

# Data Mining

*Torture the data, and it will confess to anything.*  
— Ronald Coase, economist, Nobel Prize Laureate.



# Definición

- ▶ Data mining es el proceso computacional de **explorar** y **descubrir patrones** en grandes datasets.
- ▶ Es un subcampo de las ciencias de la computación que combina variadas técnicas de la estadística, data science, teoría de bases de datos y aprendizaje automático.

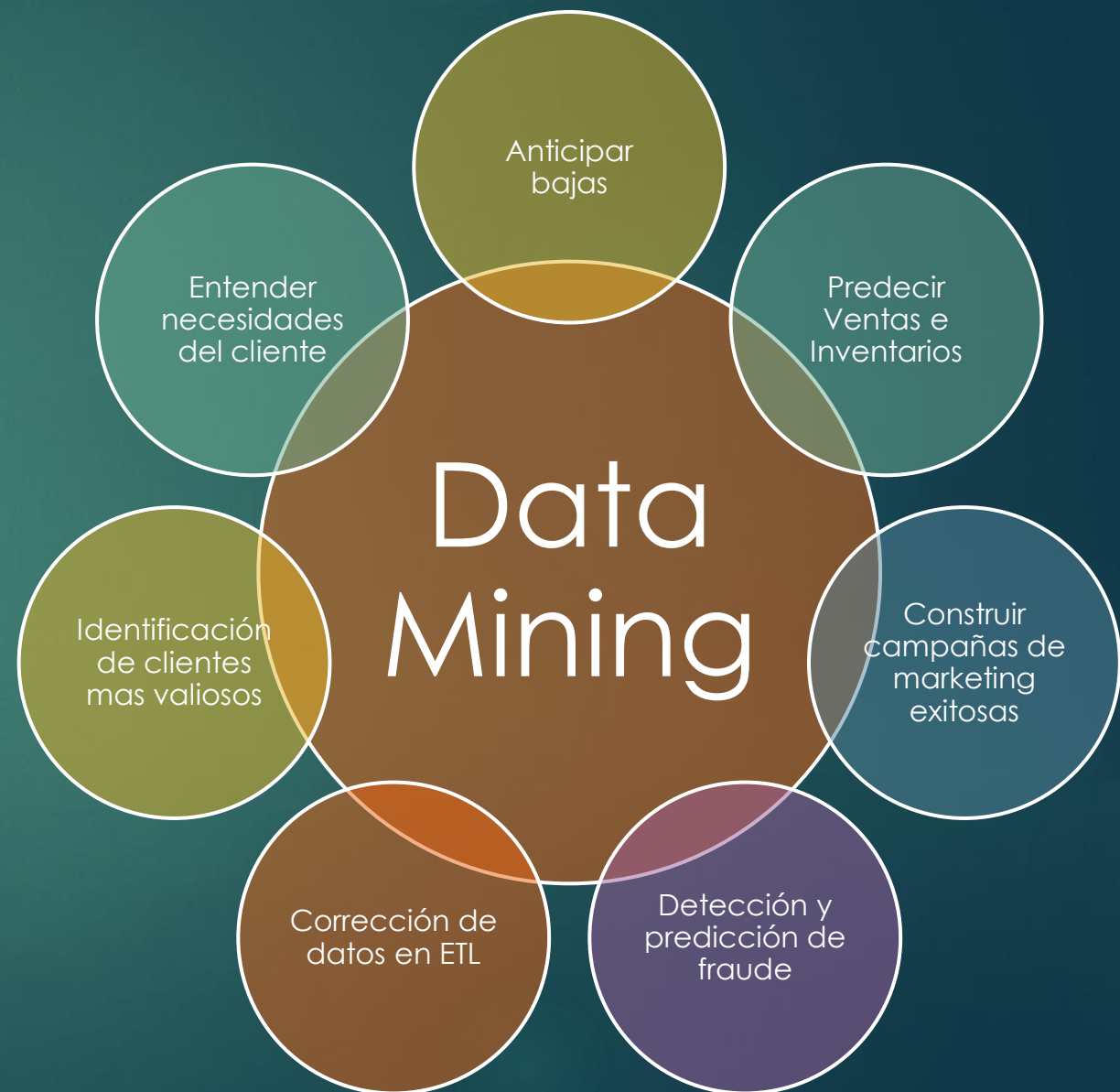
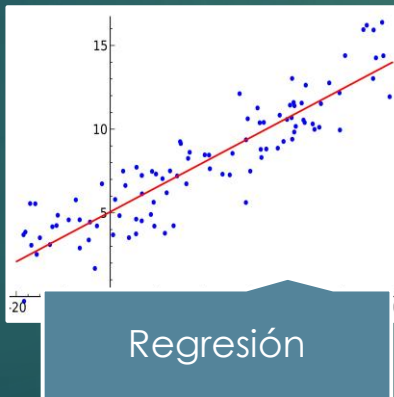
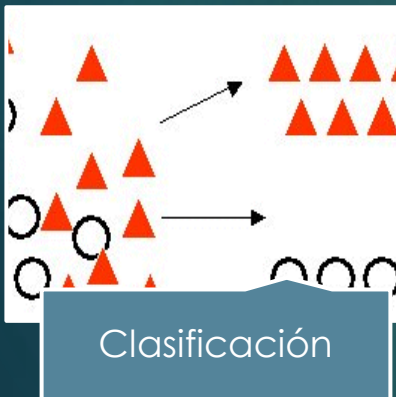
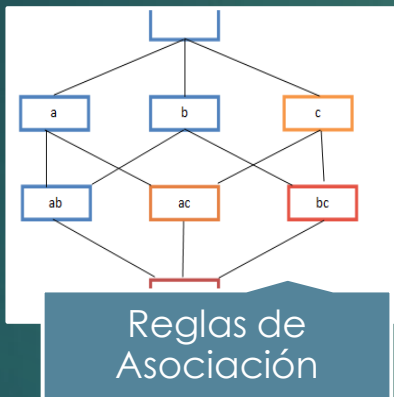
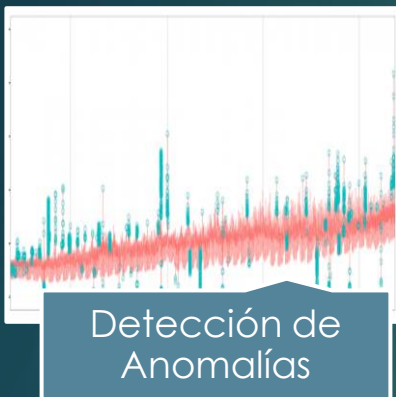


