

Supervised Learning

Amir Hosein Hadian

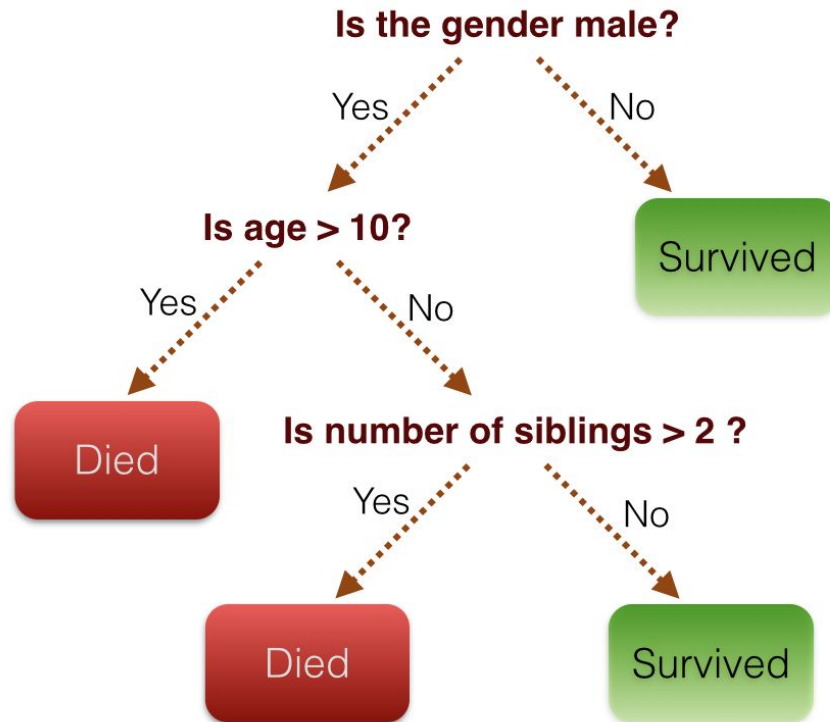
Outline

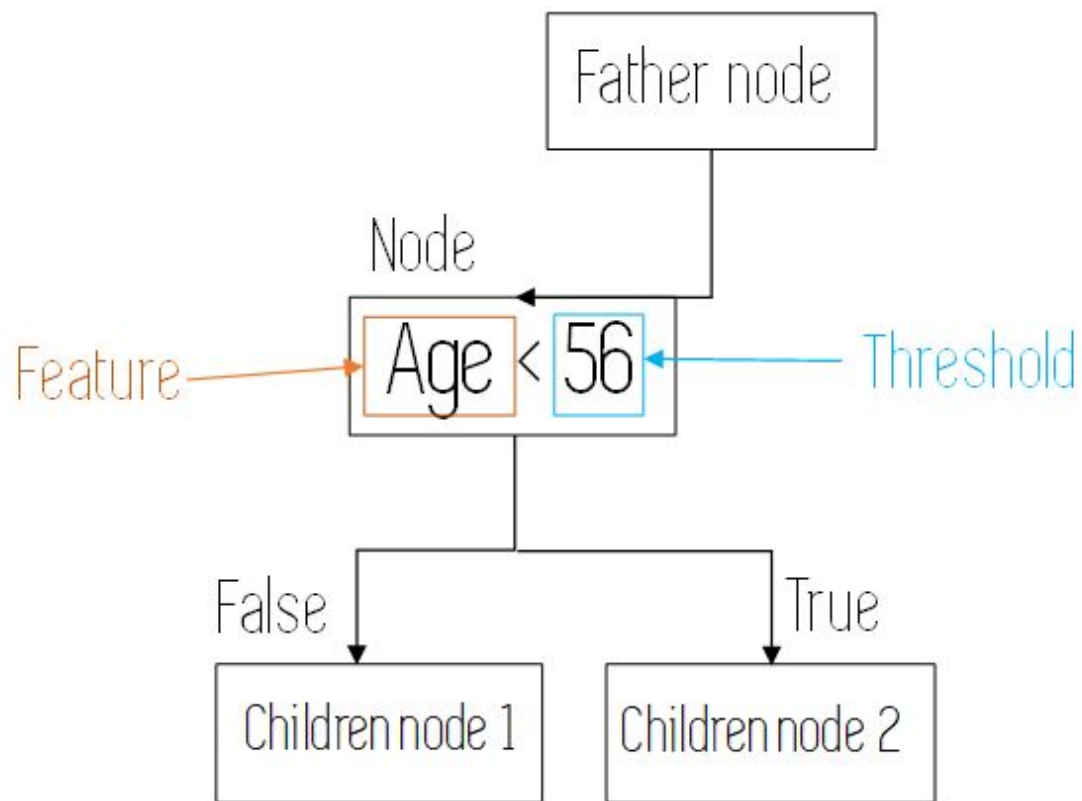
Decision tree

