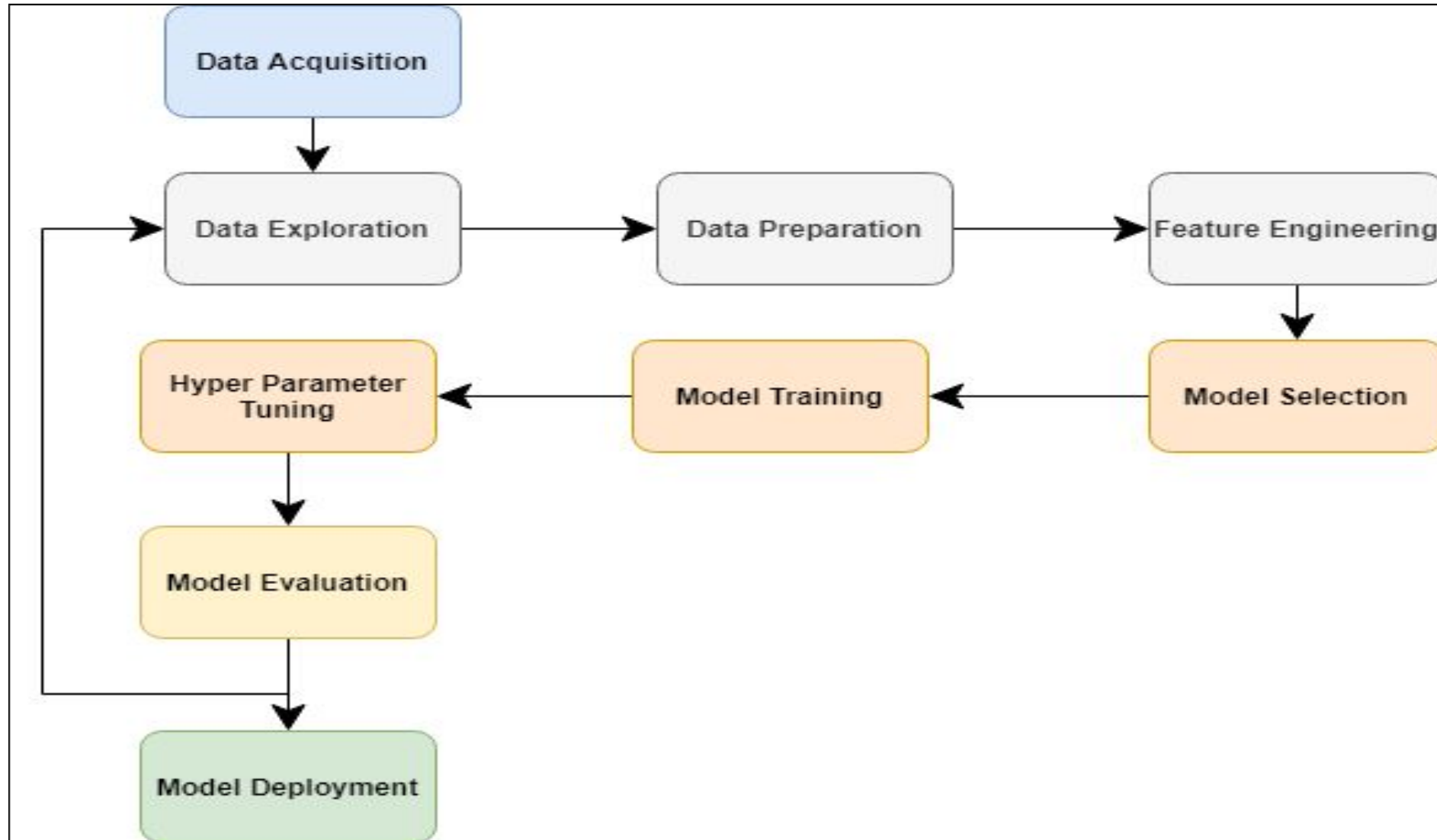
