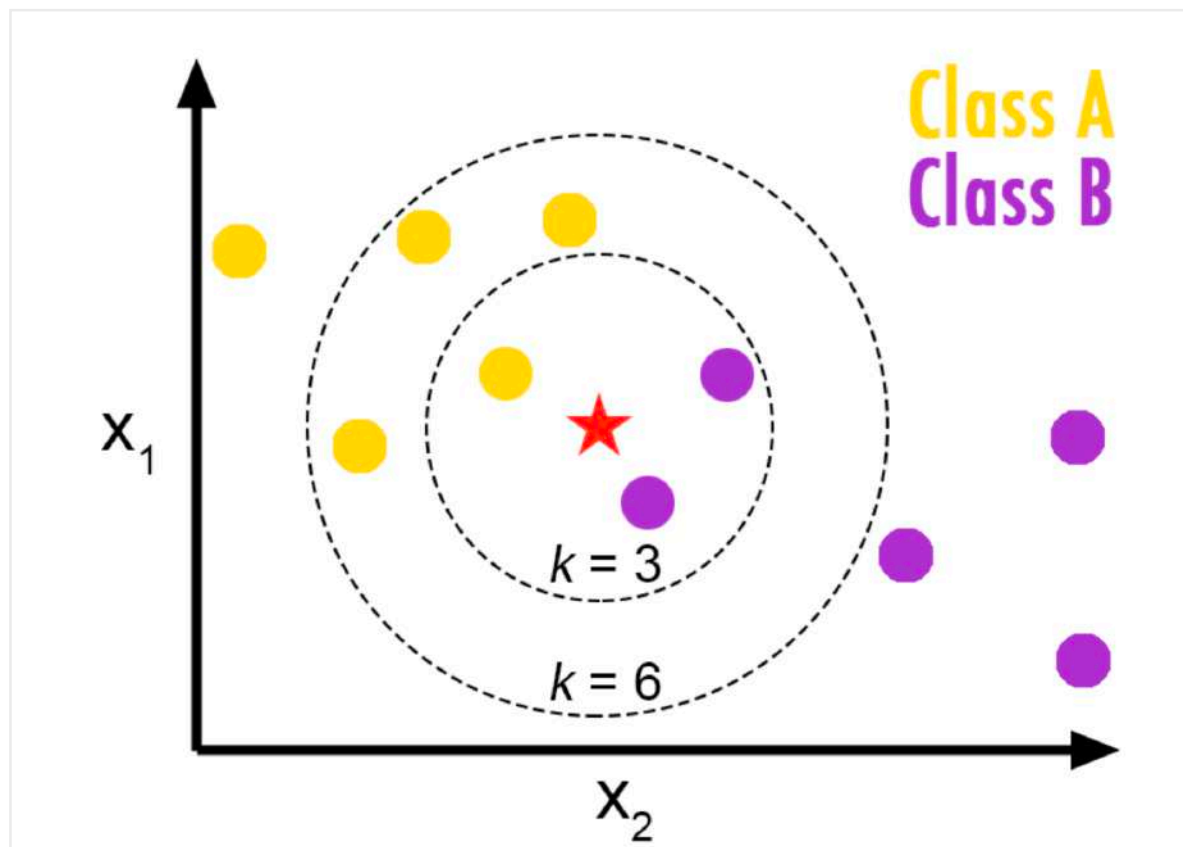


🧑 IBM - DATA SCIENCE - MACHINE LEARNING WITH PYTHON 🐍 🧠

In this Lab you will load a customer dataset, fit the data, and use K-Nearest Neighbors to predict a data point. But what is **K-Nearest Neighbors**?

K-Nearest Neighbors is a supervised learning algorithm. Where the data is 'trained' with data points corresponding to their classification. To predict the class of a given data point, it takes into account the classes of the 'K' nearest data points and chooses the class in which the majority of the 'K' nearest data points belong to as the predicted class.



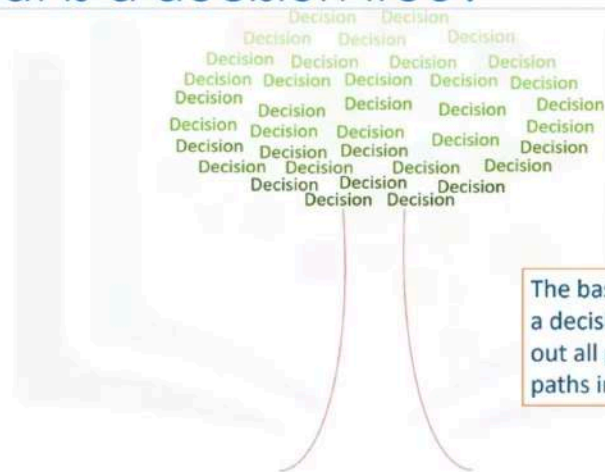
In this case, we have data points of Class A and B. We want to predict what the star (test data point) is. If we consider a k value of 3 (3 nearest data points), we will obtain a prediction of Class B. Yet if we consider a k value of 6, we will obtain a prediction of Class A.

In this sense, it is important to consider the value of k . Hopefully from this diagram, you should get a sense of what the K-Nearest Neighbors algorithm is. It considers the 'K' Nearest Neighbors (data points) when it predicts the classification of the test point.