# Historia





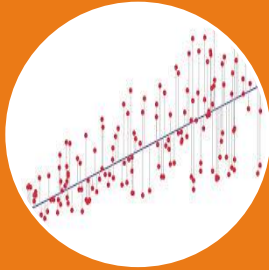
# Usos Habituales



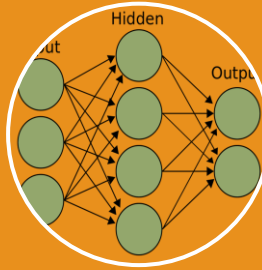
# Algoritmos empleados



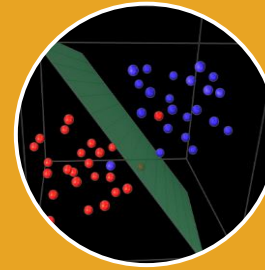
Arboles de  
Decisión



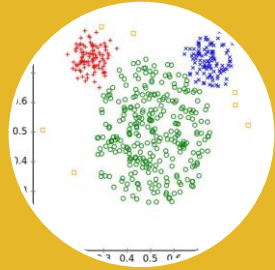
Regresión



Redes  
Neurales



Support  
Vector  
Machine



K-means  
Clustering

**Ensembles / Conjuntos**

Para un mismo modelo se puede usar uno o más de un algoritmo. Esto puede realizarse en secuencia o como una combinación de modelos donde el ensemble (o combinación) puede superar el mejor de los modelos.



# Proceso habitual

## Modelo SEMMA

Nuevos datos son incorporados y el modelo se recalibra.



### SAMPLE

**Muestrear** la información creando un data set lo suficientemente completo para incluir la información relevante.



### EXPLORE

**Explorar** los datos buscando relaciones anticipadas, y tendencias y anomalías no anticipadas - para obtener un entendimiento profundo e ideas para el análisis posterior.



### MODIFY

**Modificar** los valores observados por la creación, selección y transformación de variables. Puede mejorar la precisión del modelo.



### MODEL

**Modelar** la información usando herramientas analíticas para encontrar la combinación de valores que certeramente predice un resultado final.



### ASSESS

**verificar y contrastar** los modelos obtenidos para evaluar la usabilidad y confianza de los descubrimientos del proceso de data mining.

