

Classification:

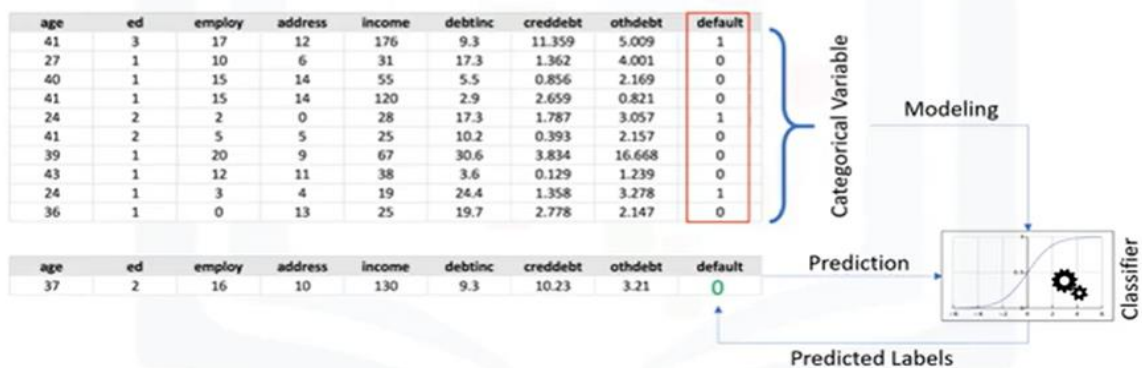
What is classification?

- A supervised learning approach.
- Categorizing some unknown items into a discrete set of categories or "classes".
- The target attribute is a categorical variable.

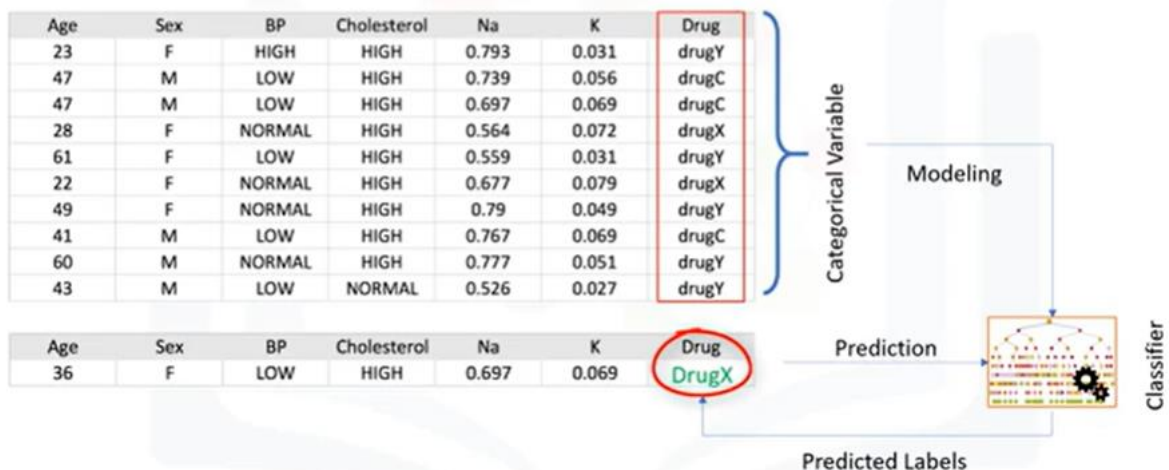
Example of a binary classification.

How does classification work?

Classification determines the class label for an unlabeled test case.



Example of multi-class classification

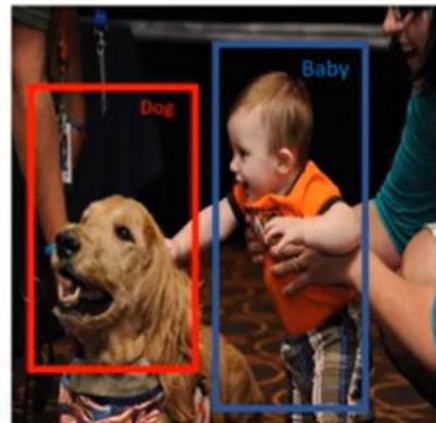


Classification use cases

	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	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	?

- Which category a customer belongs to?
- Whether a customer switches to another provider/brand?
- Whether a customer responds to a particular advertising campaign?

