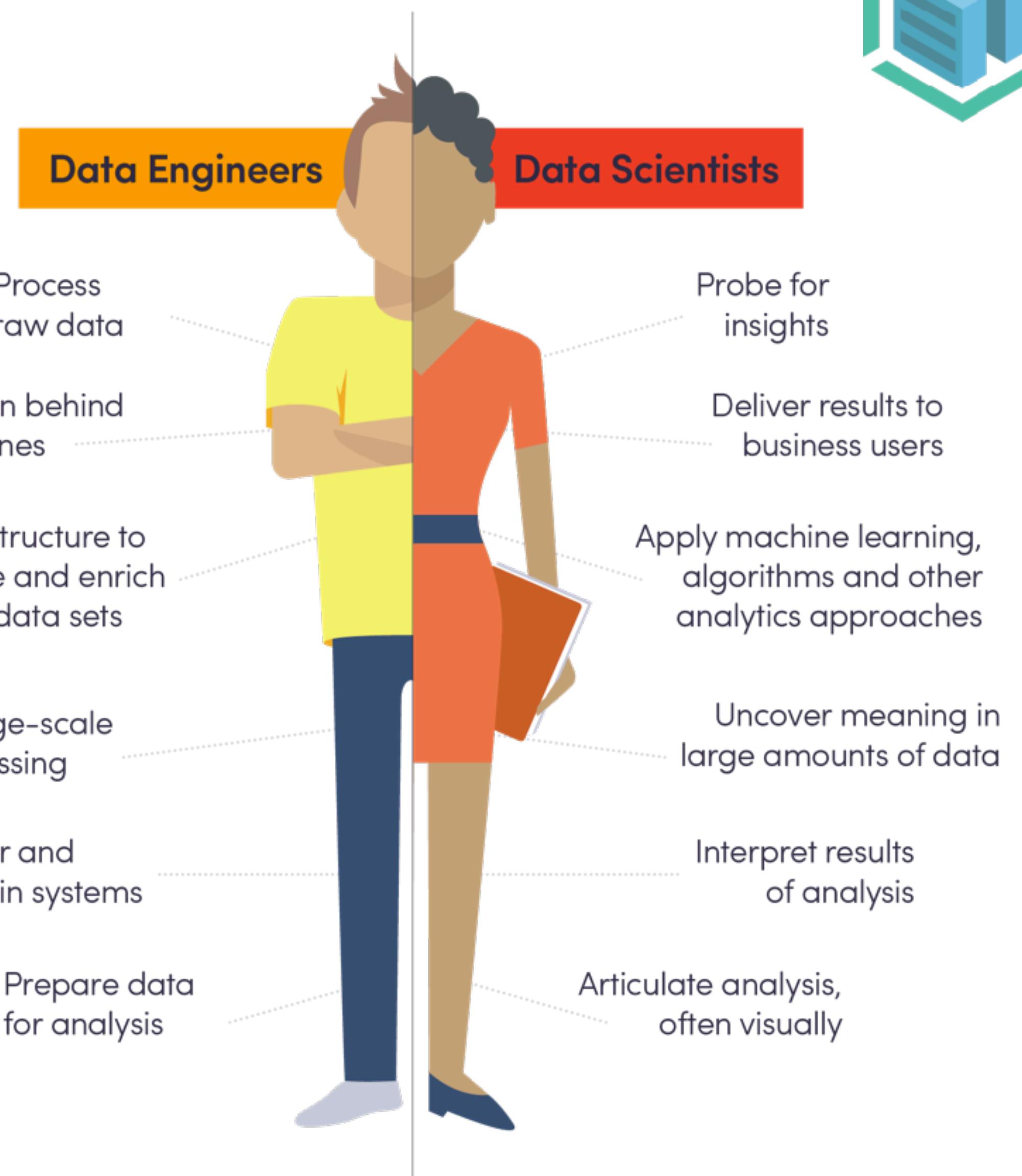
Data What?

Data Engineer

Al principio de la cadena de valor

Data Scientist

Más cerca del valor final



El Ecosistema



Personas



Procesos



Tecnología



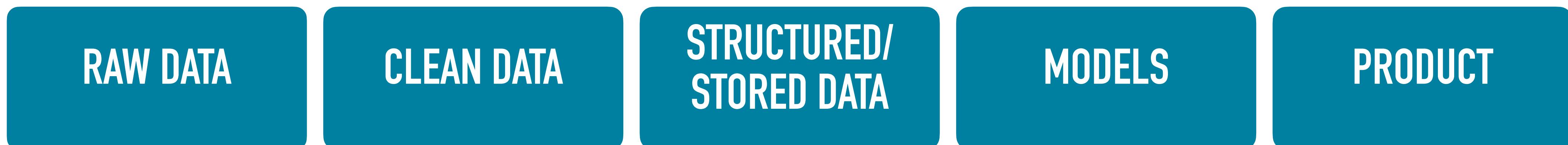
El ciclo de valor de los datos





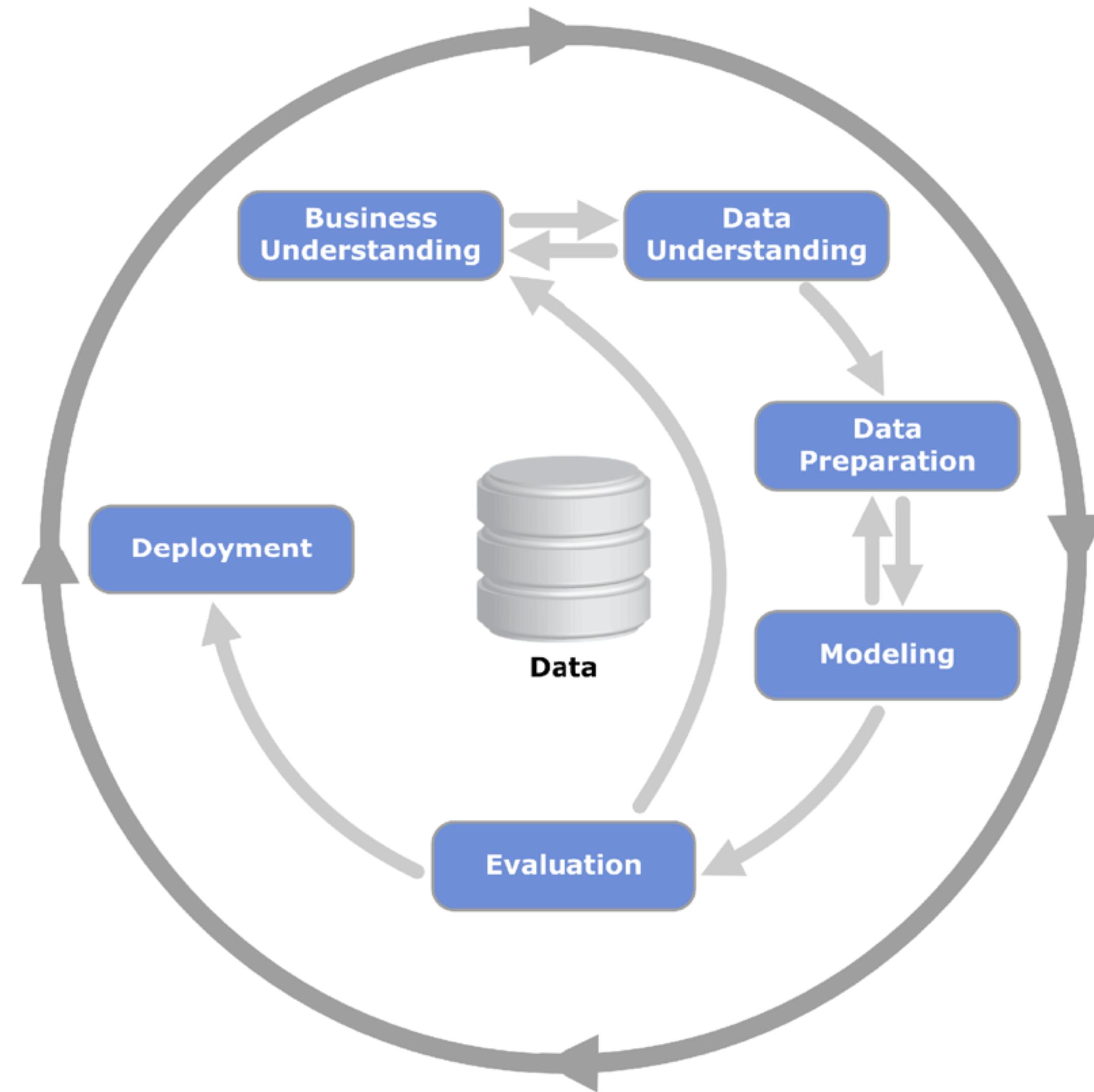
La ciencia de datos es un proceso de principio a fin

- Orientado a producto/cliente: Validado en escenarios reales
- Iterativo: Experimental y orgánico, con mejoras y adaptaciones continuas
- Data Driven: Basado en el estado de los datos





Metodologías Data Science: CRISP-DM





La “Ciencia” de la Ciencia de Datos



EL MÉTODO CIENTÍFICO



La Datificación (Transformación Digital) de las empresas

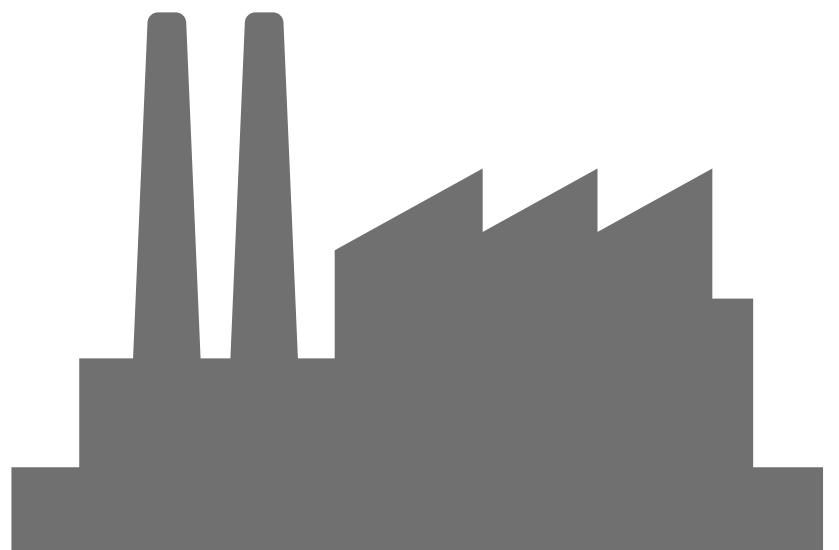
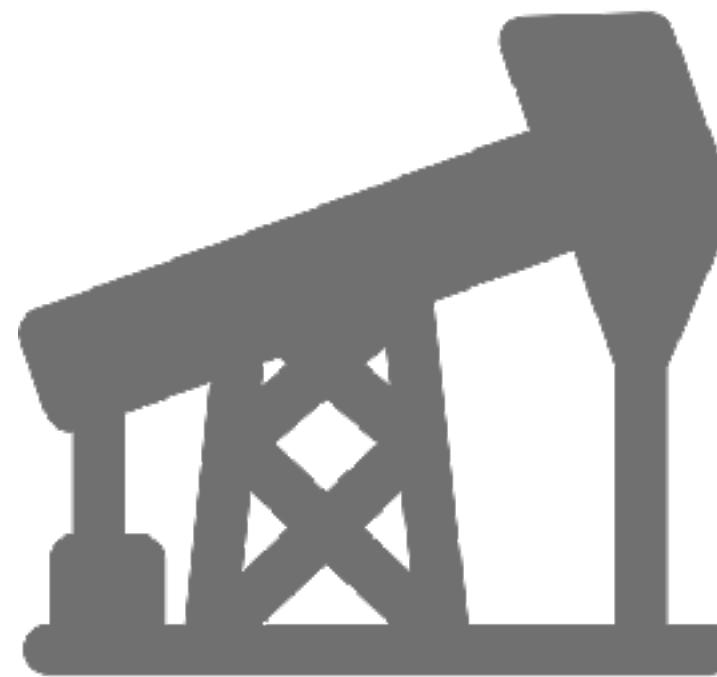
“DATA IS THE NEW OIL”





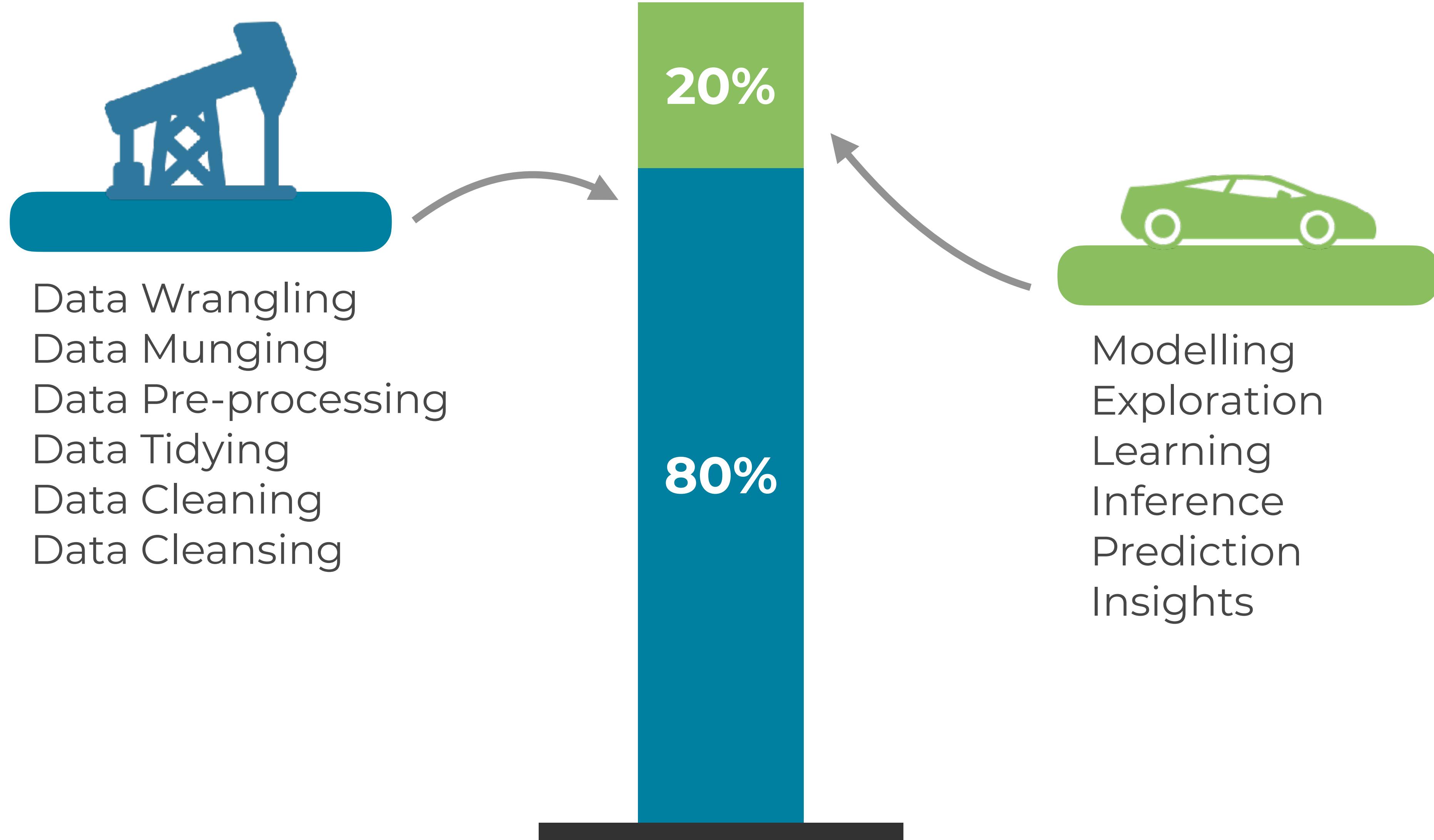
El ciclo de vida de los datos

- Extracción
- Preprocesamiento
- Almacenamiento
- Exploración
- Modelado
- Visualización





El “Trabajo Sucio”





Arquitecturas de datos

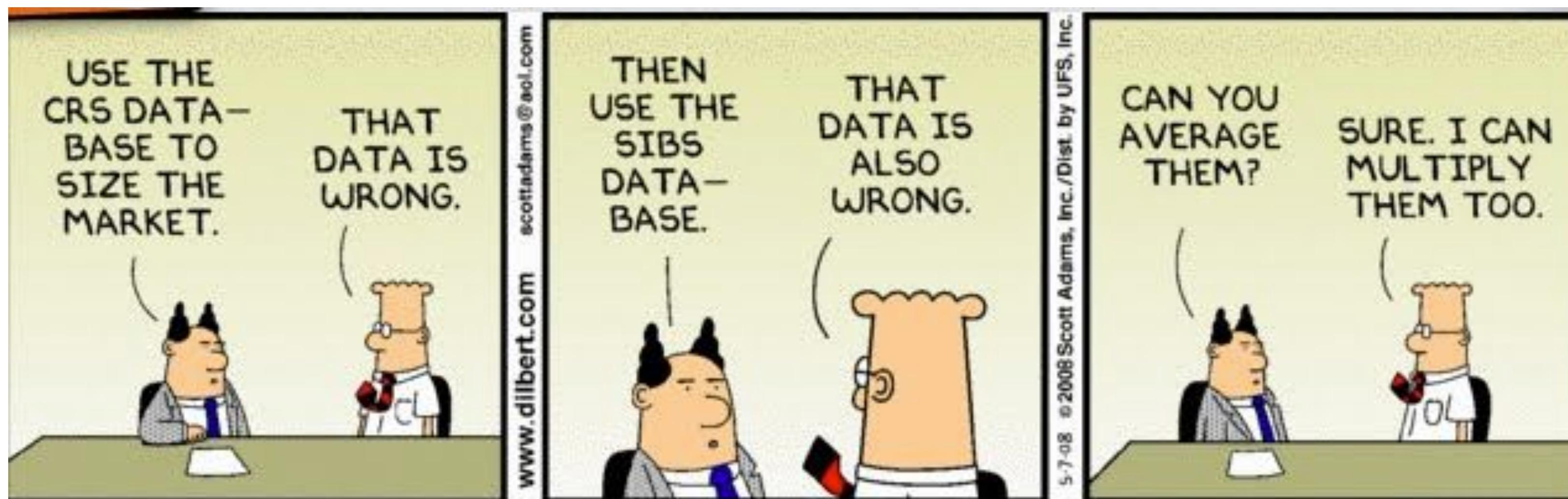
- El 80% de lo consideramos como análisis de datos no lo es realmente
- Pasamos la mayor parte del tiempo transformando y trasladando los datos
- Un formato o un sistema nunca es óptimo para todos los procesos
- Estos procesos requieren distintos perfiles en cada etapa





