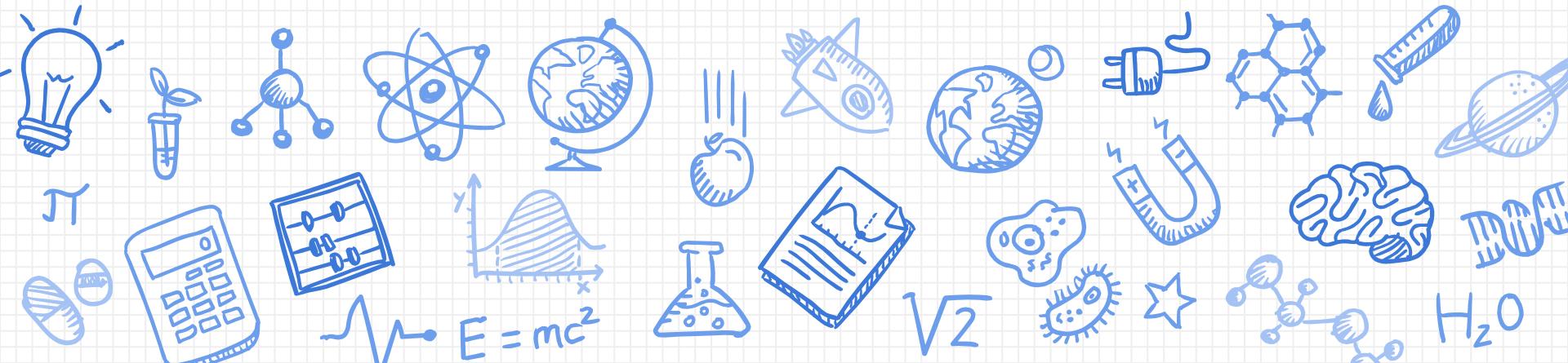


Data Science en la Industria



Soy Cristhian Boujon

- Ingeniero en Sistemas de Información.
- Sr. Software Engineer en ecommerce.
- Profesor.



Desarrollo web hace 10 años atrás...

HTTP

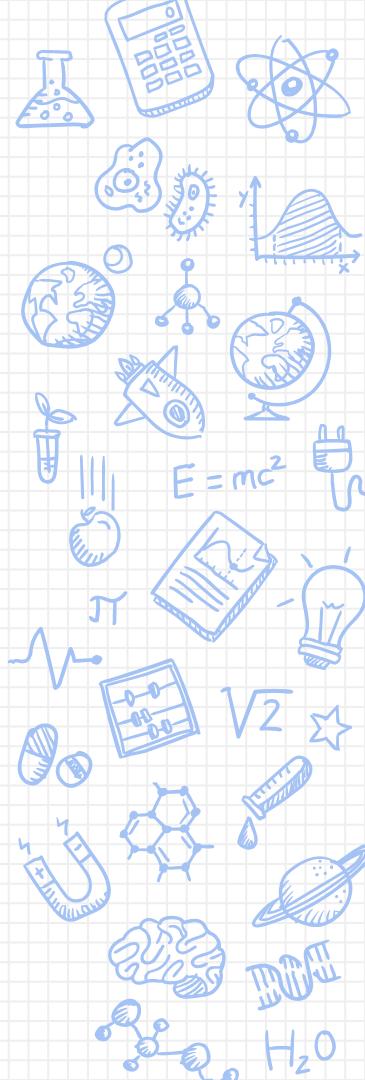
HTML

CSS

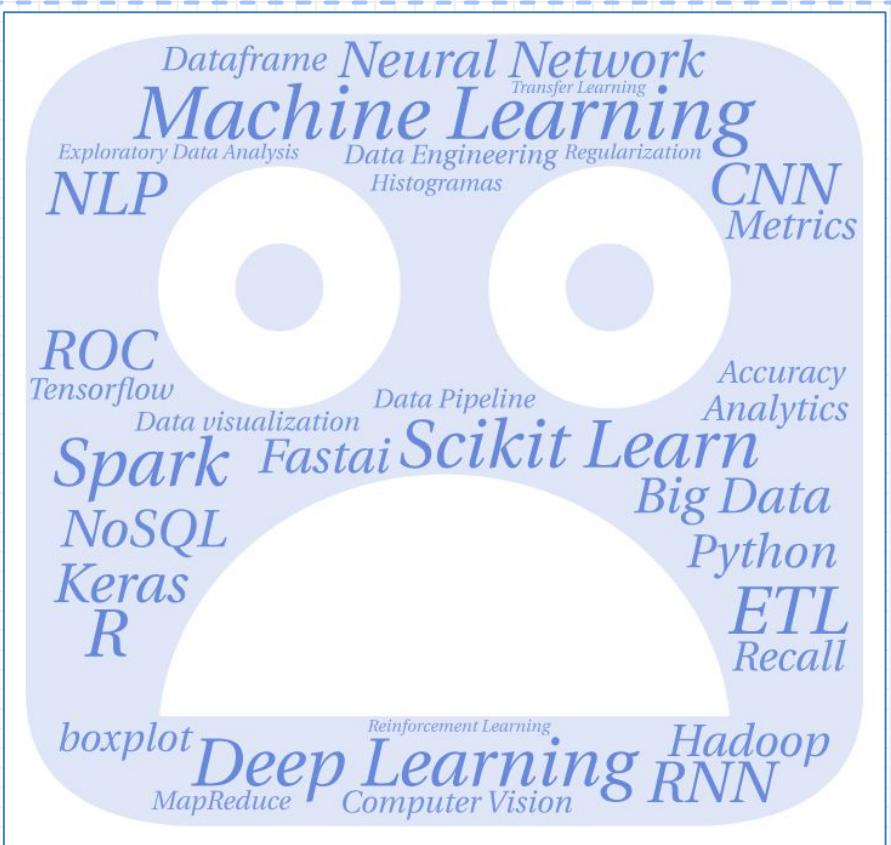
Javascript

PHP

MySQL

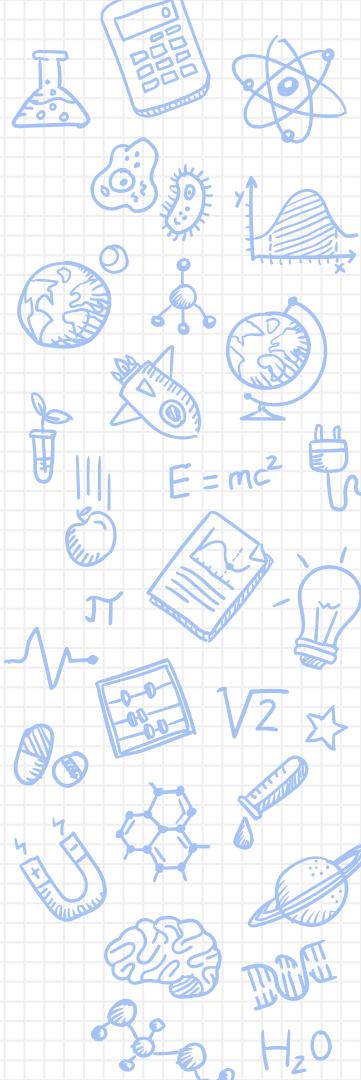


Data Science



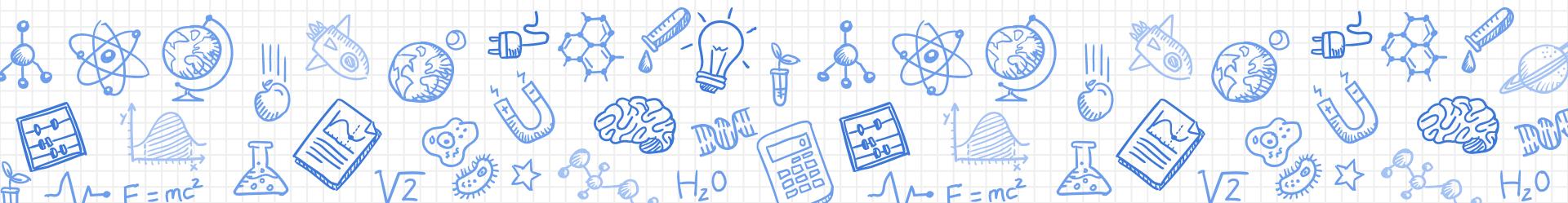
Data Science en la Industria

- Análisis de datos / Data Analysis
 - Modelado / Modeling
 - Ingeniería de datos / Data Engineering