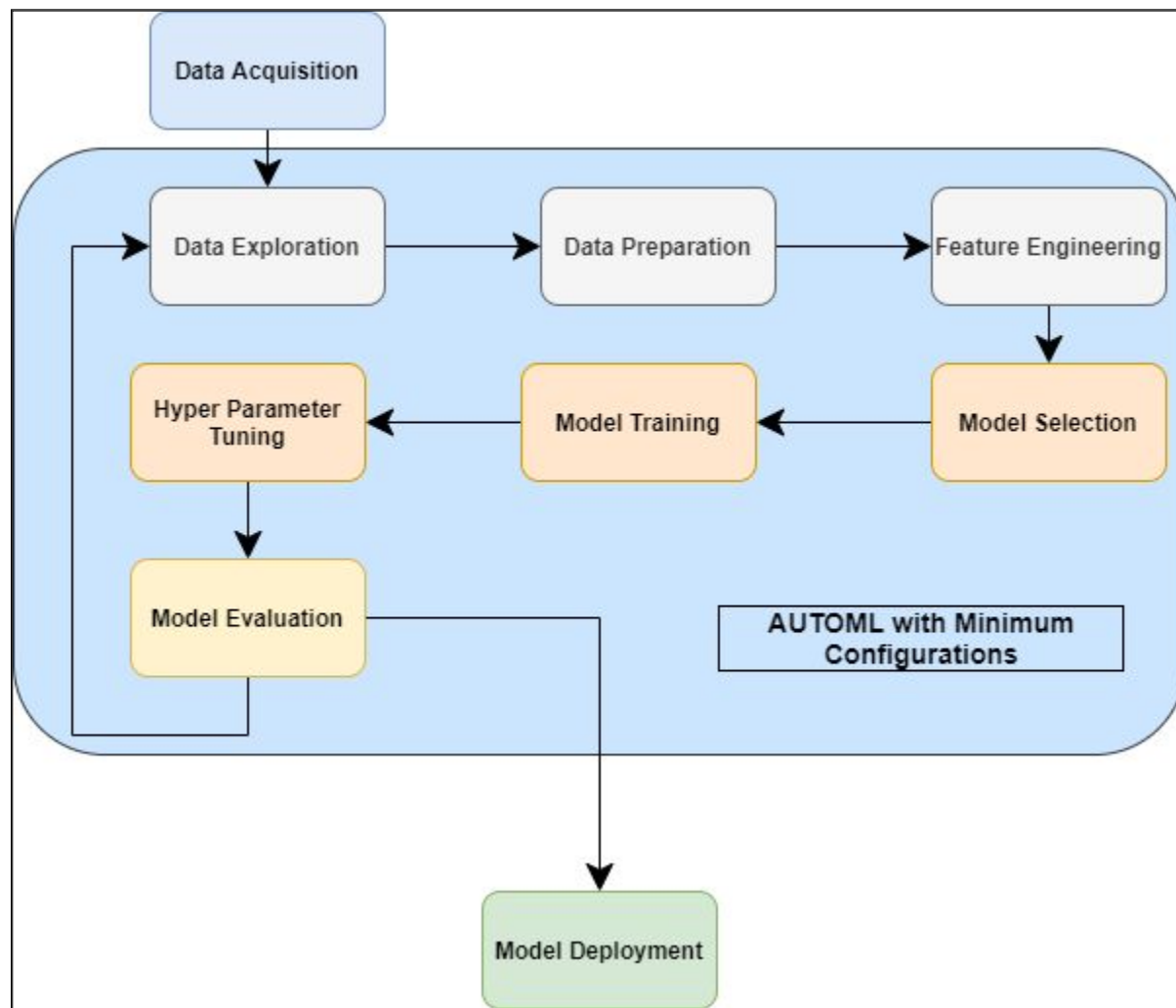