Calidad del dato

- Los fallos de calidad del datos cuestan aproximadamente \$600 mil millones en EEUU anualmente
- “Better data will always beat better models”
- “Garbage in = Garbage Out”





Limpieza/Preprocesado/Transformación

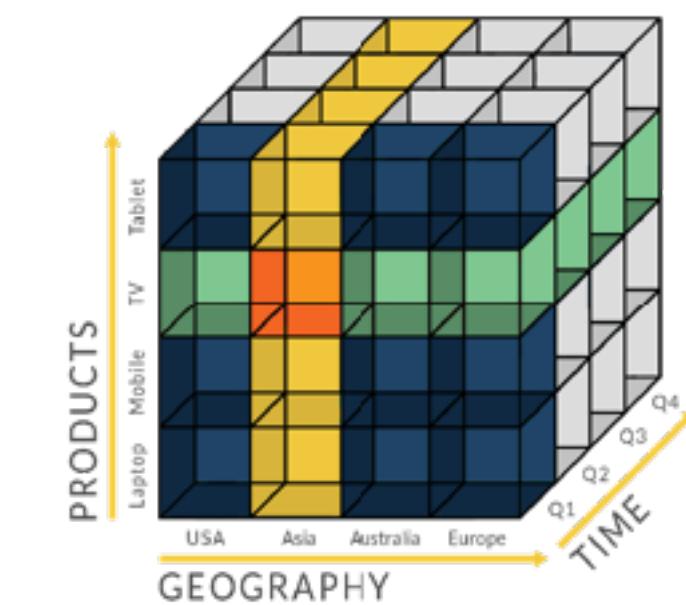
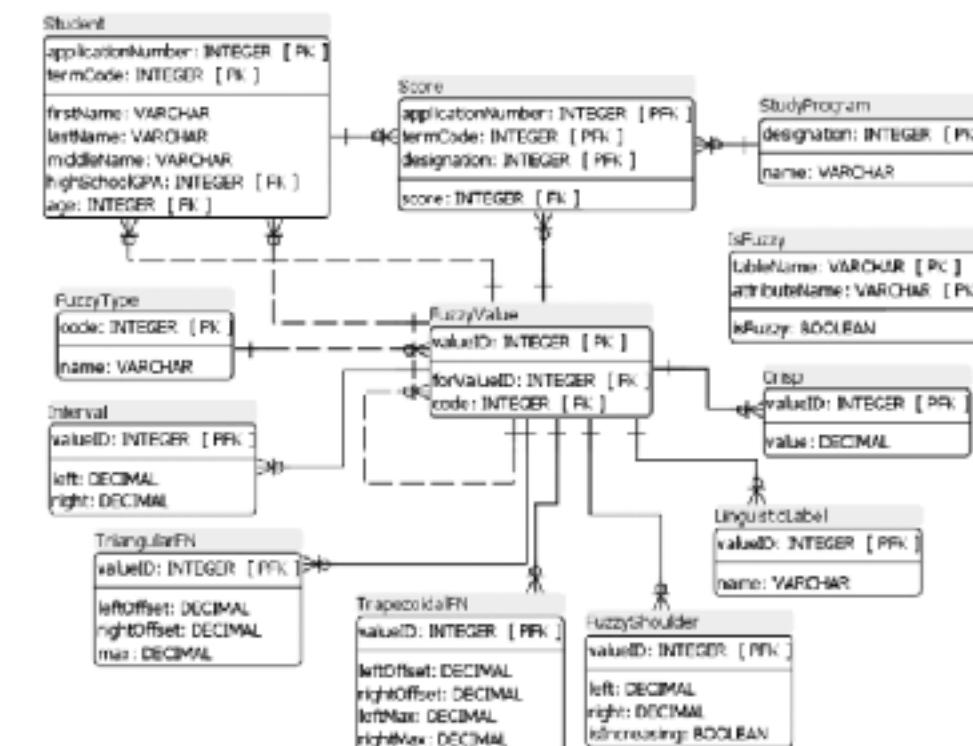
- Responsabilidad compartida en la organización:
 - Problemas de captura
 - Fallos de Hardware
 - Bugs en el Software
 - Errores humanos en la entrada
 - Errores experimentales





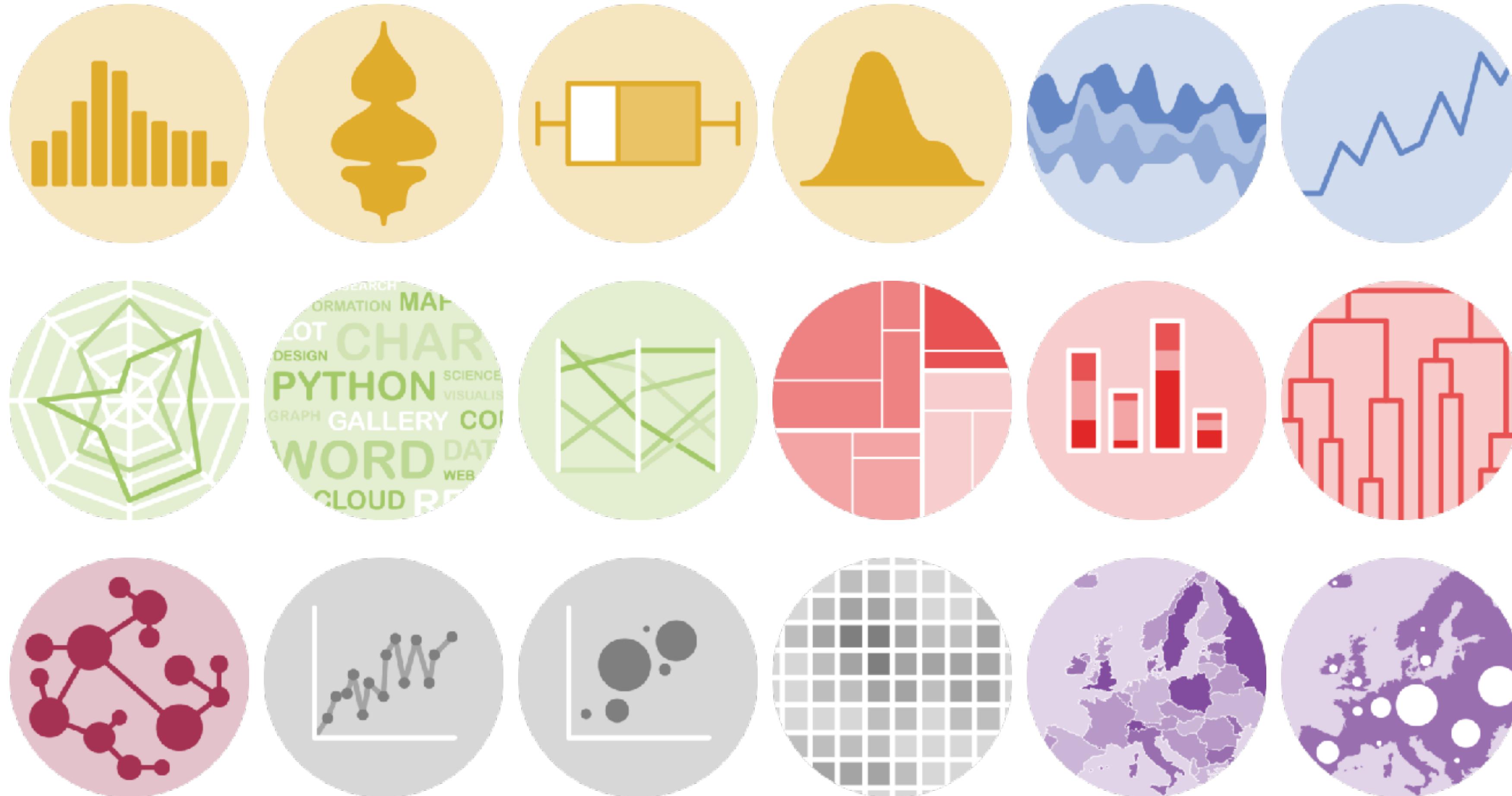
Almacenamiento

- Los datos se necesitan en distintos formatos
- La complejidad de las aplicaciones impone tecnologías concretas (SQL vs NoSQL)
- Distintas tareas requieren consultas y modelos de datos diferentes, ej: Machine Learning vs Business Intelligence





Visualización de datos

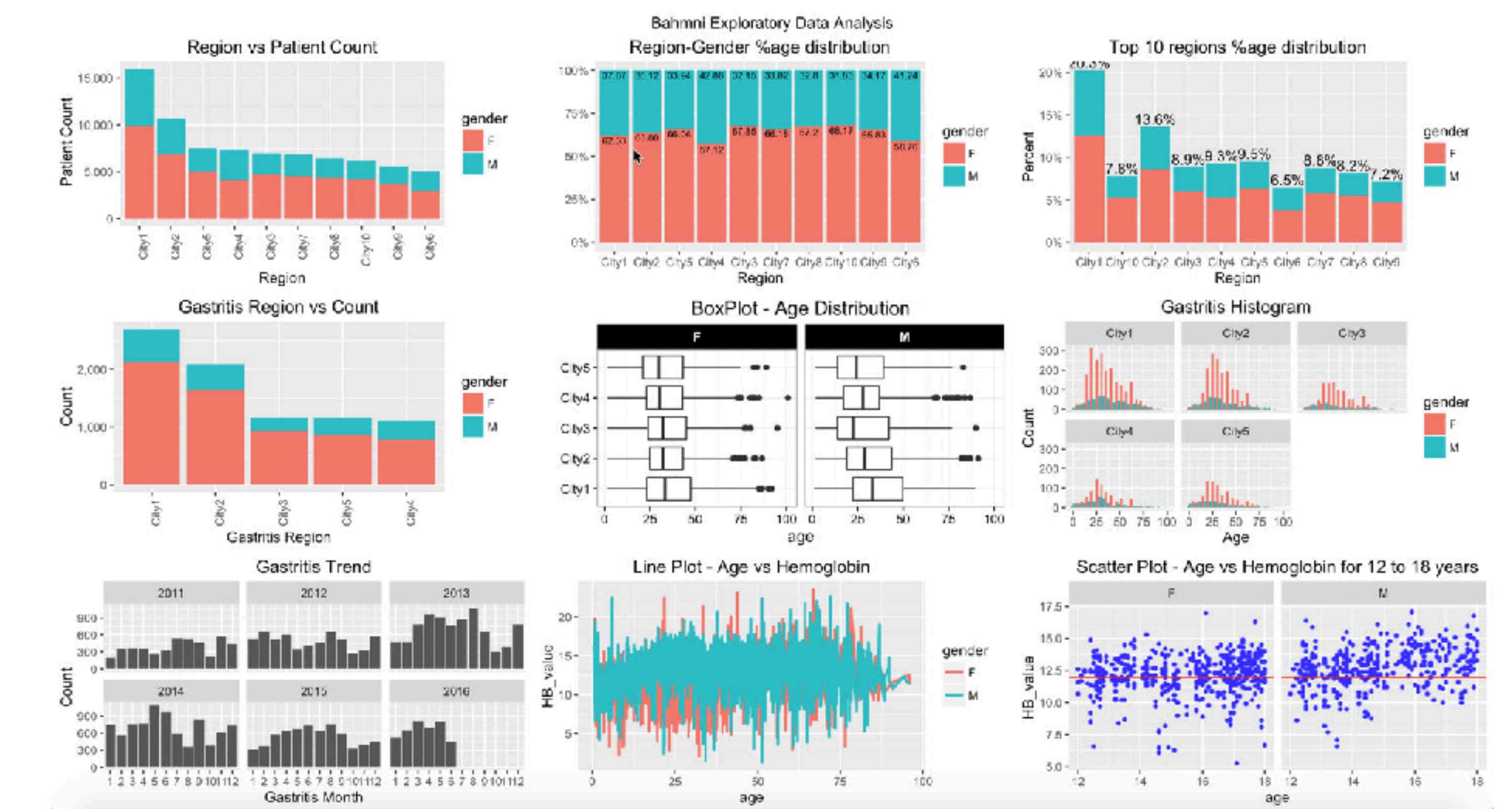


"Don't trust everything you see, even salt looks like sugar"



Análisis exploratorio (EDA)

- Realizado desde en etapas iniciales del proceso
- Aporta familiaridad y detecta problemas
- Centrado en conceptos estadísticos





Reporting y Business Intelligence

- Realizado en las etapas finales del proyecto
- Centrado en el valor final del datos en la monitorización y control de objetivos de negocio
- Realizado por expertos de negocio menos técnicos, utilizando herramientas especializadas





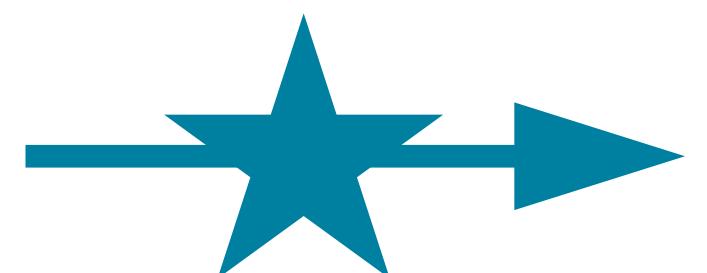
Modelado e Inteligencia

- El mundo real es **complejo, incierto y genera un flujo continuo de datos**
- Podemos capturar una **muestra de los datos generados**
- Pretendemos **explicar el proceso original usando esa muestra**
- Construiremos un **modelo que abstraiga el proceso tras los datos**

Proceso Generativo Muestra de Datos

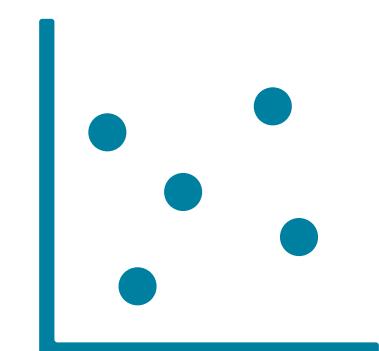


Algoritmo



Modelo

$$f(\mathbf{x})$$



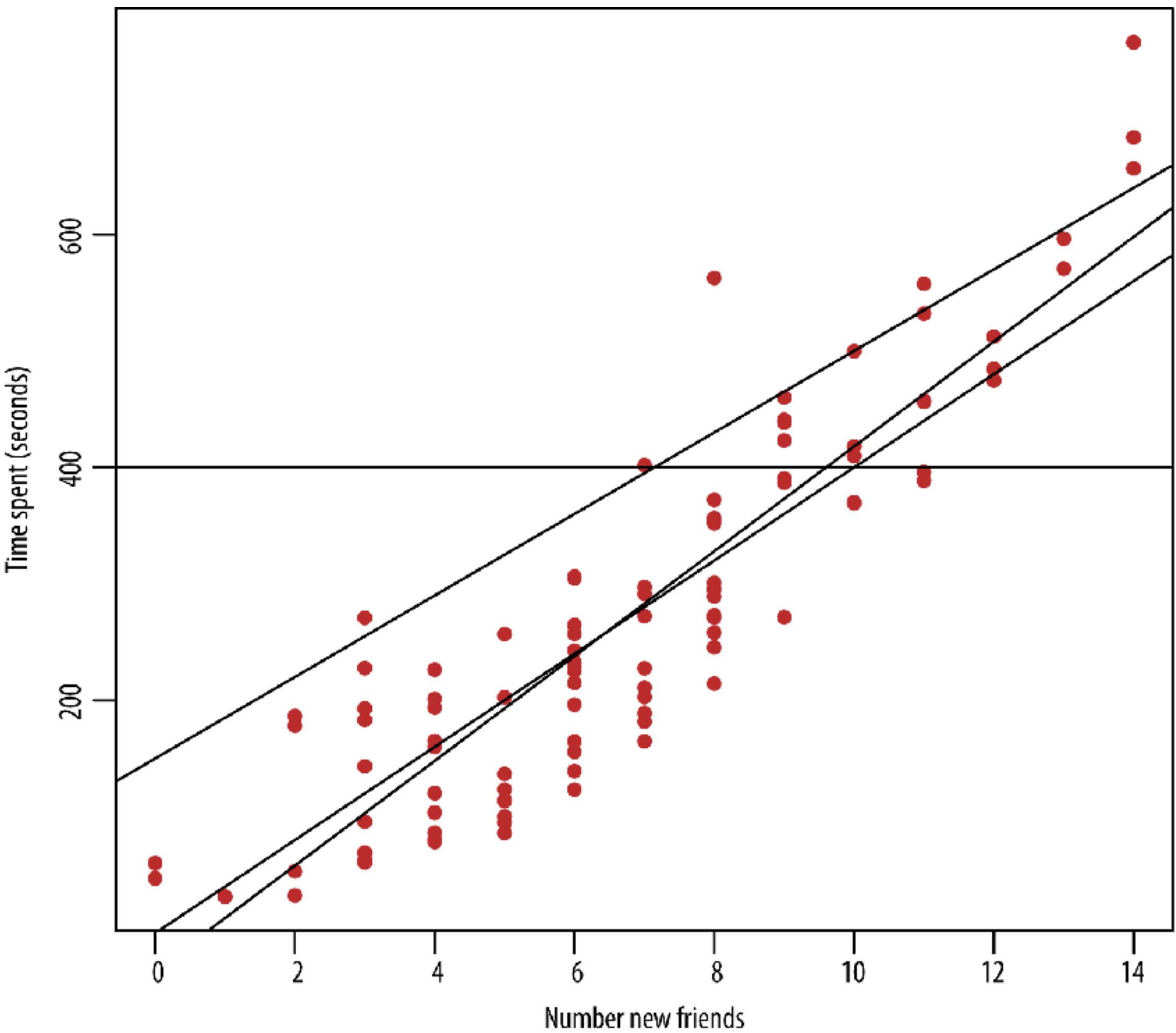


Modelado Estadístico

- Ajustar un modelo dada una muestra de datos

$$y = \alpha + x\beta$$

- Aprender los parámetros que mejor se ajusten (que menos se equivoquen)