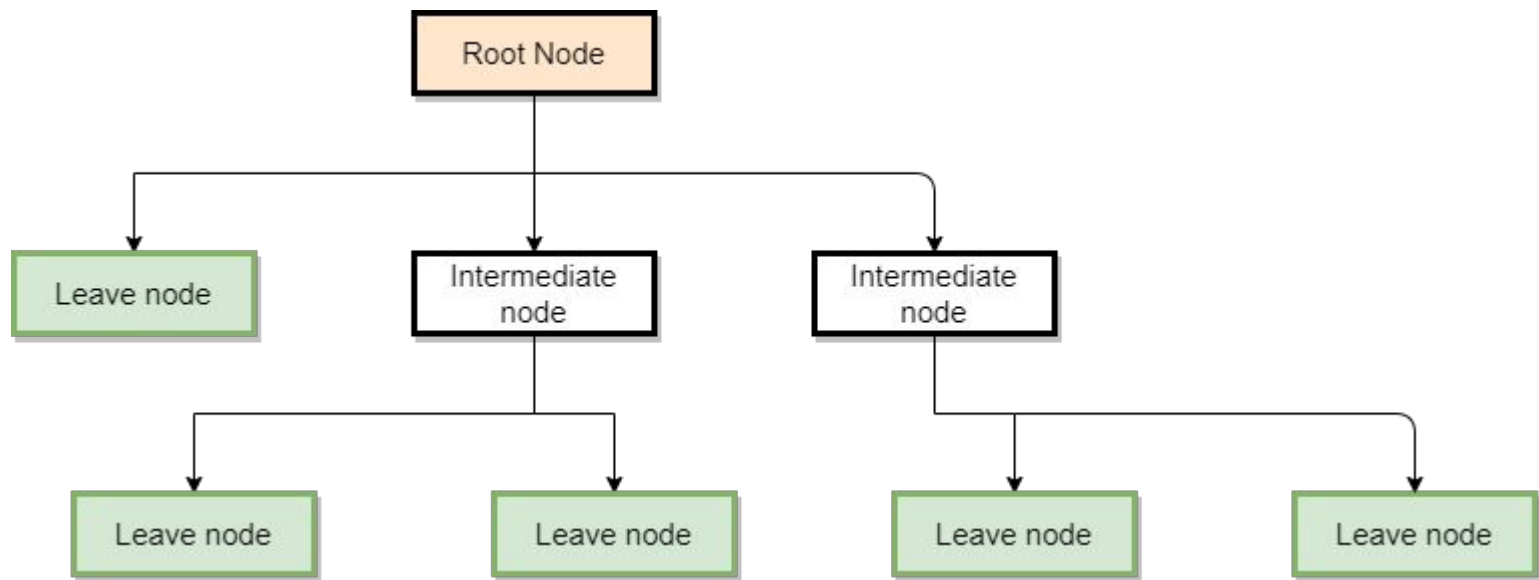
Random Forest

Model evaluation









- How are decision trees used for classification?
- Why are decision tree classifiers so popular?