Supervised Learning

Data Mining Process



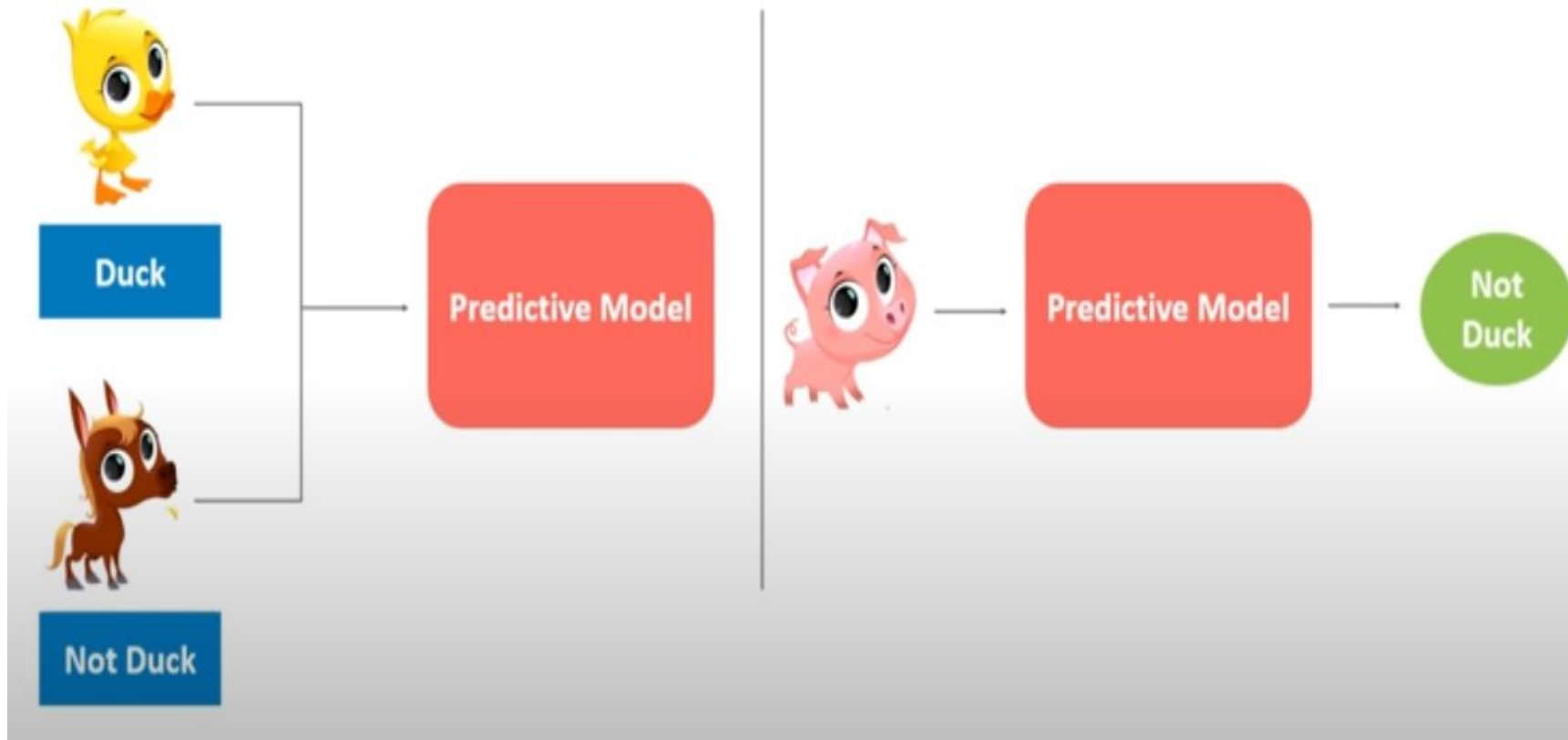
AutoML



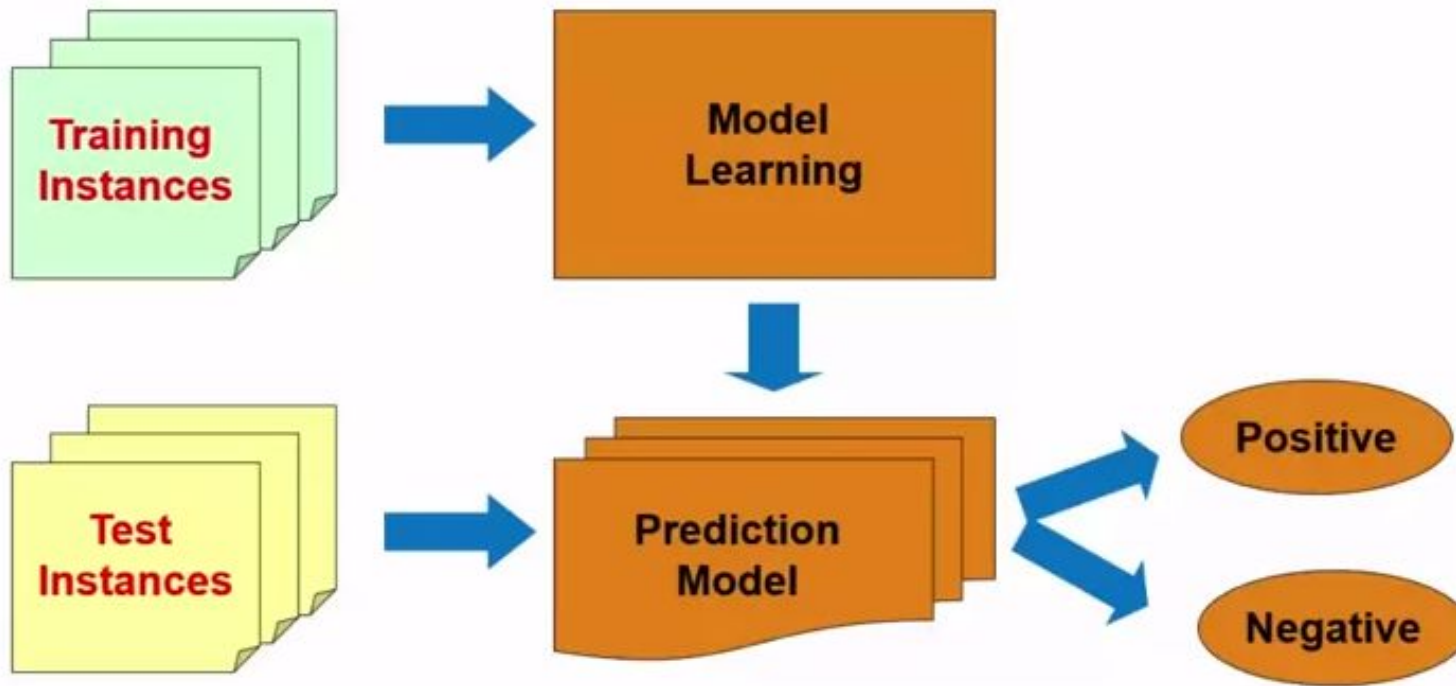
Supervised Learning

- The computer is provided with example inputs that are labeled with their desired outputs.
- Supervised learning uses patterns to predict label values on additional unlabeled data.
- An algorithm may be fed data with images of sharks labeled as fish and images of oceans labeled as water.
- By being trained on this data, the supervised learning algorithm should be able to later identify unlabeled shark images as fish and unlabeled ocean images as water.