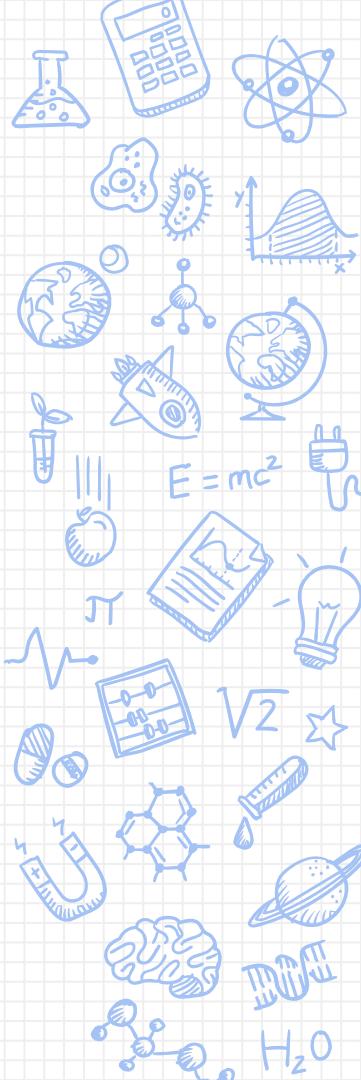
Análisis de Datos

Se examinan los datos en bruto y así sacar conclusiones sobre la información.



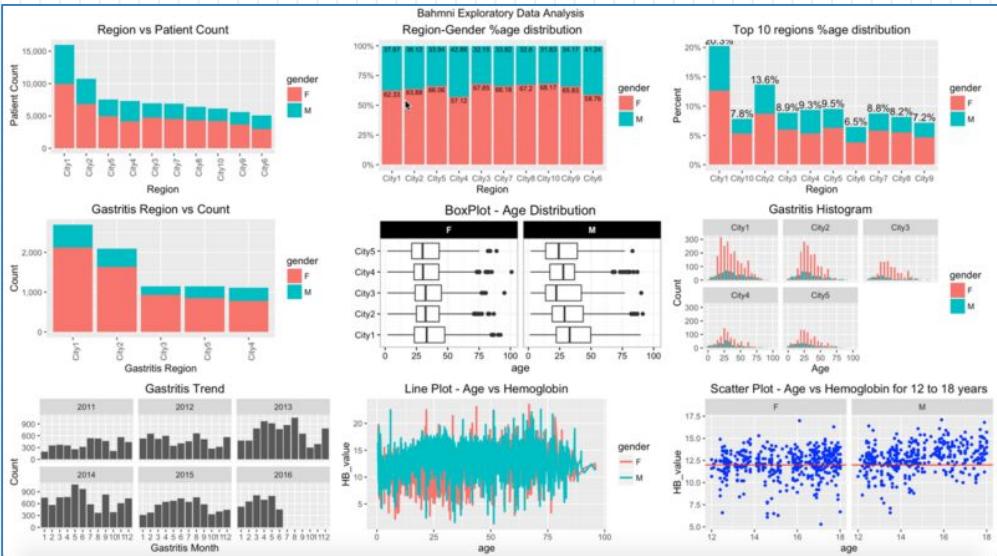
Análisis de Datos

- Conocimientos de Negocio.
 - Conocimientos en estadística

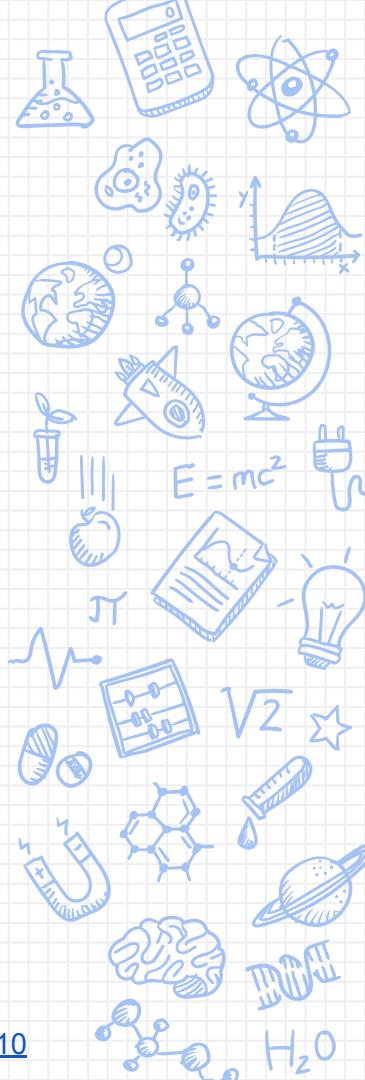


Análisis de Datos – Exploración

Consiste en entender los datos con los que contamos.



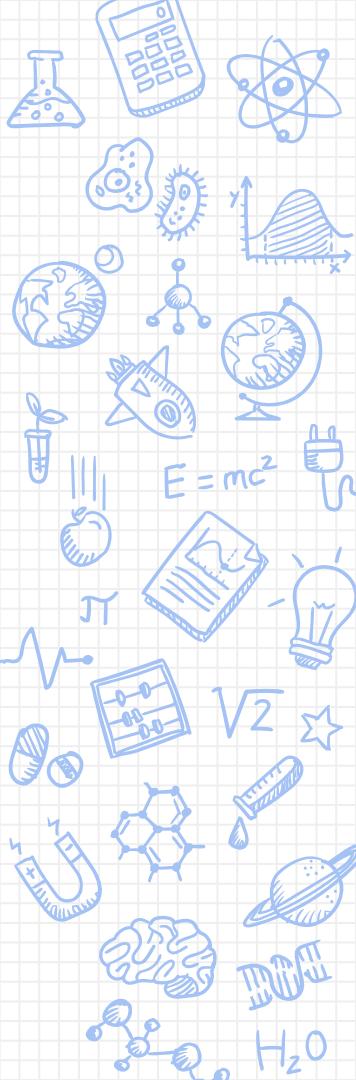
Fuente: <https://medium.com/bahmni-blog/introduction-to-exploratory-data-analysis-of-bahmni-using-r-6c186fd6f010>



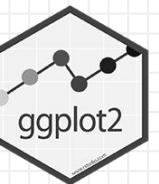
Análisis de Datos – Interpretación

Capturar el significado real de los datos en el negocio.