Algorithm: Generate_decision_tree. Generate a decision tree from the training tuples of data partition, D .

Input:

- Data partition, D , which is a set of training tuples and their associated class labels;
- *attribute_list*, the set of candidate attributes;
- *Attribute_selection_method*, a procedure to determine the splitting criterion that “best” partitions the data tuples into individual classes. This criterion consists of a *splitting_attribute* and, possibly, either a *split-point* or *splitting_subset*.

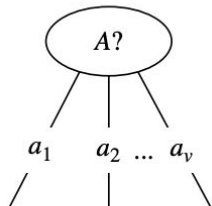
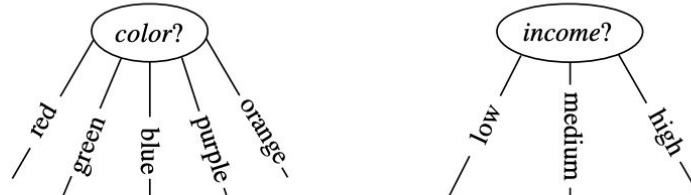
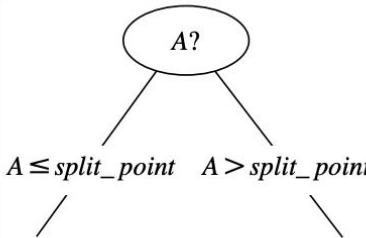
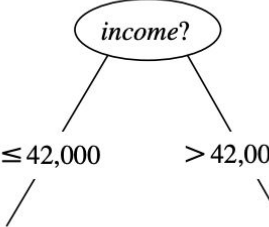
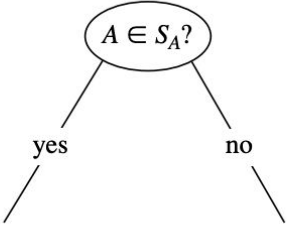
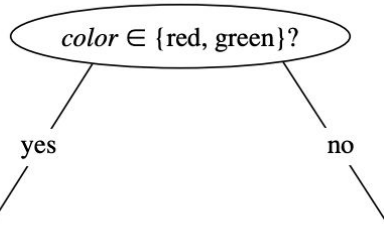
Output: A decision tree.

Method:

- (1) create a node N ;
- (2) **if** tuples in D are all of the same class, C , **then**
- (3) return N as a leaf node labeled with the class C ;
- (4) **if** *attribute_list* is empty **then**
- (5) return N as a leaf node labeled with the majority class in D ; // majority voting
- (6) apply **Attribute_selection_method**(D , *attribute_list*) to **find** the “best” *splitting_criterion*;
- (7) label node N with *splitting_criterion*;
- (8) **if** *splitting_attribute* is discrete-valued **and**
 multiway splits allowed **then** // not restricted to binary trees
- (9) *attribute_list* \leftarrow *attribute_list* – *splitting_attribute*; // remove *splitting_attribute*
- (10) **for each** outcome j of *splitting_criterion*
 // partition the tuples and grow subtrees for each partition
- (11) let D_j be the set of data tuples in D satisfying outcome j ; // a partition
- (12) **if** D_j is empty **then**
- (13) attach a leaf labeled with the majority class in D to node N ;
- (14) **else** attach the node returned by **Generate_decision_tree**(D_j , *attribute_list*) to node N ;
- endfor**
- (15) return N ;

Partitioning scenarios

Examples

<p>(a)</p>  <pre> graph TD A((A?)) --> a1[a1] A --> a2["a2 ..."] A --> av[av] </pre>	 <pre> graph TD color((color?)) --> red[red] color --> green[green] color --> blue[blue] color --> purple[purple] color --> orange[orange] income((income?)) --> low[low] income --> medium[medium] income --> high[high] </pre>
<p>(b)</p>  <pre> graph TD A((A?)) --> B["A ≤ split_point"] A --> C["A > split_point"] </pre>	 <pre> graph TD income((income?)) --> B["≤ 42,000"] income --> C[">> 42,000"] </pre>
<p>(c)</p>  <pre> graph TD A("A ∈ S_A?") --> yes[yes] A --> no[no] </pre>	 <pre> graph TD A("color ∈ {red, green}?") --> yes[yes] A --> no[no] </pre>