DECISION TREES

Introduction to Decision Trees

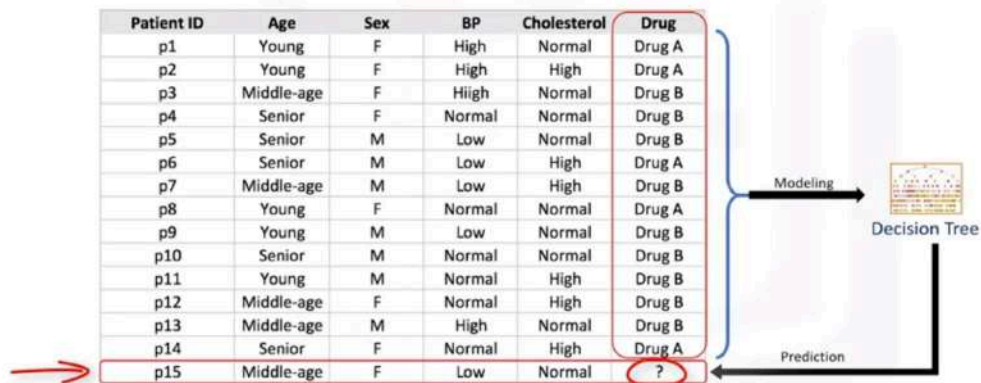
What is a decision tree?



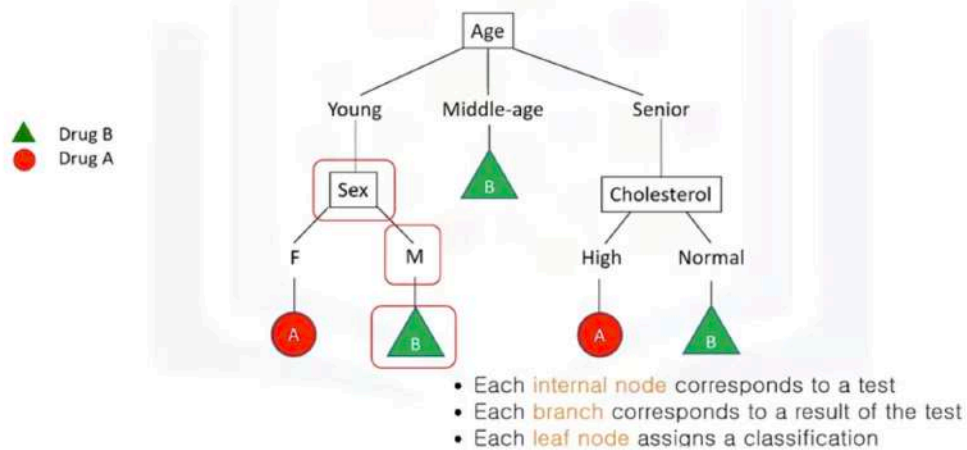
The basic intuition behind a decision tree is to map out all possible decision paths in the form of a tree.

[Narendra Nath Joshi](#)

How to build a decision tree?



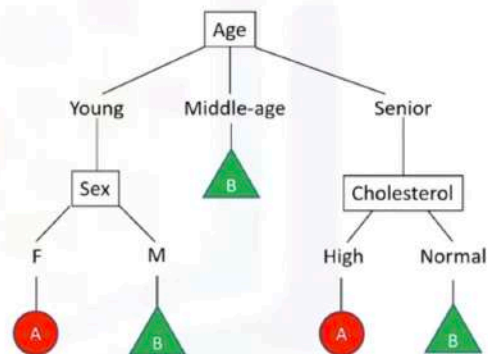
Building a decision tree with the training set



Ok muy rico todo, pero como construimos uno?

Decision tree learning algorithm

1. Choose an attribute from your dataset.
2. Calculate the significance of attribute in splitting of data.
3. Split data based on the value of the best attribute.
4. Go to step 1.

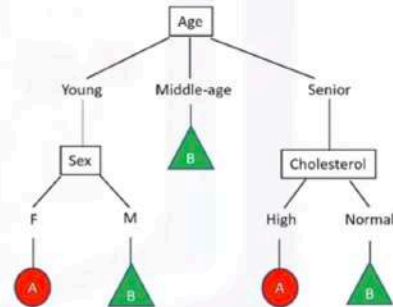


En el siguiente video vemos cómo calcular la significancia de un atributo. La joda es ir priorizándolos por este indicador e iterar así por cada nodo.

Building Decision Trees

How do you build a decision tree?

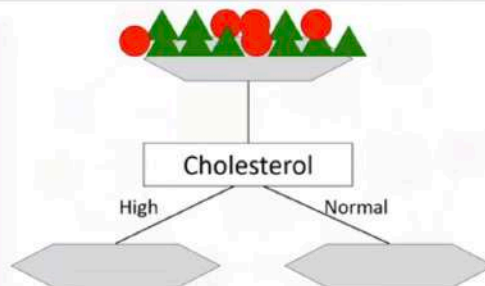
Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	Hiigh	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A
p15	Middle-age	F	Low	Normal	?



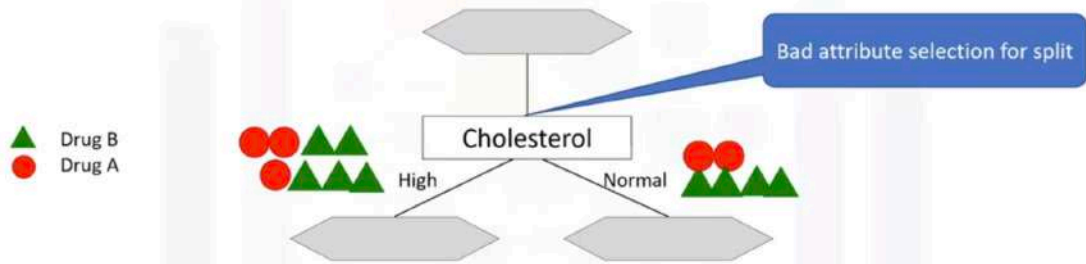
Se construyen haciendo lo que se conoce como **recursive partitioning**.

Which attribute is the best ?

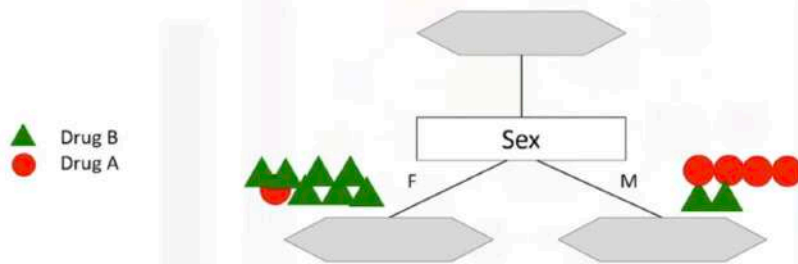
▲ Drug B
● Drug A



Which attribute is the best ?

