
CLASSIFICATION

INTRODUCTION

CLASSIFICATION!

	Categorical	Continuous
Supervised	Classification	Regression
Unsupervised	Clustering	Dimension Reduction

CLASSIFICATION

WHEN... CLASSIFICATION?

Discuss in your groups:

- ☐ What are some examples of classifications we encounter in our lives?
- ☐ How might we measure the success of a classification?

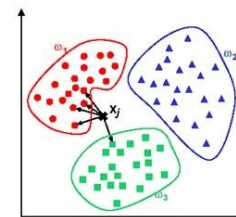
CLASSIFICATION

WHICH... CLASSIFICATION?

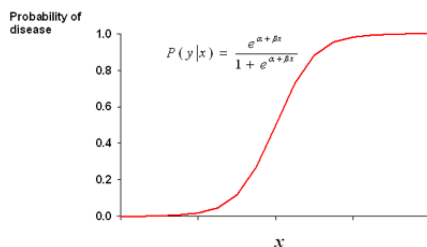
NAÏVE BAYES

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

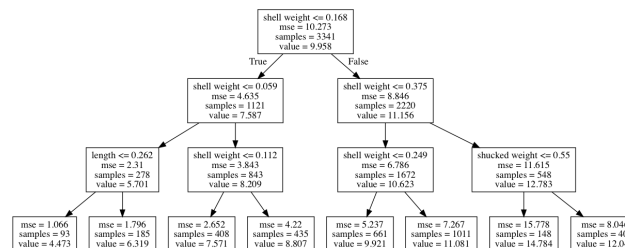
KNN



LOGISTIC REGRESSION



DECISION TREES



CLASSIFICATION

HOW... MEASURE PERFORMANCE?

Confusion Matrix: table to describe the performance of a classifier

n=165	Predicted: NO	Predicted: YES
Actual: NO	50	10
Actual: YES	5	100

Example: Test for presence of disease

NO = negative test = False = 0

YES = positive test = True = 1

- *How many classes are there?*
- *How many patients?*
- *How many times is disease predicted?*
- *How many patients actually have the disease?*

CLASSIFICATION

CONFUSION MATRIX

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)*

Accuracy:

- *Overall, how often is it **correct**?*
- *$(TP + TN) / \text{total} = 150 / 165 = 0.91$*

Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

→ *Accuracy:*

- *Overall, how often is it **correct**?*
- *$(TP + TN) / \text{total} = 150 / 165 = 0.91$*

Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

→ *Accuracy:*

- *Overall, how often is it **correct**?*
- $(TP + TN) / \text{total} = 150 / 165 = 0.91$

Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

→ *Accuracy:*

- *Overall, how often is it **correct**?*
- $(TP + TN) / \underline{\text{total}} = 150/165 = 0.91$

Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

→ *Accuracy:*

- *Overall, how often is it **correct**?*
- $(TP + TN) / total = 150 / 165 = 0.91$

Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

Accuracy:

- *Overall, how often is it **correct**?*
- *$(TP + TN) / total = 150 / 165 = 0.91$*

→ Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

Accuracy:

- *Overall, how often is it **correct**?*
- $(TP + TN) / \text{total} = 150 / 165 = 0.91$

→ Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)*

Accuracy:

- *Overall, how often is it **correct**?*
- $(TP + TN) / \text{total} = 150 / 165 = 0.91$

→ Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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)

Accuracy:

- Overall, how often is it **correct**?
- $(TP + TN) / \text{total} = 150 / 165 = 0.91$

→ Misclassification Rate (Error Rate):

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

CLASSIFICATION

CONFUSION MATRIX

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

False Positive Rate:

- *When actual value is **negative**, how often is prediction **wrong**?*
- *$FP / \text{actual no} = 10/60 = 0.17$*

→ Sensitivity:

- *When actual value is **positive**, how often is prediction **correct**?*
- *$TP / \text{actual yes} = 100/105 = 0.95$*
- *“True Positive Rate” or “Recall”*

Specificity:

- *When actual value is **negative**, how often is prediction **correct**?*
- *$TN / \text{actual no} = 50/60 = 0.83$*

CLASSIFICATION

CONFUSION MATRIX

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

False Positive Rate:

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

→ Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- “True Positive Rate” or “Recall”

Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

CONFUSION MATRIX

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

False Positive Rate:

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

→ Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- “True Positive Rate” or “Recall”

Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

CONFUSION MATRIX

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

→ False Positive Rate:

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- “True Positive Rate” or “Recall”

Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

CONFUSION MATRIX

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

→ **False Positive Rate:**

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- “True Positive Rate” or “Recall”

Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

CONFUSION MATRIX

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

False Positive Rate:

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- "True Positive Rate" or "Recall"



Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

CONFUSION MATRIX

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

False Positive Rate:

- When actual value is **negative**, how often is prediction **wrong**?
- $FP / \text{actual no} = 10/60 = 0.17$

Sensitivity:

- When actual value is **positive**, how often is prediction **correct**?
- $TP / \text{actual yes} = 100/105 = 0.95$
- "True Positive Rate" or "Recall"



Specificity:

- When actual value is **negative**, how often is prediction **correct**?
- $TN / \text{actual no} = 50/60 = 0.83$

CLASSIFICATION

RECEIVING OPERATOR CHARACTERISTIC (ROC) CURVE

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Every email is assigned a “spamminess” score by our classification algorithm. To actually make our predictions, we choose a numeric cutoff for classifying as spam.

An ROC Curve will help us to visualize how well our classifier is doing without having to choose a cutoff!

CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

The ROC plots the True Positive Rate (TRP) on the y-axis against the False Positive Rate (FPR) on the x-axis.

*TPR: When actual value is **spam**, how often is prediction **correct**?*

*FPR: When actual value is **ham**, how often is prediction **wrong**?*

CLASSIFICATION

ROC CURVE / AUC

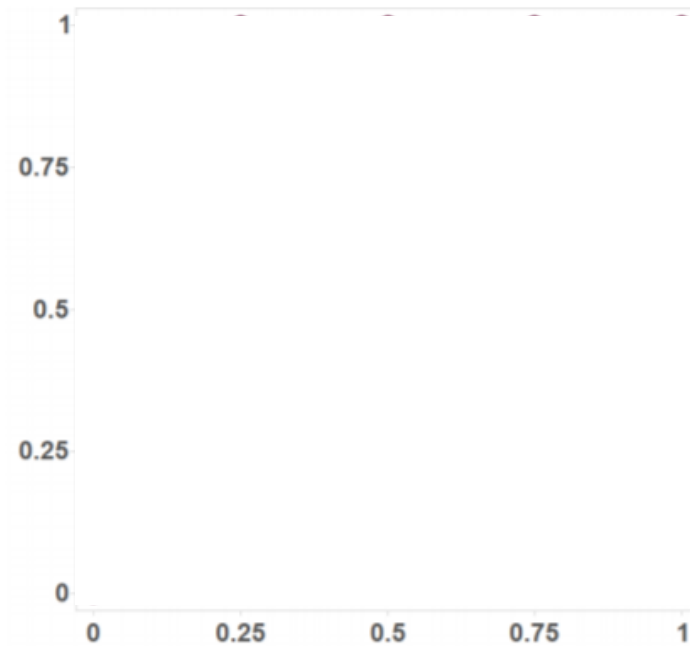
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	?	?	0.50		
0.05			0.65		
0.15			0.85		
0.25			1		

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

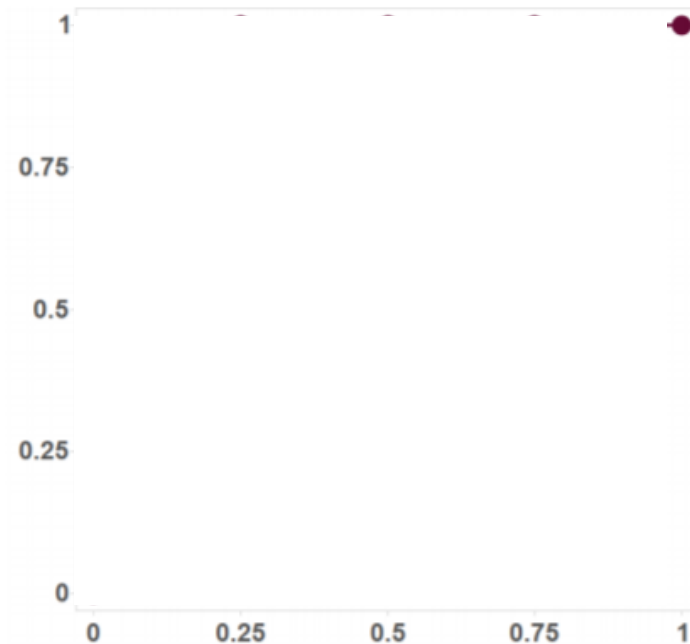
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	0	4	4
Actual: YES	0	4	4
	0	8	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50		
0.05			0.65		
0.15			0.85		
0.25			1		