Classification applications



Classification algorithms in machine learning



- Decision Trees (ID3, C4.5, C5.0)
- Naïve Bayes
- Linear Discriminant Analysis
- k -Nearest Neighbor
- Logistic Regression
- Neural Networks
- Support Vector Machines (SVM)

K-Nearest Neighbours

Intro to KNN

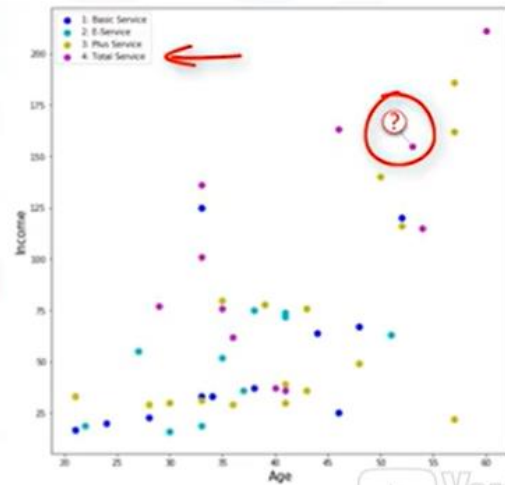
X: Independent variable											Y: Dependent variable
	region	age	marital	address	income	ed	employ	retire	gender	reside	custcat
0	2	44	1	9	64	4	5	0	0	2	1
1	3	33	1	7	136	5	5	0	0	6	4
2	3	52	1	24	116	1	29	0	1	2	3
3	2	33	0	12	33	2	0	0	1	1	1
4	2	30	1	9	30	1	2	0	0	4	3
5	2	39	0	17	78	2	16	0	1	1	3
6	3	22	1	2	19	2	4	0	1	5	2
7	2	35	0	5	76	2	10	0	0	3	4
8	3	50	1	7	166	4	31	0	0	5	?

Value	Label
1	Basic Service
2	E-Service
3	Plus Service
4	Total Service

Determining the class using 1st KNN

	region	age	marital	address	income	ed	employ	retire	gender	reside	custcat
0	2	44	1	9	64	4	5	0	0	2	1
1	3	33	1	7	136	5	5	0	0	6	4
2	3	52	1	24	116	1	29	0	1	2	3
3	2	33	0	12	33	2	0	0	1	1	1
4	2	30	1	9	30	1	2	0	0	4	3
5	2	39	0	17	78	2	16	0	1	1	3
6	3	22	1	2	19	2	4	0	1	5	2
7	2	35	0	5	76	2	10	0	0	3	4
8	3	50	1	7	166	4	31	0	0	5	?

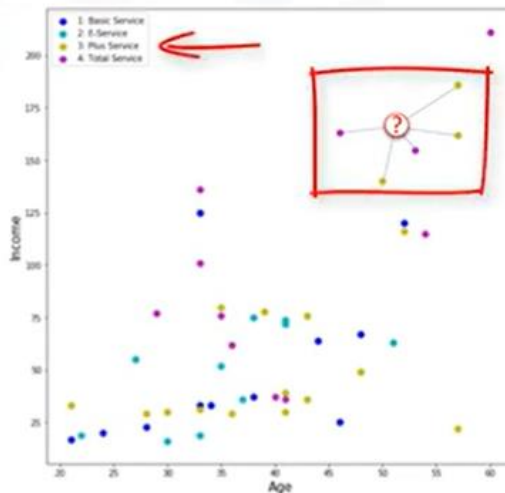
1-NN → 4: Total Service



Determining the class using the 5 KNNs

	region	age	marital	address	income	ed	employ	retire	gender	reside	custcat
0	2	44	1	9	64	4	5	0	0	2	1
1	3	33	1	7	136	5	5	0	0	6	4
2	3	52	1	24	116	1	29	0	1	2	3
3	2	33	0	12	33	2	0	0	1	1	1
4	2	30	1	9	30	1	2	0	0	4	3
5	2	39	0	17	78	2	16	0	1	1	3
6	3	22	1	2	19	2	4	0	1	5	2
7	2	35	0	5	76	2	10	0	0	3	4
8	3	50	1	7	166	4	31	0	0	5	?

5-NN → 3: Plus Service



What is K-Nearest Neighbour (or KNN)?

- A method for classifying cases based on their similarity to other cases.
- Cases that are near each other are said to be "neighbours"
- Based on similar cases with same class labels are near each other.

The K-Nearest Neighbors algorithm

1. Pick a value for K.
2. Calculate the distance of unknown case from all cases.
3. Select the K-observations in the training data that are “nearest” to the unknown data point.
4. Predict the response of the unknown data point using the most popular response value from the K-nearest neighbors.