Information Gain

$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i),$$

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j).$$

$$Gain(A) = Info(D) - Info_A(D).$$

Class-Labeled Training Tuples from the *AllElectronics* Customer Database

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

$$Info(D) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.940$$

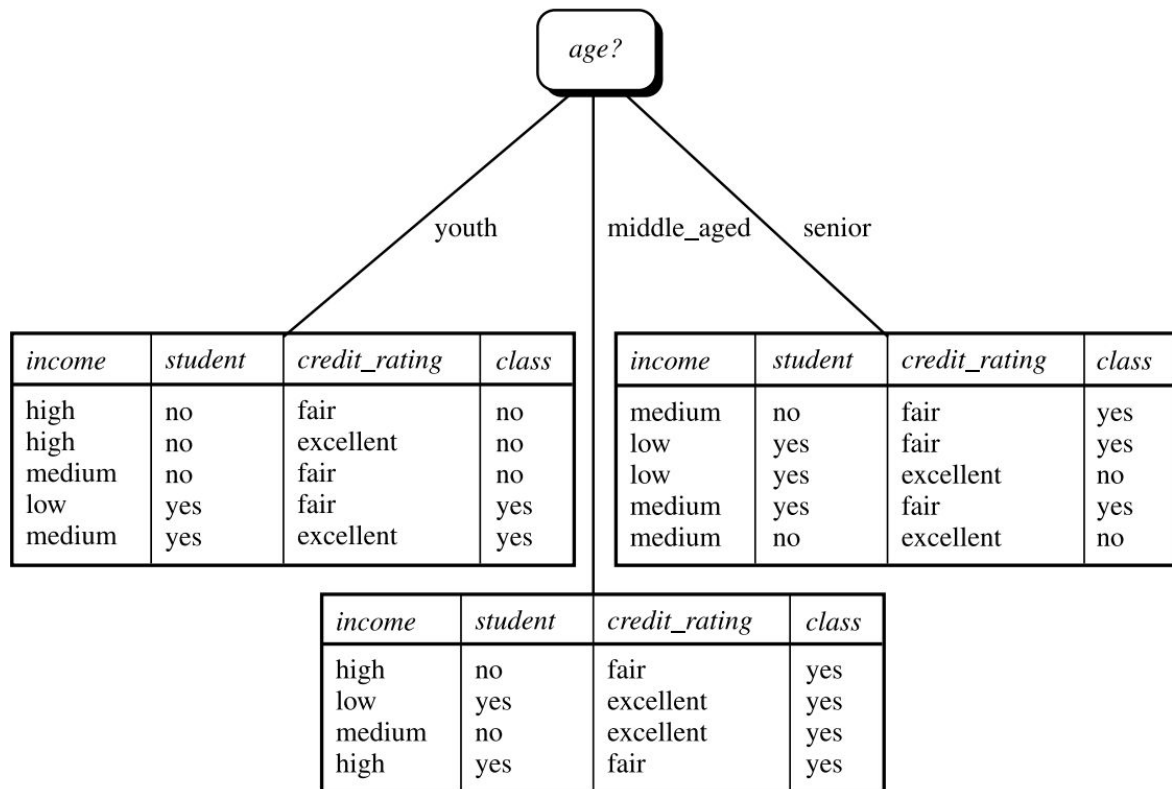
$$\begin{aligned} Info_{age}(D) &= \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) \\ &\quad + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\ &= 0.694 \end{aligned}$$

$$Gain(age) = Info(D) - Info_{age}(D) = 0.940 - 0.694 = 0.246$$

$$\textit{Gain}(\textit{income}) = 0.029$$

$$\textit{Gain}(\textit{student}) = 0.151$$

$$\textit{Gain}(\textit{credit_rating}) = 0.048$$



$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right).$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo_A(D)}.$$

