

# COMP338: ARTIFICIAL INTELLIGENCE

---

Accuracy Measure

Dr. Radi Jarrar  
Department of Computer Science  
Birzeit University



# Model accuracy

- How does a generated model,  $m$ , perform on data from domain  $D$ ?
- Which of the generated models, in means of accuracy is best to select given some data from domain  $D$ ?
- How do models produced by some learning algorithm,  $\mathcal{A}$ , perform on data from domain  $D$ ?

## Model accuracy (2)

- There is a number of approaches that are used to measure the effectiveness of a classification algorithms
- These metrics are useful for evaluating experimental scenarios

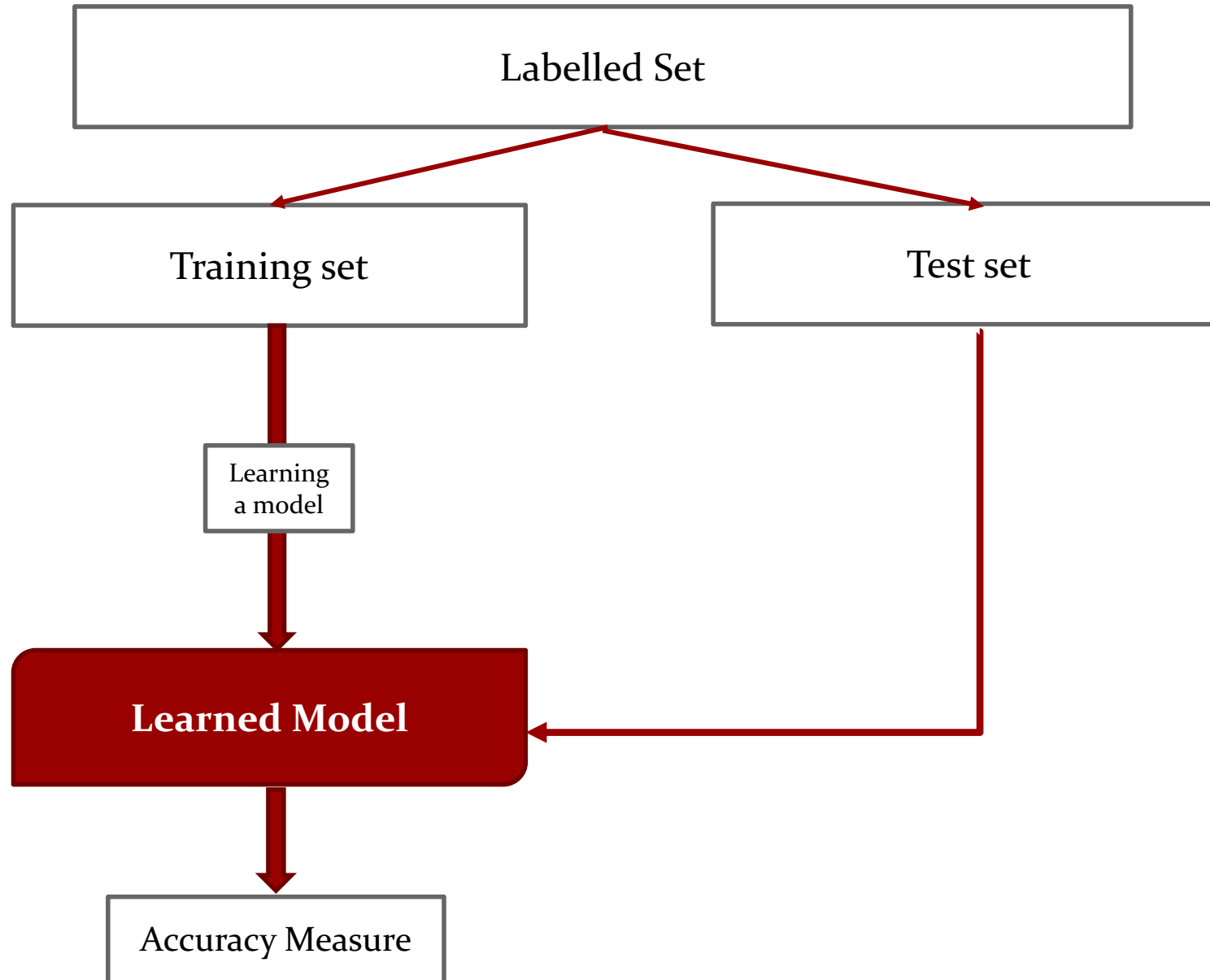
## Model accuracy (3)