- Reportes
 - Insights
 - Visualización de datos
 - Rutina en EEUU
 - Marcas



Tools

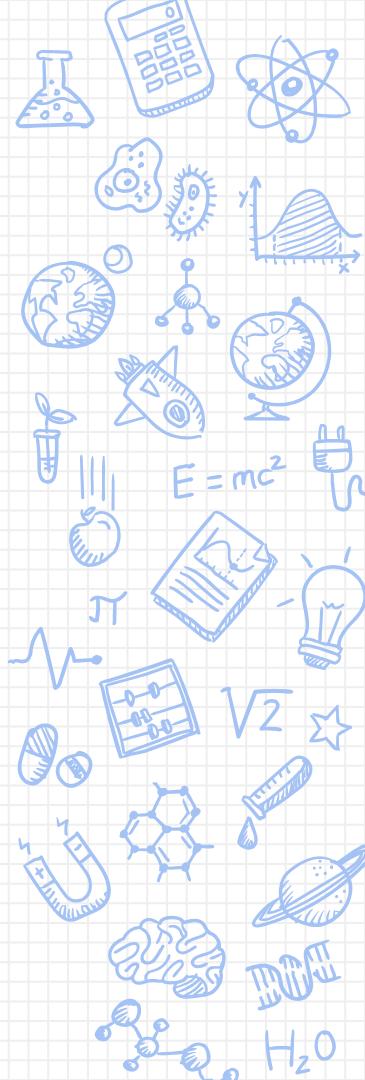


matplotlib



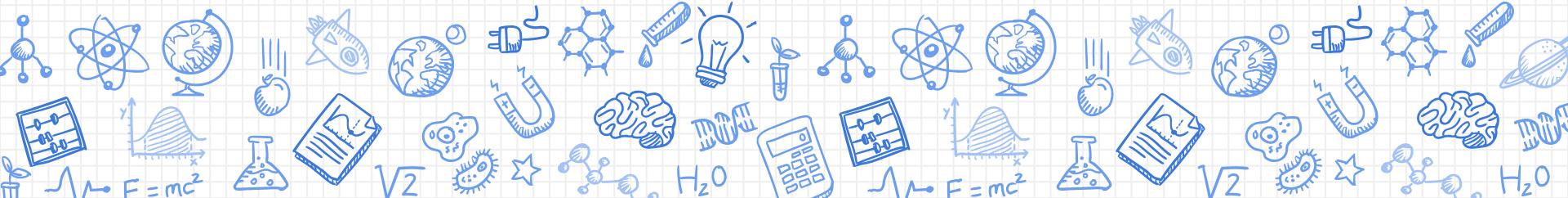
pandas

$$y_i t = \beta' x_{it} + \mu_i + \epsilon_{it}$$

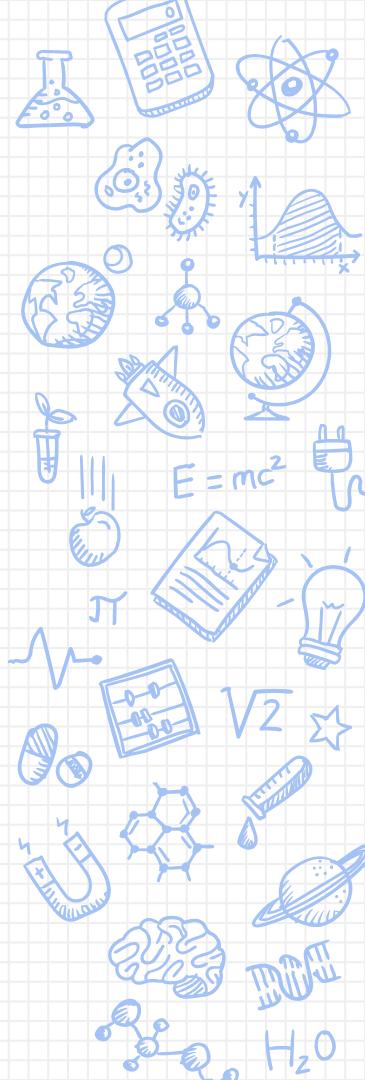
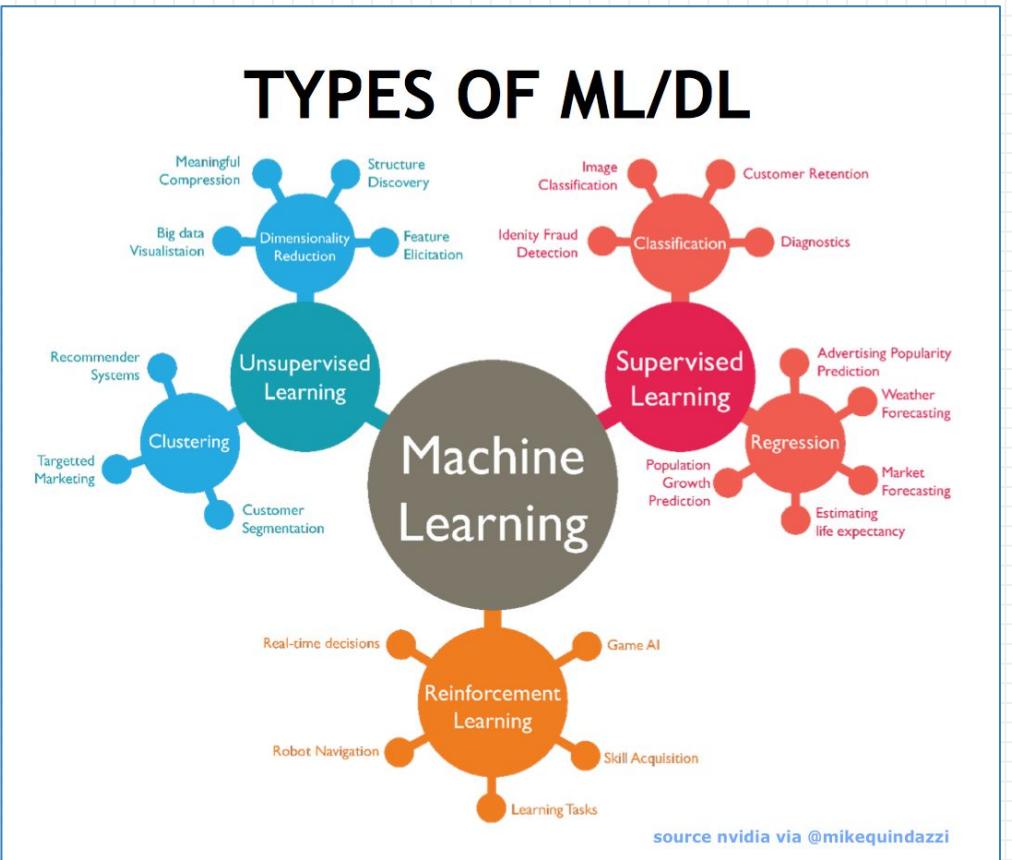


