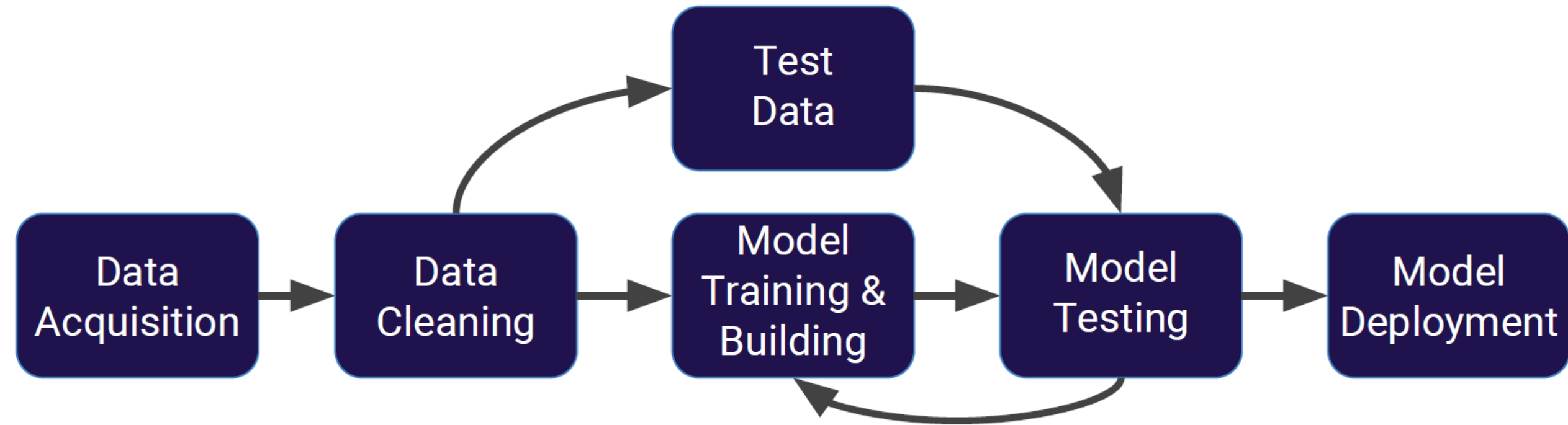


# **Introduction to Machine Learning**

# Machine Learning Process



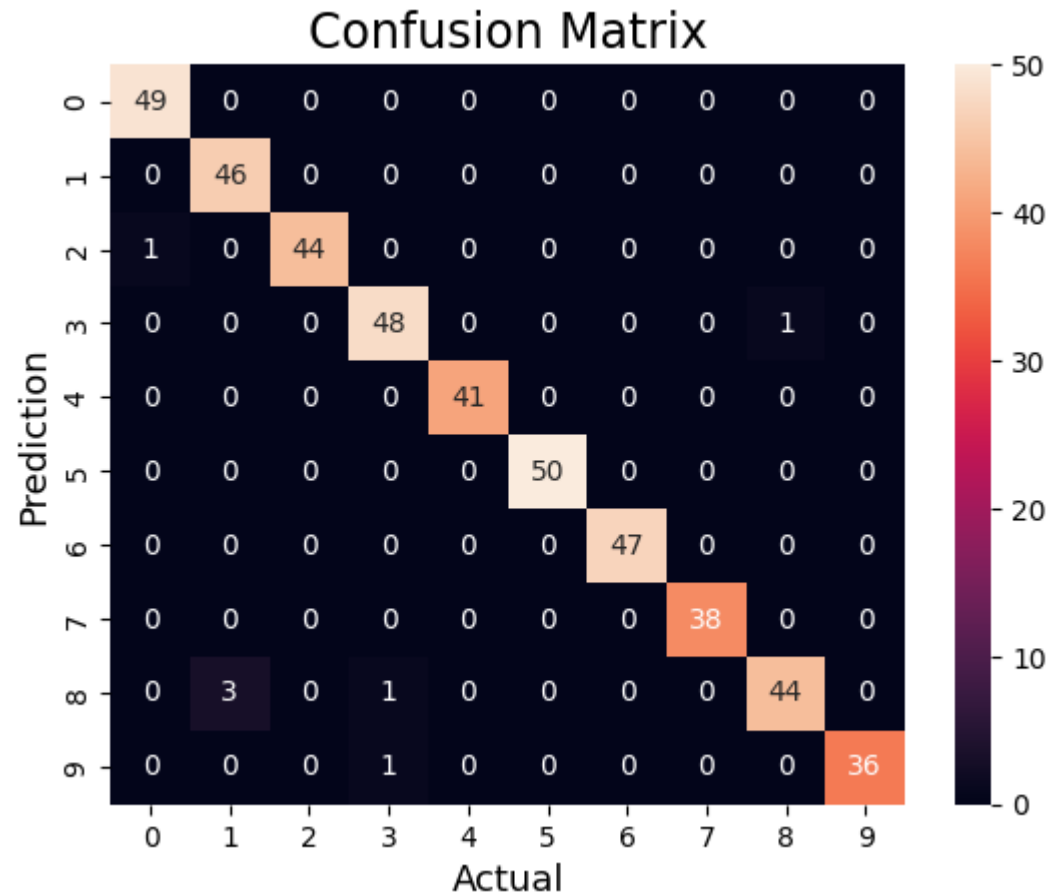
# **Model Evaluation**

# Confusion Matrix

For binary classification	Predicted <b>0</b>	Predicted <b>1</b>
Actual <b>0</b>	TN	FP
Actual <b>1</b>	FN	TP

# Confusion Matrix

For Multi-class classification



## Confusion Matrix and ROC Curve

		Predicted Class	
		No	Yes
Observed Class	No	TN	FP
	Yes	FN	TP

TN      True Negative  
FP      False Positive  
FN      False Negative  
TP      True Positive

## Model Performance

Accuracy

$$= (TN+TP)/(TN+FP+FN+TP)$$

Precision

$$= TP/(FP+TP)$$