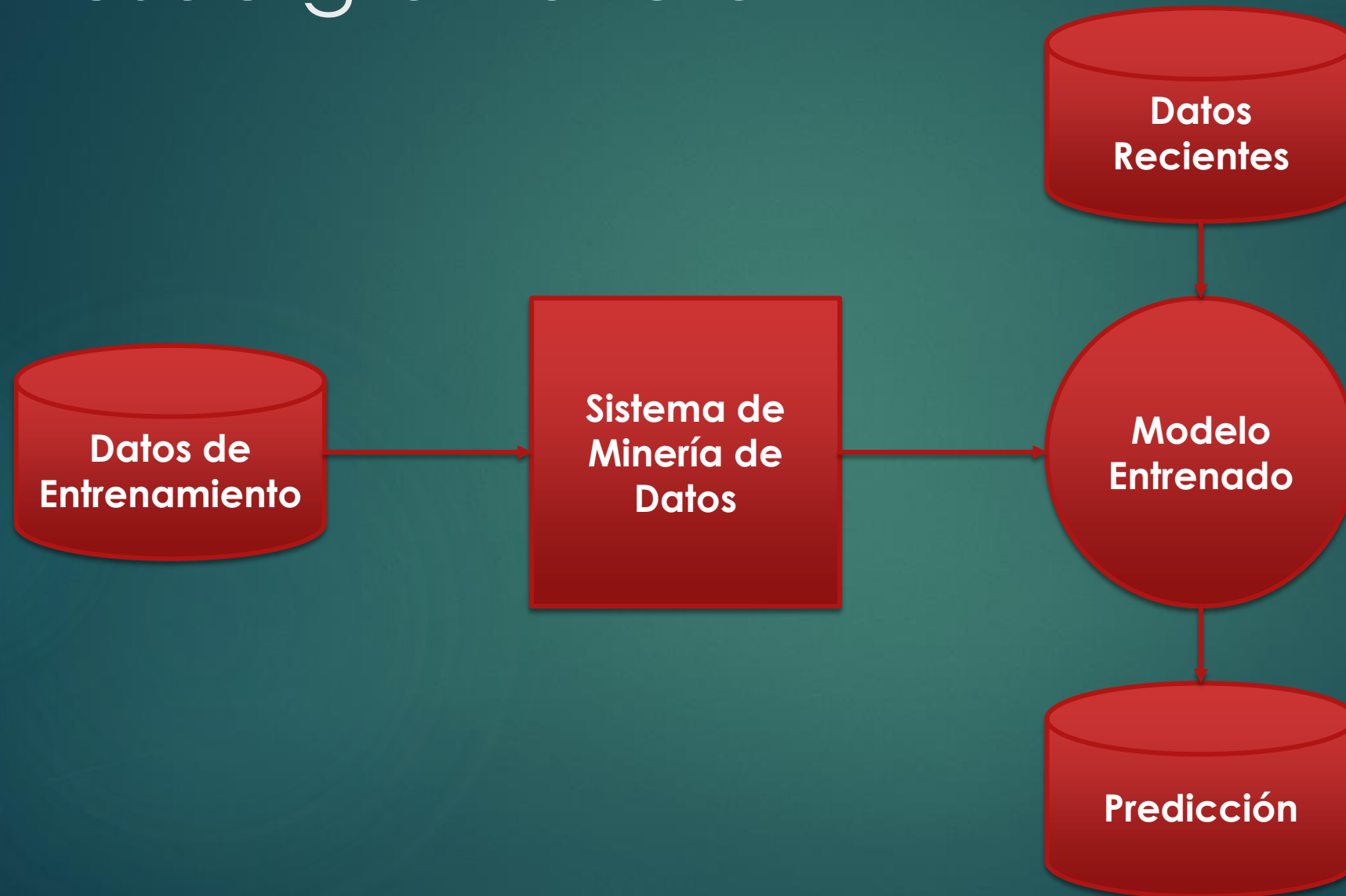
Retroalimentación



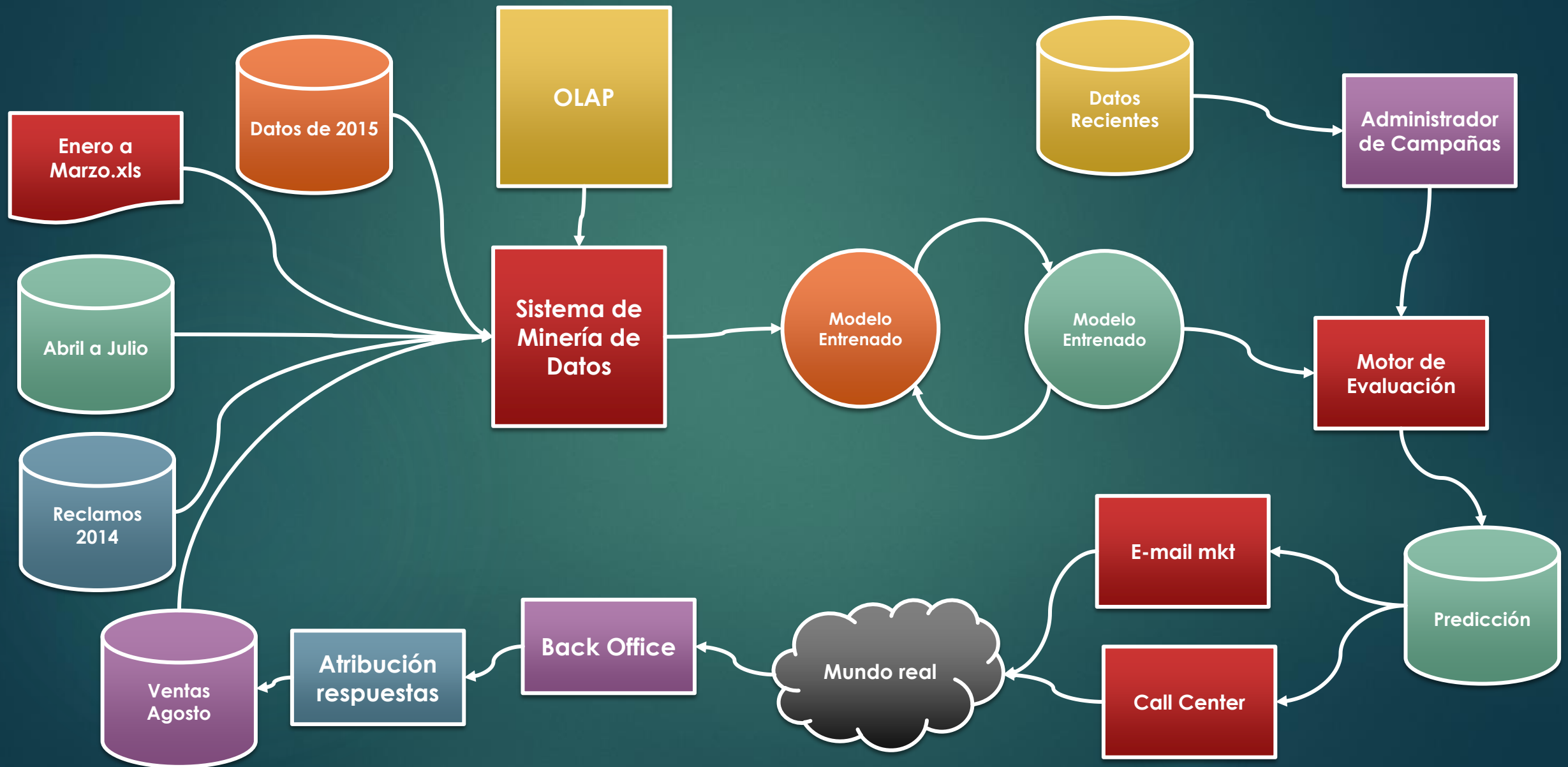
*If we have data, let's look at data. If all we have are  
opinions, let's go with mine  
— Jim Barksdale, former Netscape CEO.*

Cómo se construye y como se usa  
un modelo predictivo?

# Desde gran altura



...y en el mundo real (cuando las cosas son simples)



REG 3-5 (LN)	137.64	-0.38
REG 3-5 (FP)	93.95	0.00
REG 3-5 (LN)	137.60	0.75
REG 7-10 (FP)	138.29	0.00
REG 7-10 (LN)	143.27	0.56
REG 7-10 (GR)	98.35	0.00
REG 7-10 (IM)	144.71	1.05



*In God we trust, all others must bring data.*  
— *W. Edwards Deming, engineer, statistician,*  
*professor, author, lecturer, and management consultant.*

# Como interpretar un modelo

# Introducción – matriz de confusión

		Predicho	
		Positivo	Negativo
Real	Positivo	VP	FN
	Negativo	FP	VN

Un modelo de clasificación binaria produce 4 resultados:

- **Verdaderos Positivos:** Digo que se dan de baja y realmente se dan de baja.
- **Falsos Positivos:** Digo que se dan de baja y no lo hacen.
- **Falso Negativo:** Digo que continúan con la empresa y se dan de baja.
- **Verdadero Negativo:** Digo que continúan con la empresa y lo hacen.