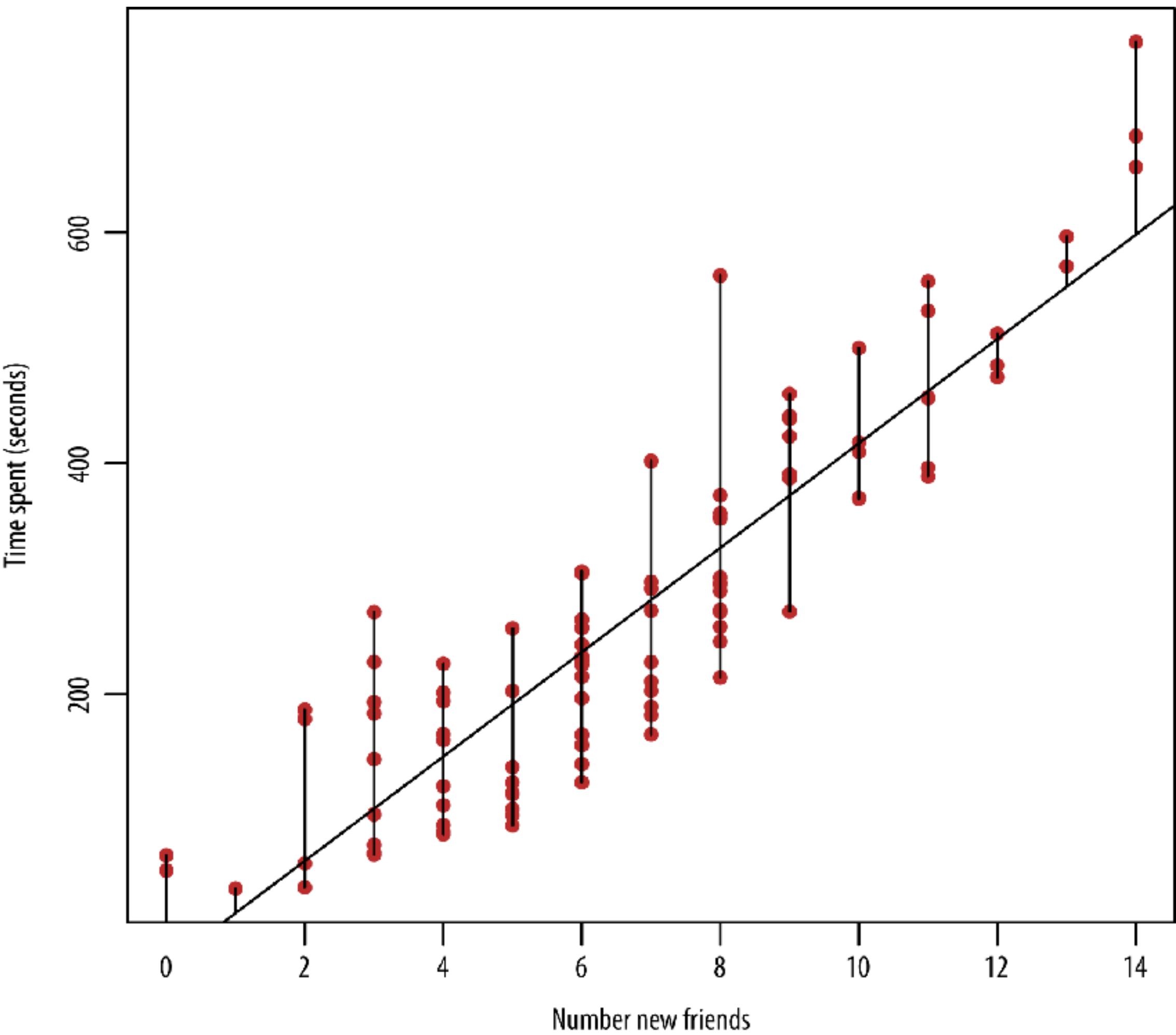


Modelado Estadístico

- Ajustar un modelo dada una muestra de datos

$$y = \alpha + x\beta$$

- Aprender los parámetros que mejor se ajusten (que menos se equivoquen)





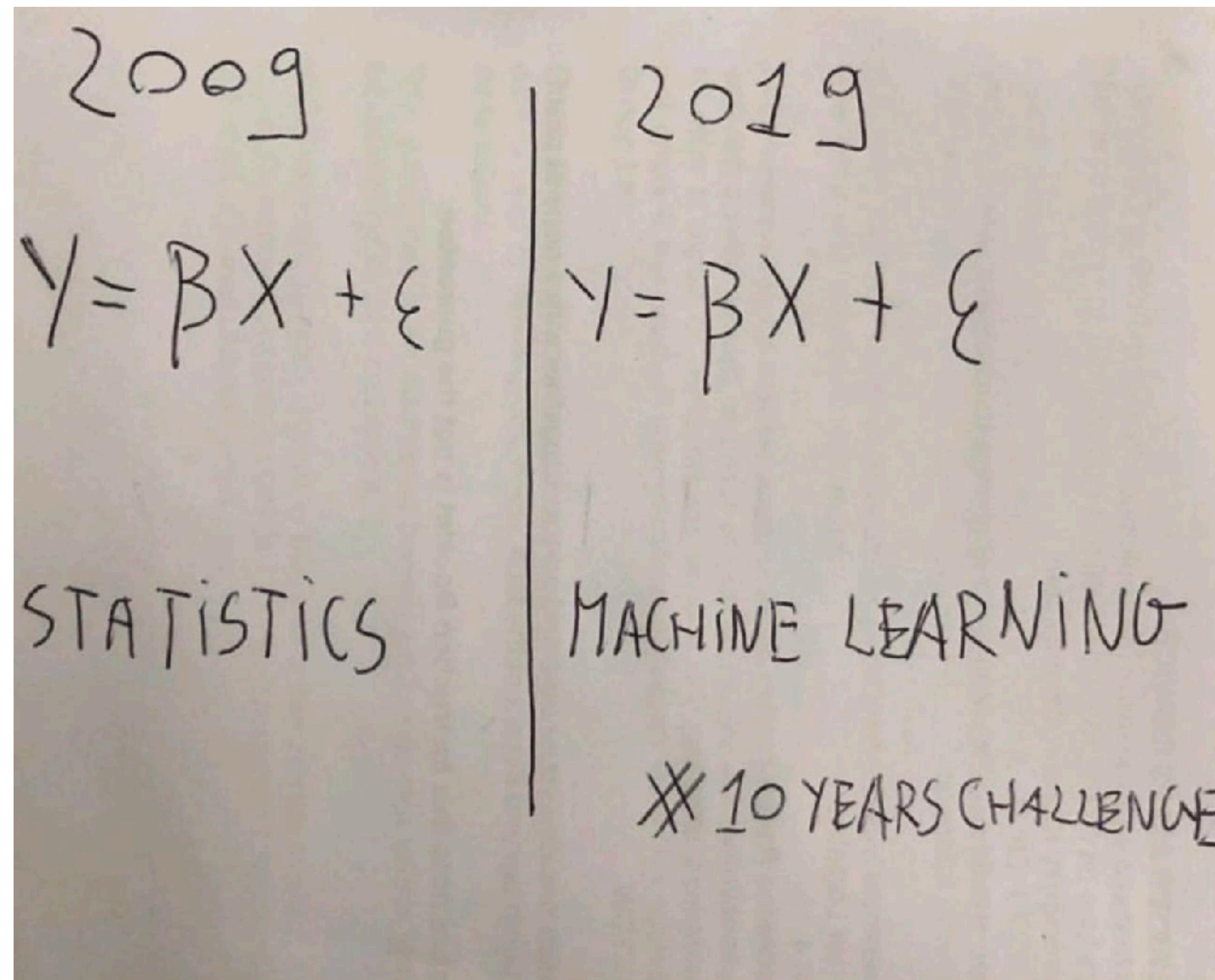
Machine Learning

- Aprendizaje de **patrones** a partir de los datos utilizando **algoritmos estadísticos automáticos**
- Usar los modelos para realizar procesos de **inferencia** tales como predicción, clasificación...
- ¿Cuál es la diferencia entre machine learning y modelado estadístico?





Isn't machine learning just glorified statistics?





Machine Learnins VS Modelado Estadístico

STATISTICAL MODELING

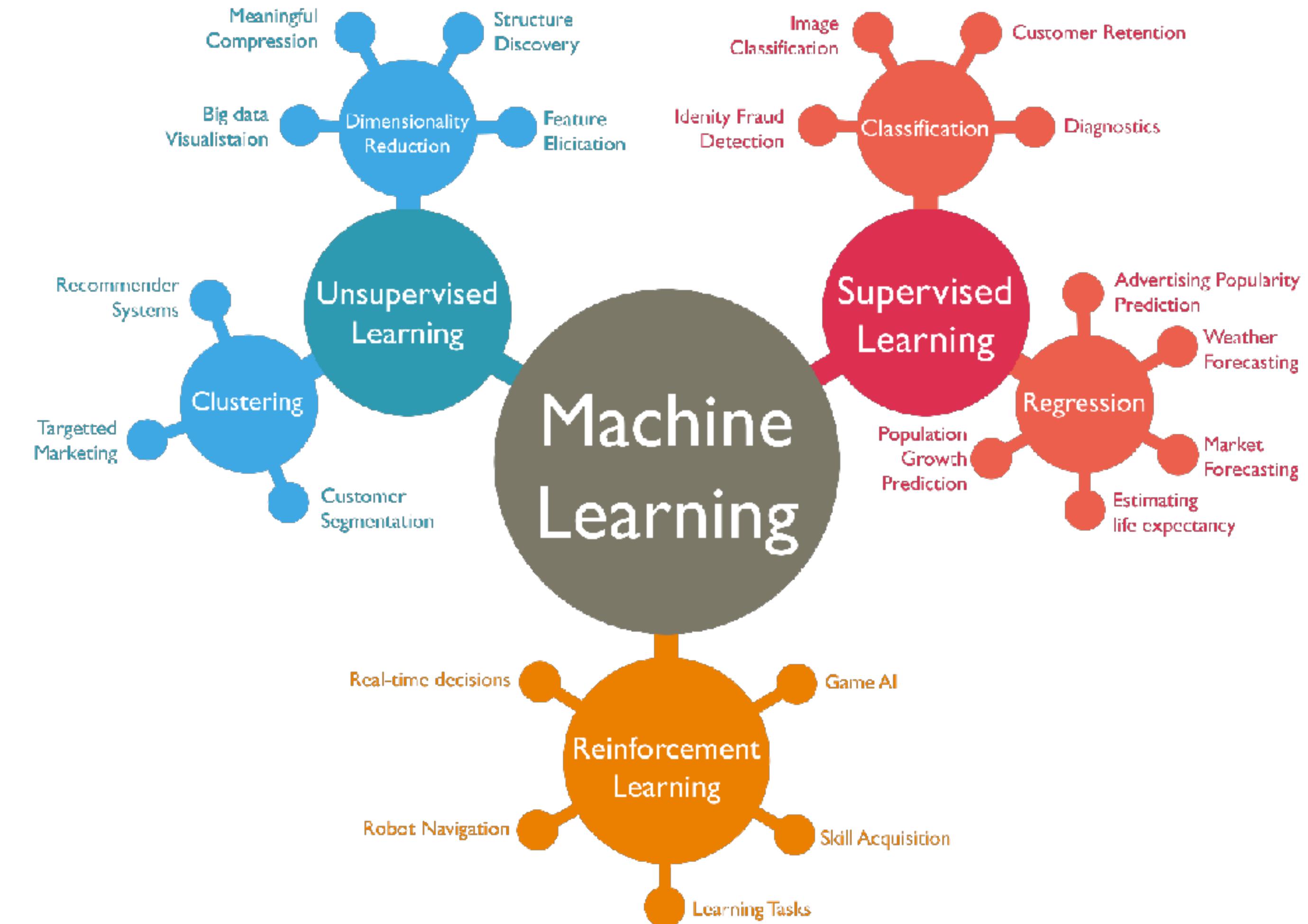
- Interpretabilidad de los parámetros
- Explicación de los resultados
- Obtener la mejor aproximación del modelo generativo

MACHINE LEARNING

- Centrado en los resultados
- Suele producir modelos “caja negra”
- Busca la automatización e integración



Tipos de Machine Learning





Hacer las preguntas correctas

Preguntas

Directas

Medibles

Clasificación

“¿Esto es un perro o un gato?”

Regresión

“¿Cuánto valdrá esta vivienda?”

Anomalías

“¿Es esta la temperatura habitual?”

Clustering

“¿Cuantos tipos de bacteria hay?”

Recomendación

“¿Qué película preferirá este usuario?”

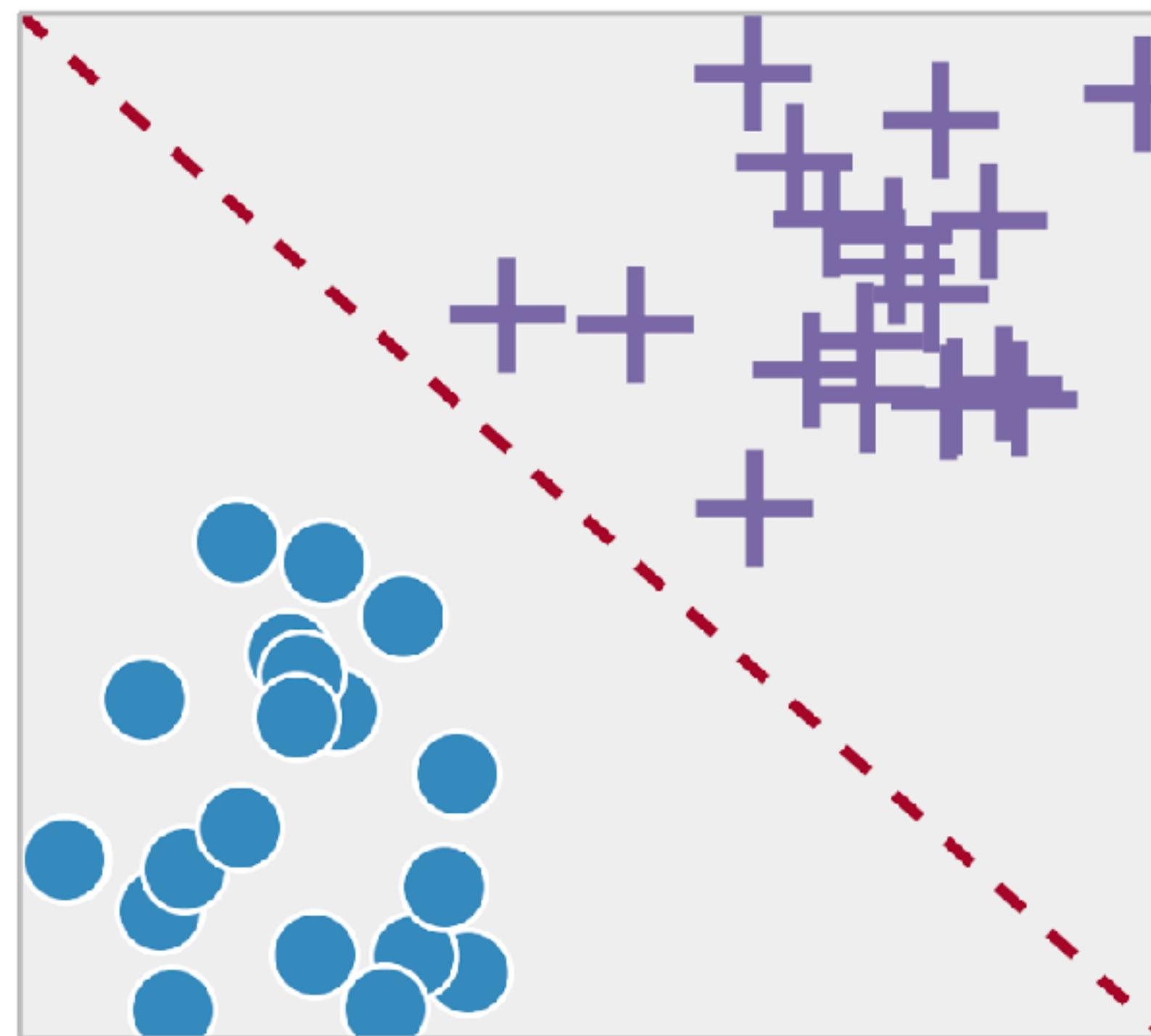
Refuerzo

“¿Tendré un accidente si cambio de carril?”