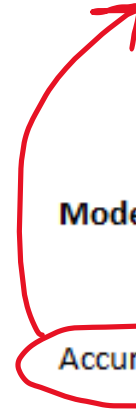
Sensitivity

$$= TP/(TP+FN)$$

Specificity

$$= TN/(TN+FP)$$

Not a good metric for unbalanced data



n=165		Predicted: NO	Predicted: YES	
Actual: NO		TN = 50	FP = 10	60
Actual: YES		FN = 5	TP = 100	105
		55	110	

## Basic Terminology:

- True Positives (TP)
- True Negatives (TN)
- False Positives (FP)
- False Negatives (FN)

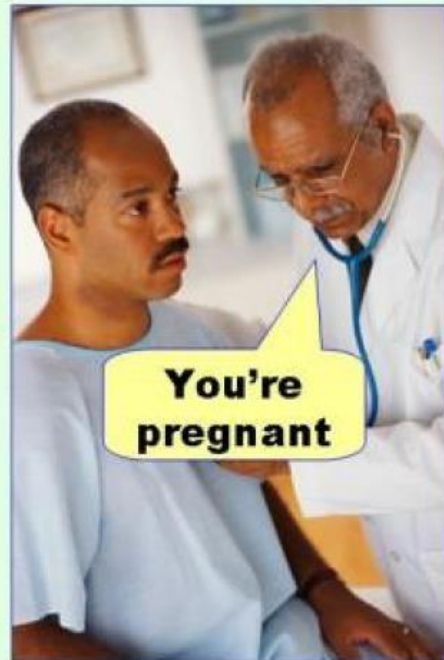
n=165	Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60
Actual: YES	FN = 5	TP = 100	105
	55	110	

Misclassification Rate  
(Error Rate):

- Overall, how often is it **wrong**?
- $(FP + FN) / \text{total} = 15/165 = 0.09$