Calculating the similarity/distance in a multi-dimensional space



Customer 1		
Age	Income	Education
34	190	3

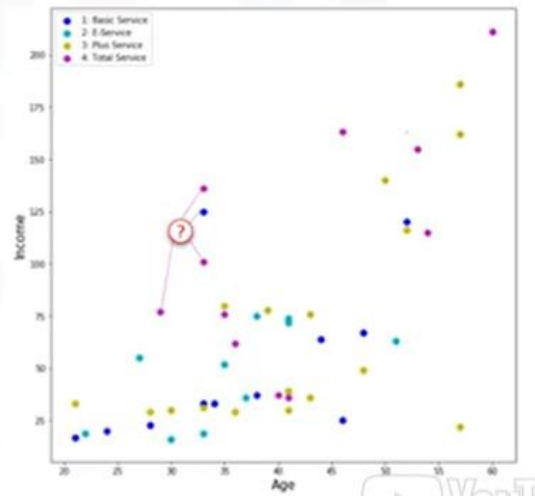


Customer 2		
Age	Income	Education
30	200	8

$$\begin{aligned}\text{Dis}(x_1, x_2) &= \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2} \\ &= \sqrt{(34 - 30)^2 + (190 - 200)^2 + (3 - 8)^2} = 11.87\end{aligned}$$

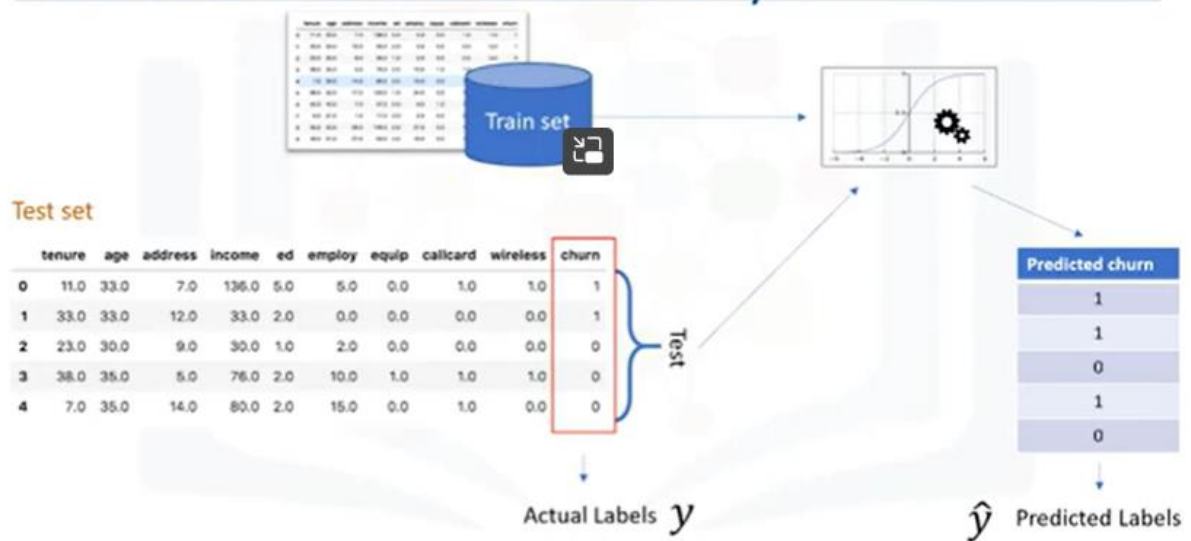
What is the best value of K for KNN?

- K = 1 class 1
- K = 20 ?



Evaluation Metrics in Classification

Classification accuracy



Jaccard index

y : Actual labels

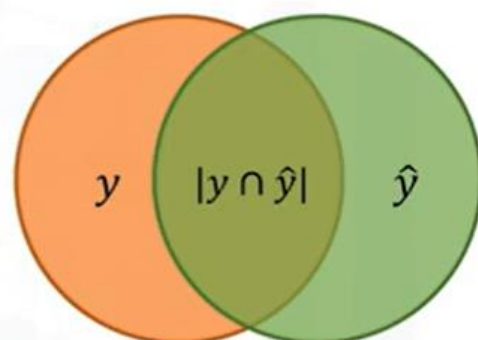
\hat{y} : Predicted labels

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$

y : [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]

\hat{y} : [1, 1, 0, 0, 0, 1, 1, 1, 1, 1]

$$J(y, \hat{y}) = \frac{8}{10+10-8} = 0.66$$



$$J(y, \hat{y}) = 0.0$$

$$J(y, \hat{y}) = 1.0$$

Higher Accuracy

F1-score

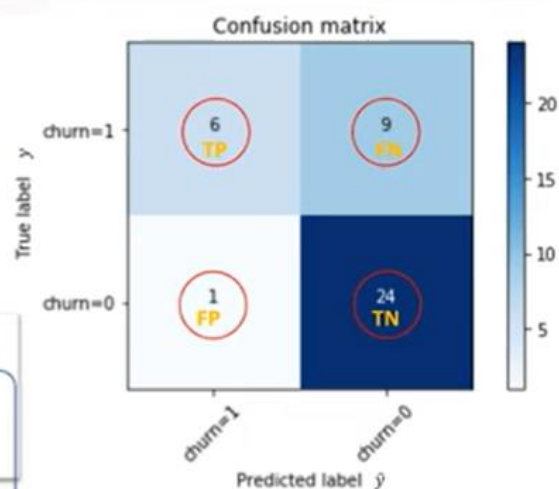
- Precision = $TP / (TP + FP)$
- Recall = $TP / (TP + FN)$
- F1-score = $2 \times (prc \times rec) / (prc + rec)$