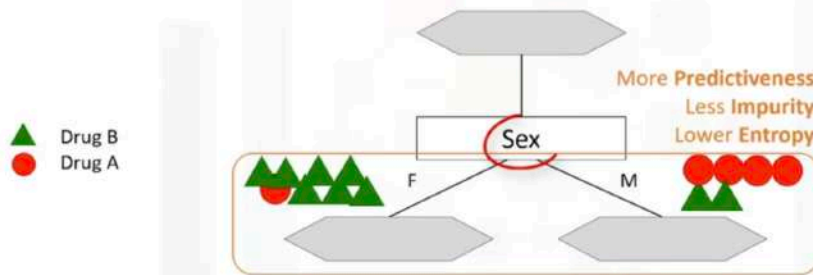


Which attribute is the best ?



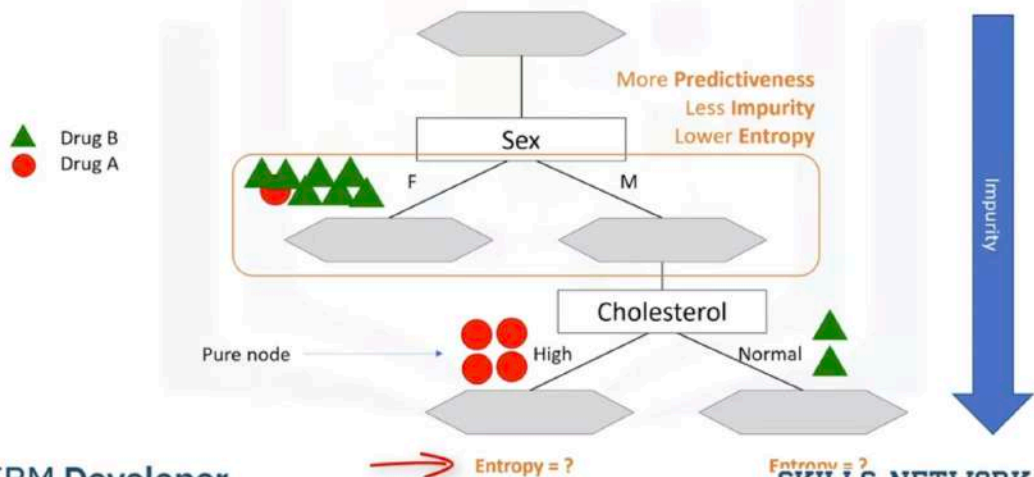
Si bien este segundo es bueno para F no lo es para M. Sin embargo sigue siendo una mejor selección de atributo que colesterol, ya que separo mejor la data, los resultados son mas puro

Which attribute is the best ?



La joda es seguir desarmando en otros atributos a medida que bajas por el árbol así obtienes resultados mas puros.

Which attribute is the best ?



Puro \longleftrightarrow 100% mismo resultado.

La impureza es calculada segun la entropia de la data en los nodos....

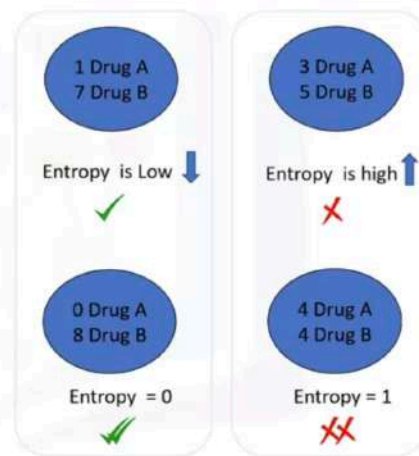
Como calculamos la entropia?

Entropy

- Measure of randomness or uncertainty

$$\text{Entropy} = -p(A)\log(p(A)) - p(B)\log(p(B))$$

The lower the Entropy, the less uniform the distribution, the purer the node.



Como ejemplo, calculemos la entropia de mi dataset para construir un arbolito:

Which attribute is the best one to use?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A

S: [9 B, 5 A]

$$E = -p(B)\log(p(B)) - p(A)\log(p(A))$$

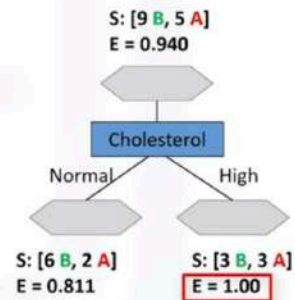
$$E = -(9/14)\log(9/14) - (5/14)\log(5/14)$$

E = 0.940

Vemos que la entropia de los datos **antes** de separar es 0.94. ahora podemos ir probando distintos atributos para encontrar el que tenga mas productividad, osea el que reduzca mas la entropia.

Is 'Cholesterol' the best attribute?

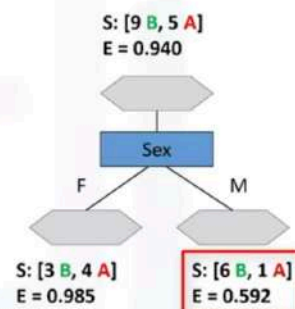
Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



Asi deberíamos ir por todos los atributos...

What about 'Sex'?

Patient ID	Age	Sex	BP	Cholesterol	Drug
p1	Young	F	High	Normal	Drug A
p2	Young	F	High	High	Drug A
p3	Middle-age	F	High	Normal	Drug B
p4	Senior	F	Normal	Normal	Drug B
p5	Senior	M	Low	Normal	Drug B
p6	Senior	M	Low	High	Drug A
p7	Middle-age	M	Low	High	Drug B
p8	Young	F	Normal	Normal	Drug A
p9	Young	M	Low	Normal	Drug B
p10	Senior	M	Normal	Normal	Drug B
p11	Young	M	Normal	High	Drug B
p12	Middle-age	F	Normal	High	Drug B
p13	Middle-age	M	High	Normal	Drug B
p14	Senior	F	Normal	High	Drug A



Which attribute is the best?