Métricas: (Positivo = Baja;  
Negativo = Continúa)

- Exactitud:
  - Clasificaciones correctas / Casos  $Ex = \frac{VP + VN}{VP + FN + FP + VN}$
- Precisión:
  - Indica proporción de bajas reales sobre bajas predichas  $Pr = \frac{VP}{VP + FP}$
- Sensibilidad:
  - Indica proporción de las bajas reales que el modelo captura  $S = \frac{VP}{VP + FN}$

# Introducción – matriz de confusión

		Predicho	
		Positivo	Negativo
Real	Positivo	VP	FN
	Negativo	FP	VN

## Exactitud

$$Ex = \frac{VP + VN}{VP + FN + FP + VN}$$

VP	FN
FP	VN

---

VP	FN
FP	VN

## Precisión

$$Pr = \frac{VP}{VP + FP}$$

VP	FN
FP	VN

---

VP	FN
FP	VN

## Sensibilidad

$$S = \frac{VP}{VP + FN}$$

VP	FN
FP	VN

---

VP	FN
FP	VN

# Introducción – comparación de modelos

...pero...

...de acuerdo al problema me puede interesar una de las métricas más que las otras. Puedo ajustar el modelo para que se ajuste a este criterio.

**El modelo me permite establecer el valor de corte (threshold).**

**Este será el grado de seguridad que le exijo antes de catalogar un caso como “positivo”.**

Esto genera una matriz de confusión para cada valor de corte posible.

## Salud:

Prefiero falsas alarmas a una lesión o muerte por falta de previsión.

**Sensibilidad**

## Justicia:

...inocente hasta que se demuestre lo contrario.

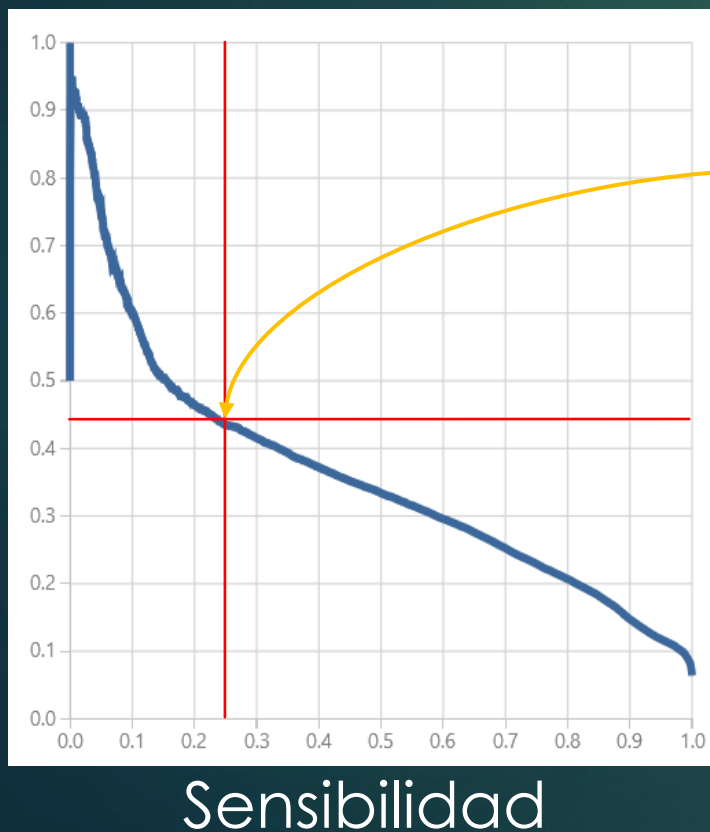
**Precisión**



# Resultados del modelo predictivo

Se presentan algunas maneras de visualizar la performance de un modelo predictivo.

## Precisión - Sensibilidad



- Precisión: Proporción de bajas reales sobre bajas predichas.
- Sensibilidad: Proporción de bajas predichas sobre total de bajas.

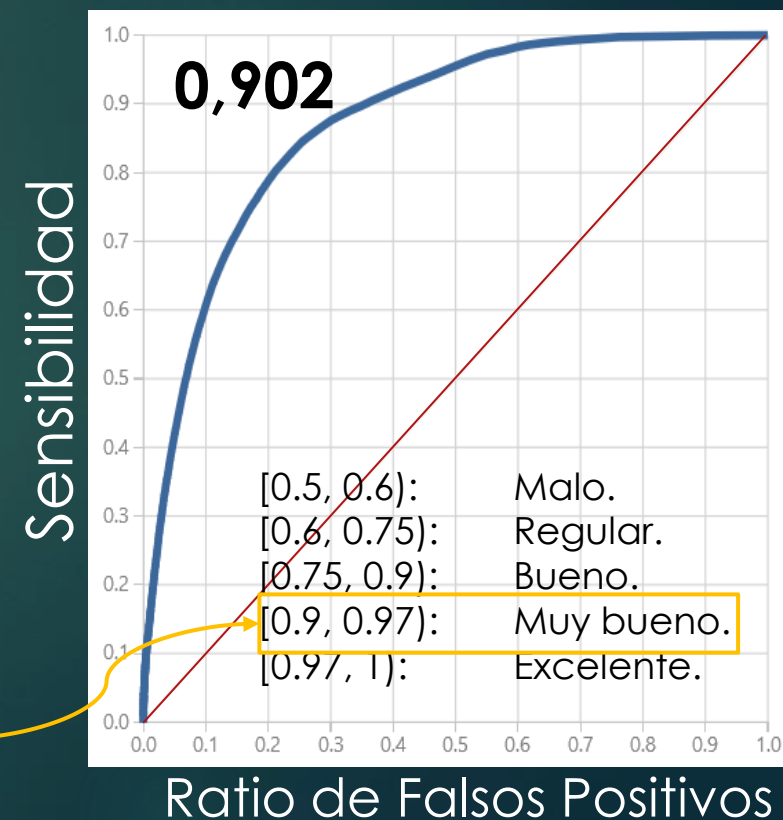
**Puedo capturar el 25% de las bajas acertando en el 45% de mis predicciones. Resultado 8 veces superior al azar.**

La curva ROC permite evaluar modelos predictivos.