F1-score: 0.00 ... 0.20 ... 0.55 ... 0.83 ... 1.00

Higher Accuracy

	precision	recall	f1-score
Churn = 0	0.73	0.96	0.83
Churn = 1	0.86	0.40	0.55

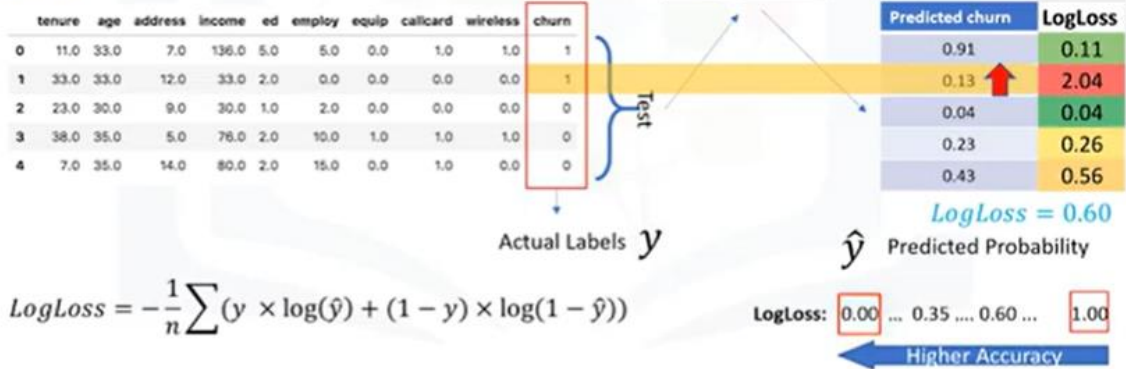
Avg Accuracy = 0.72



Log loss

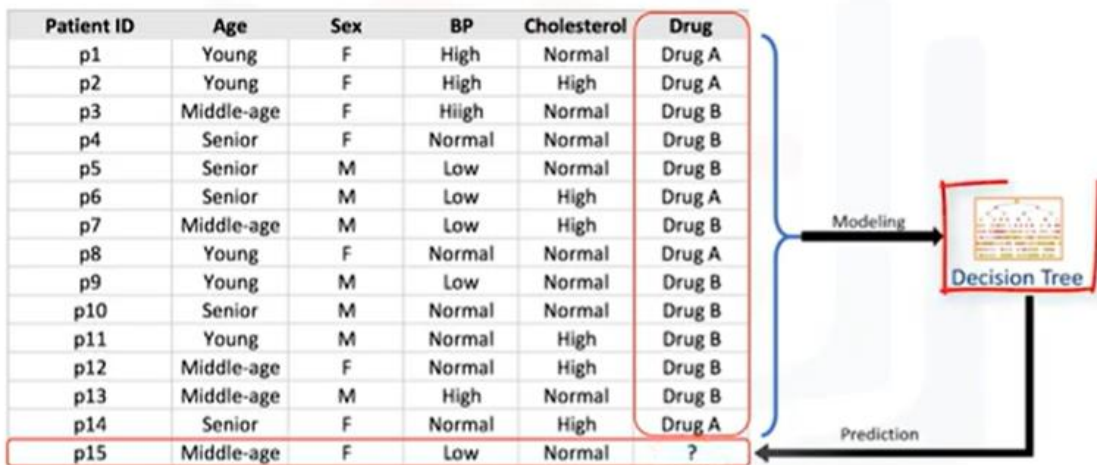
Performance of a classifier where the predicted output is a probability value between 0 and 1.