Which attribute is the best?

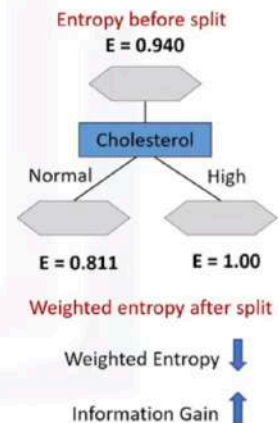


The tree with the higher **Information Gain** after splitting.

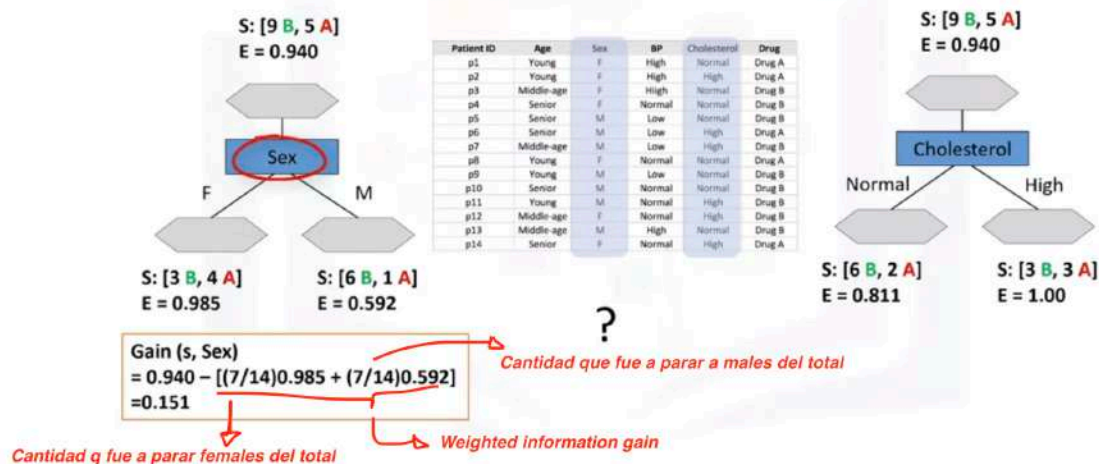
What is information gain?

Information gain is the information that can increase the level of certainty after splitting.

$$\text{Information Gain} = (\text{Entropy before split}) - (\text{weighted entropy after split})$$



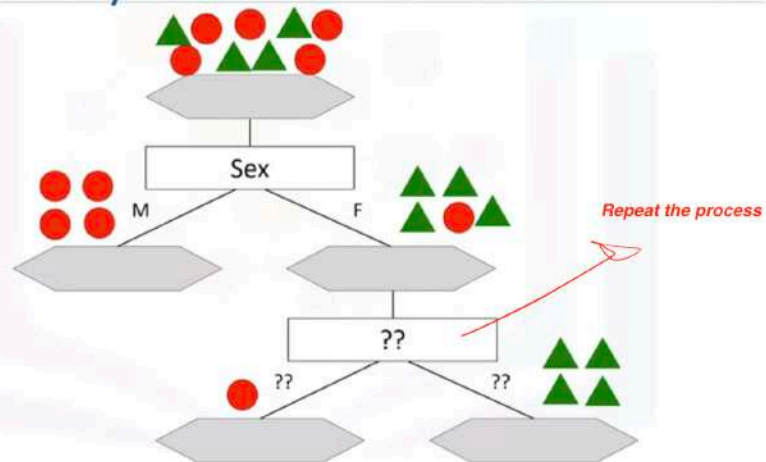
Which attribute is the best?



$$\begin{aligned} \text{Gain (s, Cholesterol)} \\ &= 0.940 - [(8/14).811 + (6/14)1.0] \\ &= 0.048 \end{aligned}$$

Por lo tanto, claramente gano el del sex. Tiene muchísimo mas Information gain.

Correct way to build a decision tree



LOGISTIC REGRESSION

Intro to Logistic Regression

Se usa para classification. Cuando lo usamos? Que es?...

What is logistic regression?

Logistic regression is a classification algorithm for categorical variables.

Independent variables										Dependent variable
tenure	age	address	income	ed	employ	equip	callicard	wireless	churn	
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	Yes
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	Yes
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	No
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	No
4	7.0	35.0	14.0	80.0	2.0	15.0	0.0	1.0	0.0	?

Continuous/Categorical variables

Categorical Vari

Es análogo a la regresión lineal, con la diferencia de que no predice valores numéricos sino categóricos. La variable dependiente que queremos predecir tiene que ser binaria (positivo negativo). Las variables independientes tienen que ser continuas. Si son categóricas las tenés que convertir transformando la data. (Dummy o indicators). Logistic regression en realidad puede ser usada para binar classification o multiclass classification, pero en este video solo vemos variables dependiente binarias.

Ejemplos de aplicación:

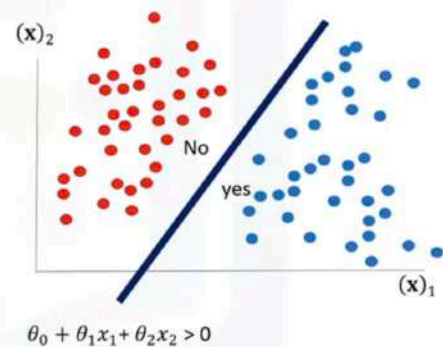
Logistic regression applications

- Predicting the probability of a person having a heart attack
- Predicting the mortality in injured patients
- Predicting a customer's propensity to purchase a product or halt a subscription
- Predicting the probability of failure of a given process or product
- Predicting the likelihood of a homeowner defaulting on a mortgage

Notar que además se entrega la probabilidad de que se pertenezca a la clase.