El valor obtenido por el modelo es 0,902

**Se puede considerar como un modelo Muy Bueno**

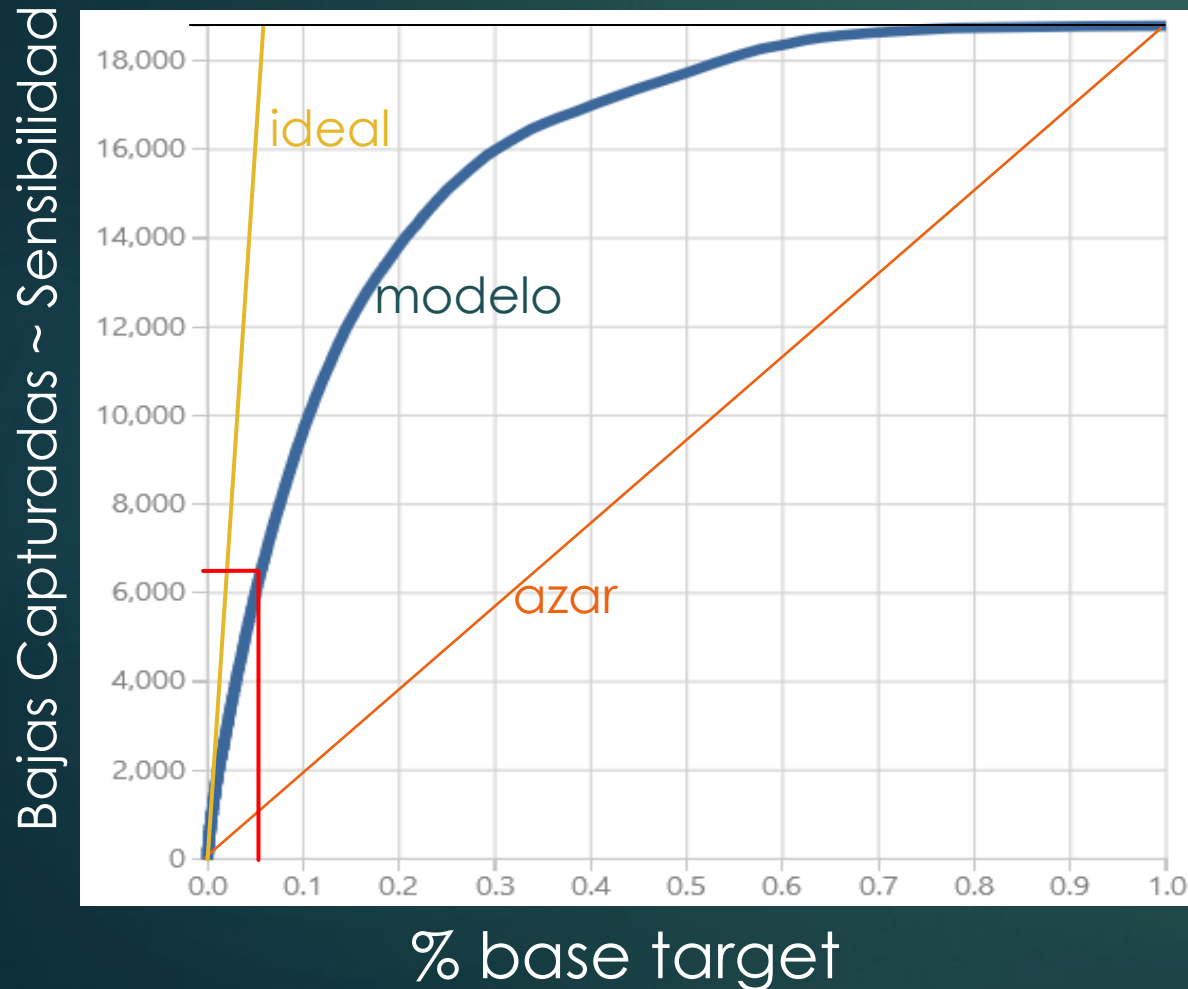
## Curva ROC



# Resultados del modelo predictivo

Otra manera de ver la performance es la Curva Lift

Curva Lift



La curva Lift muestra el número de bajas reales capturadas de acuerdo al número de clientes contactados en una acción de retención.

En este caso observamos que con este modelo contactando al **5%** de la base (15.000 de 300.000 clientes) se obtiene un **32%** de las bajas reales (6.100 de 19.200 bajas del período).

**Para conseguir esa cantidad de bajas en una muestra aleatoria hubiera sido necesario contactar 95.300 clientes.**

El modelo define una curva, pero somos nosotros los que debemos definir un punto específico en la curva.

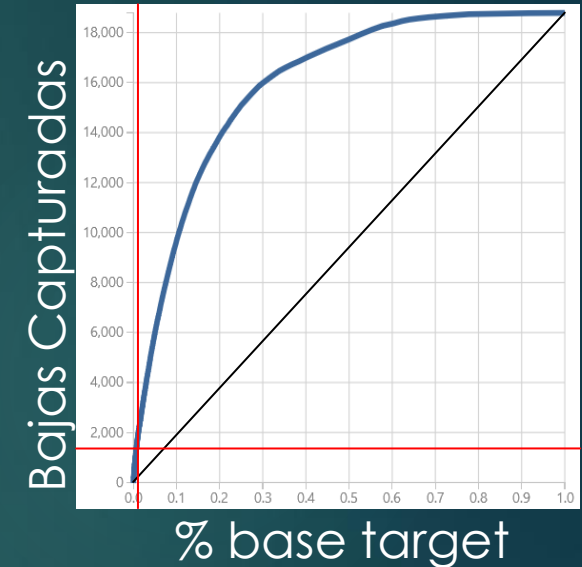
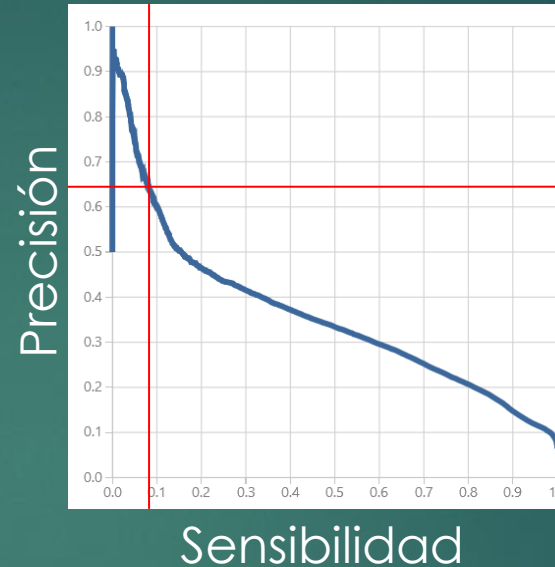
# Resultados con distintos valores de corte