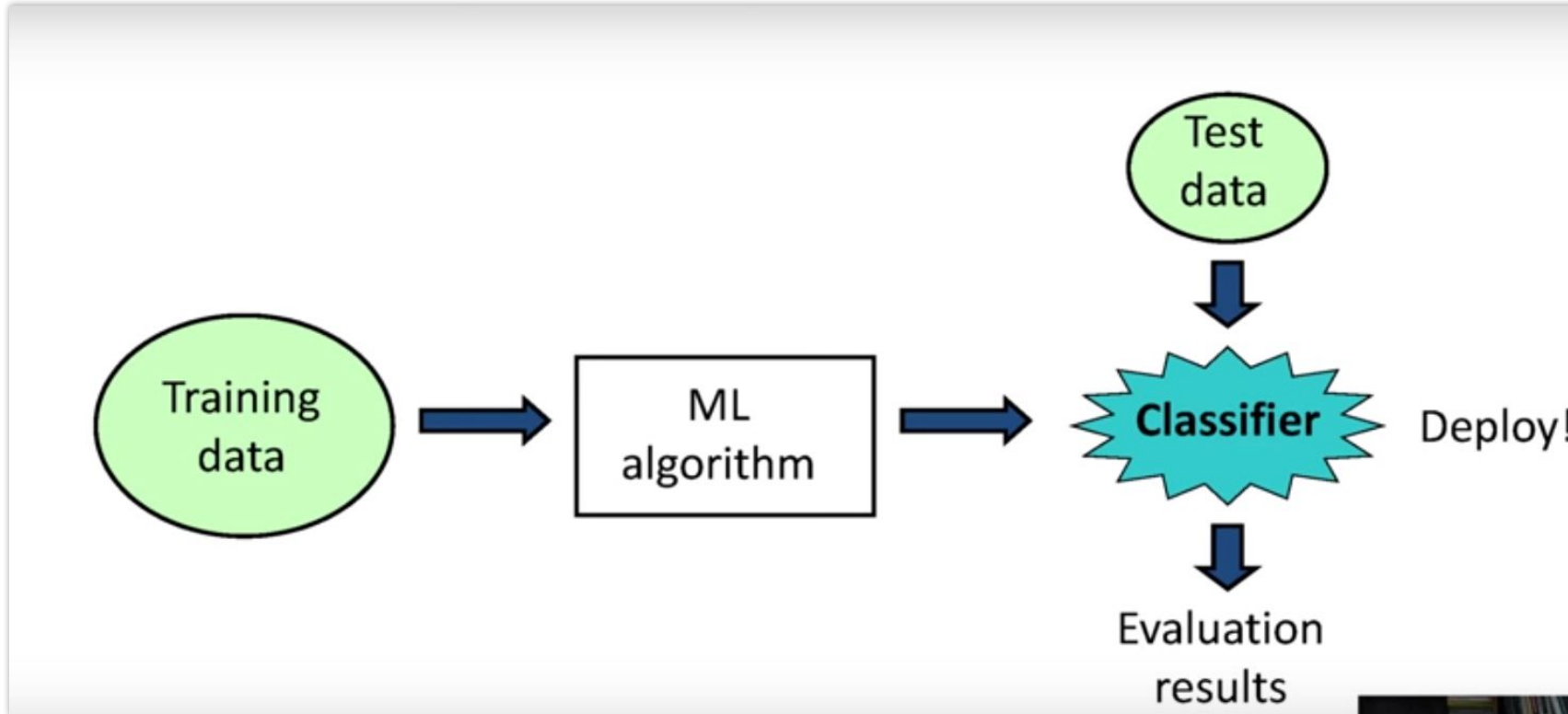
Supervised Learning



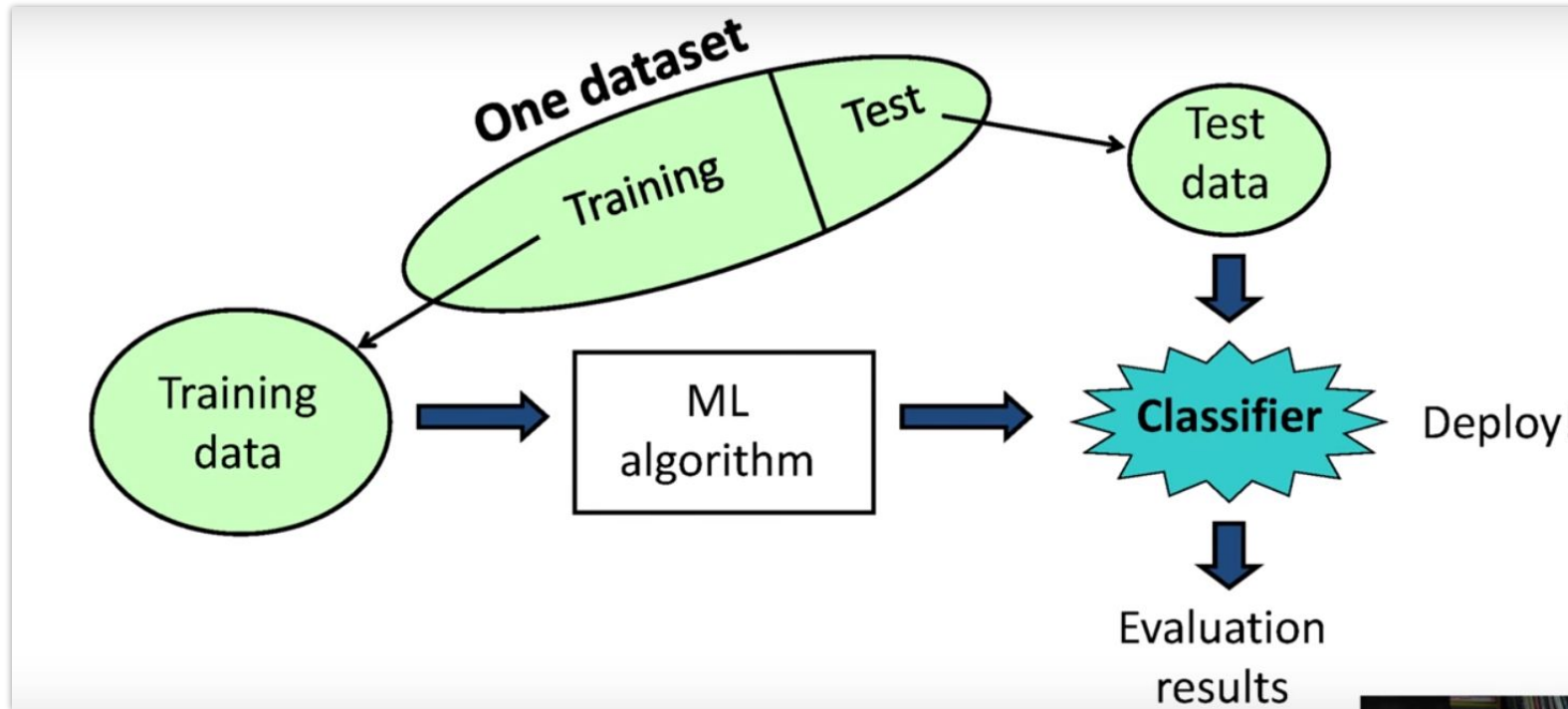
What is Classification



Training & Testing

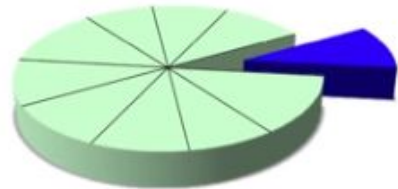
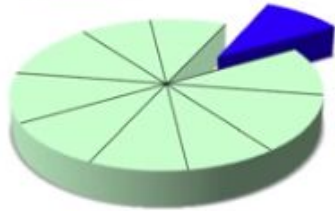
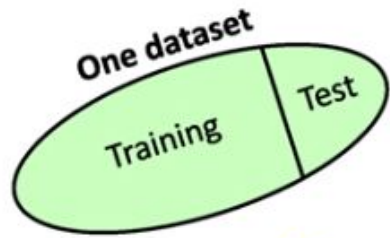


Percentage Split



Cross Validation

- ❖ Repeated holdout
(in Lesson 2.3, hold out 10% for testing, repeat 10 times)

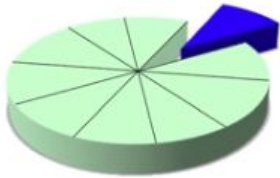


(repeat 10 times)