When is logistic regression suitable?

- If your data is binary
 - 0/1, YES/NO, True/False
- If you need probabilistic results
- When you need a linear decision boundary
- If you need to understand the impact of a feature



Building a model for customer churn

	X									y
	tenure	age	address	income	ed	employ	equip	callcard	wireless	churn
0	11.0	33.0	7.0	136.0	5.0	5.0	0.0	1.0	1.0	1.0
1	33.0	33.0	12.0	33.0	2.0	0.0	0.0	0.0	0.0	1.0
2	23.0	30.0	9.0	30.0	1.0	2.0	0.0	0.0	0.0	0.0
3	38.0	35.0	5.0	76.0	2.0	10.0	1.0	1.0	1.0	0.0

$$X \in \mathbb{R}^{m \times n}$$

$$y \in \{0,1\}$$

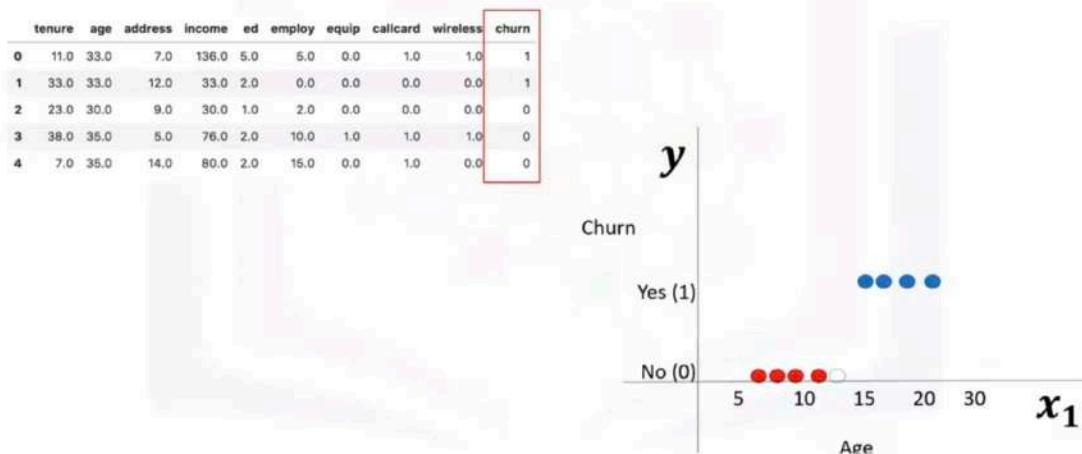
$$\hat{y} = P(y=1|x)$$

$$P(y=0|x) = 1 - P(y=1|x)$$

Logistic Regression vs Linear Regression

Vamos a ver la funcion sigmoidal.. que es "la parte principal de logístic regression"... osea que las neural networks usan logístic.

Predicting churn using linear regression



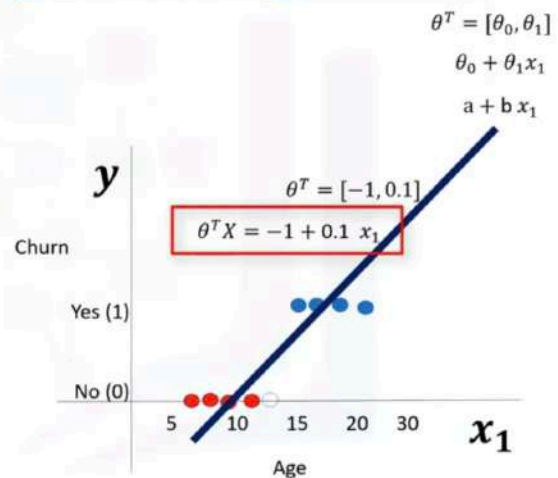
Que pasa si intentamos solucionar este problema con la misma técnica que usábamos con LR?

Predicting churn using linear regression

$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$\theta^T X = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

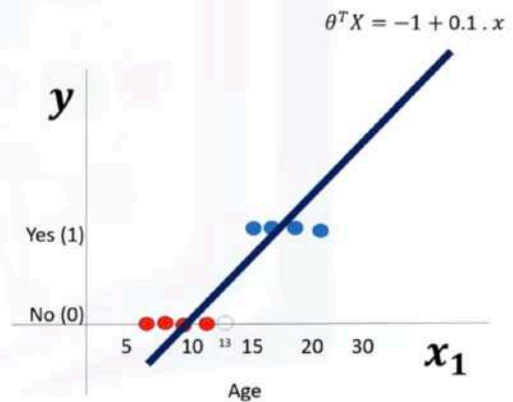


Que pasa si queremos usar esta linea de regresión obtenida para predecir el churn de mis nuevos customers?

Linear regression in classification problems?

$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$p_1 = [13] \rightarrow \begin{aligned} \theta^T X &= -1 + 0.1 \cdot x_1 \\ &= -1 + 0.1 \times 13 \\ &= 0.3 \end{aligned}$$



Pones un valor frontera que te permita decidir a que clase pertenece.