Test set

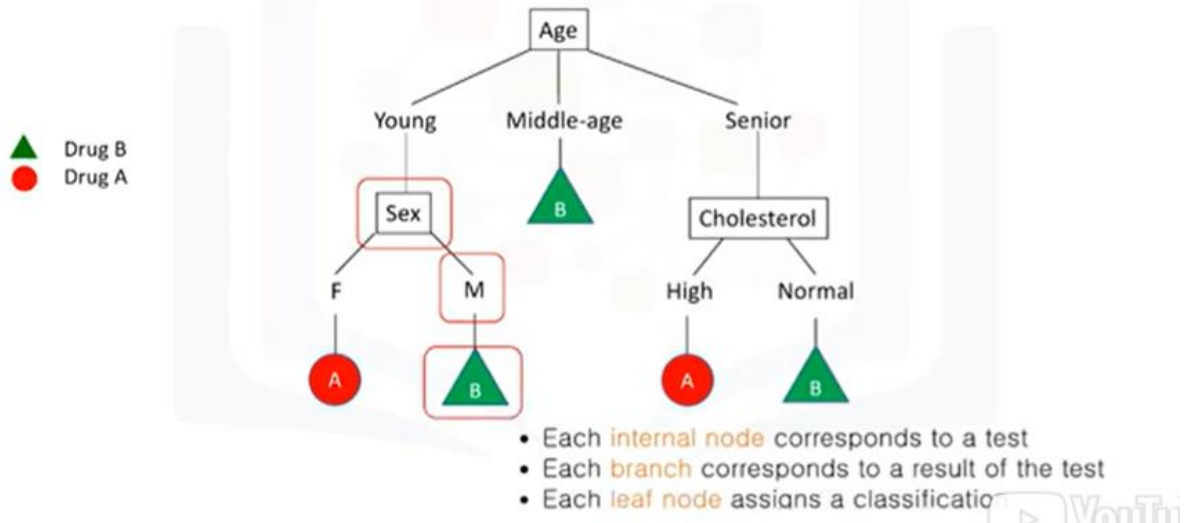


Introduction to Decision Tree.

How to build a decision tree?

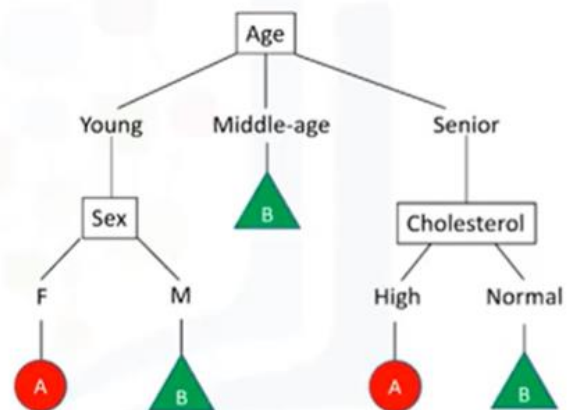


Building a decision tree with the training set



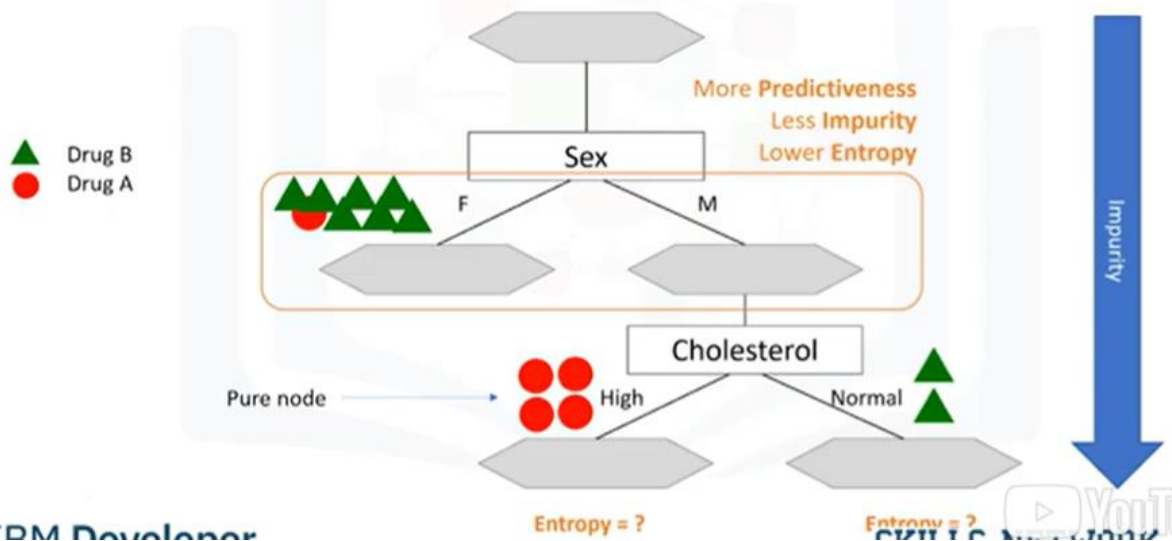
Decision tree learning algorithm

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.



Building Decision Tree:

Which attribute is the best ?