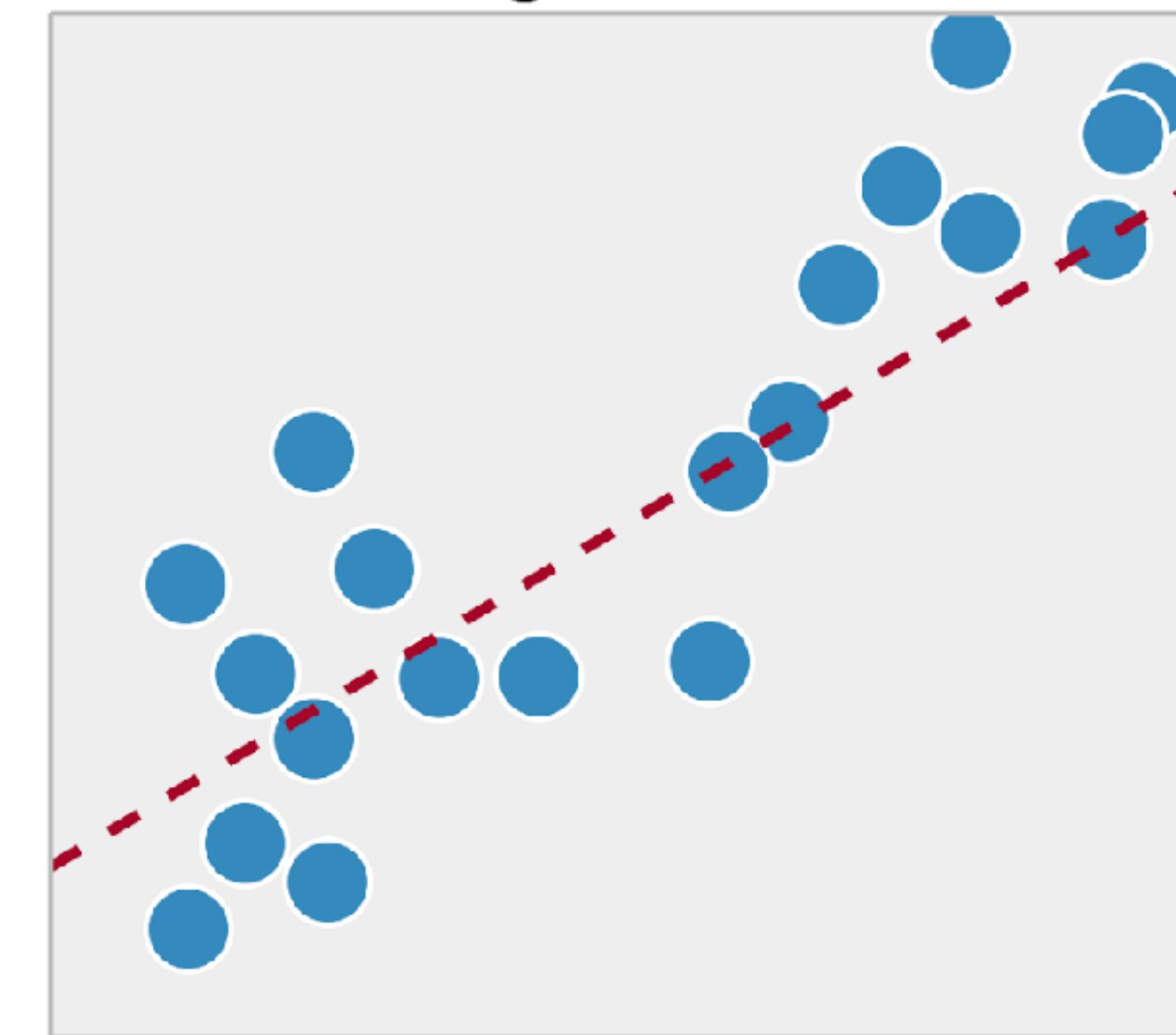


Aprendizaje Supervisado

Classification

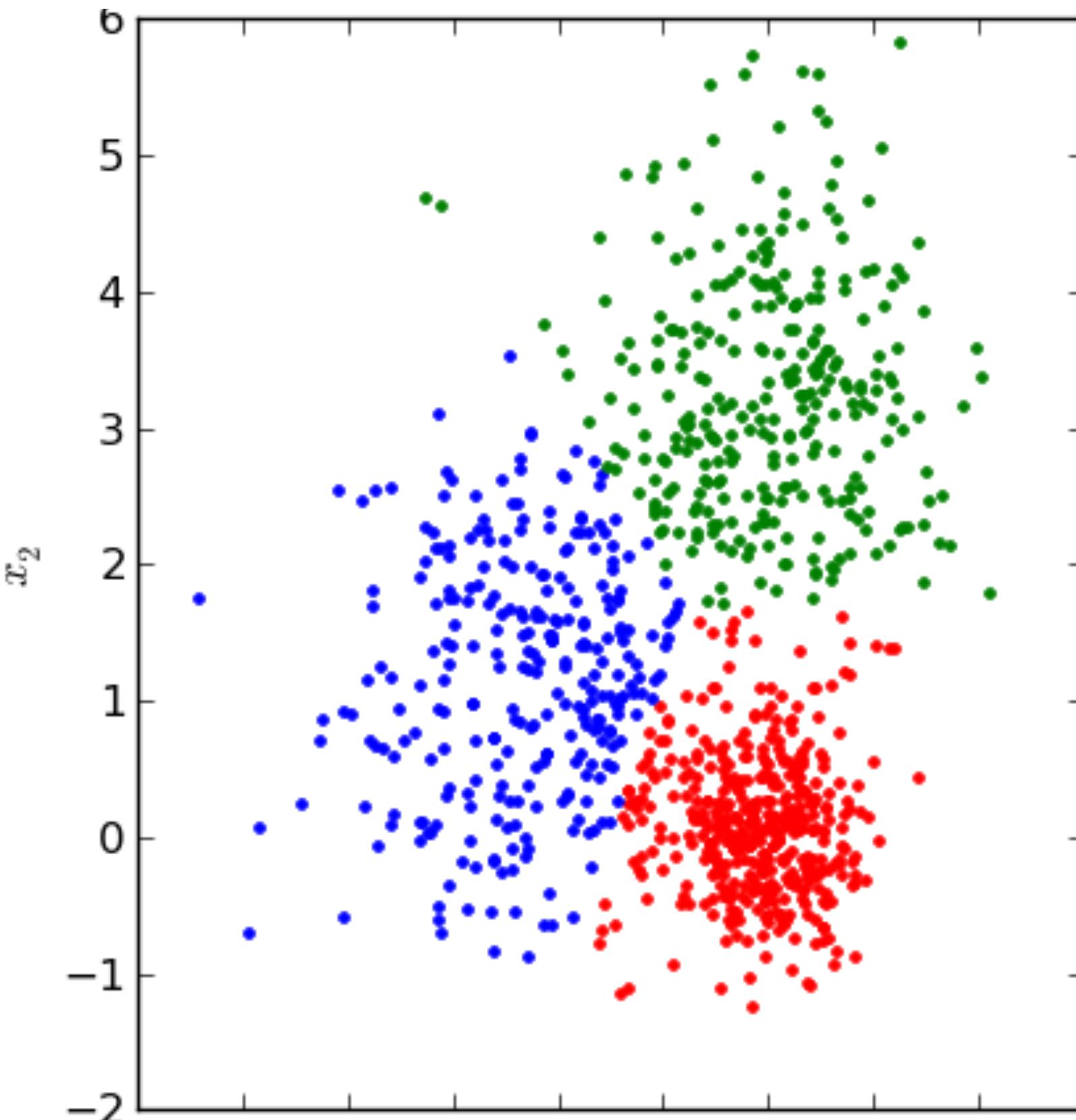
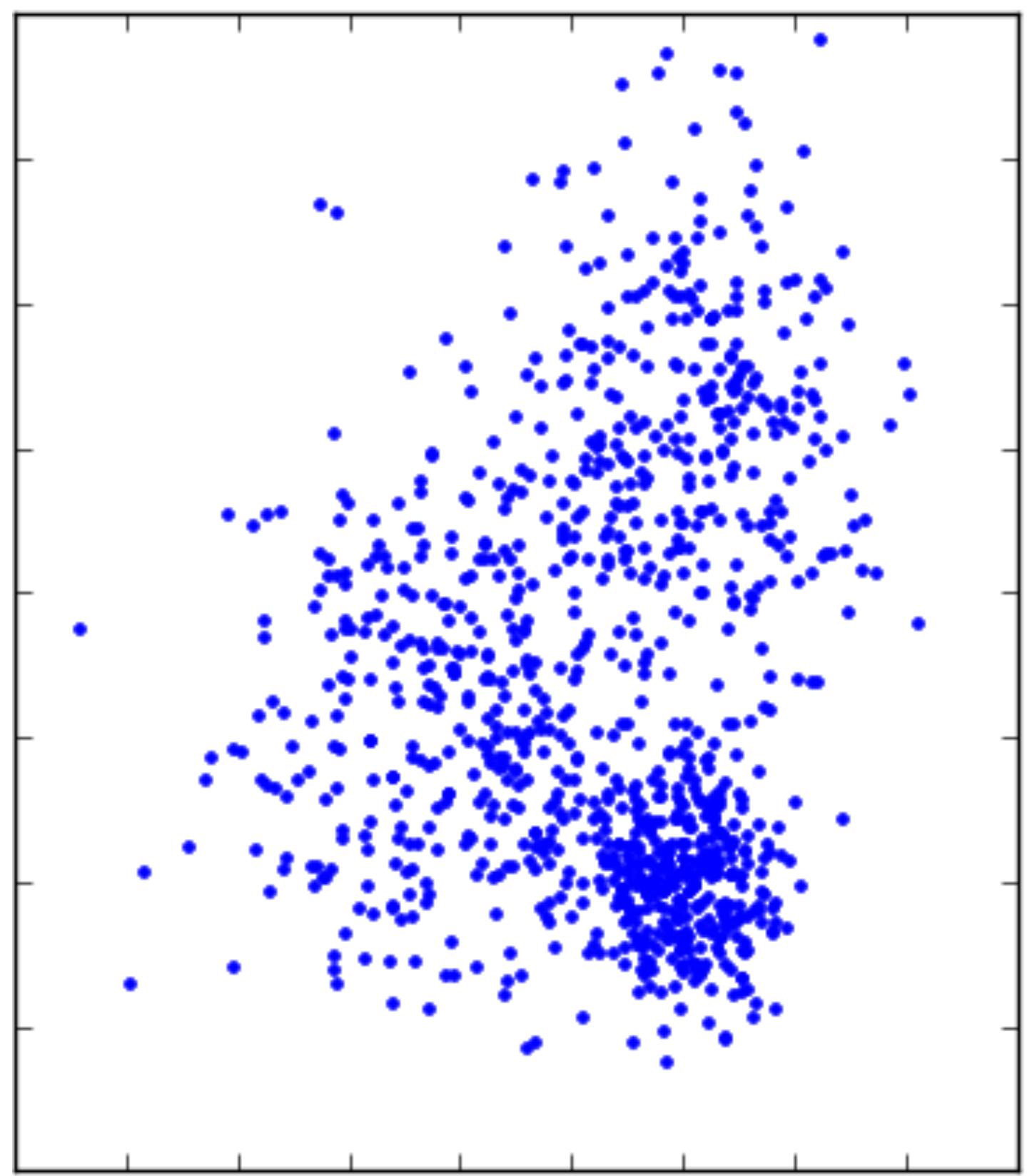