### Regression

- Root Mean Squared Error (RMSE)
- R-Square

### Classification

- Accuracy
- Precision/Recall/F-score
- Learning Curve



# EVALUATING REGRESSION MODELS

---

# Evaluating Regression - RMSE

- Root Mean Square Error
- The sample standard deviation of the differences between predicted values and the actual outputs (i.e., the residuals)
- $RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (h(x)_i - y_i)^2}$
- The best metric for predicting accuracy for regression
- Simple and present as a default metric for most model

# R-Squared

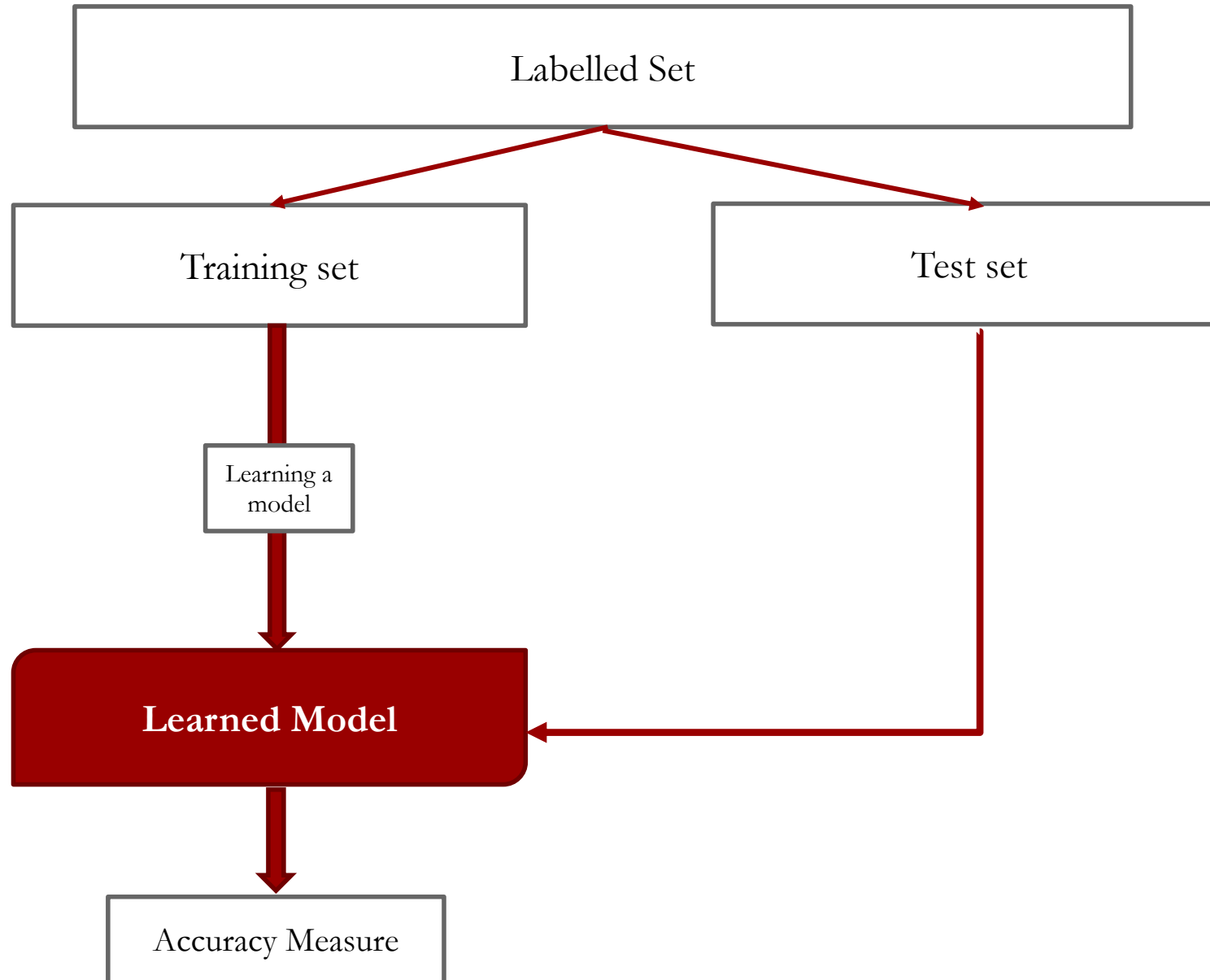
- Shows how well features fit a curve or line
- It represents the correlation between the observed outcomes and the predicted outcome values