Modelado

Creación de modelos de Machine Learning

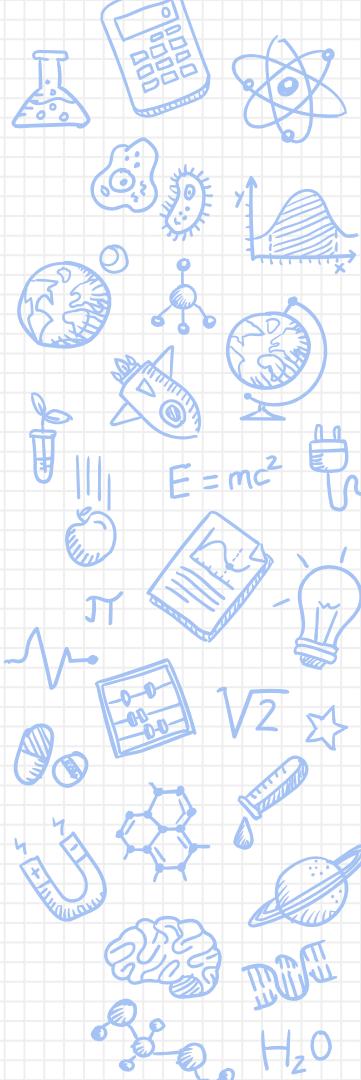


Modelado – Machine Learning Types



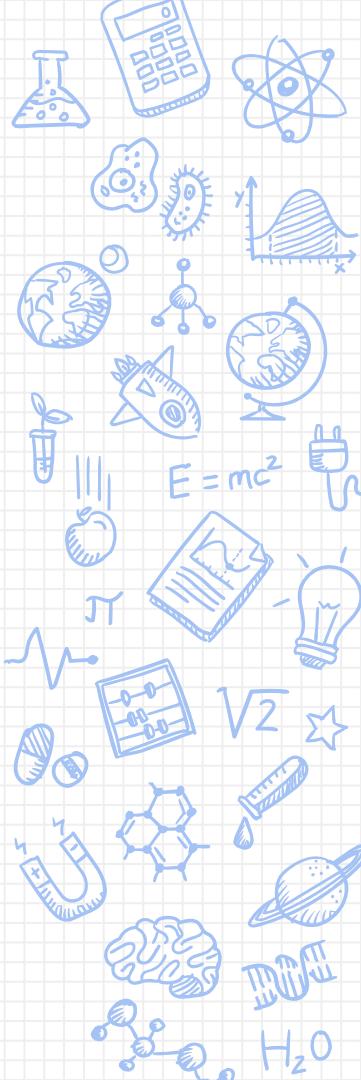
Modelado – Aprendizaje Supervisado

El modelo aprende en base a ejemplos que fueron previamente etiquetados por humanos.



Modelado – Clasificación

Si el target es discreto, hablamos de un modelo de clasificación.



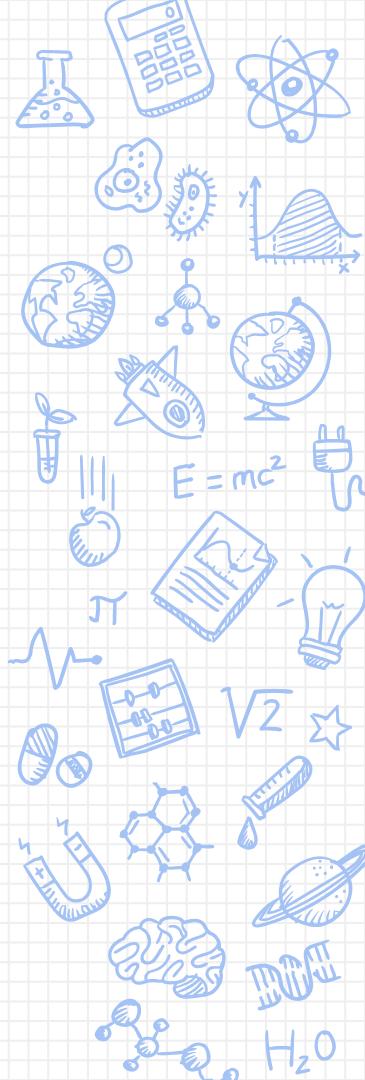
Modelado – Ejemplo de Clasificación

Features

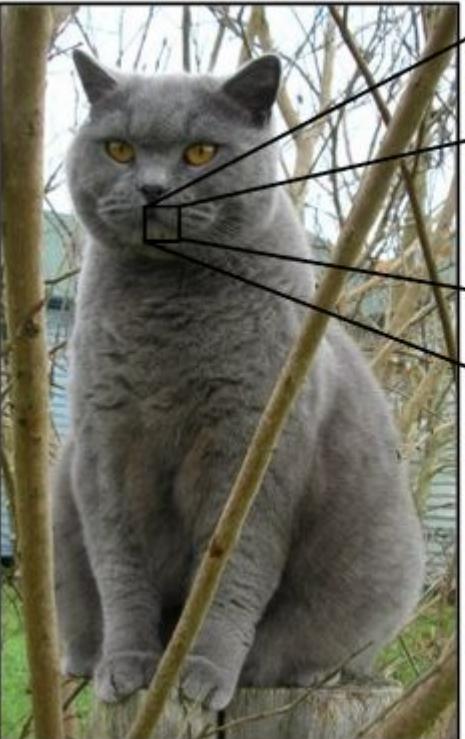
Target

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1
8	125	96	0	0	0	0.232	54	1
4	110	92	0	0	37.6	0.191	30	0
10	168	74	0	0	38	0.537	34	1
10	139	80	0	0	27.1	1.441	57	0

Fuente: <https://github.com/susanli2016/Machine-Learning-with-Python/blob/master/diabetes.csv>



Modelado – Ejemplo de Clasificación



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	80	00
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	45	54	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	53	08	30	03	49	13	36	65
52	70	95	23	04	60	11	42	62	11	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	03	59	41	92	36	54	22	40	40	28	66	33	13	80
24	47	19	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
52	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
20	34	68	87	57	62	20	72	03	46	35	67	46	55	12	32	63	93	53	69
04	42	16	73	55	39	11	24	94	72	18	08	46	29	32	40	62	76	36	20
69	36	41	72	30	23	88	31	63	93	69	82	67	59	85	74	04	36	16	20
20	73	35	29	78	31	90	01	74	31	49	73	48	56	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	21	67	46

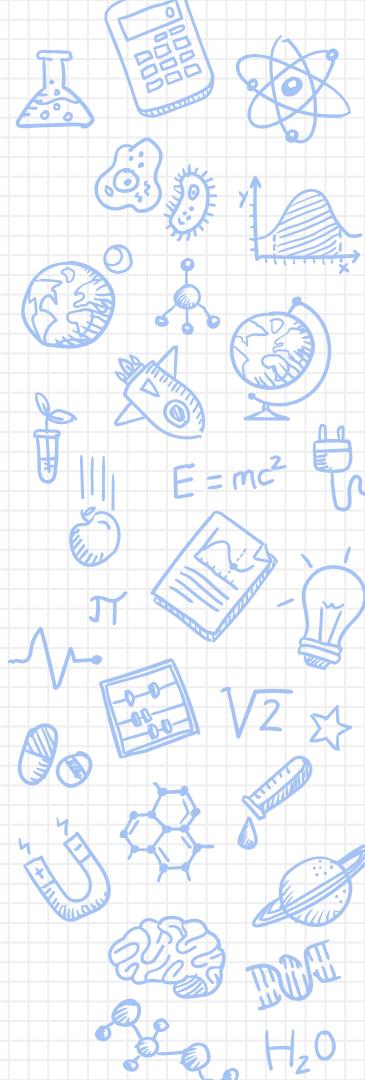
What the computer sees

image classification

82% cat
15% dog
2% hat
1% mug

Fuente:

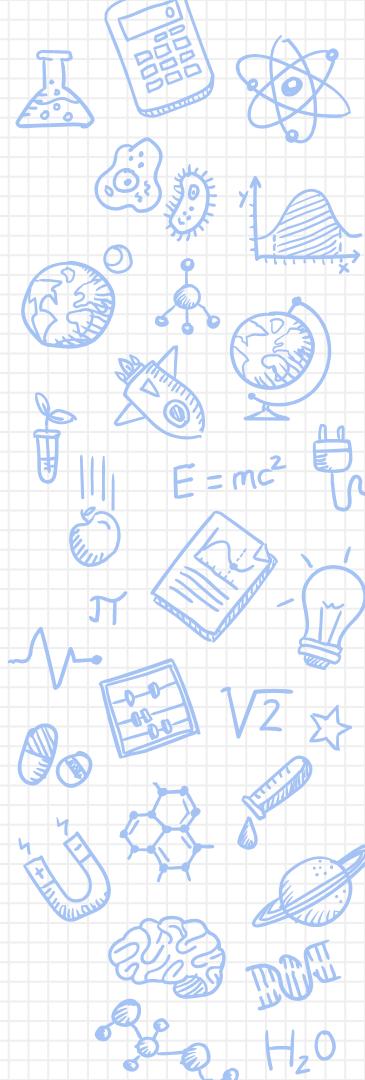
<https://medium.com/the-data-experience/building-a-data-pipeline-from-scratch-32b712cfb1db>



Modelado – Ejemplo de Clasificación

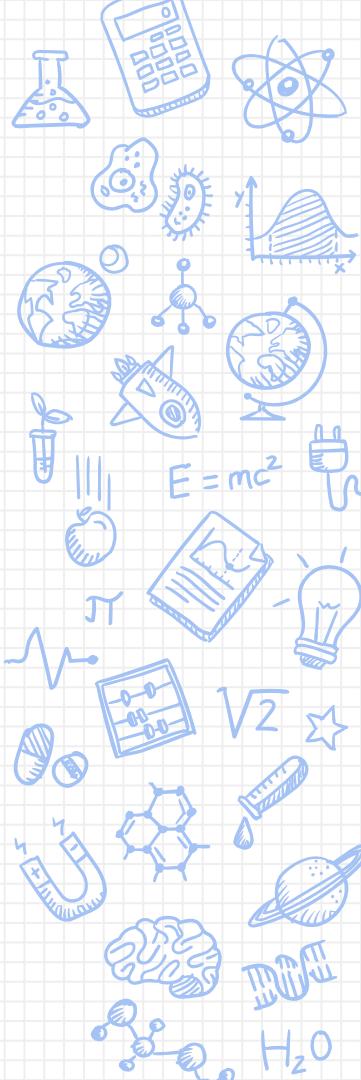
Features	Target
v2	v1
Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat... Ok lar... Joking wif u oni...	ham
Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry q... U dun say so early hor... U c already then say...	ham
Nah I don't think he goes to usf, he lives around here though	ham
FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? T... Even my brother is not like to speak with me. They treat me like aids patient.	spam
As per your request 'Melle Melle (Oru Minnaminunginte Nurunqu Vettam)' has been set as your callertune	ham
WINNER!! As a valued network customer you have been selected to receivea ♣900 prize reward! To cla...	spam
Had your mobile 11 months or more? U R entitled to Update to the latest colour mobiles with camera for 1...	ham
I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough to...	spam
SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 0...	ham
URGENT! You have won a 1 week FREE membership in our ♣100,000 Prize Jackpot! Txt the word: CLA...	spam
I've been searching for the right words to thank you for this breather. I promise i wont take your help for g...	spam
I HAVE A DATE ON SUNDAY WITH WILL!!	ham

Fuente: <https://www.kaggle.com/uciml/sms-spam-collection-dataset>



Modelado – Regresión

Si el target son continuas, hablamos de un modelo de Regresión.



Modelado – Ejemplo de Regresión

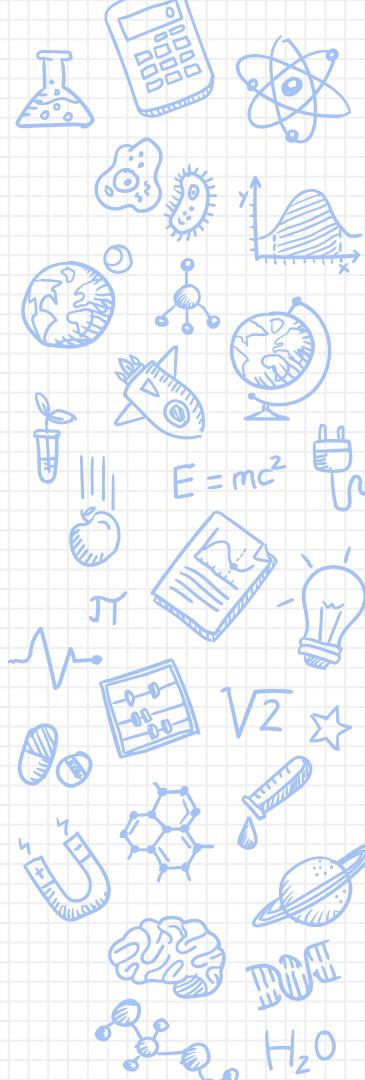
Ejemplos →

Features ↓

Target ↓

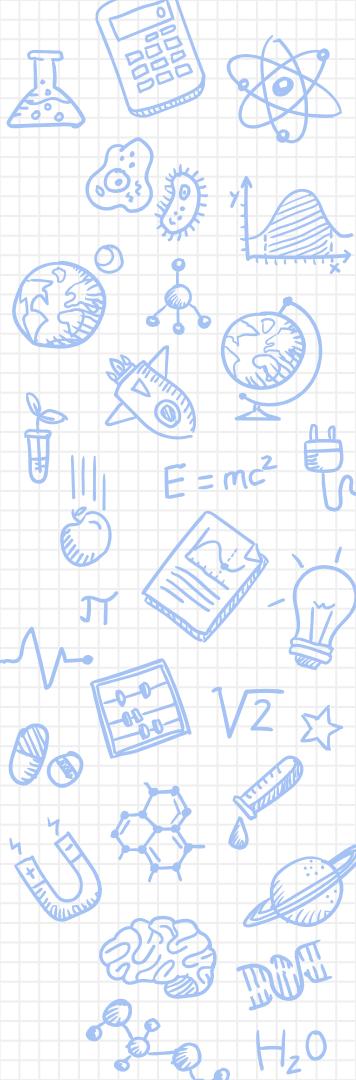
bedrooms	bathrooms	sqft_living	sqft_lot	floors	yr_built	lat	long	price
3	1	1180	5650	1	1955	47.5112	-122.257	221900
3	2	2570	7242	2	1951	47.721	-122.319	538000
2	1	770	10000	1	1933	47.7379	-122.233	180000
4	3	1960	5000	1	1965	47.5208	-122.393	604000
3	2	1680	8080	1	1987	47.6168	-122.045	510000
4	5	5420	101930	1	2001	47.6561	-122.005	1225000
3	2	1715	6819	2	1995	47.3097	-122.327	257500
3	2	1060	9711	1	1963	47.4095	-122.315	291850
3	1	1780	7470	1	1960	47.5123	-122.337	229500
3	3	1890	6560	2	2003	47.3684	-122.031	323000
3	3	3560	9796	1	1965	47.6007	-122.145	662500
2	1	1160	6000	1	1942	47.69	-122.292	468000
3	1	1430	19901	1.5	1927	47.7558	-122.229	310000