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

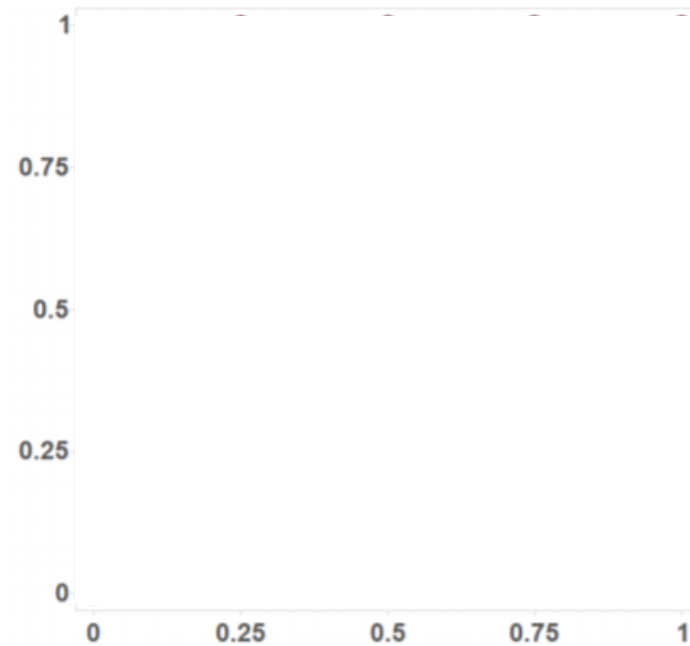
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

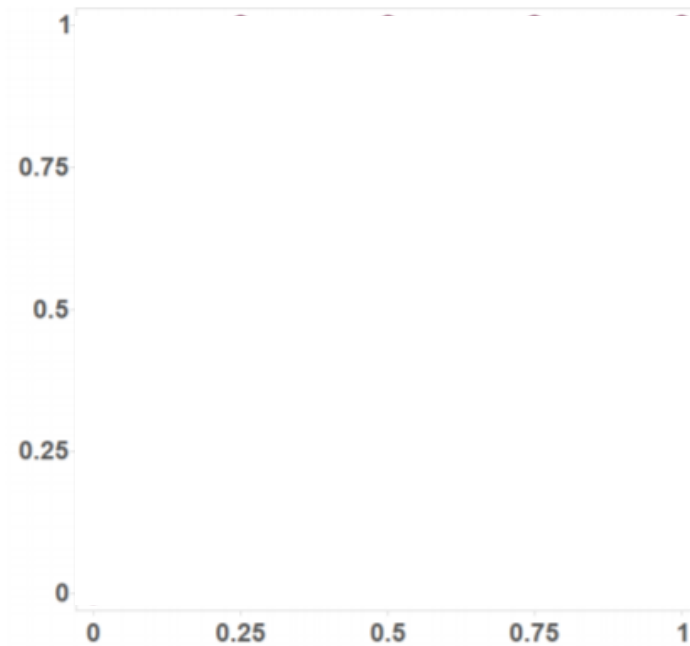
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