## Modelo exigente

## Precisión vs azar = 10

Ratio bajas = 6,4%

		Predicho	
		Positivo	Negativo
Real	Positivo	1586	17593
	Negativo	906	279915



Apunto al 0,8% de mi base y obtengo el 8,3% de las bajas

Exactitud	Precisión	Sensibilidad
94%	64% 10 veces superior al azar	8,3%

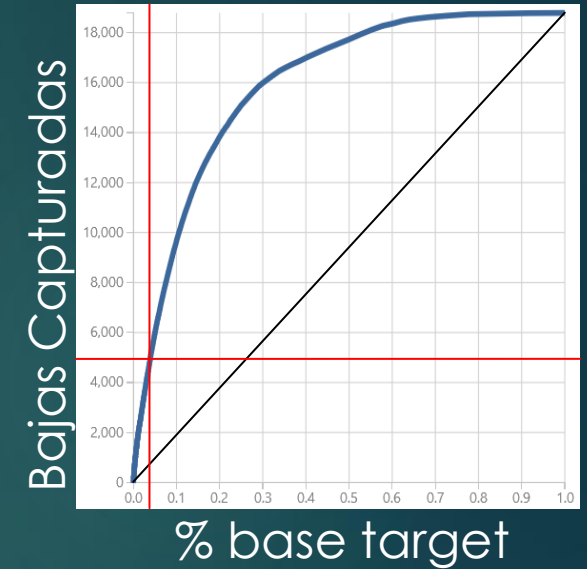
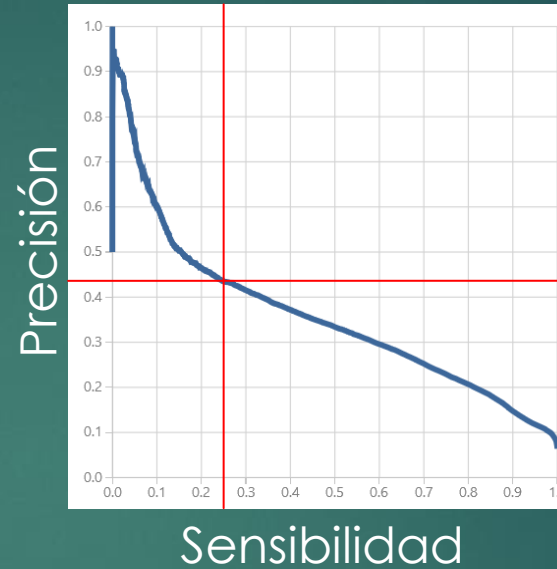
# Resultados con distintos valores de corte

## Modelo medio

## Precisión vs azar = 7

Ratio bajas = 6,4%

		Predicho	
		Positivo	Negativo
Real	Positivo	4845	14334
	Negativo	6310	274511



Apunto al 3,7% de mi base y obtengo el 25,3% de las bajas

Exactitud	Precisión	Sensibilidad
93%	43% 7 veces superior al azar	25%



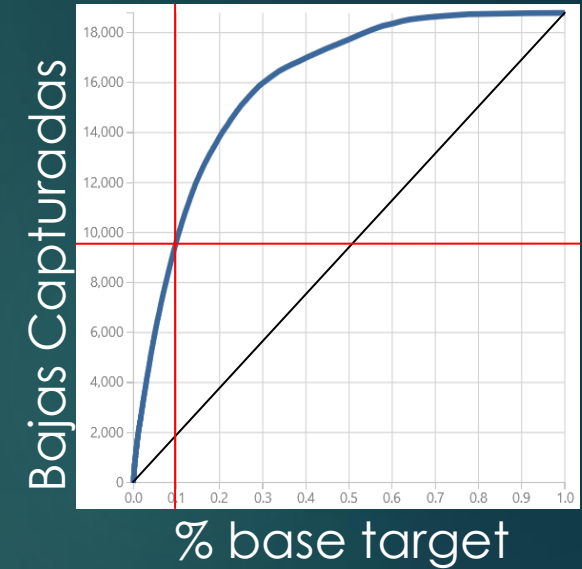
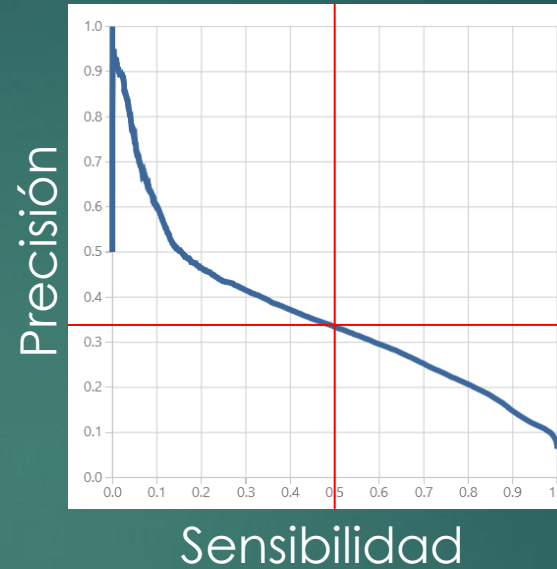
# Resultados con distintos valores de corte