Regression



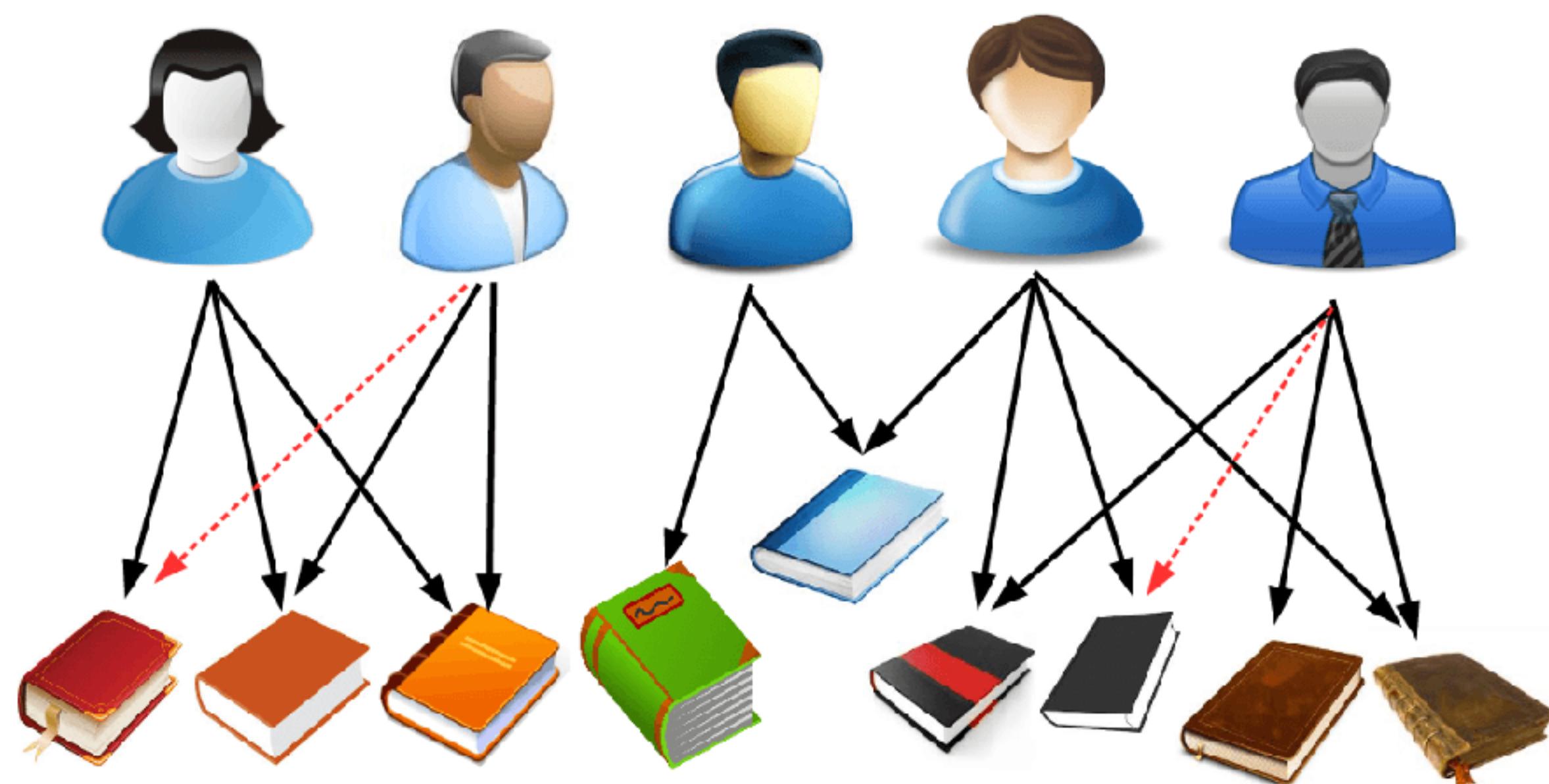


Aprendizaje no Supervisado





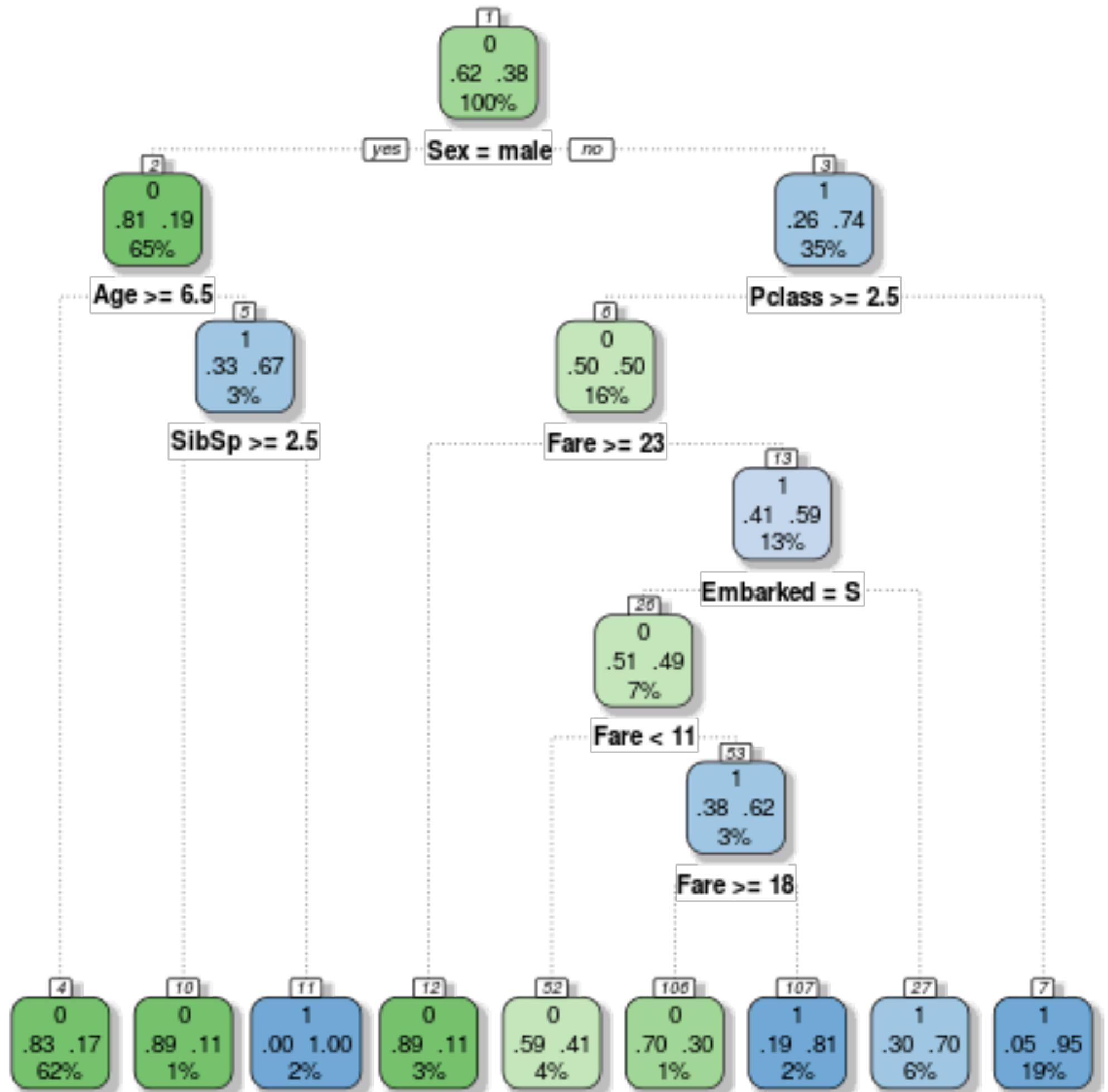
Sistemas de Recomendación





Tipos de Modelos

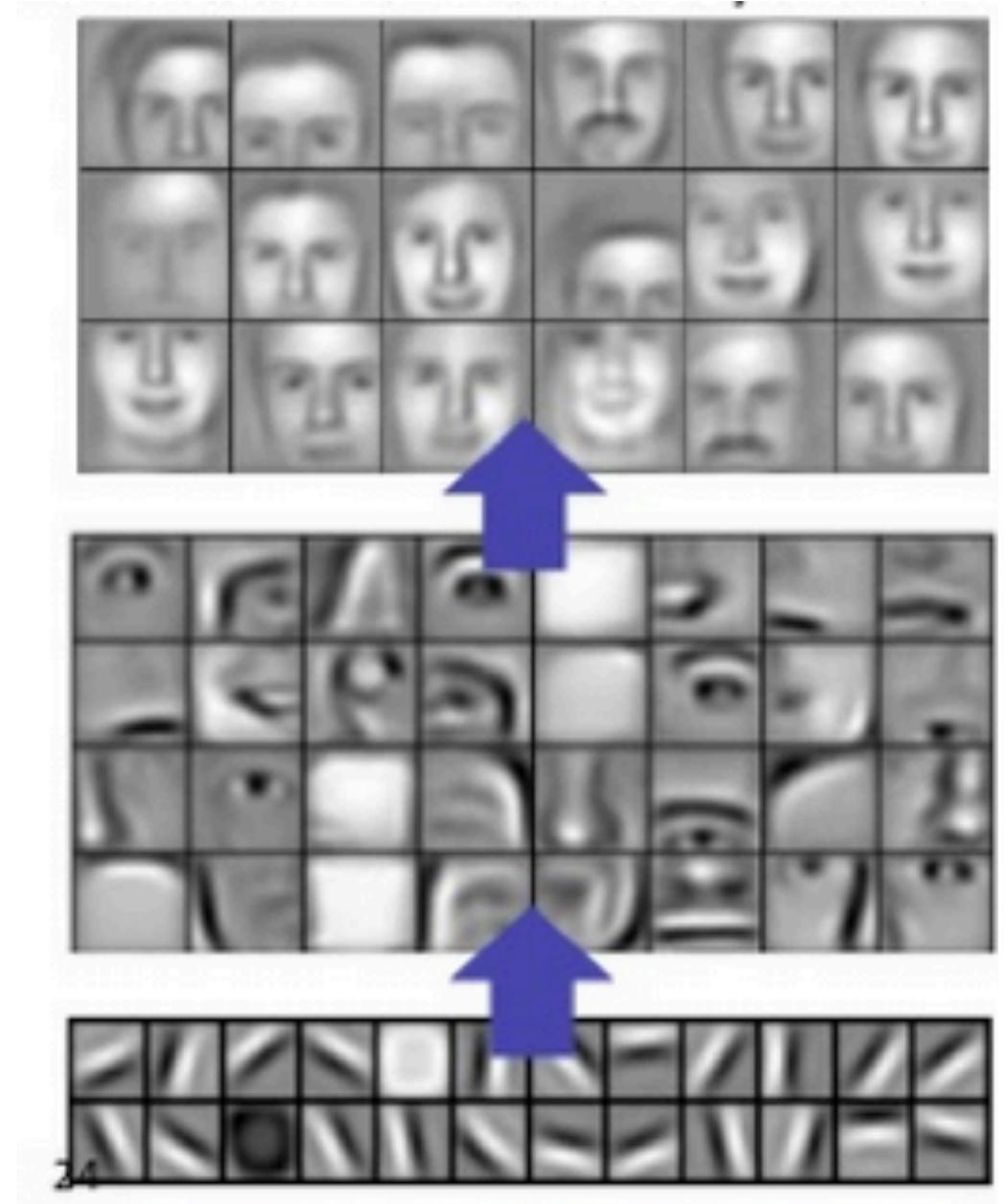
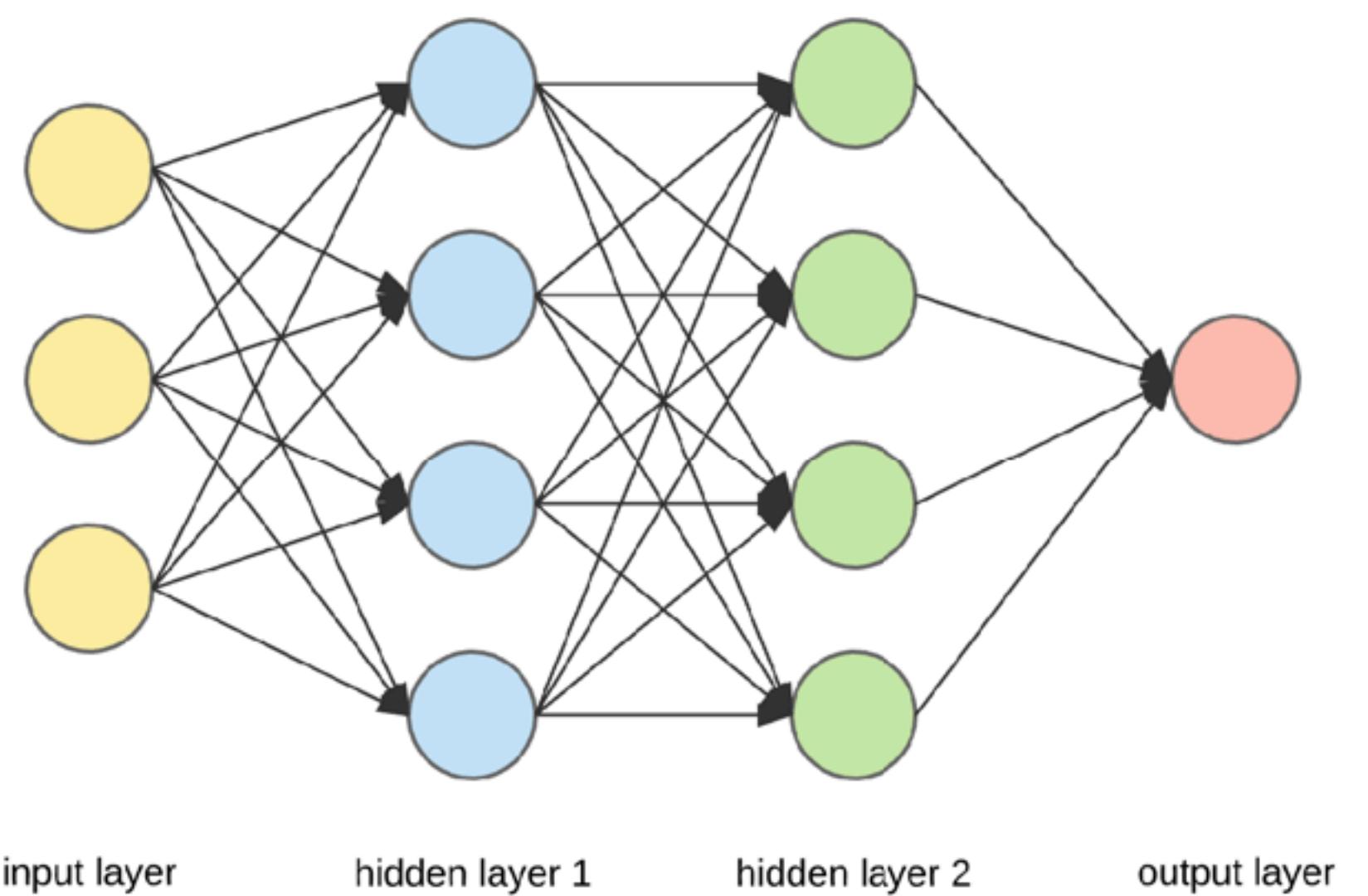
- Árbol de decisión
- Aprende un árbol binario que clasifica objetos
- Utiliza métricas de probabilidad y teoría de la información
- Puede ser fácilmente interpretado





Tipos de Modelos

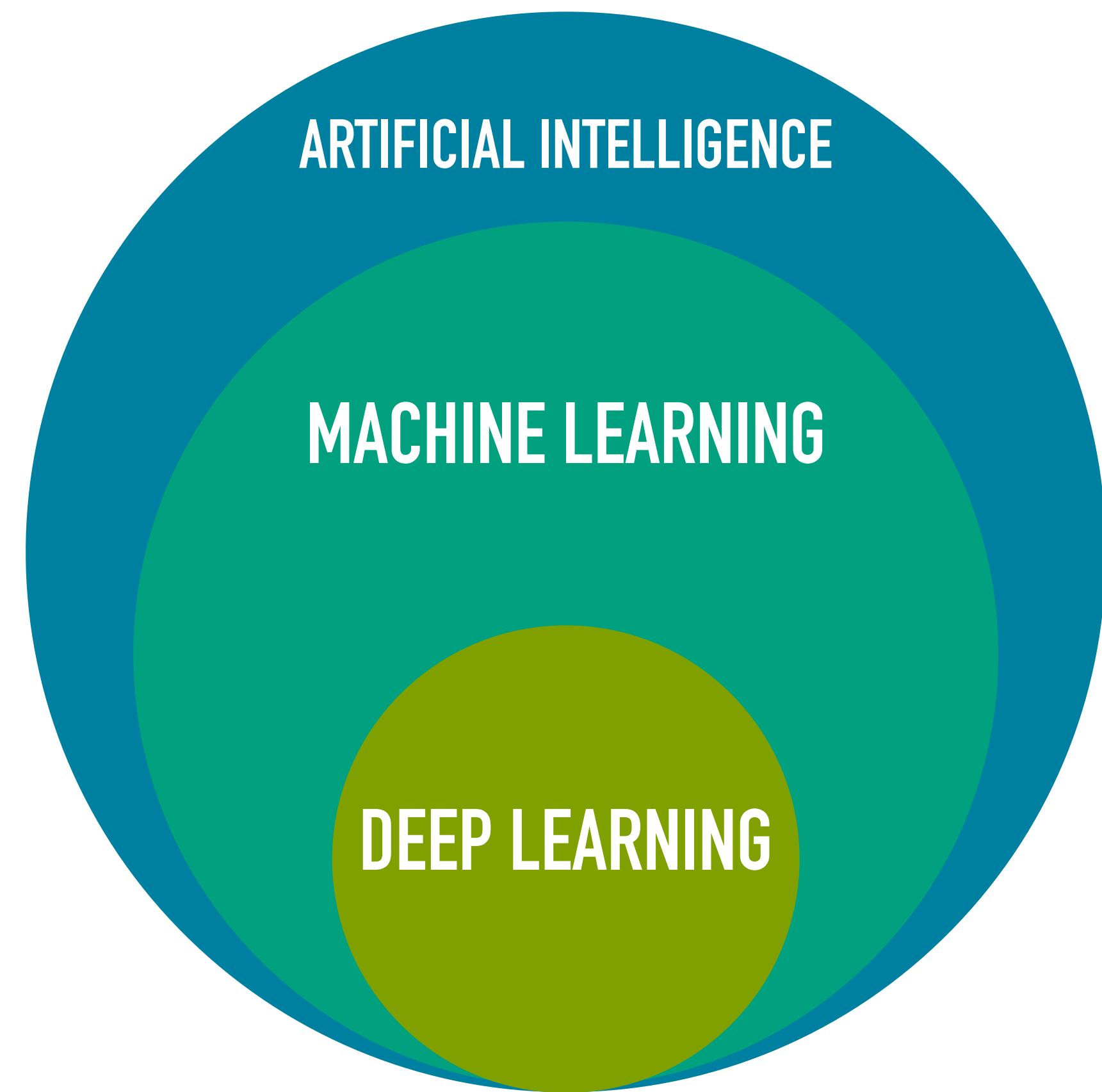
- Red Neuronal
- Modelo de caja negra
- Muy poca o nula interpretación
- Inspirado en como se conectan las neuronas entre si





Deep Learning

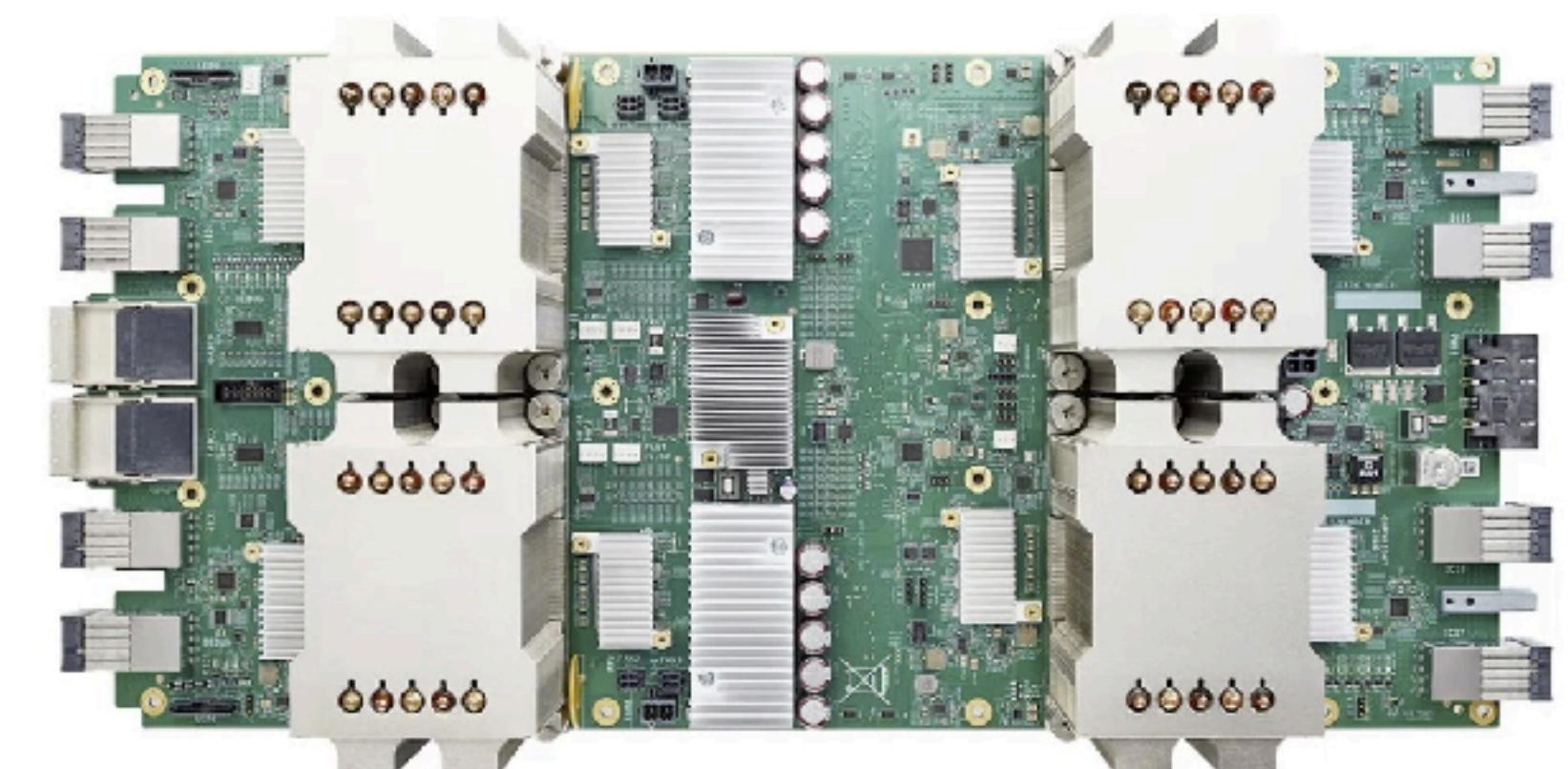
- Técnicas avanzadas que aprenden redes inmensas
- Requieren grandes cantidades de datos y capacidad de cómputo
- Los modelos no pueden ser interpretados
- En ocasiones se pueden transferir a otros dominios





Deep Learning = Datos + Cómputo Intensivo

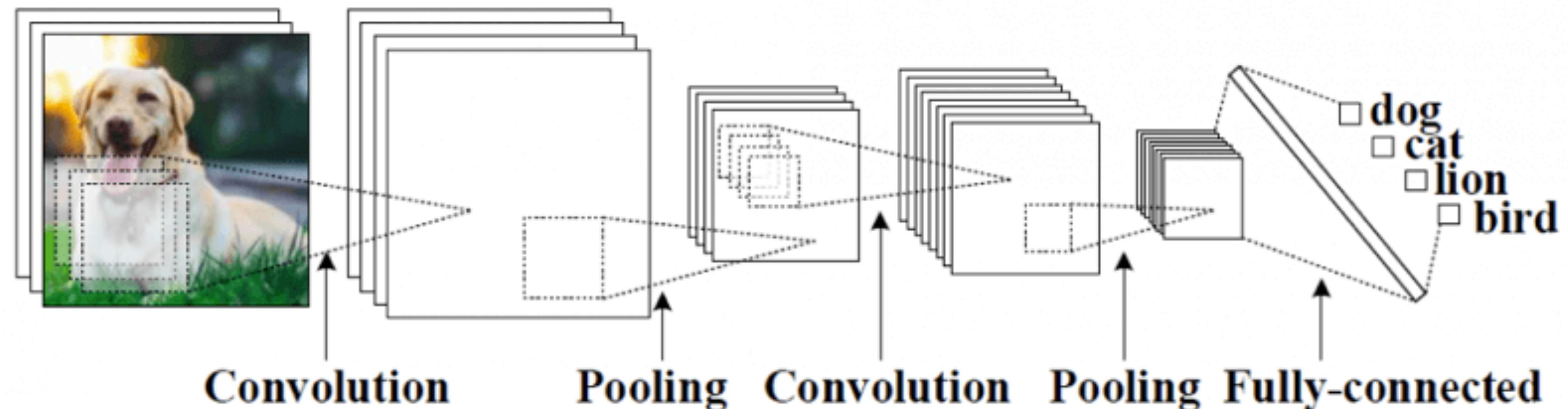
- La teoría básica data de los 50s y 80s
- Popularizado en los 2000s
- Reservado a grandes volúmenes de datos y hardware especializado (GPU y TPU)





Ejemplos de Deep Learning

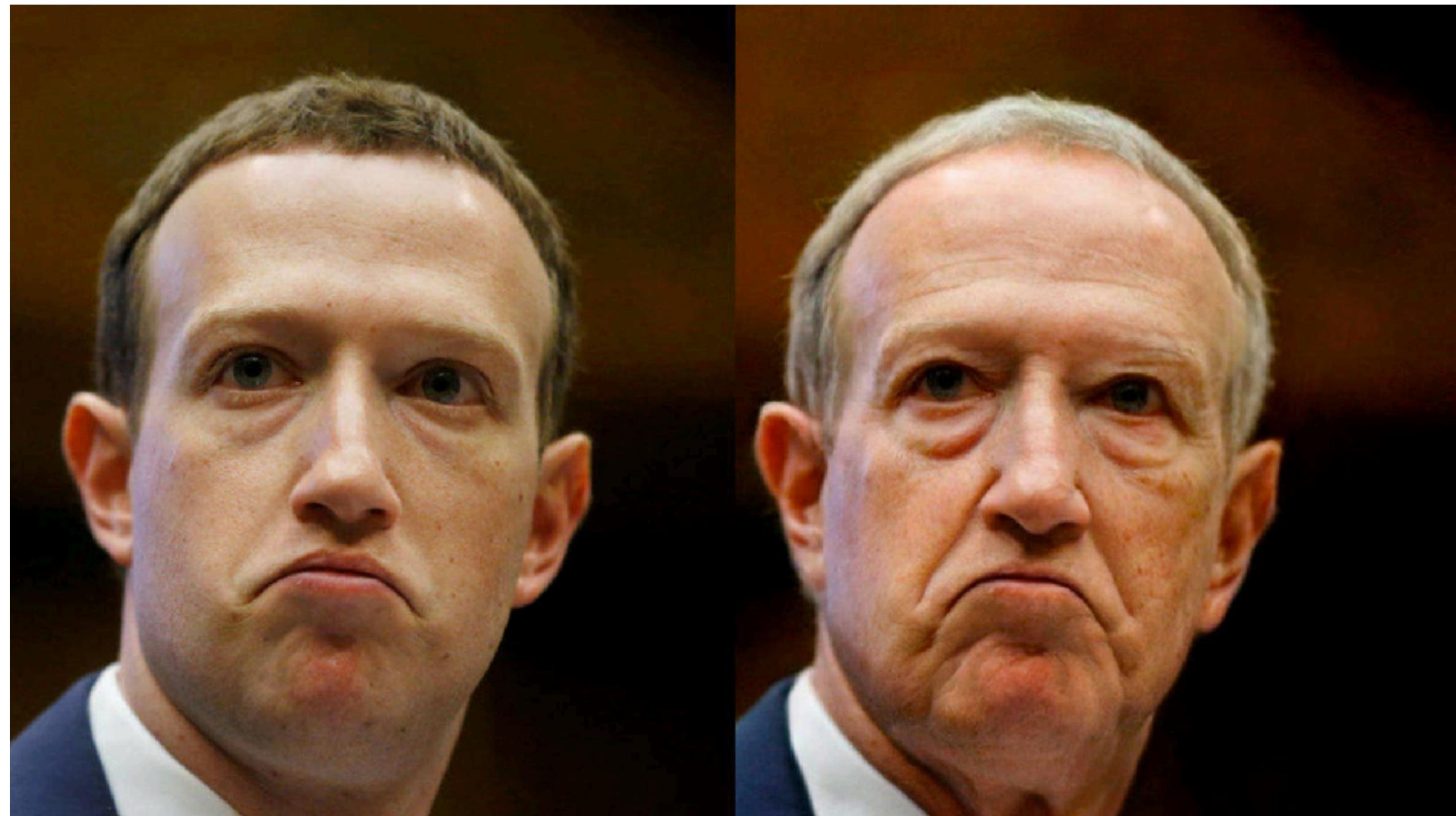
- Redes Convolucionales: Clasificación y regresión
- Reducción de la dimensionalidad sin pérdida
- Aprendizaje de muchas características intermedias