Fuente: <https://www.kaggle.com/harlfoxem/housesalesprediction>



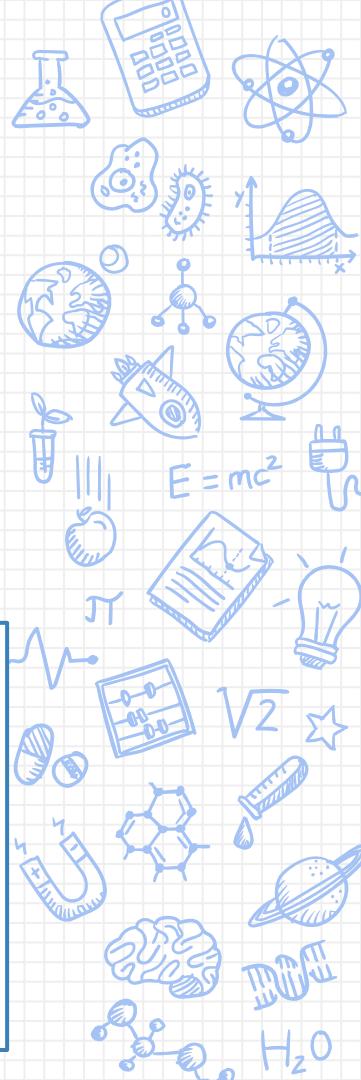
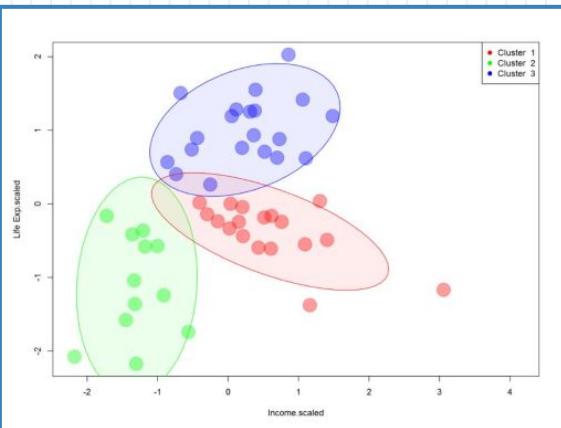
Modelado - Aprendizaje No Supervisado

El modelo aprende en base a ejemplos que sin ningún tipo de etiqueta.



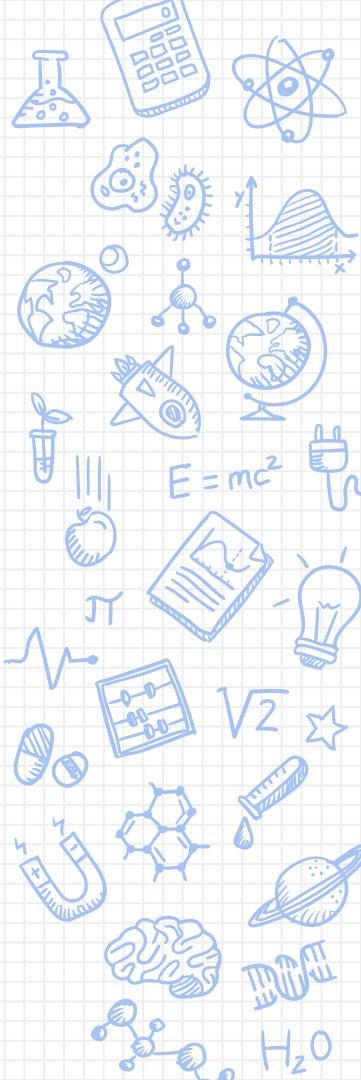
Modelado – Ejemplo de Clustering

- Segmentación de Clientes basados en sus preferencias de consumo, datos demográficos, intereses, etc.
- Clasificación en la digitalización de documentos.

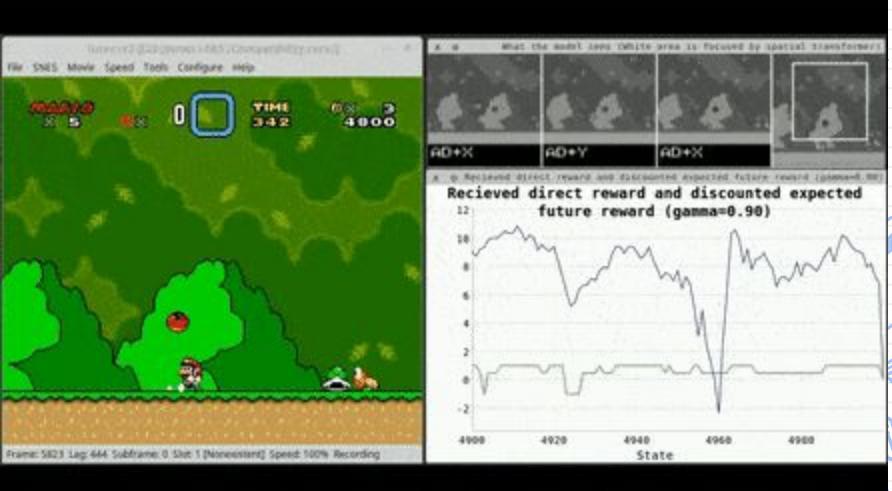
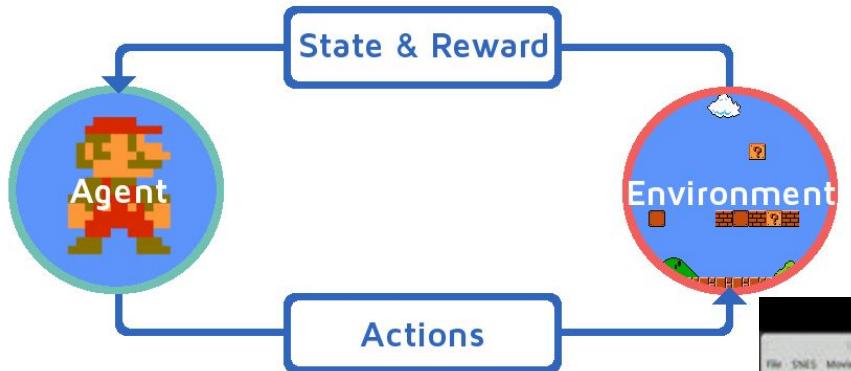