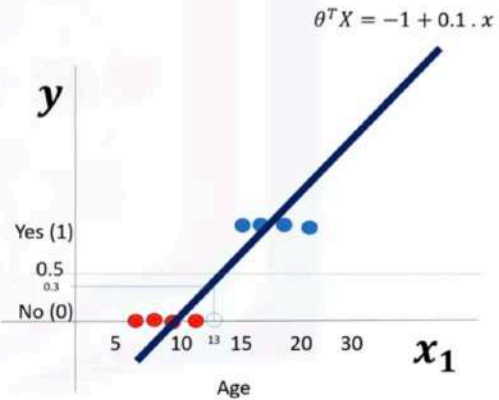
Linear regression in classification problems?

$$\theta^T X = \theta_0 + \theta_1 x_1$$

$$p_1 = [13] \rightarrow \theta^T X = -1 + 0.1 \cdot x_1 \\ = -1 + 0.1 \times 13 \\ = 0.3$$

$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

$$\theta^T X = 0.3 \\ \theta^T X < 0.5 \rightarrow \text{Class 0}$$

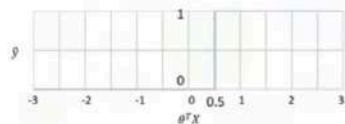


Pero tenemos un problema... cuál es la probabilidad de que ese punto de data pertenezca? Es baja.... No es la mejor forma de resolver este problema. Además hay otros problemas que verifican que linear regression no es el mejor método para clasificar (no los menciona)

Aca es donde entra la funcion sigmoidal, para cambiar la frontera

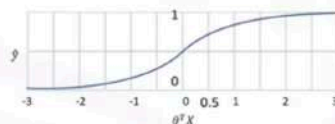
The problem with using linear regression

$$\theta^T X = \theta_0 + \theta_1 x_1 + \dots$$



$$\hat{y} = \begin{cases} 0 & \text{if } \theta^T X < 0.5 \\ 1 & \text{if } \theta^T X \geq 0.5 \end{cases}$$

$$\sigma(\theta^T X) = \sigma(\theta_0 + \theta_1 x_1 + \dots)$$



$$\hat{y} = \sigma(\theta^T X)$$

$$P(y=1|x)$$

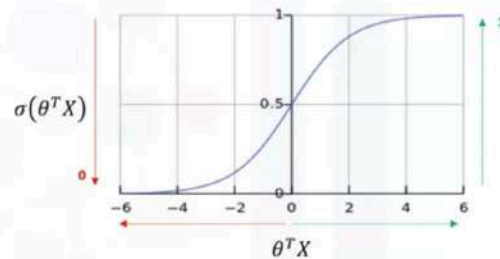
Sigmoid function in logistic regression

- Logistic Function

$$\sigma(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$$

$$\sigma(\theta^T X) = 1$$

$$\sigma(\theta^T X) = 0$$



Clarification of the customer churn model

What is the output of our model?

- $P(Y=1|X)$
- $P(y=0|X) = 1 - P(y=1|x)$
- $P(\text{Churn}=1|\text{income,age}) = 0.8$
- $P(\text{Churn}=0|\text{income,age}) = 1 - 0.8 = 0.2$

$$\sigma(\theta^T X) \rightarrow P(y=1|x)$$

$$1 - \sigma(\theta^T X) \rightarrow P(y=0|x)$$

Como logramos eso? Con el training process.

The training process

1. Initialize θ .
2. Calculate $\hat{y} = \sigma(\theta^T X)$ for a customer.
3. Compare the output of \hat{y} with actual output of customer, y , and record it as error.
4. Calculate the error for all customers.
5. Change the θ to reduce the cost.
6. Go back to step 2.

$$\sigma(\theta^T X) \rightarrow P(y=1|x)$$

$$\theta = [-1, 2]$$

$$\hat{y} = \sigma([-1, 2] \times [2, 5]) = 0.7$$

$$\text{Error} = 1 - 0.7 = 0.3$$

$$\text{Cost} = J(\theta)$$

$$\theta_{\text{new}}$$

como cambiamos los valores de theta? La forma mas popular es el **descenso del gradiente**.

La funcion sigmoideal es una constante, predefinida. No cambia en el proceso de entrenamiento.

Logistic Regression Training

- Entrenamiento de LR
- Cost function
- Tuning de parametros
- Gradient Descent

La joda de todo esto es ir cambiando los parámetros para que el modelo rinda mejor. Esto lo hacemos a través de la cost función. Usando la derivada de la función de costo (gradiente) vamos a ir poder cambiando los parámetros para reducirlo.

$$Cost(\hat{y}, y) = \frac{1}{2} (\underbrace{\sigma(\theta^T X)}_{\text{predicado}} - \underbrace{y}_{\text{valores reales}})^2$$

Se toma de forma cuadrática para eliminar negativos.

Se considera la mitad de este error como el costo

$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2 \rightarrow \text{Para un caso específico}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}, y) \rightarrow \text{Para todos los casos de mi dataset (promedio)}$$

Como calculamos cómo hacer mínimo este error? Debemos encontrar el mínimo de la función. Como es una función multivariable se convierte en un gradiente.

General cost function

$$\sigma(\theta^T X) \rightarrow P(y=1|x)$$

- Change the weight -> Reduce the cost

- Cost function

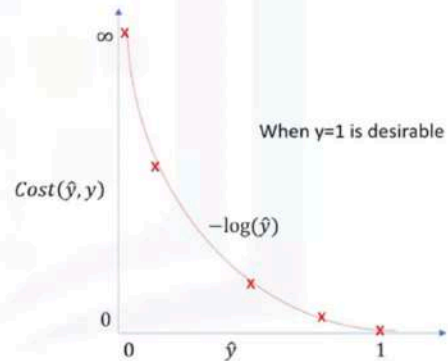
$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}, y)$$

Lo que dice el curso es que "difícil" encontrar el mínimo de esta función. Entonces lo que haces es proponer otra función de costo.