**Type I error**  
(false positive)



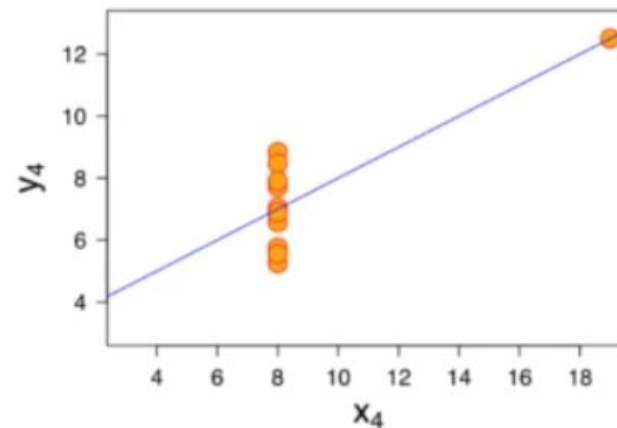
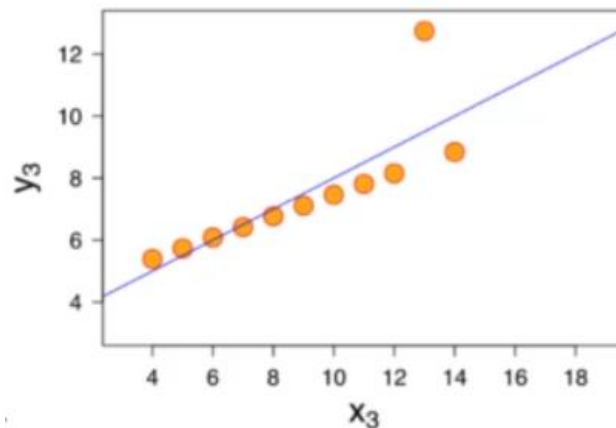
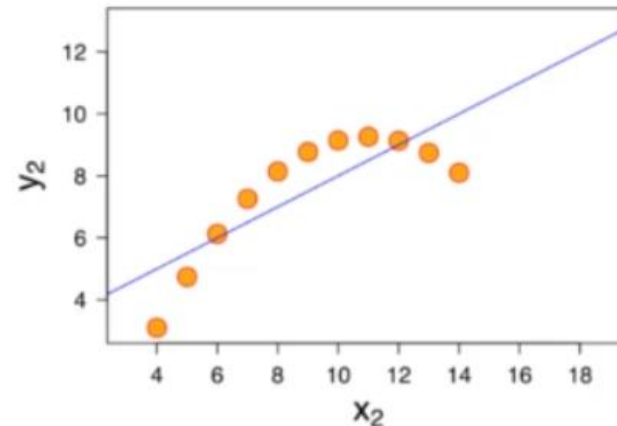
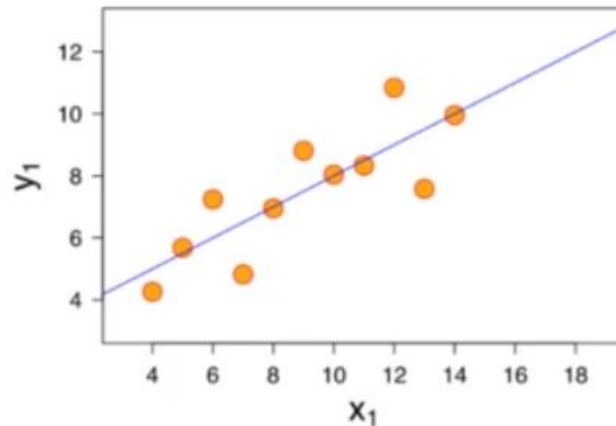
**Type II error**  
(false negative)