$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2 \left(\frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left(\frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left(\frac{4}{14} \right)$$

$$= 1.557.$$

$$Gain(income) = 0.029.$$

$$GainRatio(income) = 0.029/1.557 = 0.019.$$

Gini Index

$$Gini(D) = 1 - \sum_{i=1}^m p_i^2,$$

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2).$$

$$\Delta Gini(A) = Gini(D) - Gini_A(D).$$

$$Gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459.$$

$$Gini_{income \in \{low, medium\}}(D)$$

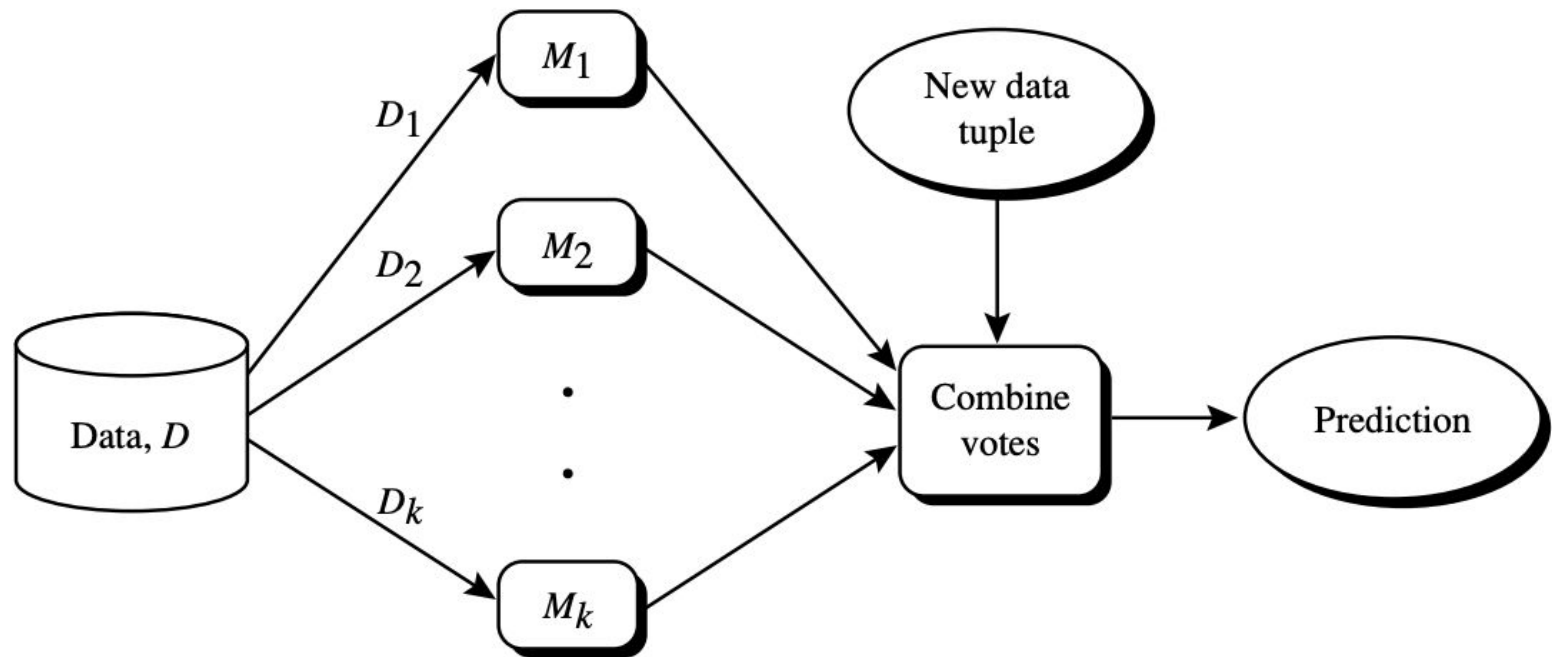
$$= \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