Plotting the cost function of the model

- Model \hat{y}
- Actual Value $y=1$ or 0
- If $Y=1$, and $\hat{y}=1 \rightarrow \text{cost} = 0$
- If $Y=1$, and $\hat{y}=0 \rightarrow \text{cost} = \text{large}$



La podemos reescribir como:

Logistic regression cost function

- So, we will replace cost function with:

$$Cost(\hat{y}, y) = \frac{1}{2} (\sigma(\theta^T X) - y)^2$$

$$Cost(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

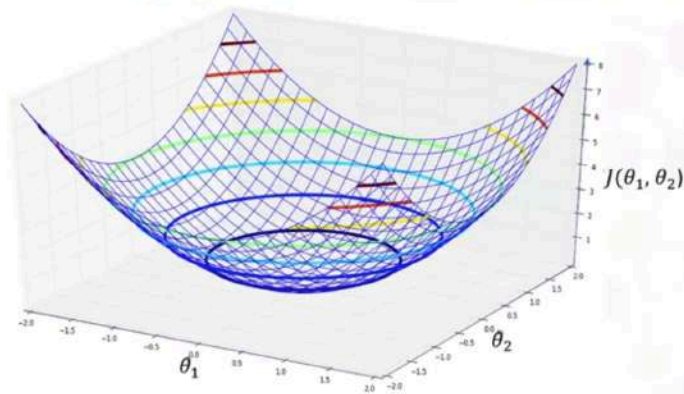
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(\hat{y}^i, y^i)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

Minimizing the cost function of the model

- How to find the best parameters for our model?
 - Minimize the cost function
- How to minimize the cost function?
 - Using Gradient Descent
- What is gradient descent?
 - A technique to use the derivative of a cost function to change the parameter values, in order to minimize the cost

Using gradient descent to minimize the cost

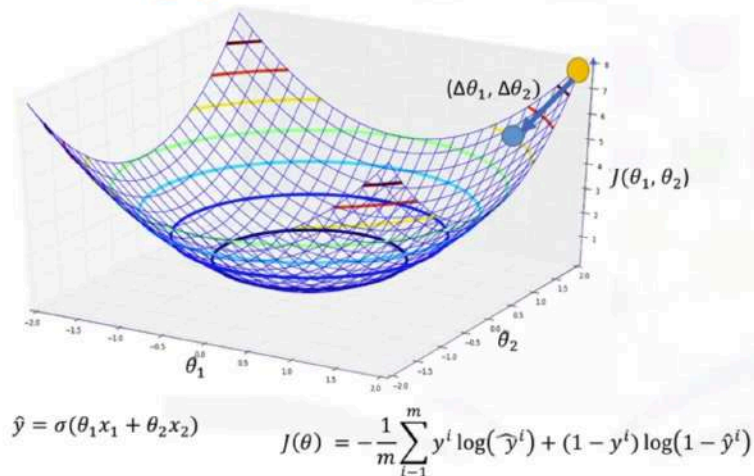


$$\hat{y} = \sigma(\theta_1 x_1 + \theta_2 x_2)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

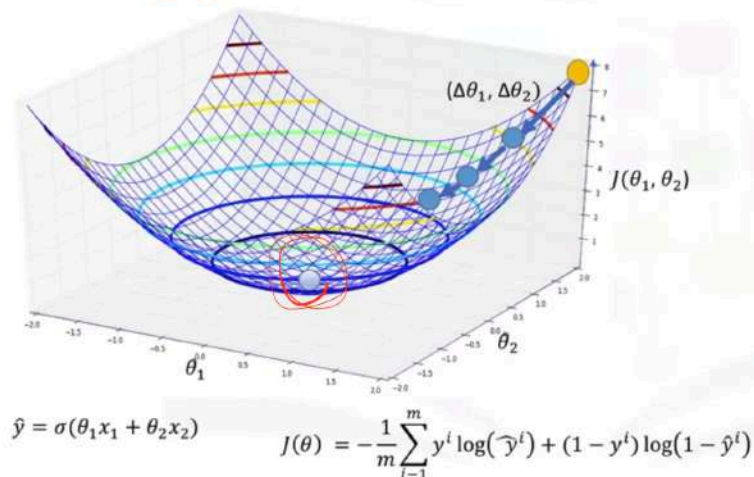
Empezas con un punto random y vas iterando cambiando de a diferenciales (deltas) de los parámetros Tita sub n:

Using gradient descent to minimize the cost



Vas caminando por la superficie y análisis los valores del gradiente.

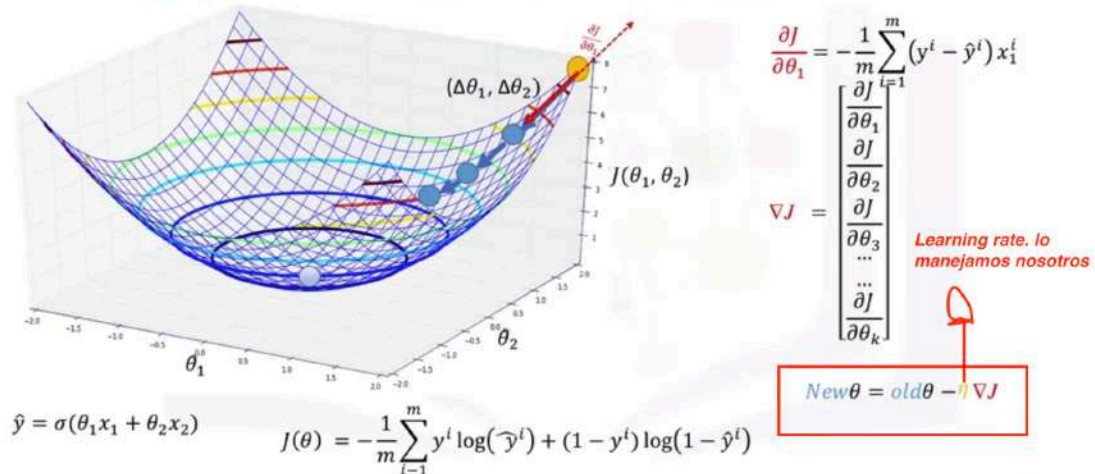
Using gradient descent to minimize the cost



Para encontrar la dirección y tamaño de estos pasos, necesitamos ir calculando el gradiente de la función de costo. El gradiente es la "pendiente" de la superficie en cada punto y la dirección del gradiente apunta en hacia donde crece la función. Nos tenemos que mover en el sentido contrario del gradiente para estar seguros de que vamos hacia donde la función tiene su mínimo!