Aprendizaje por refuerzo

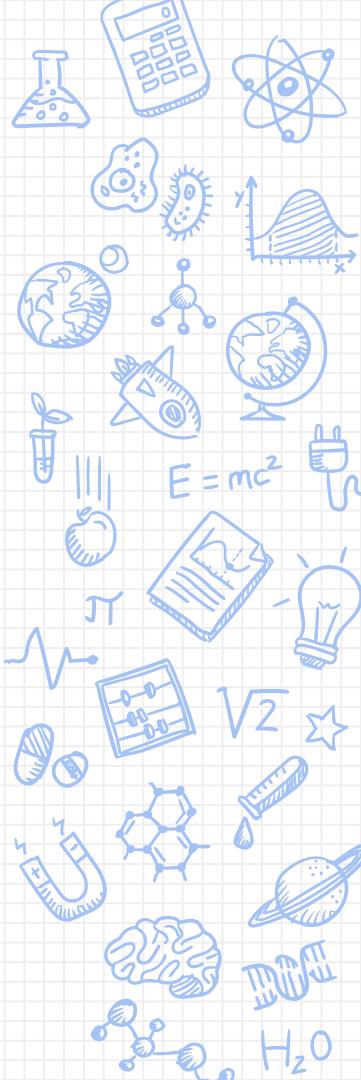
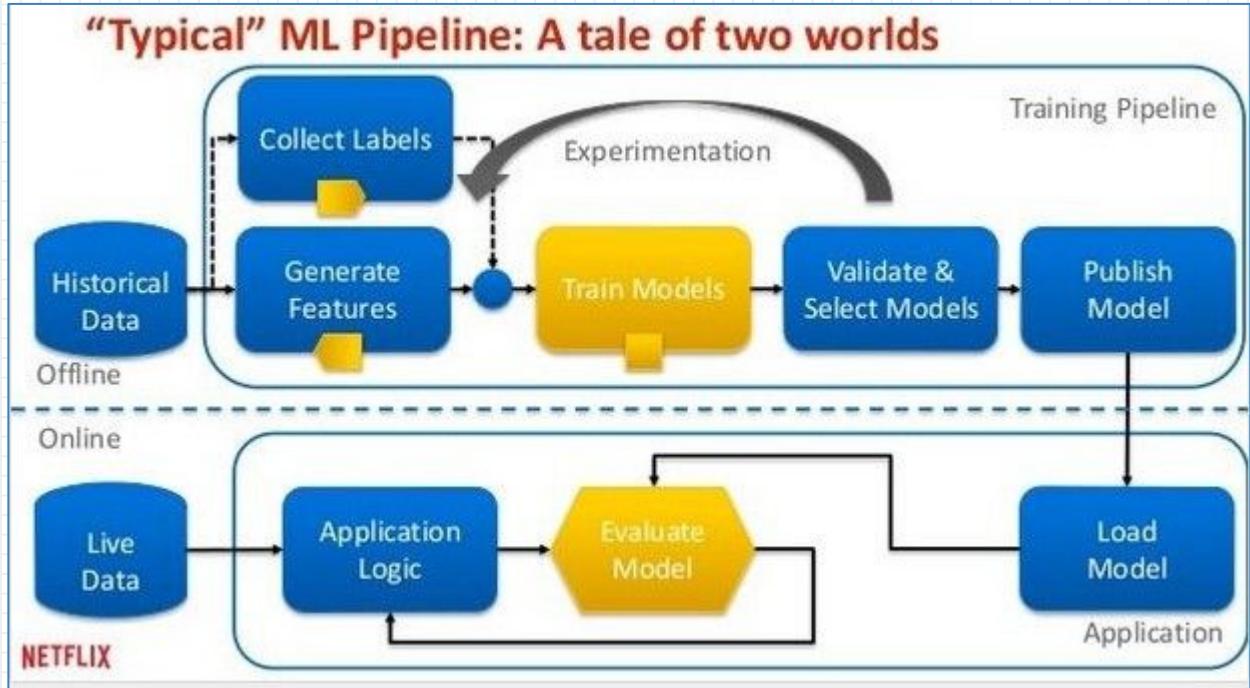
En el tercer tipo de Machine Learning, los modelos aprenden creando su propia experiencia realizando una tarea, optimizando la “recompensa” que pueda obtener ante determinada acción.



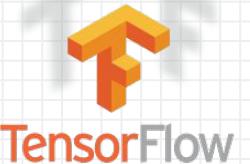
Aprendizaje por refuerzo – Ejemplos



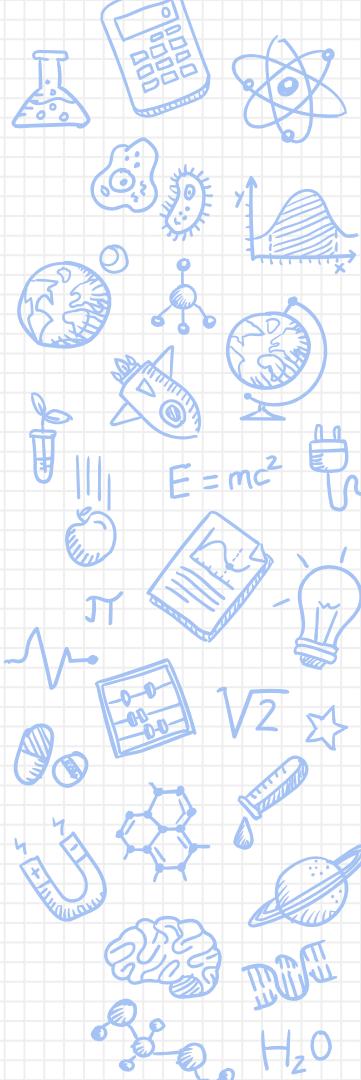
Modelado – Ciclo de Vida de ML



Tools

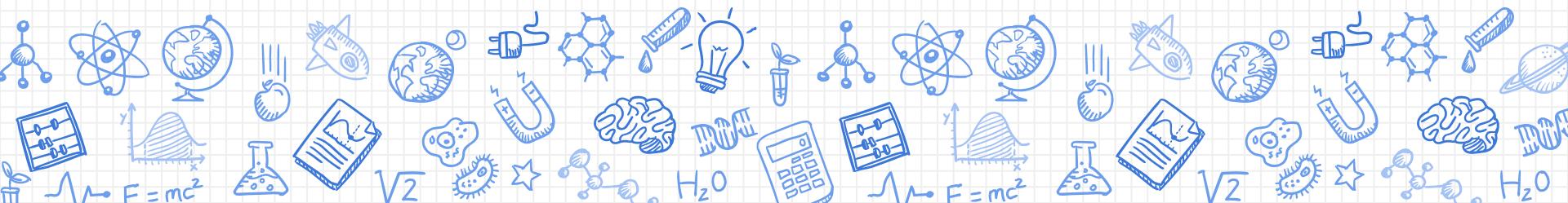


Natural Language Toolkit (NLTK)



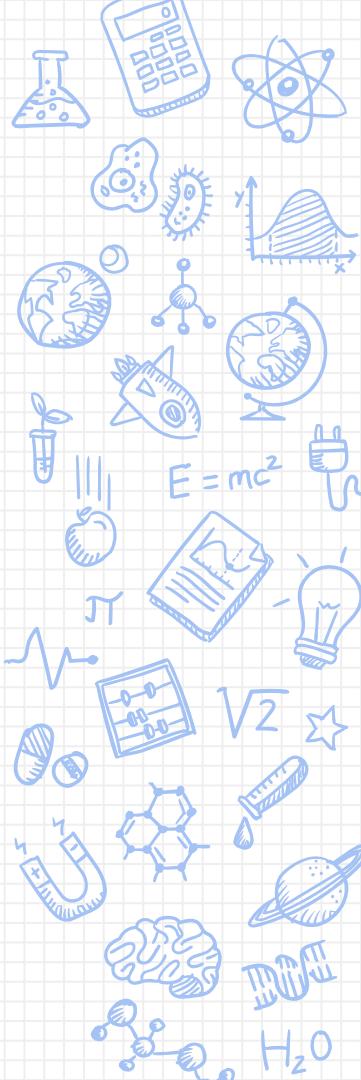
Ingeniería de Datos

Se encarga de que todo funcione
eficientemente, estable y escalable.



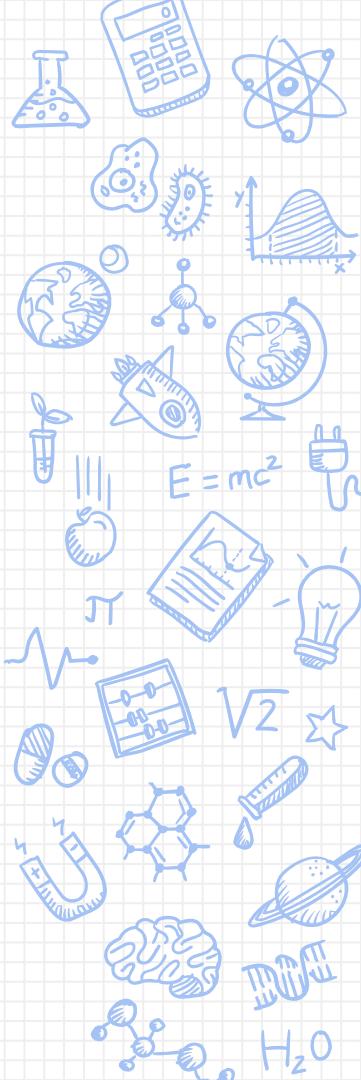
Gestión de los Datos

- Extracción de datos.
 - Transformación de datos.
 - Almacenamiento de datos.