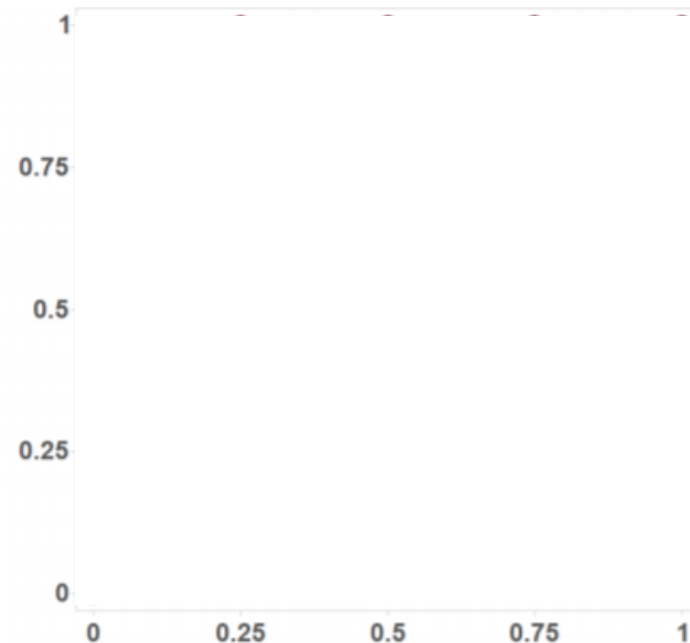
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES		1

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

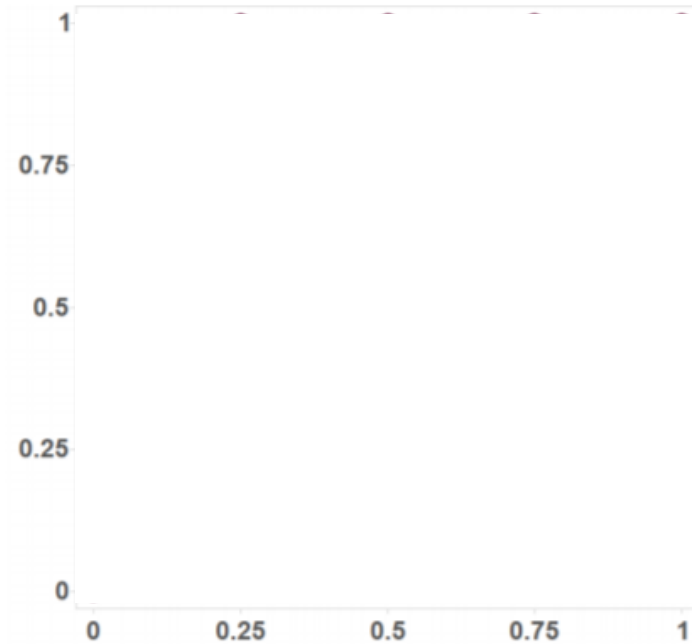
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO			
Actual: YES		2	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO			
Actual: YES		3	

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

ROC CURVE / AUC

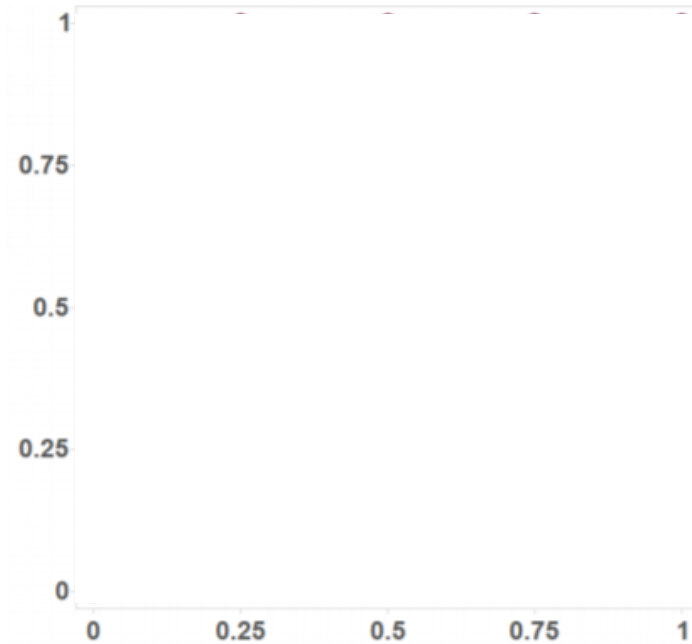
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO			
Actual: YES		4	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