El valor del gradiente (numérico) nos indica también que tamaño de paso tenemos que tomar. Este va a ir disminuyendo a medida que aumenten las iteraciones.

Using gradient descent to minimize the cost



Training algorithm recap

1. initialize the parameters randomly.
2. Feed the cost function with training set, and calculate the error.
3. Calculate the gradient of cost function.
4. Update weights with new values.
5. Go to step 2 until cost is small enough.
6. Predict the new customer X.

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots]$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

$$\nabla J = \left[\frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \frac{\partial J}{\partial \theta_3}, \dots, \frac{\partial J}{\partial \theta_k} \right]$$

$$\theta_{\text{new}} = \theta_{\text{prv}} - \eta \nabla J$$

$$P(y=1|x) = \sigma(\theta^T X)$$

SUPPORT VECTOR MACHINE (S.V.M)

Support Vector Machine

Classification with SVM



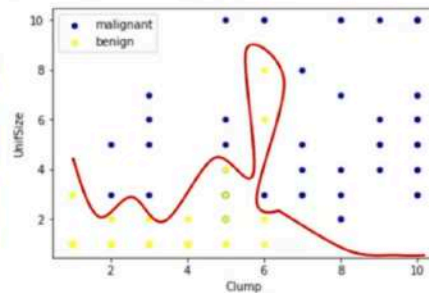
Necesitamos un separator. El tema es que la mayoría de los datasets del mundo real no aceptan separaciones lineales, sino con curvas

What is SVM?

SVM is a supervised algorithm that classifies cases by finding a separator.

1. Mapping data to a **high-dimensional** feature space
2. Finding a **separator**

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucI	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign



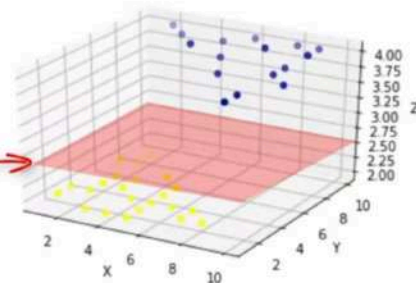
Si transferimos la data a un espacio de R3:

What is SVM?

SVM is a supervised algorithm that classifies cases by finding a separator.

1. Mapping data to a **high-dimensional** feature space
2. Finding a **separator**

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucI	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

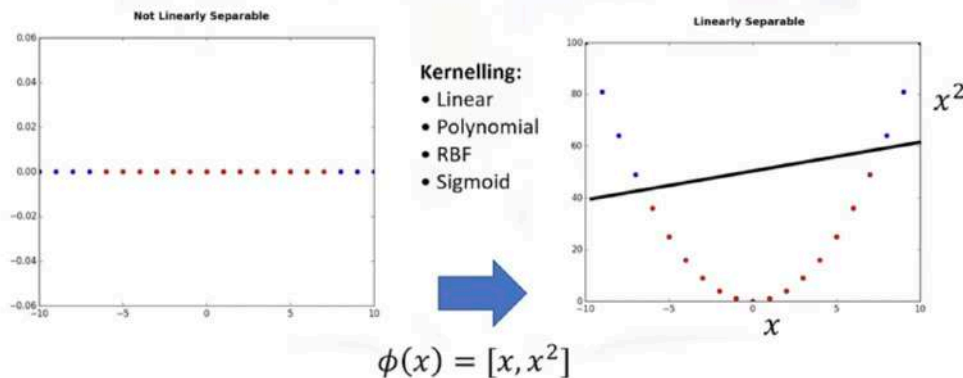


El output de SVM es un hiperplane optimizado para separar de la mejor forma los casos de la clasificación.

Como transferimos data de una manera tal que podamos dibujar un separador como un hyperplane? Y como podemos encontrar el mejor de estos planos post transformación?

Mapping data into a higher dimensional space is called Kernelling

Data transformation



El Fi de la slide de arriba es lo que se conoce como Kernell Function

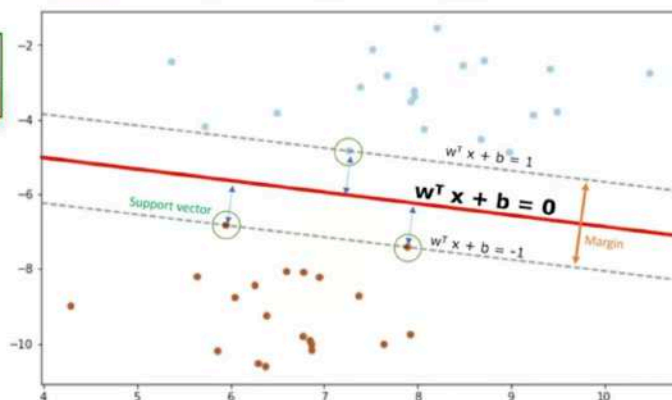
La siguiente pregunta es como encontramos el separador mas optimo después de haber transformado mis datos?

Un primer approach seria elegir el que represente el mayor margen entre las clases

Lo que esta cerca de la frontera es lo que llamamos support vector.

Using SVM to find the hyperplane

Find w and b such that
 $\Phi(w) = \frac{1}{2} w^T w$ is minimized;
and for all $\{(x_i, y_i)\}: y_i (w^T x_i + b) \geq 1$



No entra en detalle con la matemática :(

Este problema de optimización dice que también se puede resolver por descenso del gradiente, pero que no lo vemos en este video.

Como usan solamente los support vectores para el algoritmo, son eficientes. Pero para datasets chiquitos. (<1000 filas)

Lo malo es que muy común overfitteen si la cantidad de features es mayor a la cantidad de **Samples**.

Pros and cons of SVM

- Advantages:
 - Accurate in high-dimensional spaces
 - Memory efficient
- Disadvantages:
 - Prone to over-fitting
 - No probability estimation
 - Small datasets

SVM applications

- Image recognition
- Text category assignment
- Detecting spam
- Sentiment analysis
- Gene Expression Classification
- Regression, outlier detection and clustering