Ejemplos de Deep Learning

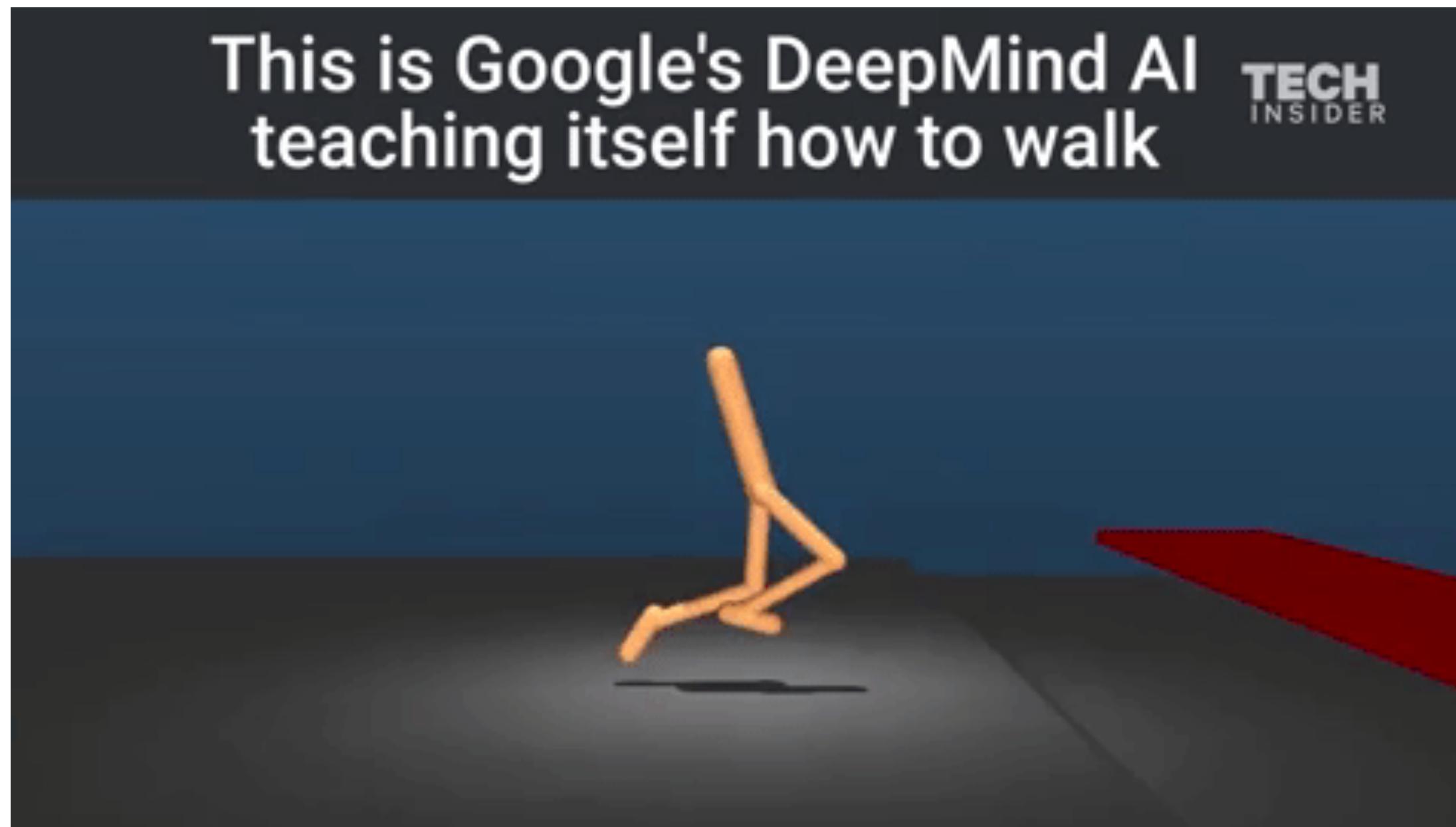
- Redes generativas adversariales (GANs): Generación de datos sintéticos





Ejemplos de Deep Learning

- Aprendizaje por refuerzo: Aprende de la experiencia, datos por simulación.
Aplicaciones en robótica, videojuegos, navegación...



El Ecosistema



Personas



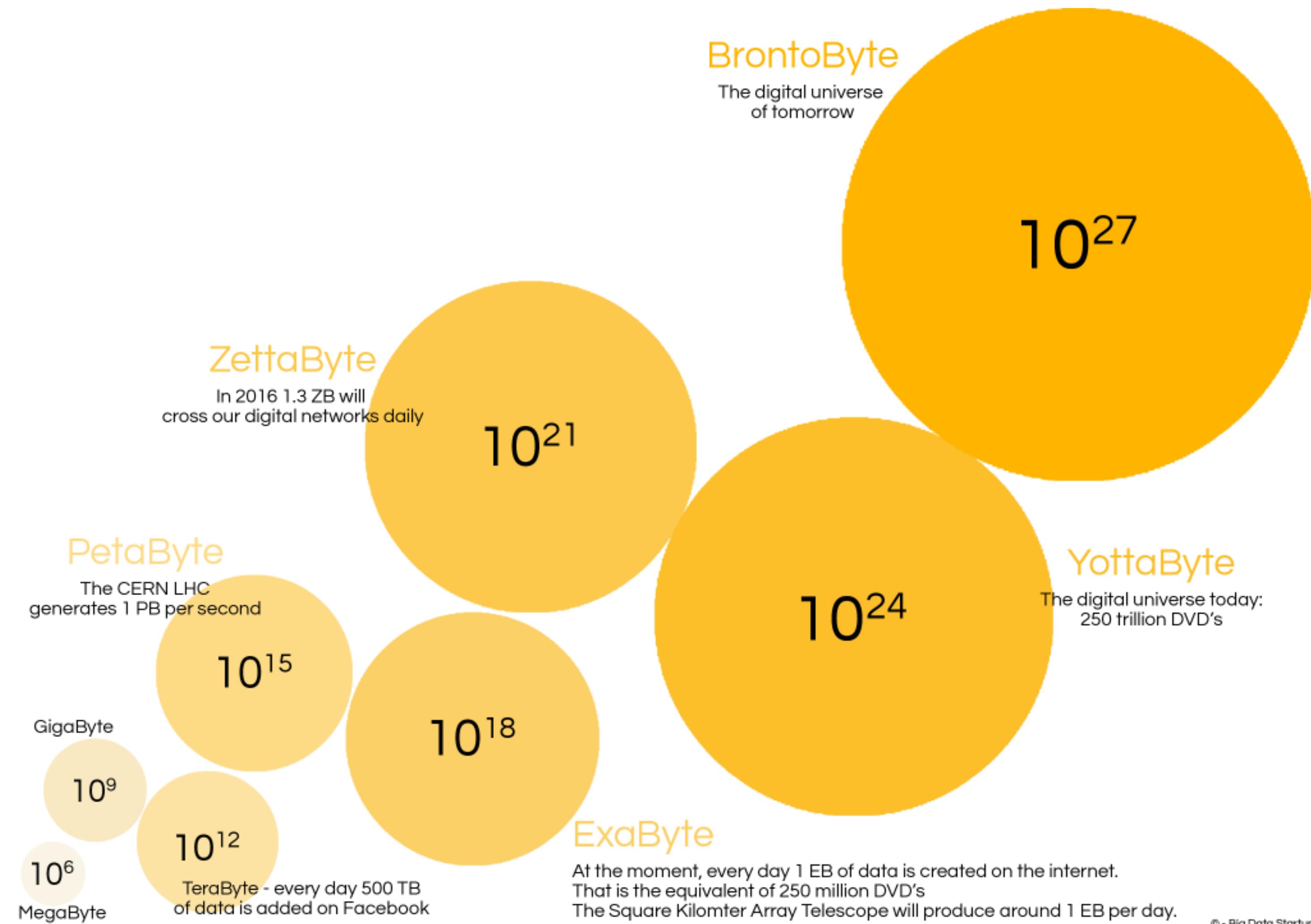
Procesos



Tecnología



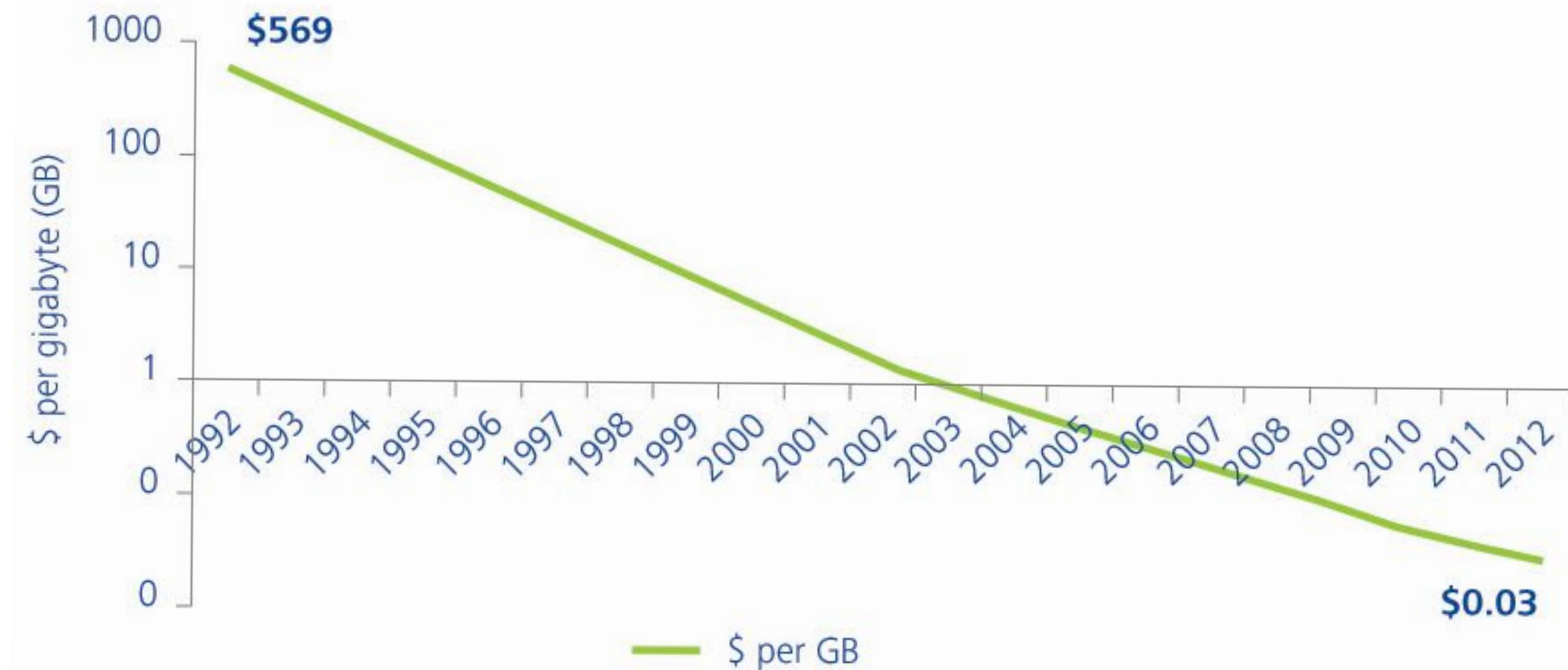
Las escalas son relativas a la tecnología





La tecnología impone límites

- “Big” es cuando no nos cabe en una sola máquina





Big Data = Big Computing Power

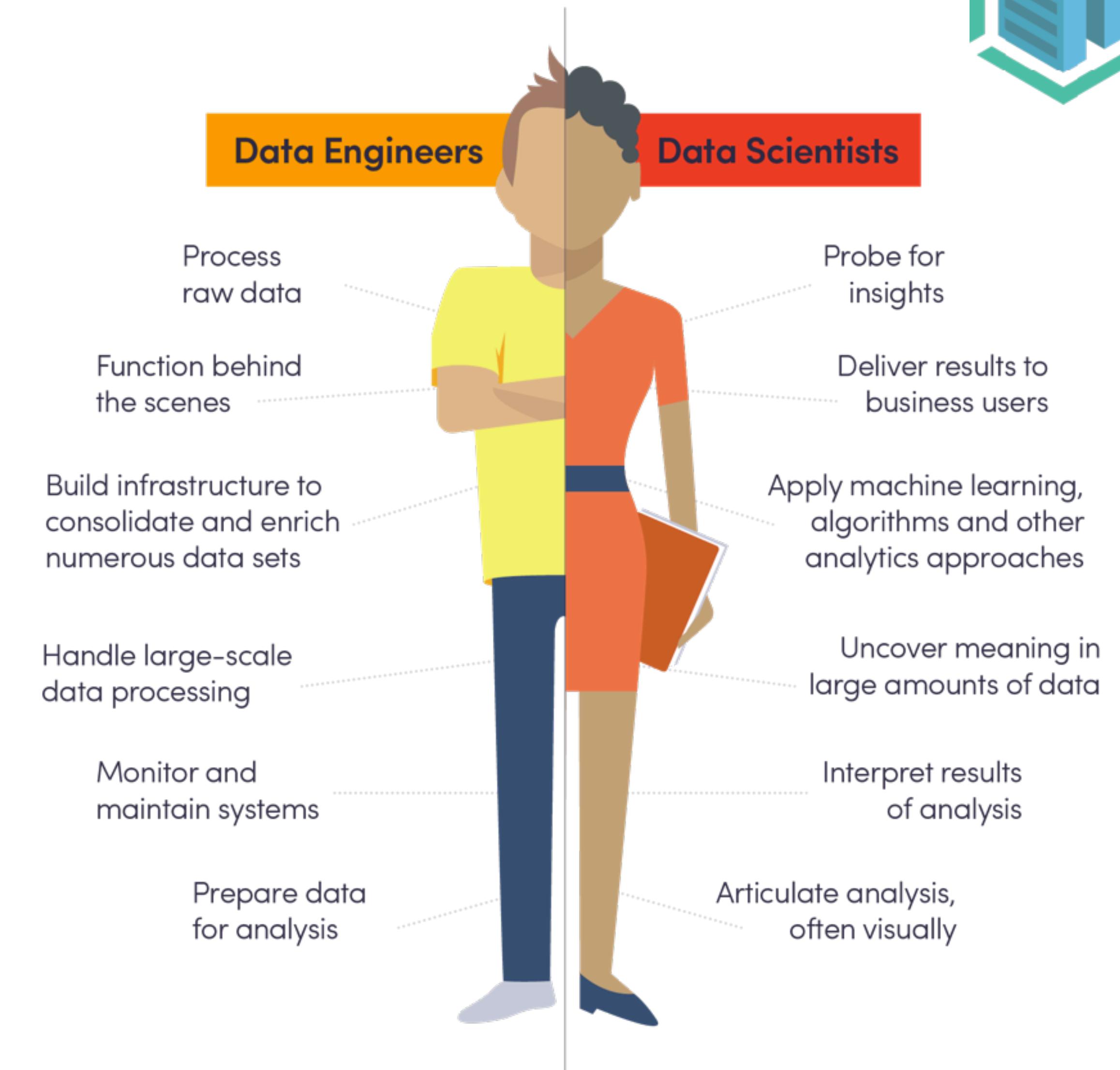
- Un set específico de tecnologías que requieres habilidades concretas
- Gestionar múltiples máquinas simultáneamente es costoso
- Necesitamos abstraer la deslocalización de los datos y la ubicuidad del cómputo





El perfil del Ingeniero de Datos

- Centrado en los datos, pero especializado en la tecnología
- Gestiona el complejo stack de big data y cloud
- Trabaja cerca del origen y la infraestructura que aloja los datos
- Clave en la captura y limpieza de los datos

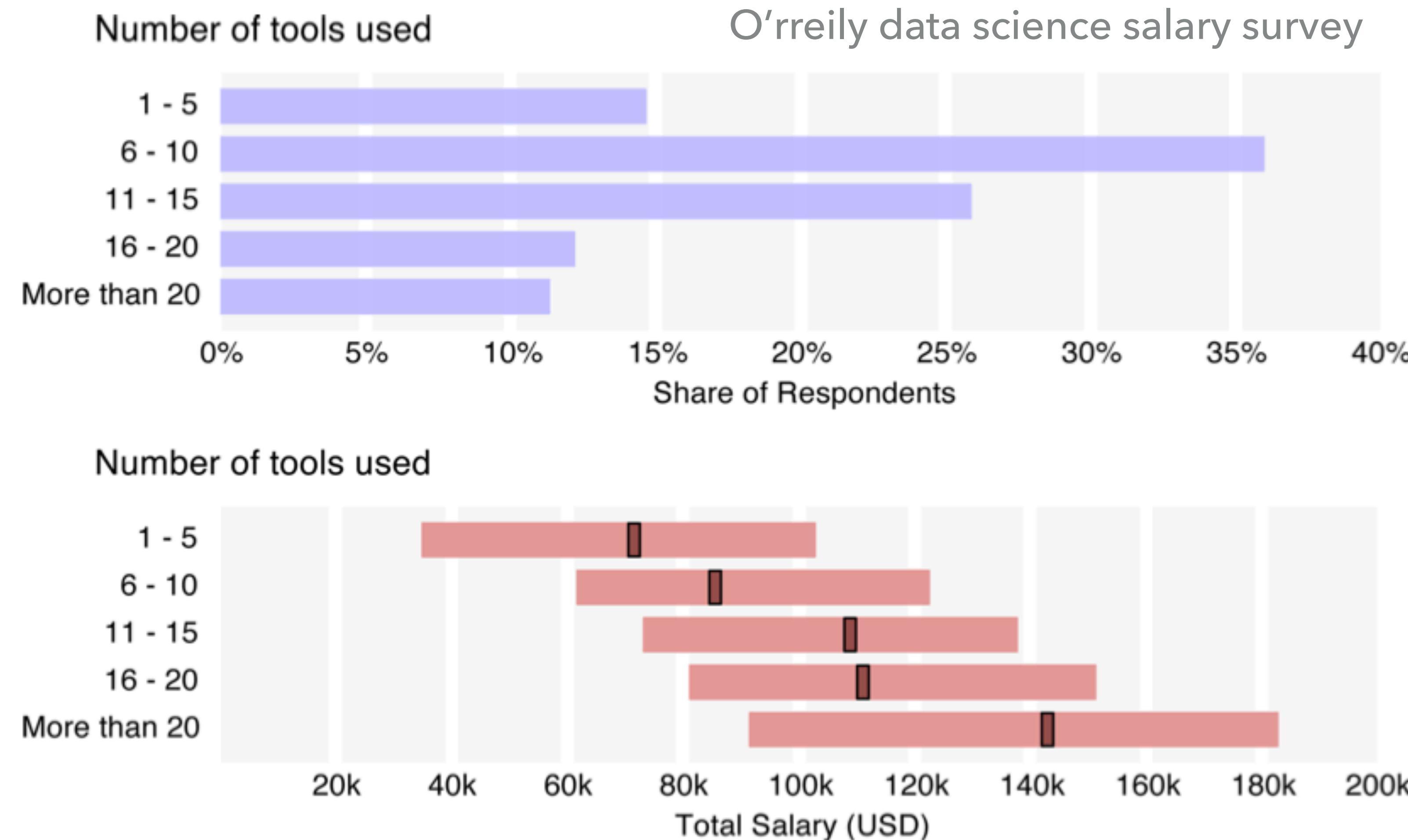


YA TENGO SUFICIENTES HERRAMIENTAS

Esto no lo dijo un Data Scientist



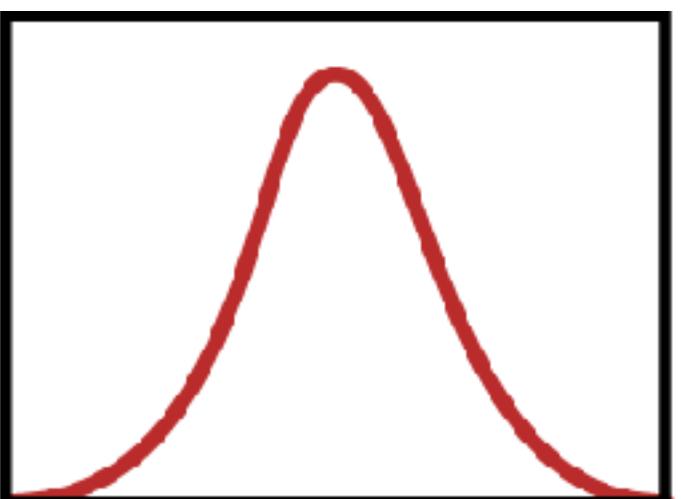
Un ecosistema altamente especializado





¡Matemáticas y Estadística!