# Mean Absolut Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

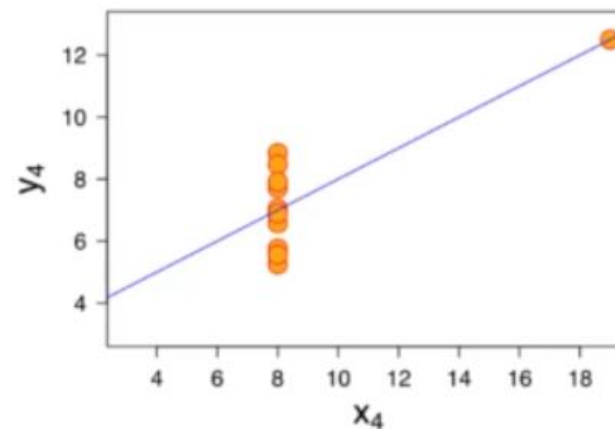
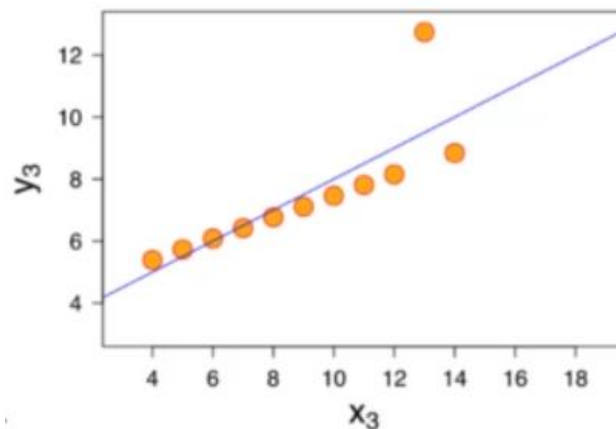
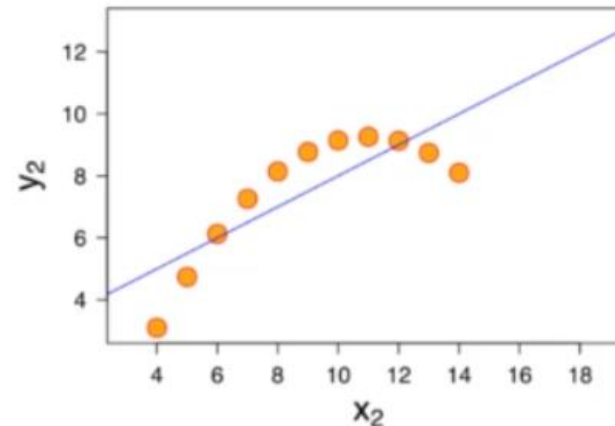
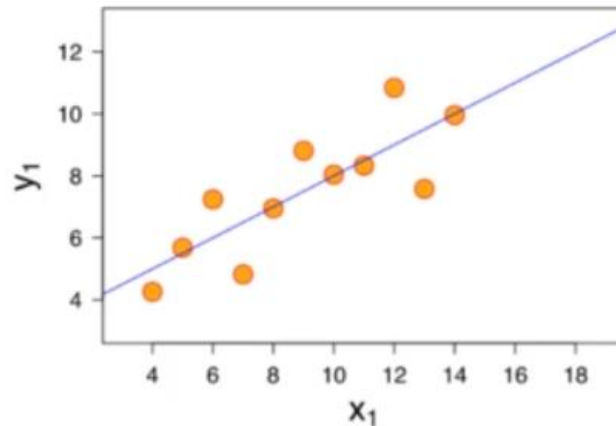