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - h(x)_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{\frac{1}{N} \sum_{i=1}^N (y_i - h(x)_i)^2}{\frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2}$$

- Notice that the numerator is MSE (i.e., average of squares of the residuals)
- The denominator is the variance in y values
- The higher the MSE the poorer the model
- The higher the  $R^2$  the better the model

# EVALUATING CLASSIFICATION MODELS

---



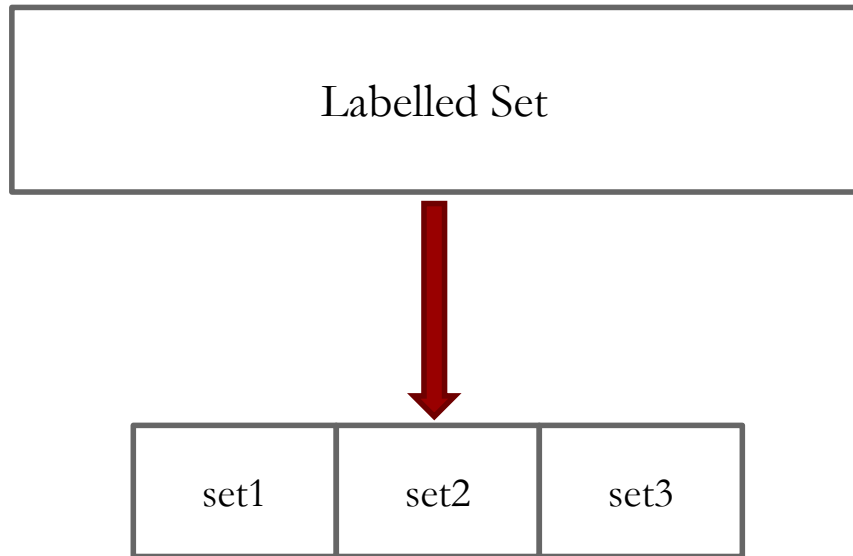
# Single Dataset?

- If there is a single dataset, or if the data is small, this will not tell how sensitive accuracy is to a particular training sample
- Larger datasets give better estimations on the accuracy of the model

# Cross Validation

- Cross-validation is a technique that is used to avoid overfitting (later in this course)
- In cross-validation, the training dataset is split into a number of folds (subsets) that are used to test the performance of the generated model while the training process is taking place
- Assume a training dataset of 900 records, It can be divided into 3-subsets each of around 300 records namely set1, set2, and set3
- 5-fold and 10-fold cross validations are widely used

## Cross Validation (2)



Iteration	Train-on	Test-on

## Cross Validation (3)

- E.g., suppose you have 90, using 3-fold cross-validation, estimate the accuracy such that:





Iteration	Train-on	Test-on	Correctly classified
1	set1, set2	set3	18/30
2	set2, set3	set1	20/30
3	set1, set3	set2	22/30

- $\text{Accuracy} = 60 / 90 = 0.66 = 66\%$

# Confusion matrix

- The confusion matrix is a well-known method for classification systems
- It contains all information about the actual (the original class label) and the predicted classification assigned by the classification method
- Columns represent predictor's output while the rows represent the actual class labels

# Confusion matrix (2)

		Predicted Class	
		A	B
Actual Class	A		
	B		

# Confusion matrix (3)

		Predicted Class	
		Pos	Neg
Actual Class	Pos	<b>TP</b> True Positive	<b>FN</b> False Negative
	Neg	<b>FP</b> False Positive	<b>TN</b> True Negative

## Confusion matrix (5)

- **TP (True positive)** is the number of correct predictions that an instance is positive (classified as class of interest)
- **TN (True negative)** is the number of correct predictions that an instance is negative (not class of interest)
- **FP (False positive)** is the number of incorrect predictions that an instance is negative (incorrectly classified as class of interest)
- **FN (False negative)** is the number of incorrect predictions that an instance is positive (incorrectly classified as not a class of interest)

# Accuracy

- The confusion matrix is used to derive a number of performance metrics
- The Accuracy metric measures the proportion of the total number of predictions that were correctly classified
- Used to measure the overall effectiveness of a classifier

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

## Accuracy (2)

- Is accuracy always good to be used?
  - It is not the best choice when data is imbalanced (i.e., there is a skew in data towards one class)
    - E.g., Is 95% is good when 90% of the data is negative?
- Cost—Getting a positive wrong costs more than getting a negative wrong
  - E.g., in medical domain, false positives results in wrong tests; however, false negative results in a failure to treat a disease