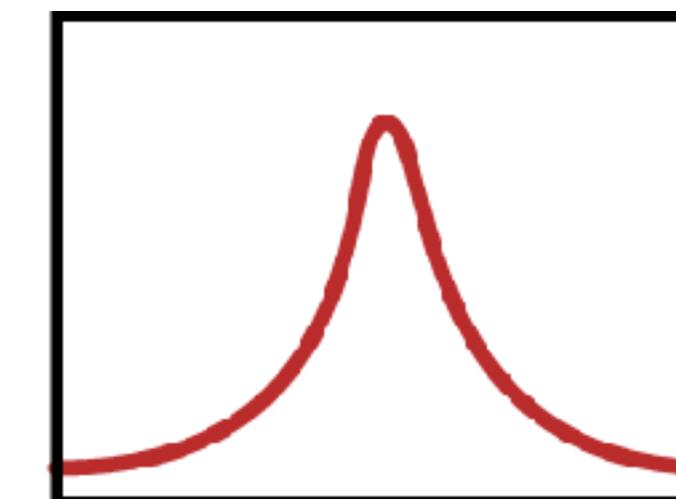
- Álgebra Lineal
- Probabilidad
- Estadística descriptiva e inferencial
- Teoría de la información
- Teoría de grafos



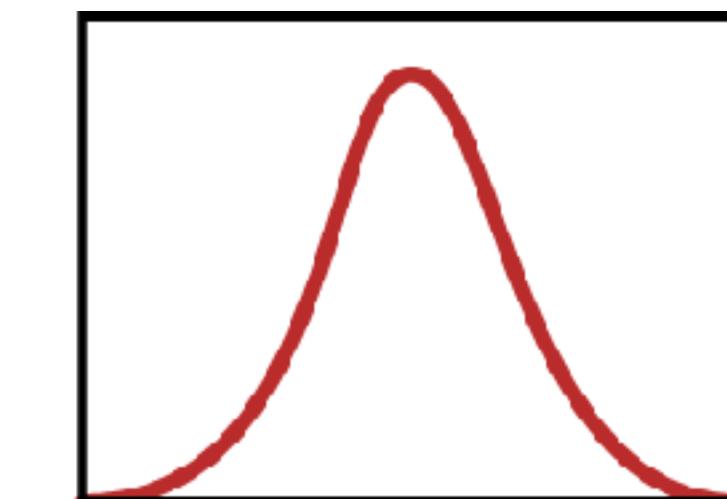
Normal Distribution



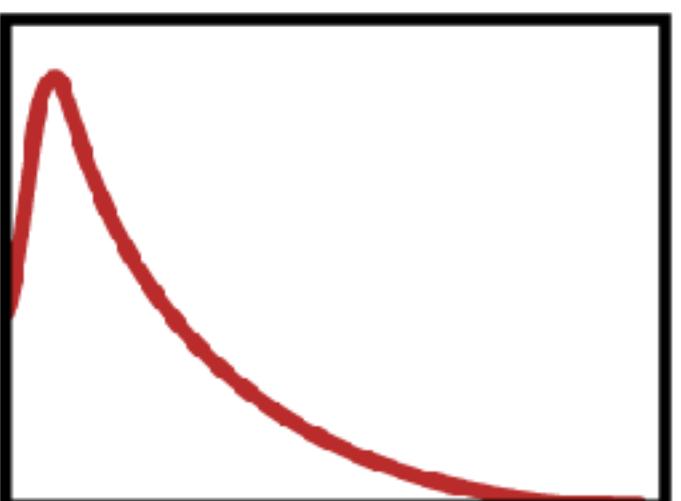
Uniform Distribution



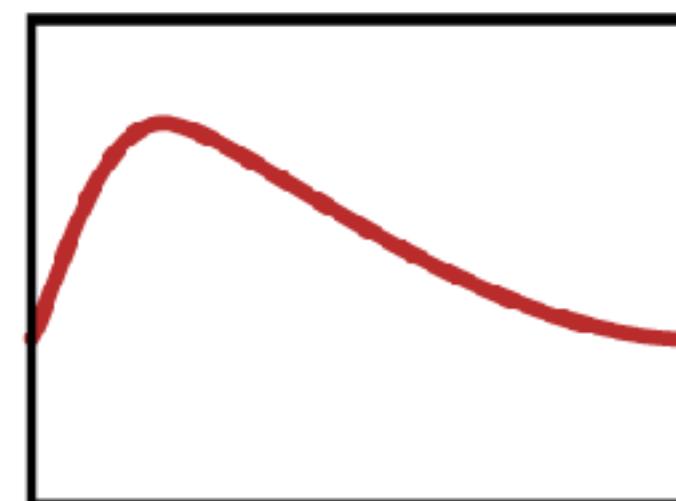
Cauchy Distribution



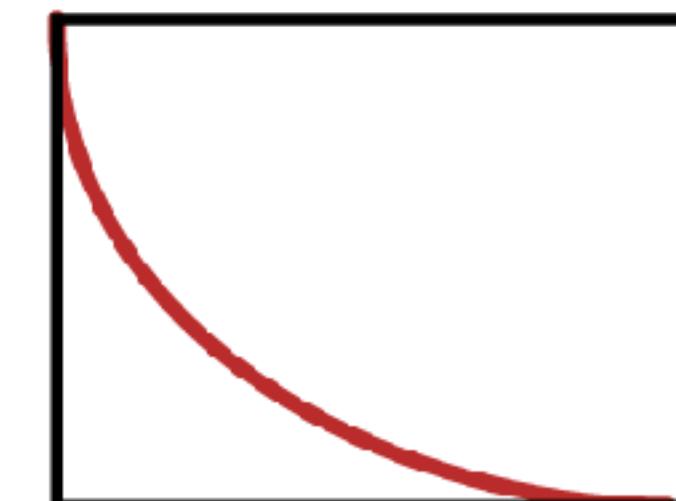
t Distribution



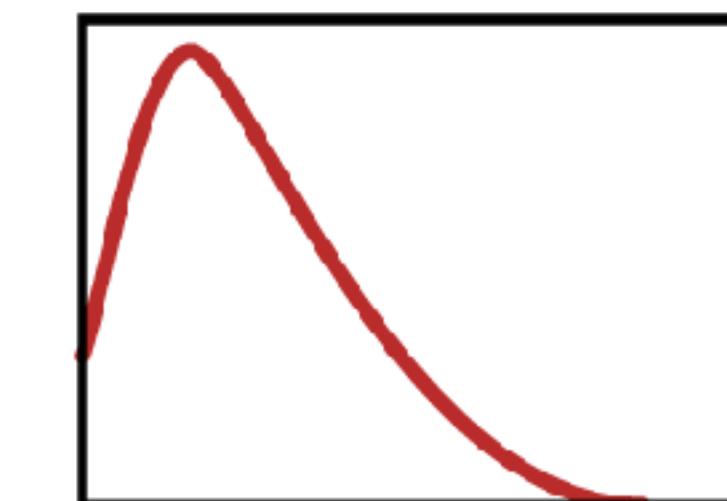
F Distribution



Chi-Square Distribution



Exponential Distribution



Weibull Distribution



Lenguajes y habilidades de programación

- Hay que diferenciar entre las habilidades de programación y el lenguaje/librerías concretos que se utilizan
- Los fundamentos son abstractos y se pueden reutilizar
- Los lenguajes son herramientas y se solapan





Lenguajes y habilidades de programación

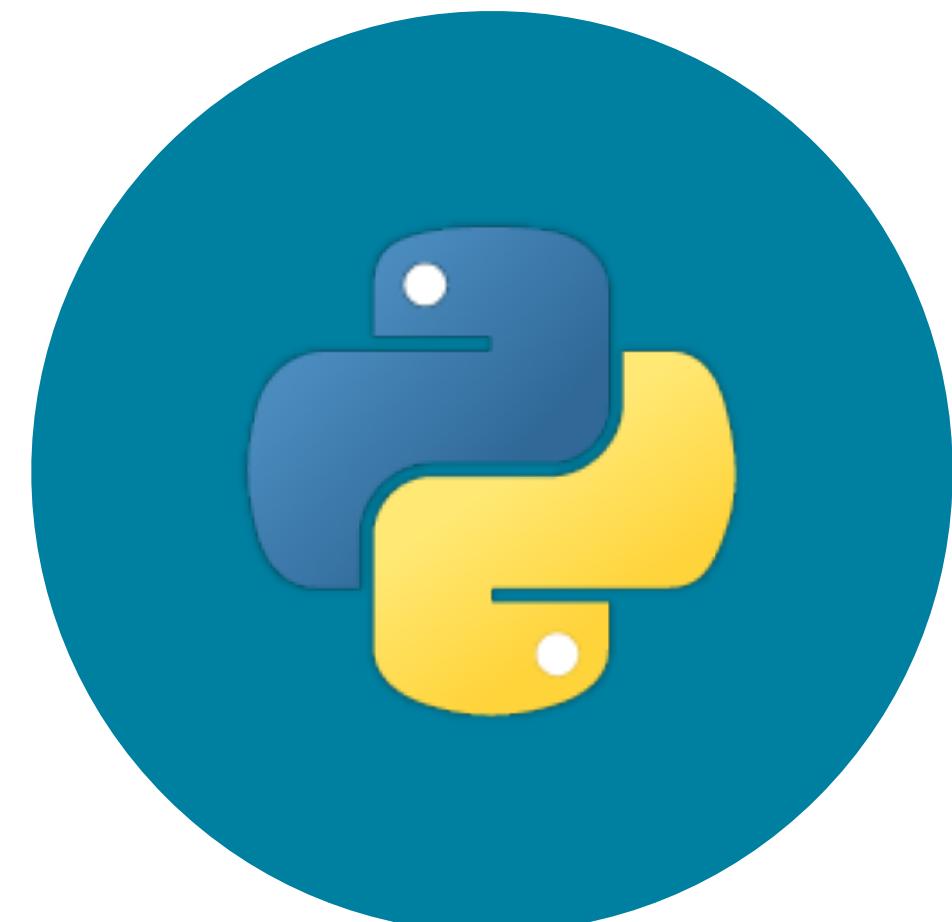
- Hay que diferenciar entre las habilidades de programación y el lenguaje/librerías concretos que se utilizan
- Los fundamentos son abstractos y se pueden reutilizar
- Los lenguajes son herramientas y se solapan





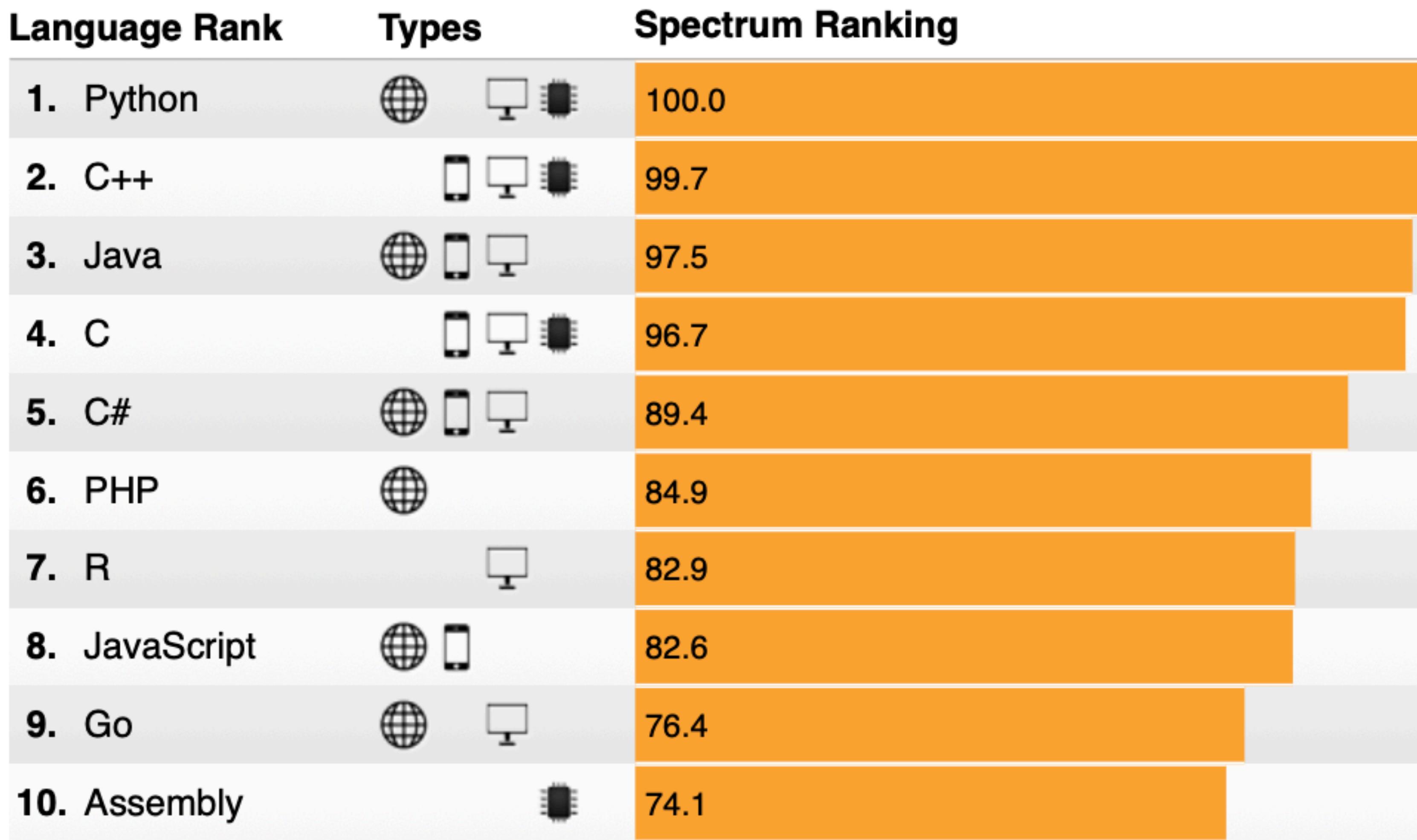
Lenguajes y habilidades de programación

- Hay que diferenciar entre las habilidades de programación y el lenguaje/librerías concretos que se utilizan
- Los fundamentos son abstractos y se pueden reutilizar
- Los lenguajes son herramientas y se solapan