# Mean Absolut Error

square

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

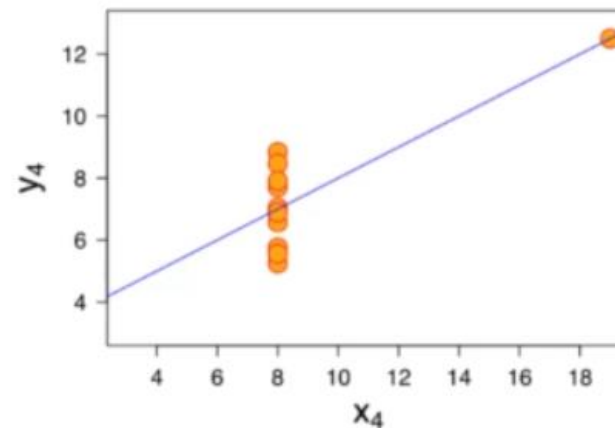
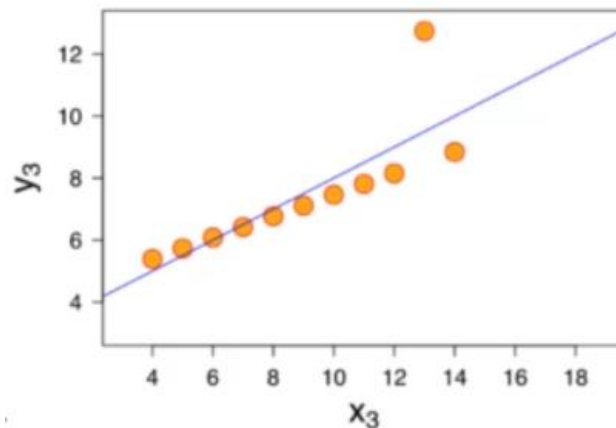
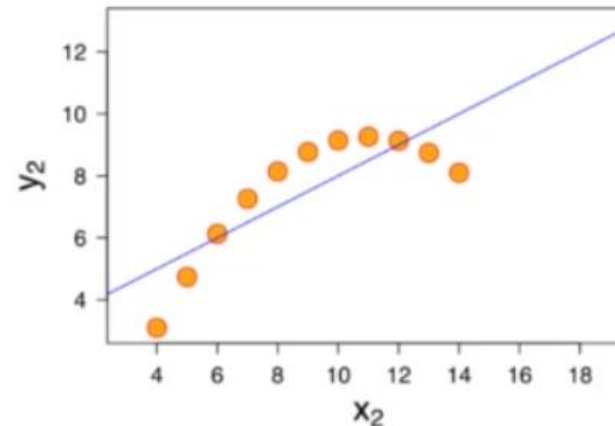
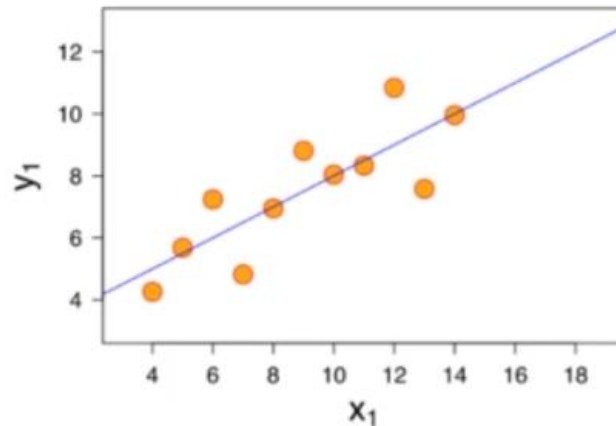
$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$



# Mean Absolut Error

$$R_{MSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$



		True condition			
Total population		Condition positive	Condition negative	$Prevalence = \frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$	$Accuracy \text{ (ACC)} = \frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$
Predicted condition	Predicted condition positive	<b>True positive</b> , Power	<b>False positive</b> , Type I error	$Positive \text{ predictive value (PPV), Precision} = \frac{\Sigma \text{ True positive}}{\Sigma \text{ Predicted condition positive}}$	$False \text{ discovery rate (FDR)} = \frac{\Sigma \text{ False positive}}{\Sigma \text{ Predicted condition positive}}$
	Predicted condition negative	<b>False negative</b> , Type II error	<b>True negative</b>	$False \text{ omission rate (FOR)} = \frac{\Sigma \text{ False negative}}{\Sigma \text{ Predicted condition negative}}$	$Negative \text{ predictive value (NPV)} = \frac{\Sigma \text{ True negative}}{\Sigma \text{ Predicted condition negative}}$
		$True \text{ positive rate (TPR), Recall, Sensitivity, probability of detection} = \frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	$False \text{ positive rate (FPR), Fall-out, probability of false alarm} = \frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	$Positive \text{ likelihood ratio (LR+)} = \frac{TPR}{FPR}$	$Diagnostic \text{ odds ratio (DOR)} = \frac{LR+}{LR-}$
		$False \text{ negative rate (FNR), Miss rate} = \frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	$Specificity \text{ (SPC), Selectivity, True negative rate (TNR)} = \frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	$Negative \text{ likelihood ratio (LR-)} = \frac{FNR}{TNR}$	
				$F_1 \text{ score} = \frac{1}{\frac{1}{Recall} + \frac{1}{Precision}}$	

- Regression

- $R^2$

$$R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$

- RMSE

- Classification

- Precision

- Recall

- Clustering

- Within Sum of Squares Error

**Install scikit learn**