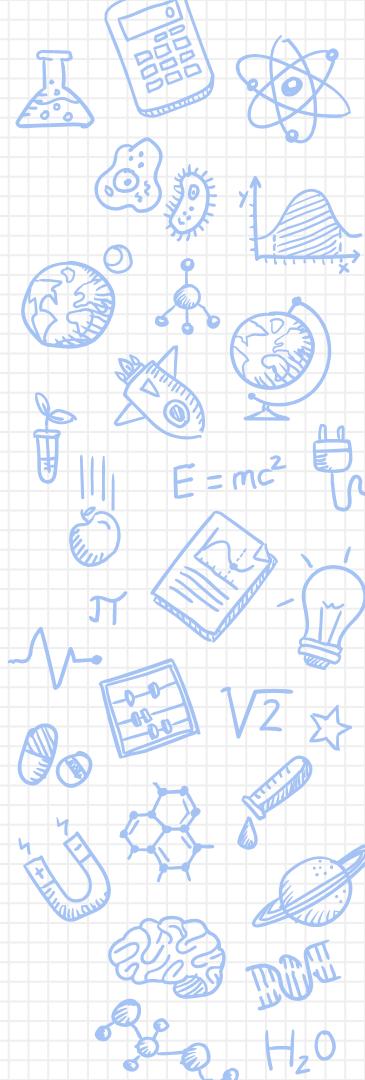
Desarrollo de Software

- Data pipelines.
 - Herramientas internas.
 - APIs.
 - Entrenamiento Automático de Modelos.
 - Monitoreo.



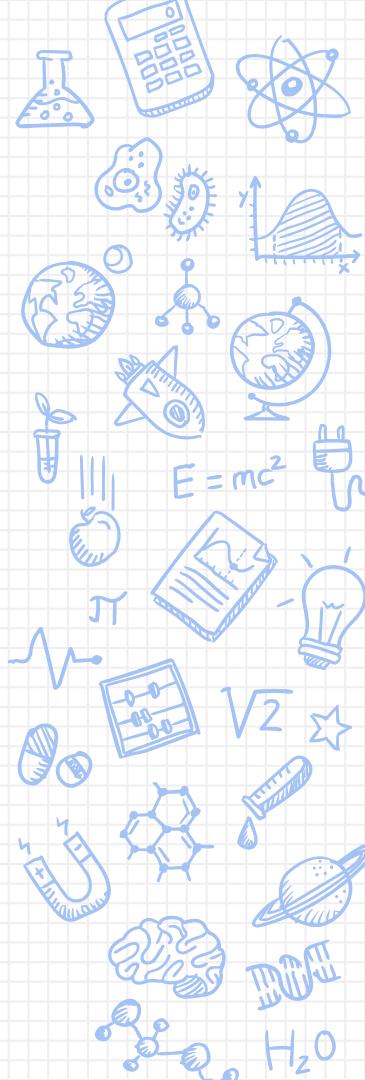
Ingeniería de Software

- Adopción de Metodologías Ágiles.
- Versionado.
- Mantenibilidad.
- Escalabilidad.
- Modularidad.

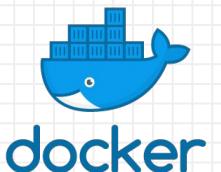
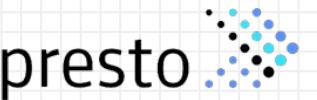


Ingeniería de Datos – Desafíos

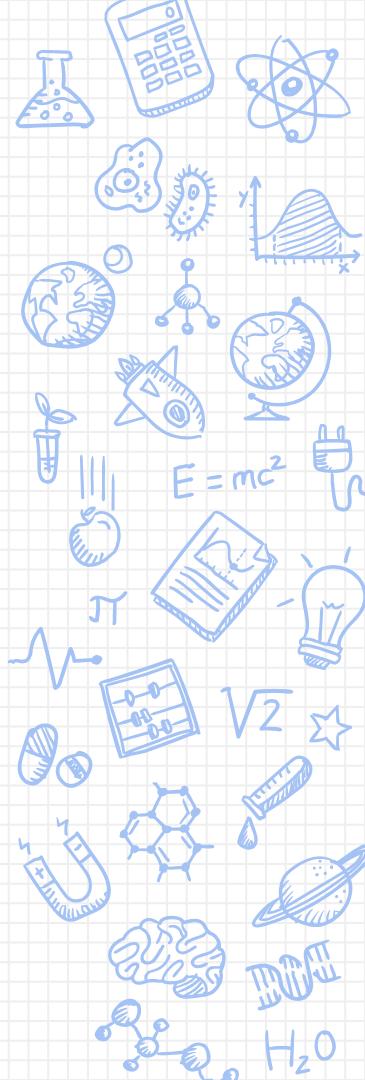
- Procesamiento Real-Time.
- Multiprocessing y Multithreading.
- Conurrencia y Asincronía.
- Tolerancia a Fallas.
- Almacenamiento Distribuido.



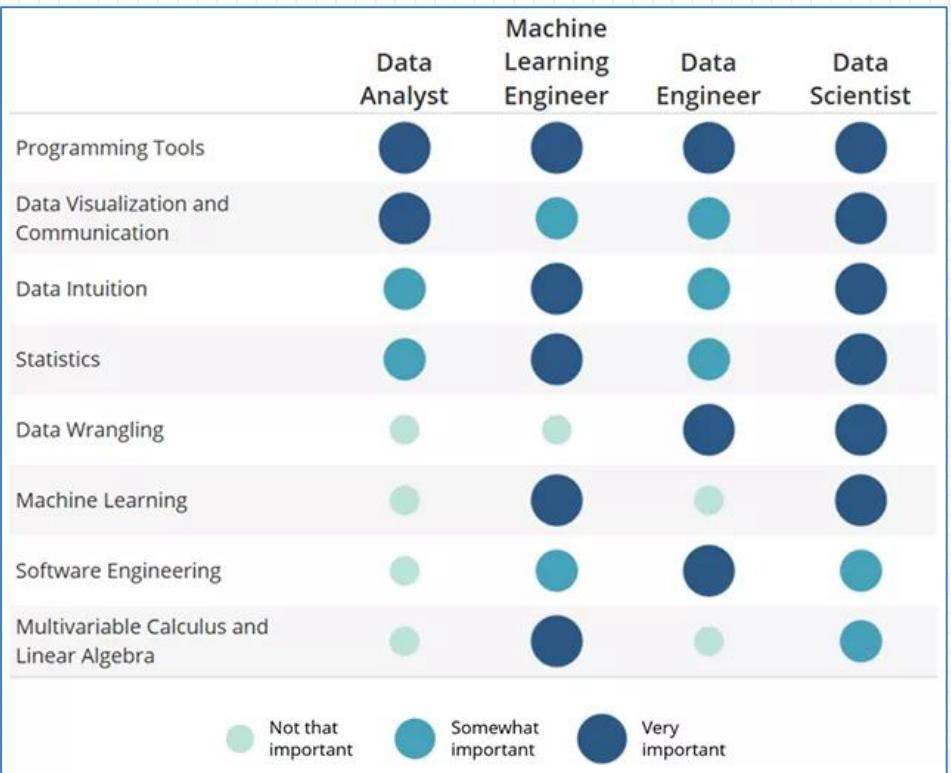
Tools



kubernetes



Skills



Fuente: <https://blog.udacity.com/2014/11/data-science-job-skills.html>