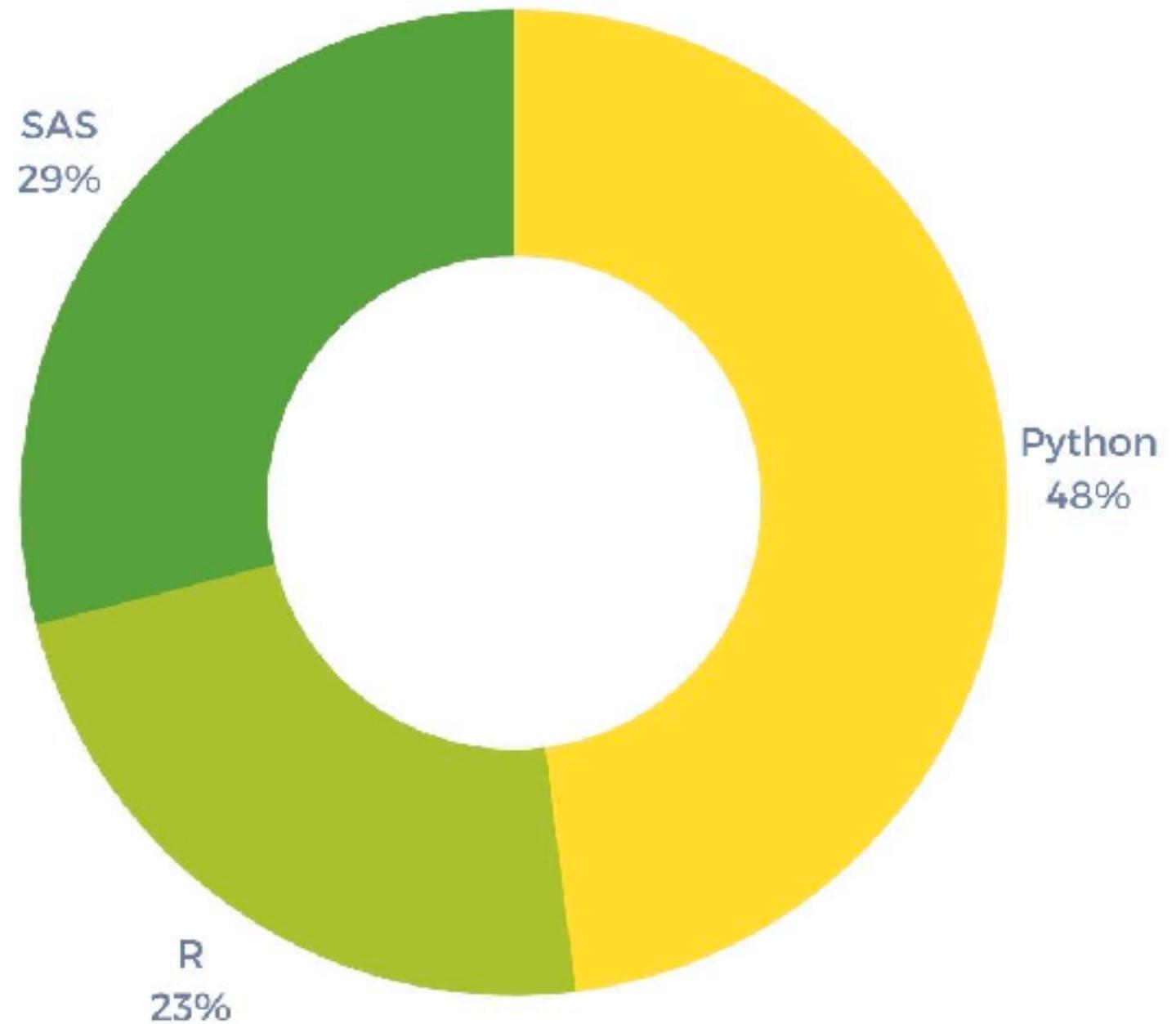
Lenguajes de programación





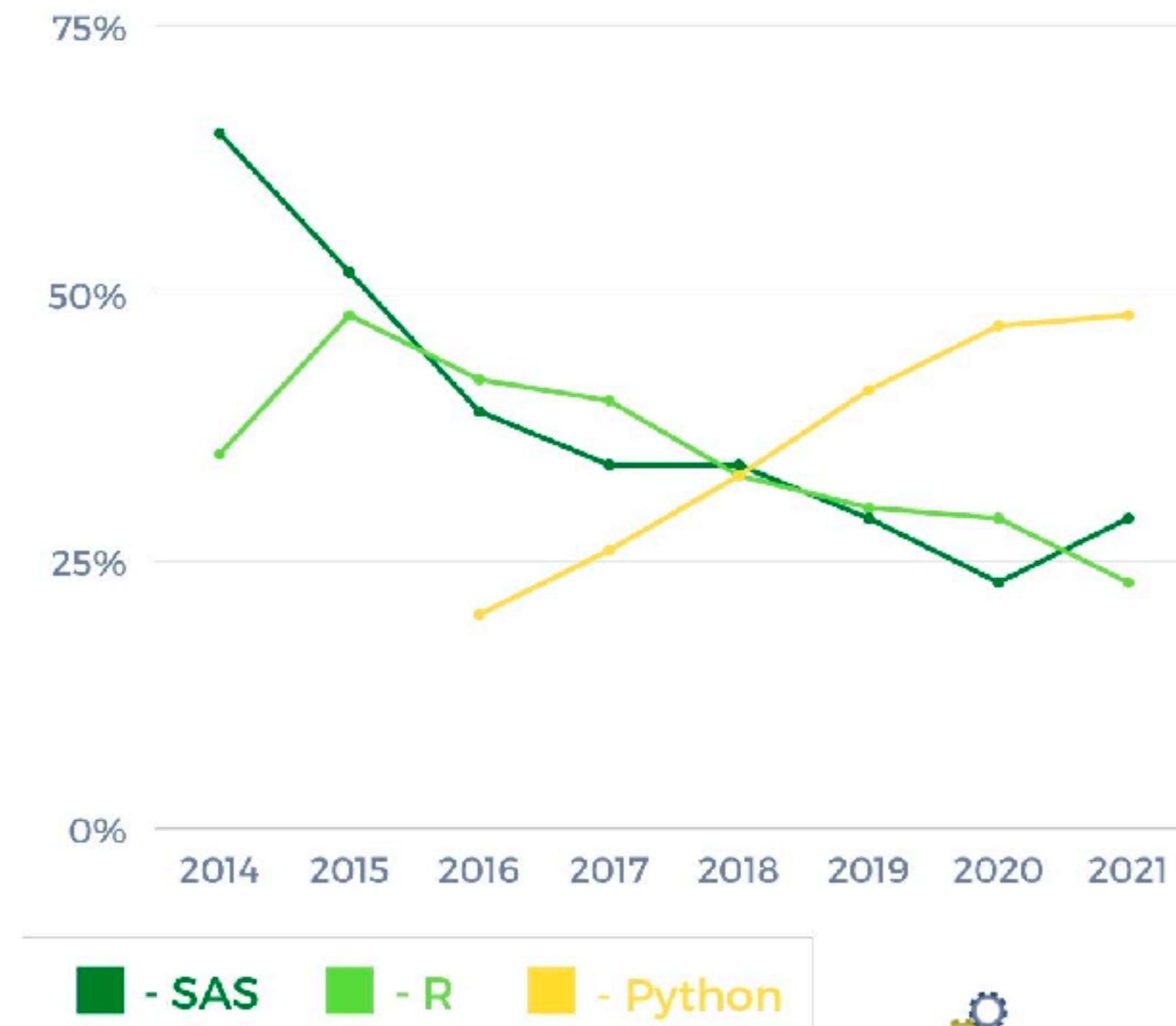
Lenguajes de programación

SAS, R, or Python 2021 Overall Results



Burtsch Works
Executive Recruiting

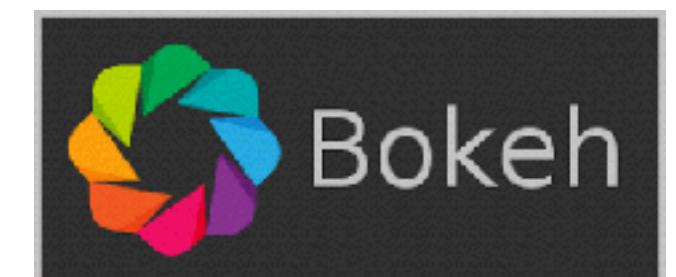
SAS, R, or Python Preference: 8-Year Trend



Burtsch Works
Executive Recruiting



Ecosistema python



Ecosistema Big Data





A hombros de gigantes: Proveedores cloud público



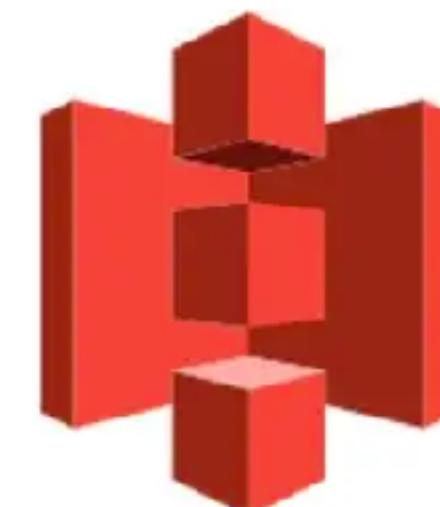
Google Cloud Platform



Microsoft Azure



A hombros de gigantes: Ecosistema AWS



amazon
S3



AWS Lambda



Amazon EC2



Amazon DynamoDB

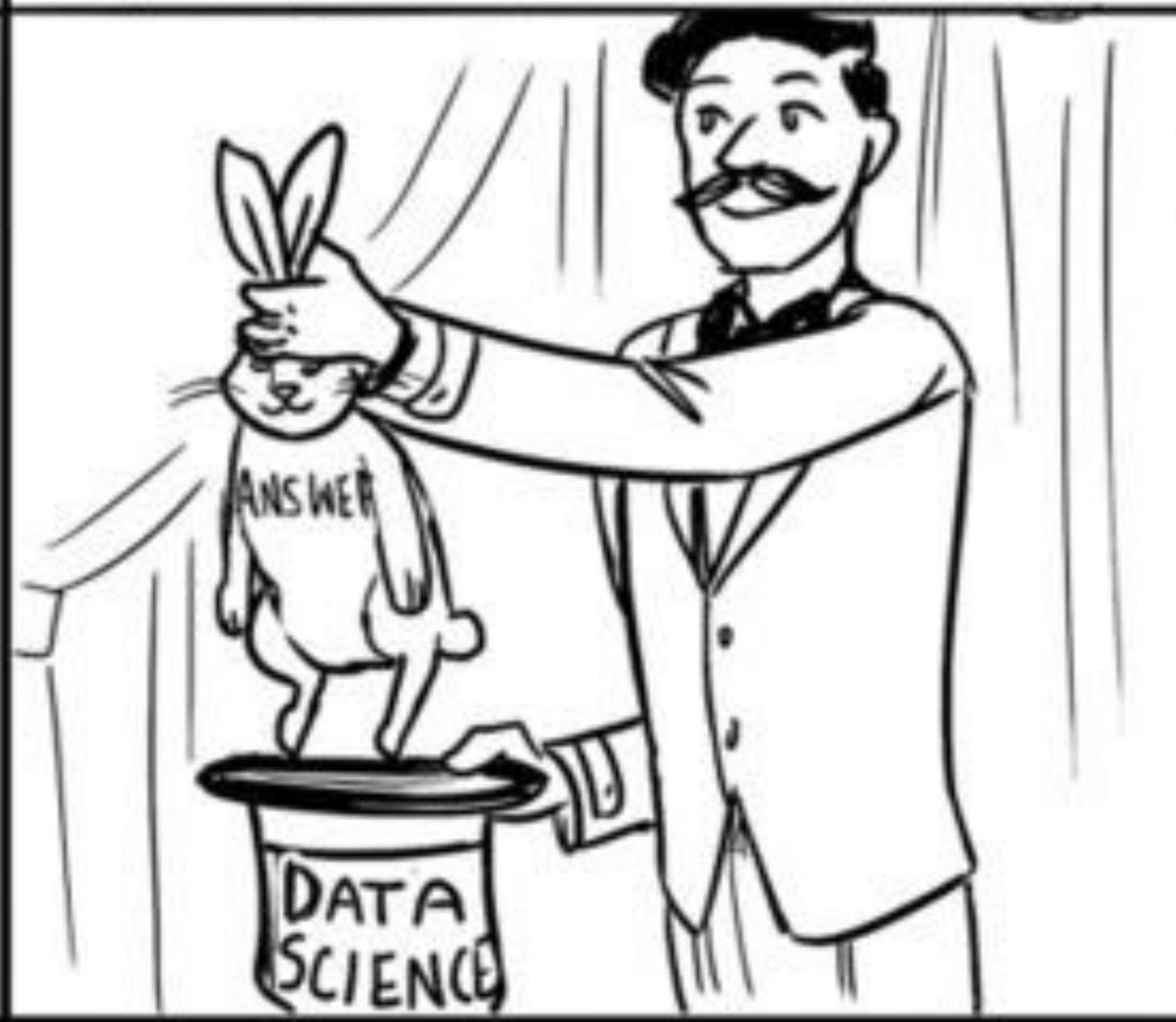


Amazon RDS

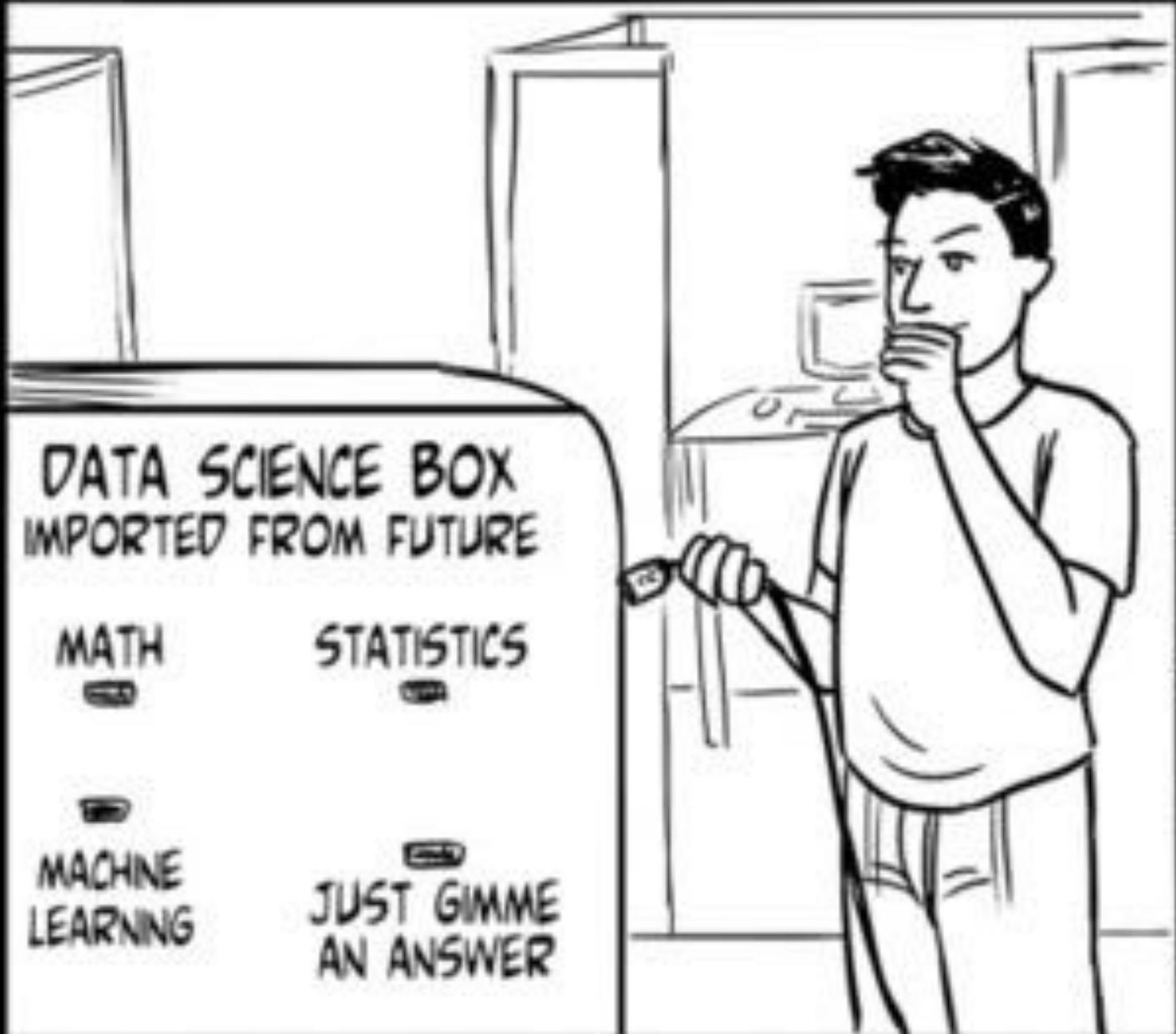
WHAT MY BOSS THINKS DATA SCIENCE IS



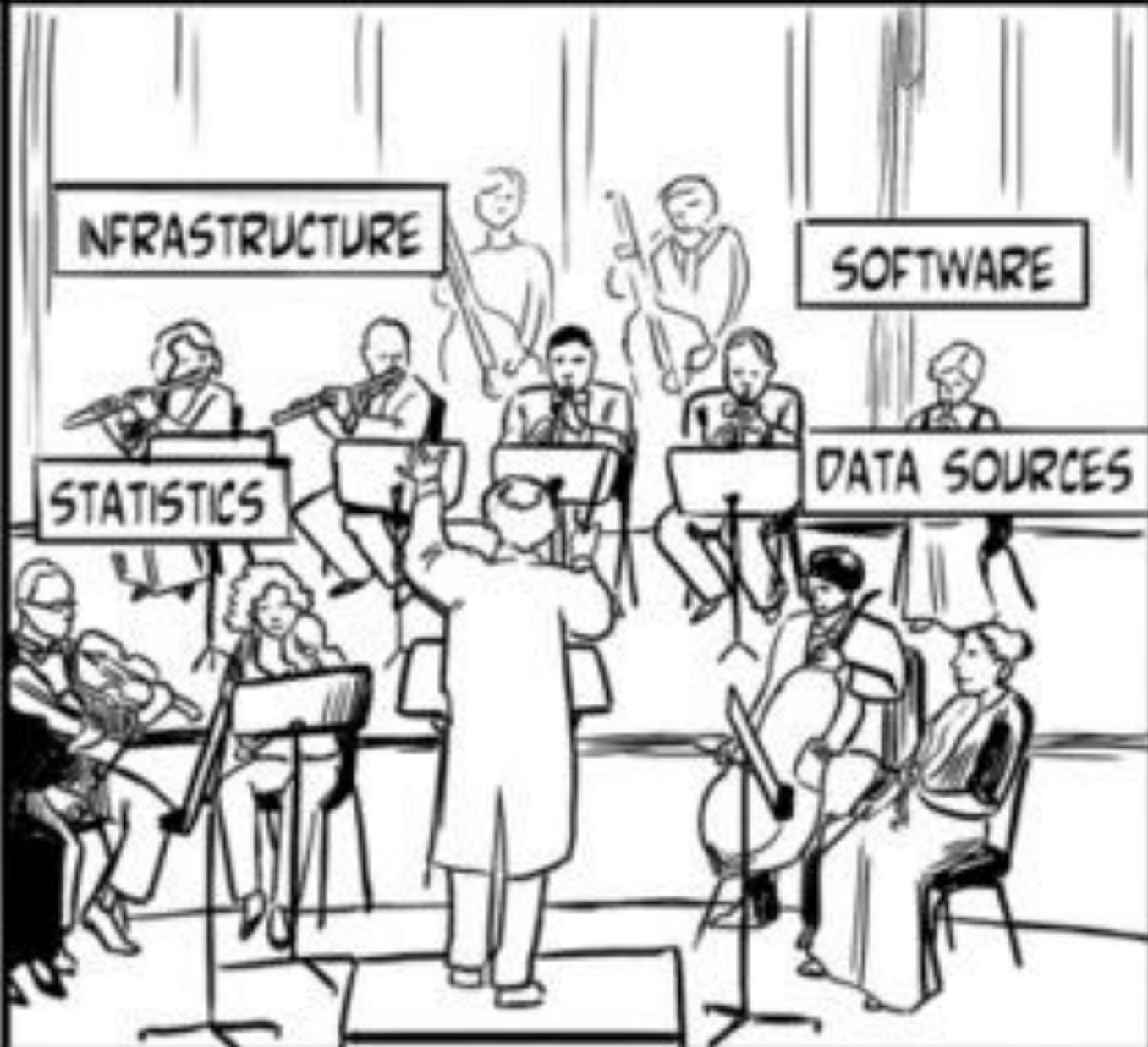
WHAT MY CUSTOMERS THINK DATA SCIENCE IS



WHAT SOFTWARE ENGINEERS THINK DATA SCIENCE IS



WHAT I THINK DATA SCIENCE IS



Módulo 1. Introducción a la ciencia de datos

Introducción a la Ciencia de Datos e Ingeniería de Datos

Jacinto Arias

Máster en Ciencia de Datos e Ingeniería de Datos en la Nube