# Error rate

- Is the proportion of incorrectly classified instances
- $\text{Error rate} = 1 - \text{Accuracy}$

# Precision

- Precision measure the accuracy such that a class has been predicted correctly
- Defines the proportion of positive examples that are correctly classified

$$\textit{Precision} = \frac{tp}{(tp + fp)}$$

# Recall

- Recall measures the completeness of the results (in this context, it is the also the true positive rate or sensitivity)
- It measures the proportion of positive examples that were correctly classified (from the dataset)
- High recall indicates a large portion of positive examples captured in the model

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

# F-score

- The F-score is a harmonic mean between the precision and recall
- It has the advantage that it combines both the precision and recall in a single value

$$F - score = \frac{2 \times tp}{2 \times tp + fp + fn}$$

# Learning Curves

- Learning curves show the effect of the datasets size and the accuracy

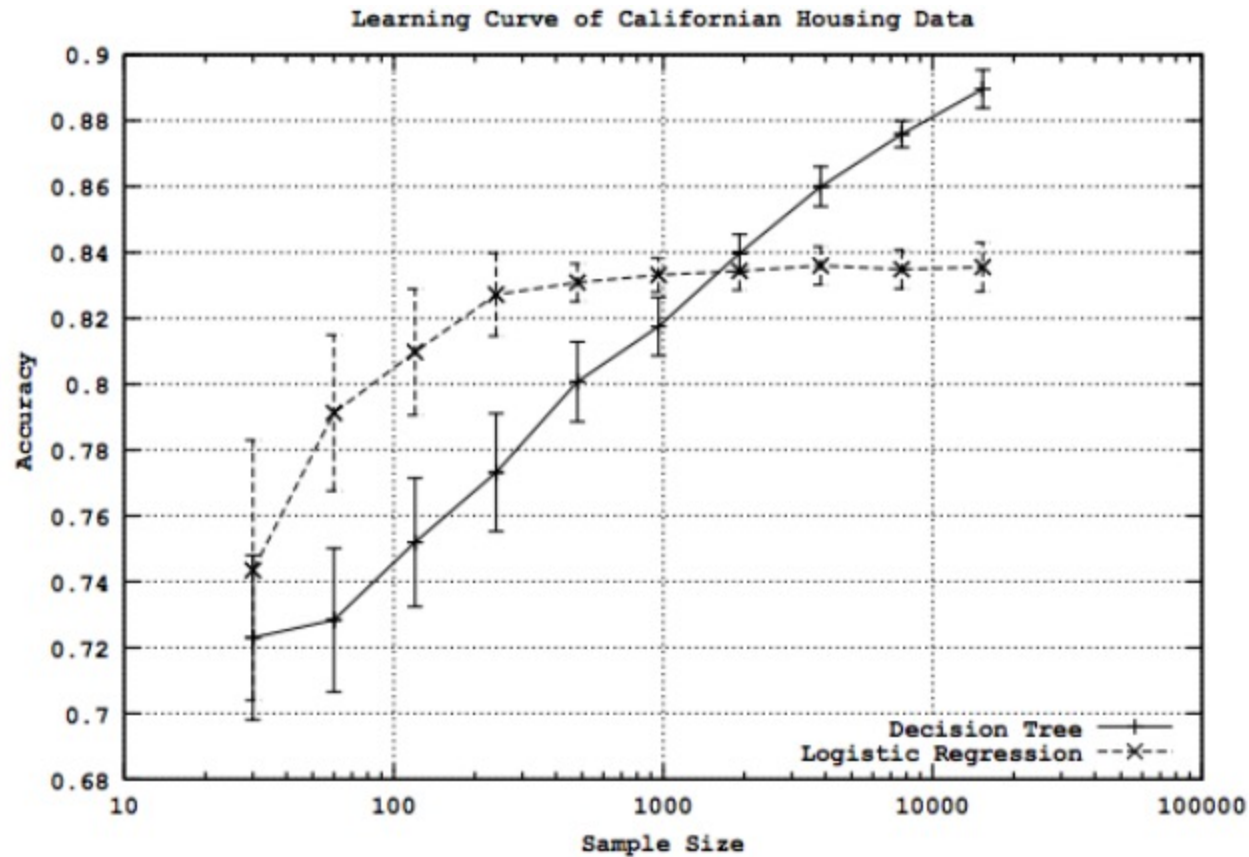


Figure from Perlich et al. *Journal of Machine Learning Research*, 2003