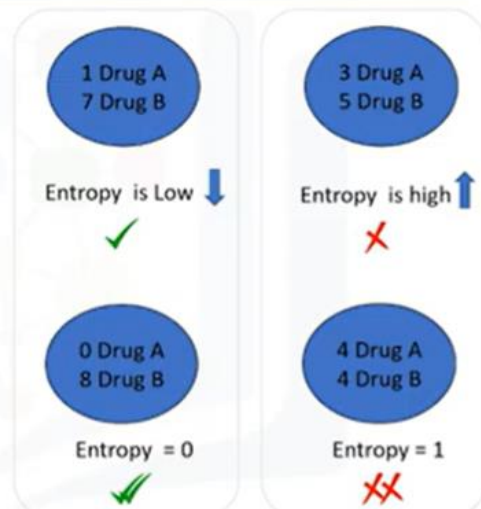


Entropy

- Measure of randomness or uncertainty

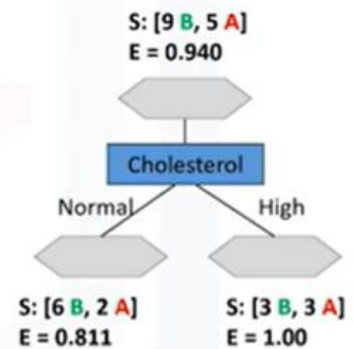
$$\text{Entropy} = -p(A)\log(p(A)) - p(B)\log(p(B))$$

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



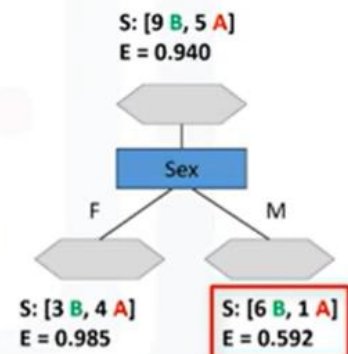
Is 'Cholesterol' the best attribute?

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

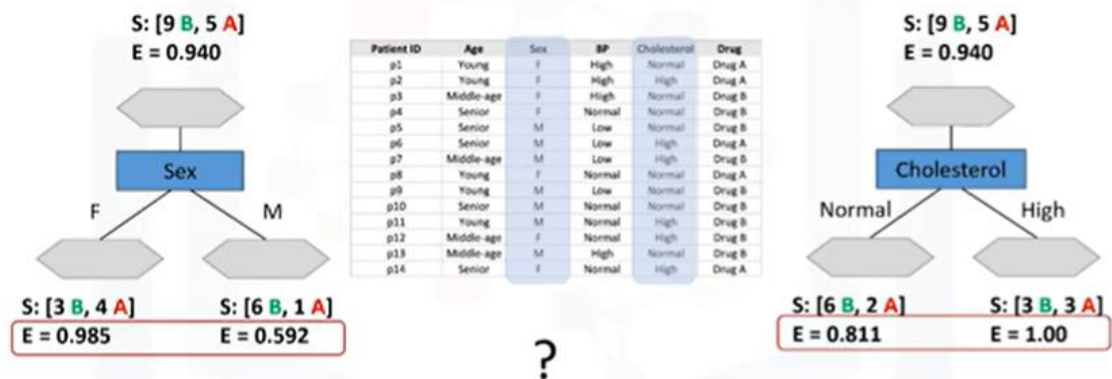


What about 'Sex'?

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?

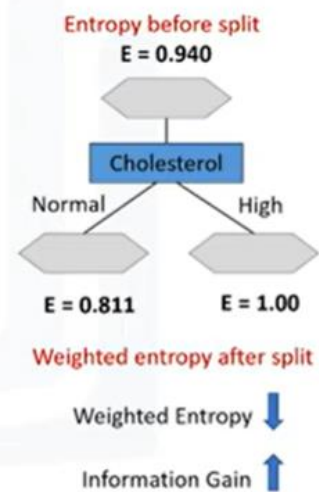


The tree with the higher **Information Gain** after splitting.

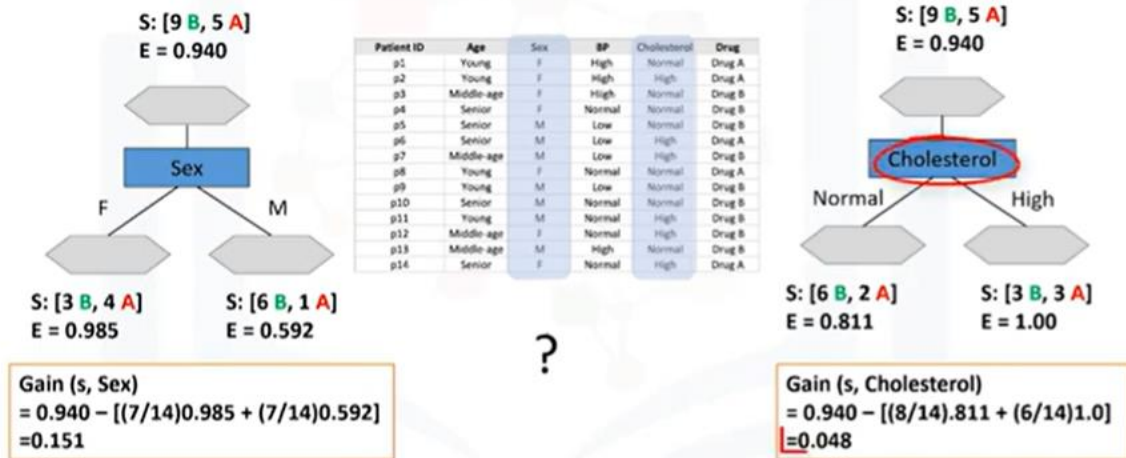
What is information gain?