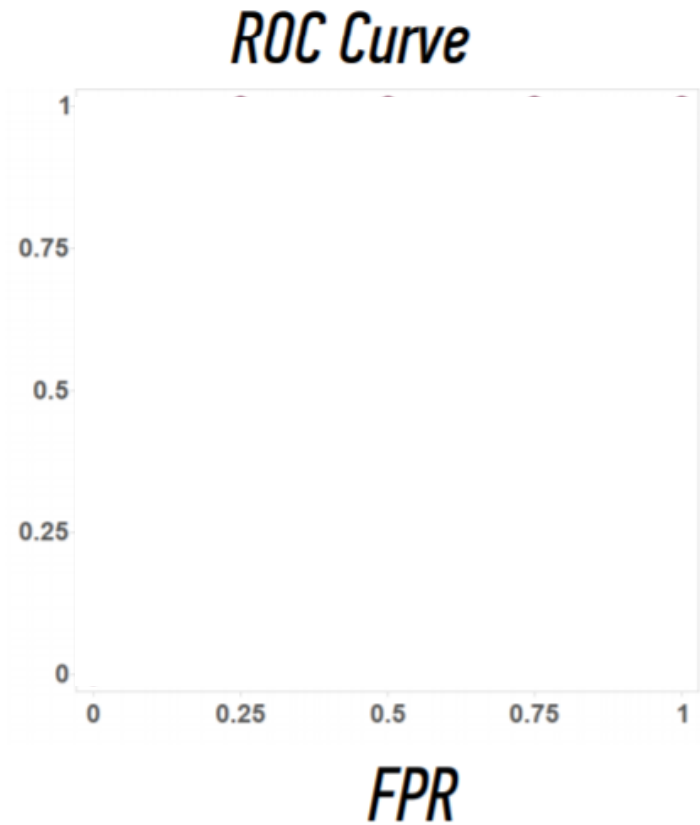
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO			
Actual: YES		5	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0



CLASSIFICATION

ROC CURVE / AUC

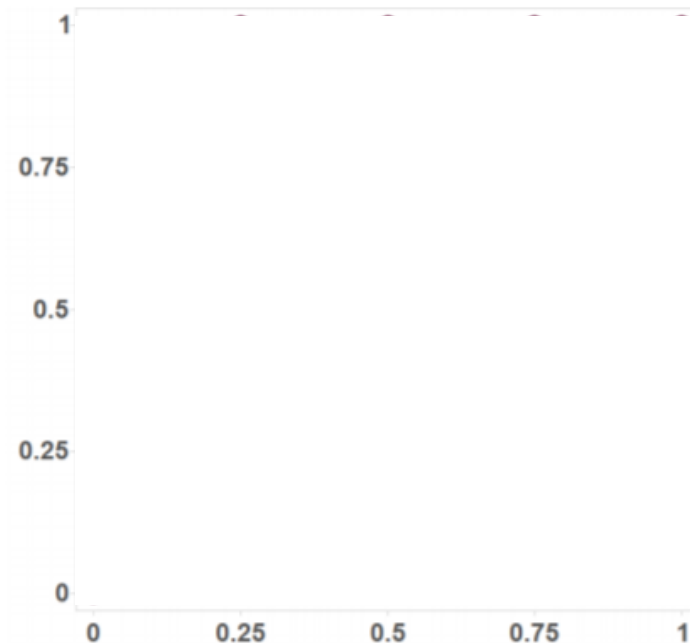
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO			
Actual: YES		7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

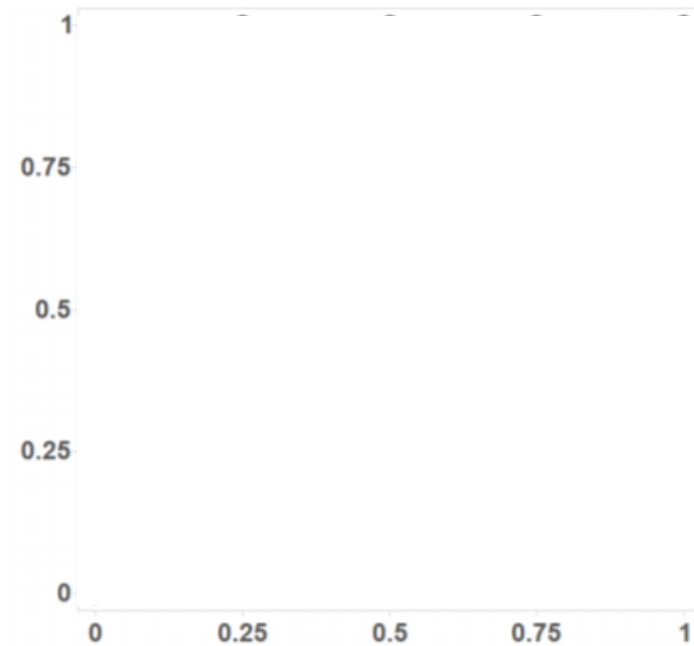
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		
Actual: YES	1	7