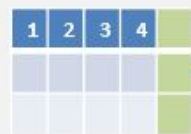
$$= \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 \right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right)$$

$$= 0.443$$

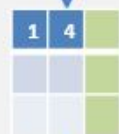
$$= Gini_{income \in \{high\}}(D).$$



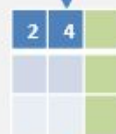
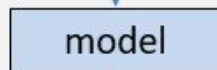




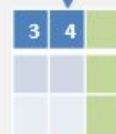
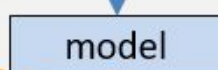
training



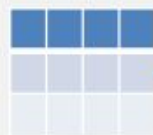
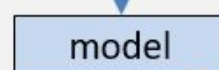
Bootstrap sample



Bootstrap sample



Bootstrap sample



new

Yes

No

Yes

Yes

Bagging

Algorithm: Bagging. The bagging algorithm—create an ensemble of classification models for a learning scheme where each model gives an equally weighted prediction.

Input:

- D , a set of d training tuples;
- k , the number of models in the ensemble;
- a classification learning scheme (decision tree algorithm, naïve Bayesian, etc.).

Output: The ensemble—a composite model, M_* .

Method:

- (1) **for** $i = 1$ to k **do** // create k models:
- (2) create bootstrap sample, D_i , by sampling D with replacement;
- (3) use D_i and the learning scheme to derive a model, M_i ;
- (4) **endfor**

To use the ensemble to classify a tuple, X :

let each of the k models classify X and return the majority vote;

Model Evaluation





		ACTUAL VALUES	
		Positive	Negative
PREDICTED VALUES	Positive	TP	FP
	Negative	FN	TN

The predicted value is positive and its positive

Type I error :
The predicted value is positive but it False

Type II error :
The predicted value is negative but its positive

The predicted value is Negative and its Negative

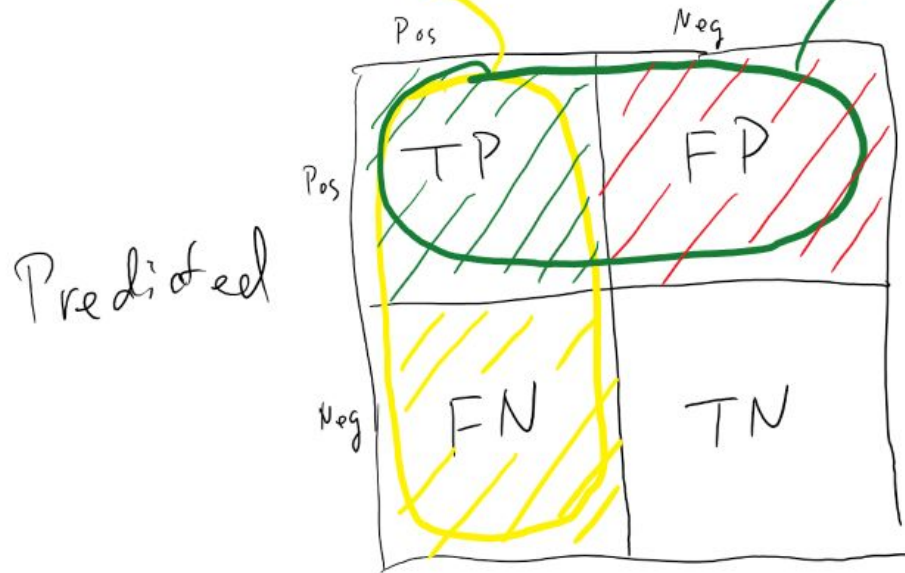
		Actual Values	
		1	0
Predicted Values	1	TRUE POSITIVE  You're pregnant	FALSE POSITIVE  You're pregnant TYPE 1 ERROR
	0	FALSE NEGATIVE  You're not pregnant TYPE 2 ERROR	TRUE NEGATIVE  You're not pregnant