## Modelo flexible

## Precisión vs azar = 5

Ratio bajas = 6,4%

		Predicho	
		Positivo	Negativo
Real	Positivo	9525	9654
	Negativo	18871	261950



Apunto al 9,5% de mi base y obtengo el 50% de las bajas

Exactitud	Precisión	Sensibilidad
90%	34% 5 veces superior al azar	49,7%

# Ejercicio Económico Supuestos

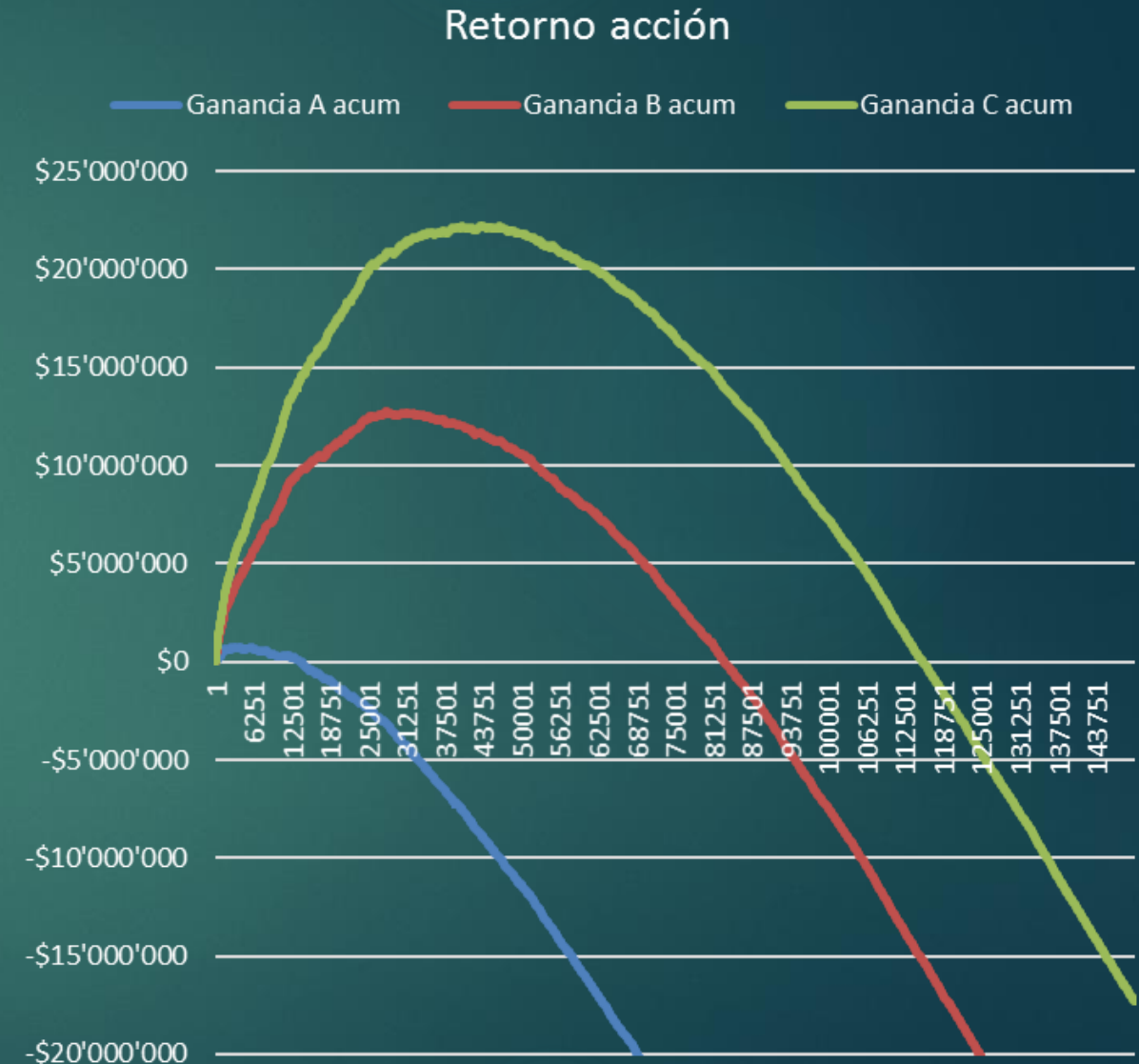
Inversión por Cliente	% Retención	% Gasto Bajas no retenidas	Ingreso Acierto	Cientes contactados	Bajas reales	Bajas evitadas	Inversión	Ingresos	Ganancia
<ul style="list-style-type: none"><li>•Costo del descuento o regalo a cada cliente incluído en la acción</li></ul>	<ul style="list-style-type: none"><li>•Porcentaje de las bajas que permanecen en la empresa por el incentivo de la acción de retención</li><li>•Se podrá obtener con muestras pequeñas sujetas a diferentes incentivos.</li></ul>	<ul style="list-style-type: none"><li>•Si el incentivo se distribuye a lo largo de varios meses, los clientes que no se logró retener no recibirán el total del incentivo, y por lo tanto el costo será menor</li></ul>	<ul style="list-style-type: none"><li>•Ingresos adicionales para la empresa por cada cliente retenido</li></ul>	<ul style="list-style-type: none"><li>•Cantidad de clientes a los que se ofrece el incentivo</li></ul>	<ul style="list-style-type: none"><li>•Bajas que hubieran abandonado la empresa dentro del grupo de Clientes contactados.</li></ul>	<ul style="list-style-type: none"><li>•Bajas que se evitaron por la acción</li></ul>	<ul style="list-style-type: none"><li>•Monto que se distribuyó a todos los clientes a los que se les ofreció el incentivo</li></ul>	<ul style="list-style-type: none"><li>•Monto que ingresa a la empresa gracias a las Bajas evitadas</li></ul>	<ul style="list-style-type: none"><li>•Diferencia entre los ingresos y el costo de los incentivos (Ingresos – Inversión)</li></ul>

# Ejercicio Económico

	Escenario A	Escenario B	Escenario C
<b>Inversión por cliente</b>	\$600	\$600	\$600
<b>% Retención</b>	10%	20%	25%
<b>% Gasto Bajas no retenidas</b>	100%	66%	33%
<b>Ingreso Acierto</b>	\$ 15'000	\$ 15'000	\$ 15'000
<b>Clientes contactados</b>	3791	27933	43241
<b>Bajas reales</b>	2164	9434	12212
<b>Bajas evitadas</b>	216	1887	3053
<b>Inversión total</b>	\$ 2'475'600	\$ 15'542'724	\$ 23'532'708
<b>Ingreso</b>	\$ 3'246'000	\$ 28'302'000	\$ 45'795'000
<b>Ganancia</b>	\$ 770'400	\$ 12'759'276	\$ 22'262'292

Suponiendo que el descuento implica un costo de \$600, que solo el 20% de las personas a las que se les ofrece el descuento deciden permanecer y que cada uno de estos genera luego \$15.000; se evitarían 1887 bajas con un resultado neto de \$12.000.000

Es posible (y recomendable) iniciar con muestras más pequeñas para validar los supuestos.