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

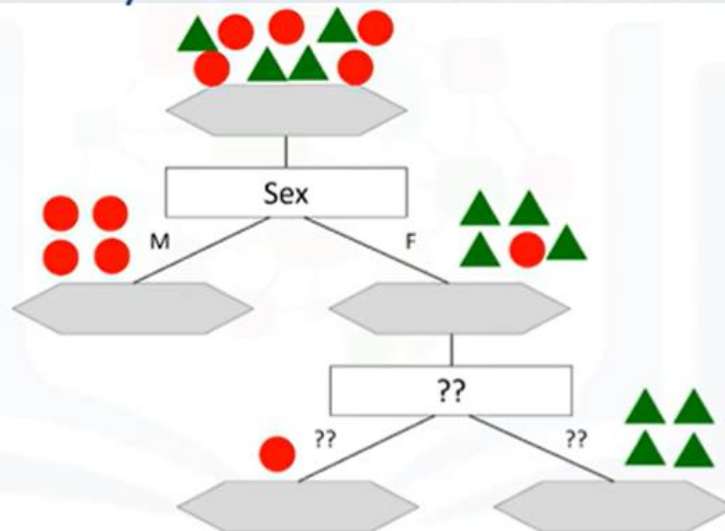
$$\text{Information Gain} = (\text{Entropy before split}) - (\text{weighted entropy after split})$$



Which attribute is the best?



Correct way to build a decision tree



Introduction to Logistic Regression:

What is logistic regression?

Logistic regression is a classification algorithm for categorical variables.

Independent variables										Dependent variable
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	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

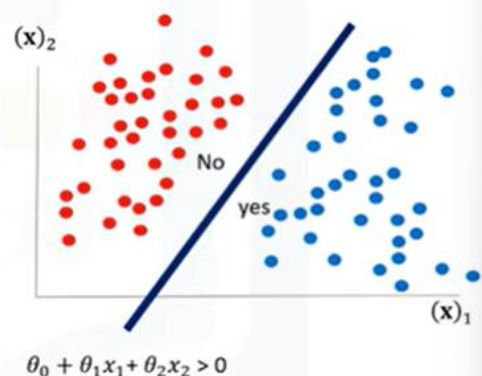
Categorical Variable

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.

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	calcard	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

Sigmoid function in logistic regression

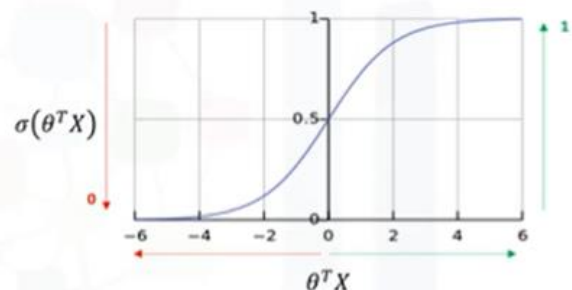
• Logistic Function

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

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

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

$$[0, 1]$$



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



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

Clarification of the customer churn model

What is the output of our model?

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

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

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

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.

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

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

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

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

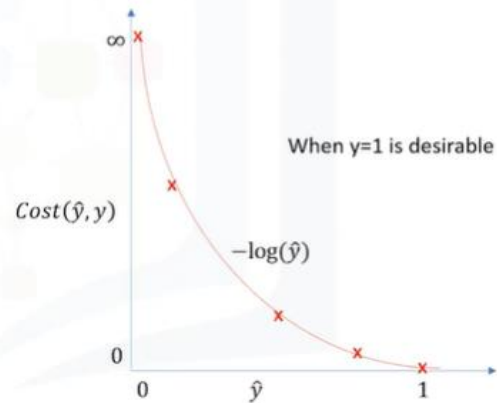
$$\text{Cost} = J(\theta)$$

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

Logistic Regression Training:

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 $\hat{y}=0 \rightarrow \text{cost} = \text{large}$