Model Evaluation





$$\text{Recall} = \frac{TP}{TP + FN} \quad \text{Actual}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Accuracy} = \frac{TP + TN}{\text{Total}}$$



Model Evaluation

		PREDICTED VALUES	
		Positive (CAT)	Negative (DOG)
ACTUAL VALUES	Positive (CAT)	 <p>TRUE POSITIVE</p> <p>6</p> <p>YOU ARE A CAT</p>	 <p>FALSE NEGATIVE</p> <p>1</p> <p>TYPE II ERROR</p> <p>YOU ARE A DOG</p>
	Negative (DOG)	 <p>FALSE POSITIVE</p> <p>2</p> <p>TYPE I ERROR</p> <p>YOU ARE A CAT</p>	 <p>TRUE NEGATIVE</p> <p>11</p> <p>YOU ARE NOT A CAT</p>

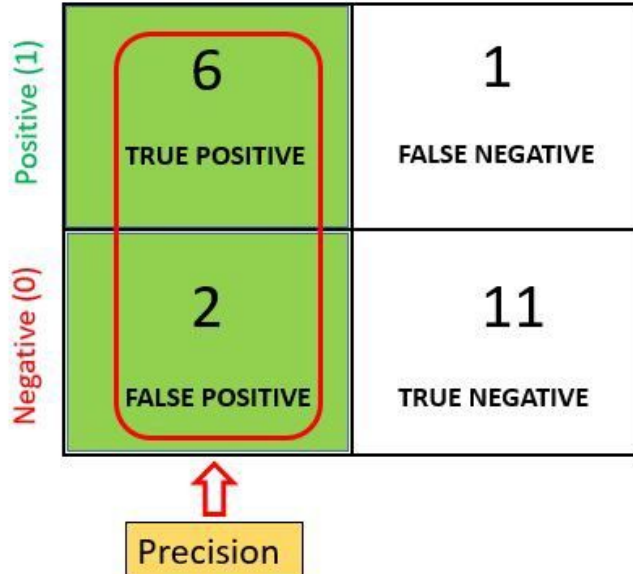
Model Evaluation

		PREDICTED VALUES	
		Positive (1)	Negative (0)
ACTUAL VALUES	Positive (1)	6 TRUE POSITIVE	1 FALSE NEGATIVE
	Negative (0)	2 FALSE POSITIVE	11 TRUE NEGATIVE

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} = \frac{6 + 11}{6 + 11 + 2 + 1} = 85\%$$

Accuracy


		PREDICTED VALUES	
		Positive (1)	Negative (0)
ACTUAL VALUES	Positive (1)	<div>6</div> <div>TRUE POSITIVE</div>	<div>1</div> <div>FALSE NEGATIVE</div>
	Negative (0)	<div>2</div> <div>FALSE POSITIVE</div>	<div>11</div> <div>TRUE NEGATIVE</div>



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{Predictions Actually Positive}}{\text{Total Predicted positive}} = \frac{6}{6 + 2} = 0.75$$

Model Evaluation

		PREDICTED VALUES	
		Positive (1)	Negative (0)
ACTUAL VALUES	Positive (1)	6 TRUE POSITIVE	1 FALSE NEGATIVE
	Negative (0)	2 FALSE POSITIVE	11 TRUE NEGATIVE



$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{Predictions Actually Positive}}{\text{Total Actual positive}} = \frac{6}{6 + 1} = 0.85$$

Model Evaluation

$$\text{F1-Score} = 2 * \frac{(\text{Recall} * \text{Precision})}{(\text{Recall} + \text{Precision})} = 2 * \frac{(0.85 * 0.75)}{(0.85 + 0.75)} = 0.79$$

Model Evaluation