# Google

Google Search

I'm Feeling L

**DATA: SEARCHING THE SEARCHES**



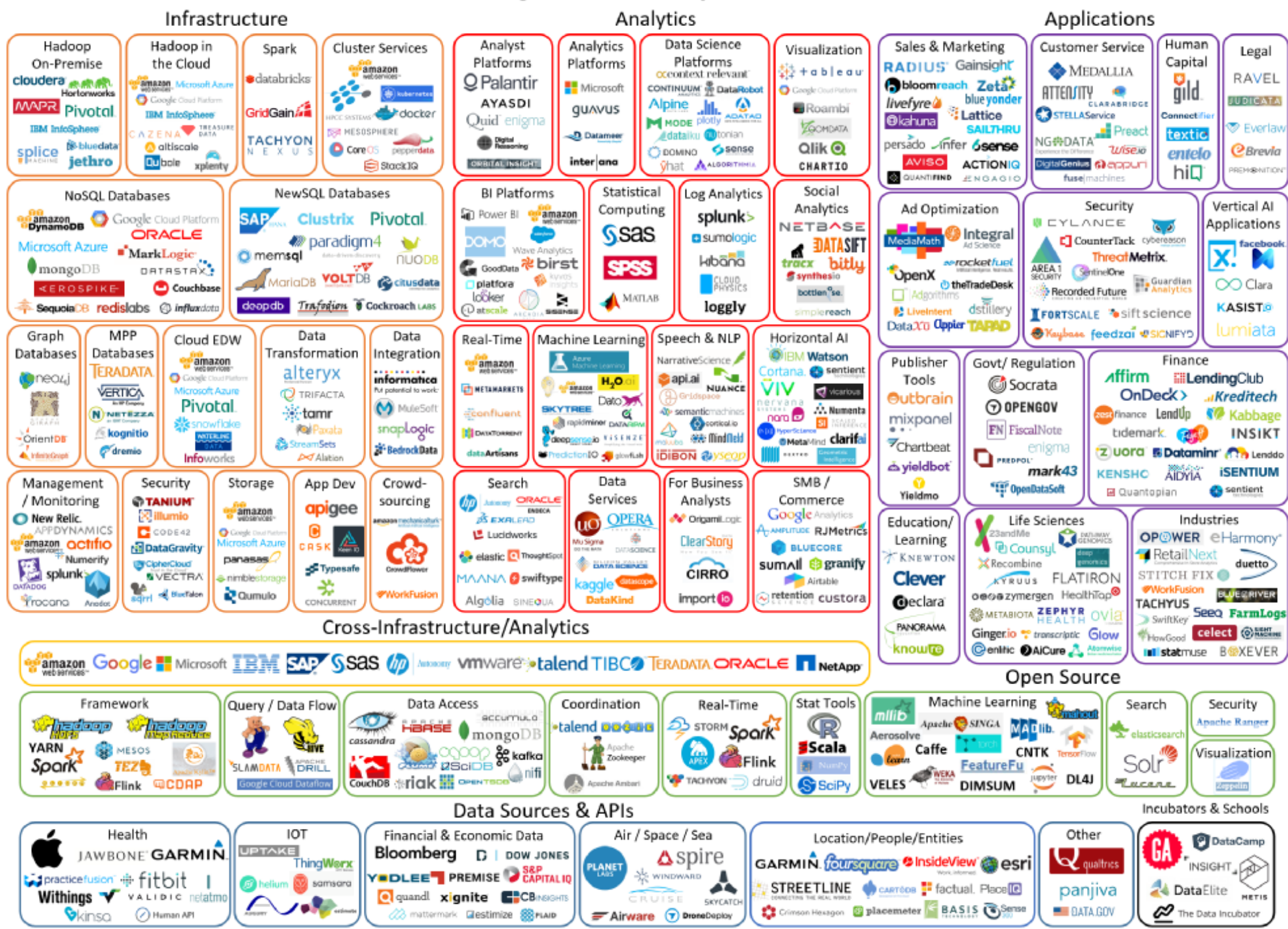
# Herramientas Disponibles

*“Data!data!data!” he cried impatiently. “I can’t make bricks without clay.”*

*— Arthur Conan Doyle, The Adventure of the Copper Beeches (Sherlock Holmes)*

# Landscape Completo

## Big Data



# Algunos productos destacados



R

- [www.r-project.org](http://www.r-project.org)



Python

- [www.python.org](http://www.python.org)
- [www.continuum.io](http://www.continuum.io)



Weka

- [www.cs.waikato.ac.nz/ml/weka/](http://www.cs.waikato.ac.nz/ml/weka/)



KNIME

- [www.knime.org](http://www.knime.org)



RapidMiner

- [rapidminer.com](http://rapidminer.com)

Open  
Source

Pagas



SAS

- [www.sas.com](http://www.sas.com)



SPSS

- [www.ibm.com/analytics/](http://www.ibm.com/analytics/)



Google Cloud

- [cloud.google.com/ml/](http://cloud.google.com/ml/)



Azure ML

- [studio.azureml.net](http://studio.azureml.net)



Amazon

- [aws.amazon.com/es/machine-learning/](http://aws.amazon.com/es/machine-learning/)

Cloud

Muchas gracias