Cross Validation

10-fold cross-validation

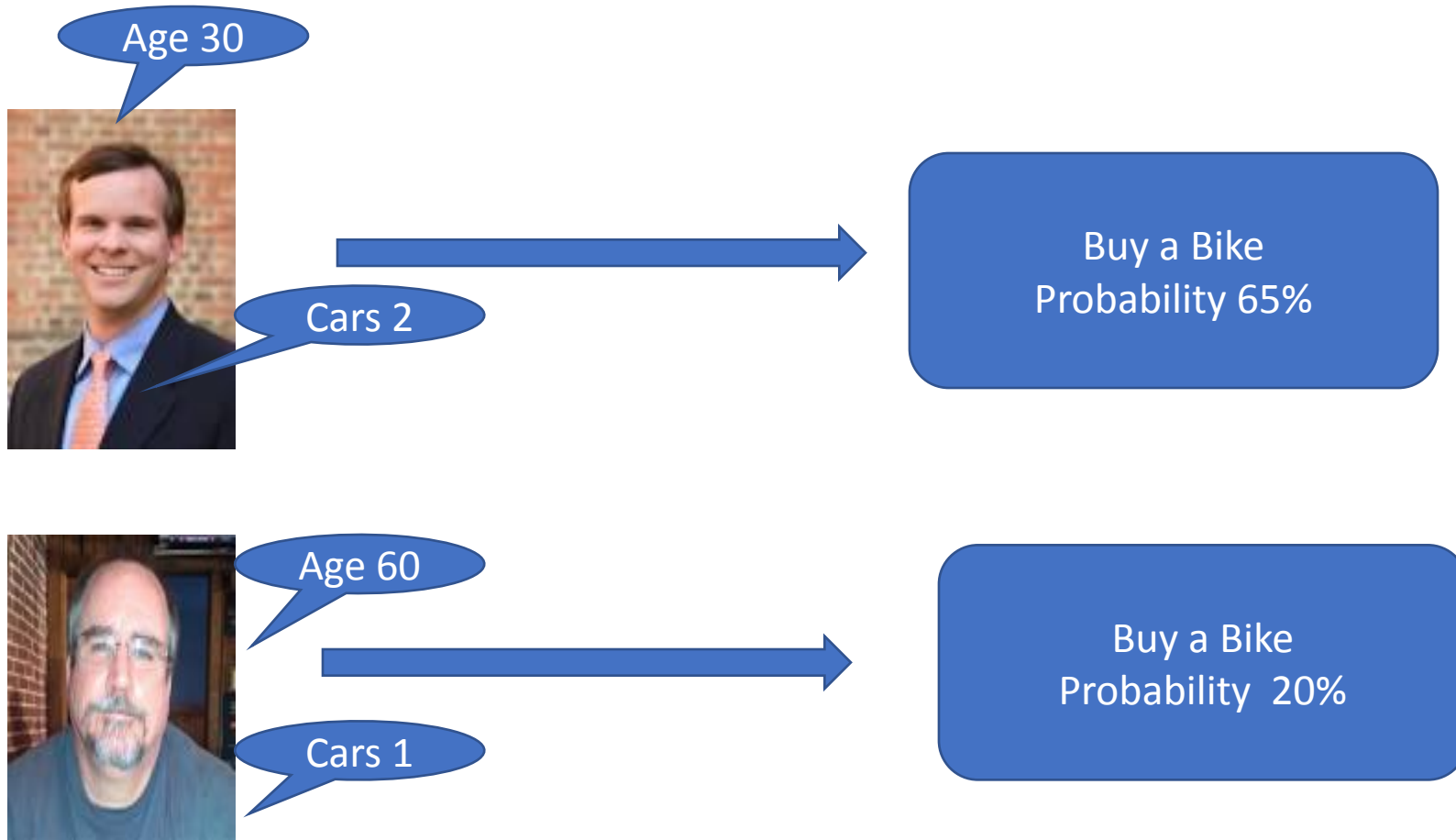
- ❖ Divide dataset into 10 parts (folds)
- ❖ Hold out each part in turn
- ❖ Average the results
- ❖ Each data point used once for testing, 9 times for training



Stratified cross-validation

- ❖ Ensure that each fold has the right proportion of each class value

In Simple Terms?





Iris-setosa



Iris-versicolor



Iris-virginica

		True condition			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$	Accuracy (ACC) = $\frac{\Sigma \text{True positive} + \Sigma \text{True negative}}{\Sigma \text{Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\Sigma \text{True positive}}{\Sigma \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{False positive}}{\Sigma \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\Sigma \text{False negative}}{\Sigma \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\Sigma \text{True negative}}{\Sigma \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\Sigma \text{True positive}}{\Sigma \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\Sigma \text{False positive}}{\Sigma \text{Condition negative}}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$ $F_1 \text{ score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
		False negative rate (FNR), Miss rate $= \frac{\Sigma \text{False negative}}{\Sigma \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\Sigma \text{True negative}}{\Sigma \text{Condition negative}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$	

Accuracy of NB

- Error Rate

- $(T - C) / T$

- T – Number of Objects

- C – correctly classified objects.

- Confusion matrix

Predicated Class	True Class		
	A	B	C
A	8	1	1
B	2	9	2
C	0	0	7

Classification Performance Measures

- Recall = $a/(a+c)$ where $a + c > 0$ (o.w. undefined).
 - Did we find all of those that belonged in the class?
- Precision = $a/(a+b)$ where $a+b>0$ (o.w. undefined).
 - Of the times we predicted it was “in class”, how often are we correct?

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

F1 Score / F1 Measure

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Matthews Correlation Coefficient

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

The MCC doesn't depend on which class is the positive one, which has the advantage over the F1 score to avoid incorrectly defining the positive class.

- False Positive (FP)

- FP cases are those that did not belong to a class but allocated to it.

- False Negative (FN)

- FN are cases that belong to a class but were not allocated to it.

	FP	FN
Class 1	2	2
Class 2	4	1
Class 3	0	3

- Sensitivity = $TP / (TP + FN)$
- Specificity = $TN / (TN + FP)$

Decision Trees

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