Logistic regression cost function

- So, we will replace cost function with:

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

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

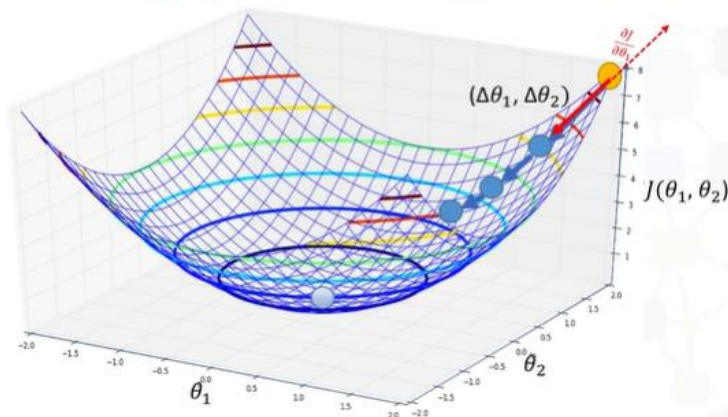
$$Cost(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

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

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



$$\frac{\partial J}{\partial \theta_1} = -\frac{1}{m} \sum_{i=1}^m (y^i - \hat{y}^i) x_1^i$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \frac{\partial J}{\partial \theta_3} \\ \dots \\ \frac{\partial J}{\partial \theta_k} \end{bmatrix}$$

$$\text{New } \theta = \text{old } \theta - \eta \nabla J$$

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

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

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.

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

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

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

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

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

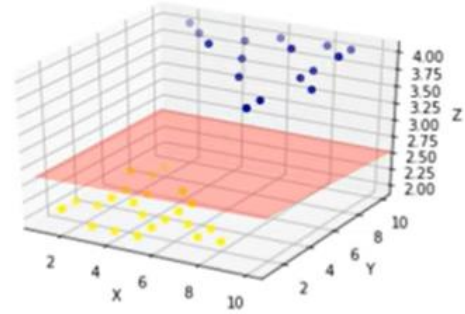
Support Vector Machines:

What is SVM?

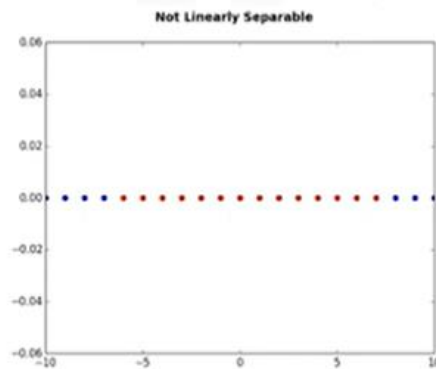
SVM is a supervised algorithm that classifies cases by finding a separator.

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

Clump	UnitSize	UnitShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	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

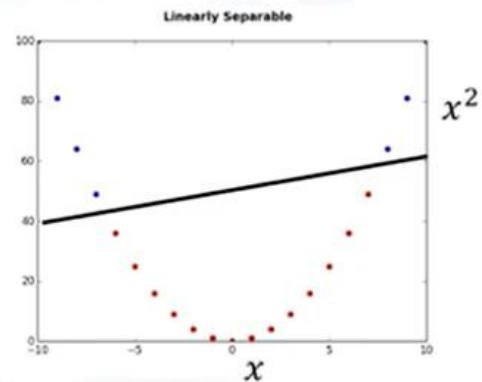


Data transformation



Kernelling:

- Linear
- Polynomial
- RBF
- Sigmoid

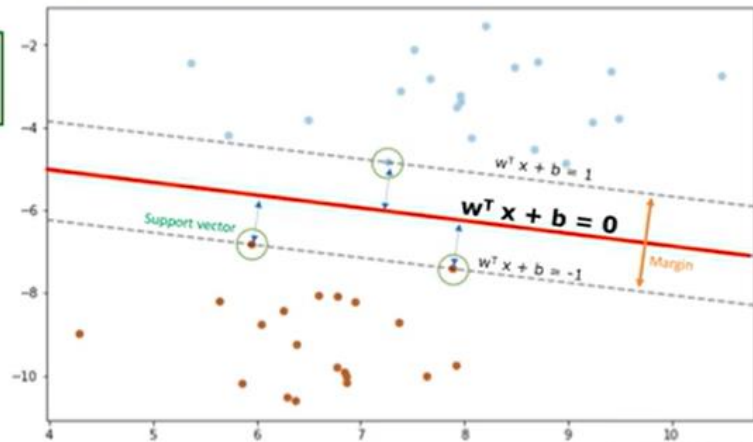


$$\phi(x) = [x, x^2]$$



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 $\{(x_i, y_i)\}$: $y_i (w^T x_i + b) \geq 1$



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