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

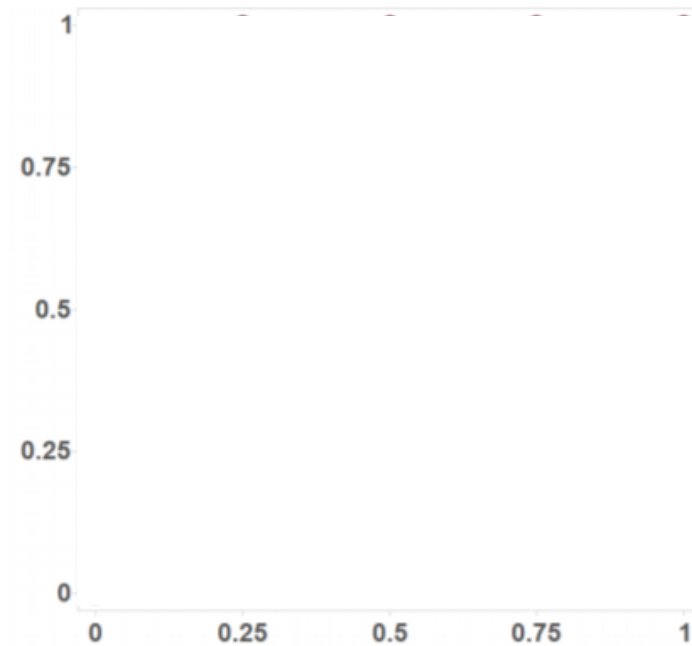
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO		3
Actual: YES	1	7

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

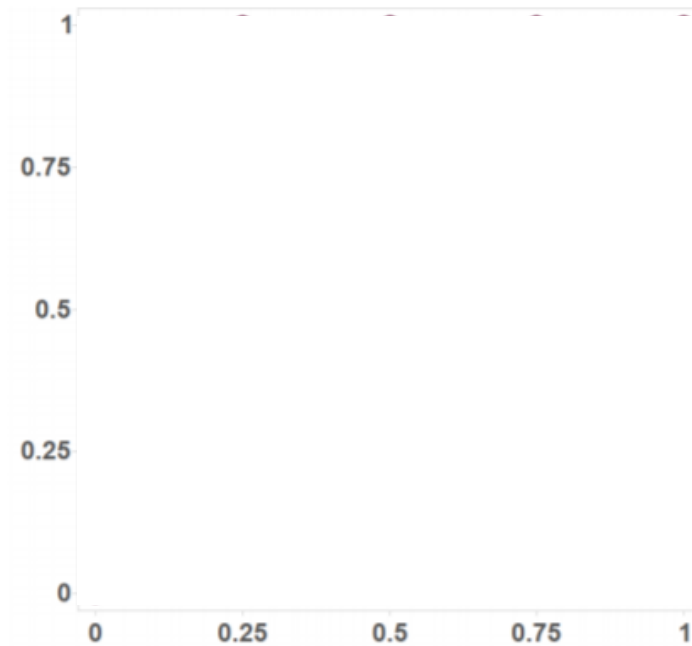
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES
Actual: NO	1	3
Actual: YES	1	7

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

ROC CURVE / AUC

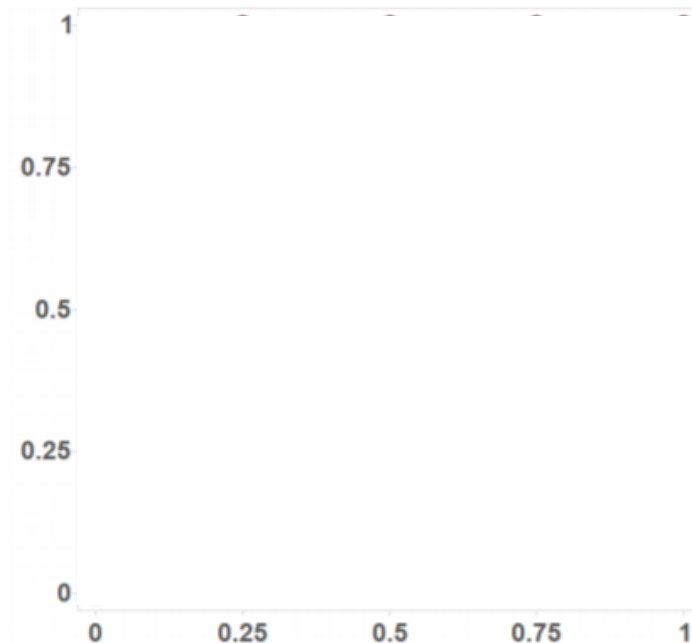
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	1	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

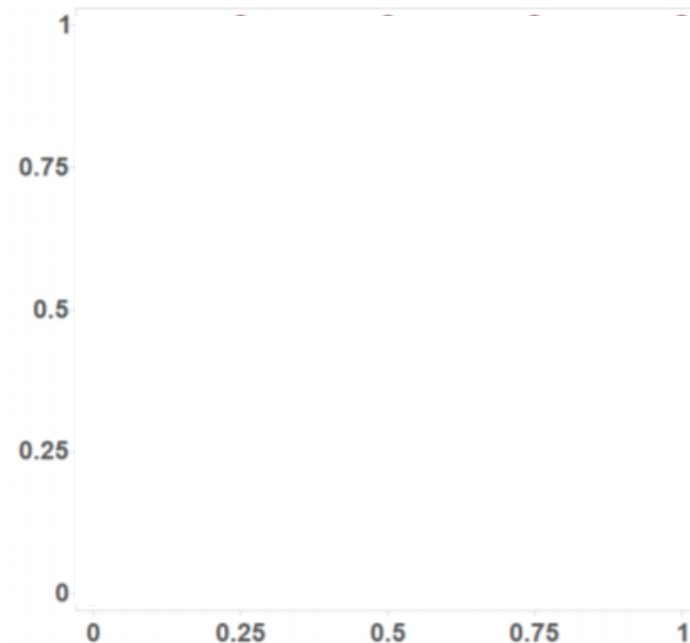
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	?		
	1	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

ROC CURVE / AUC

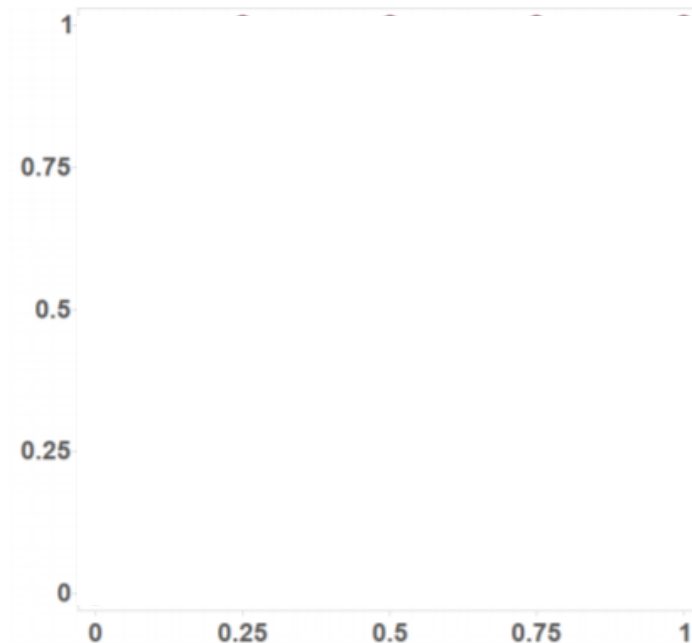
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

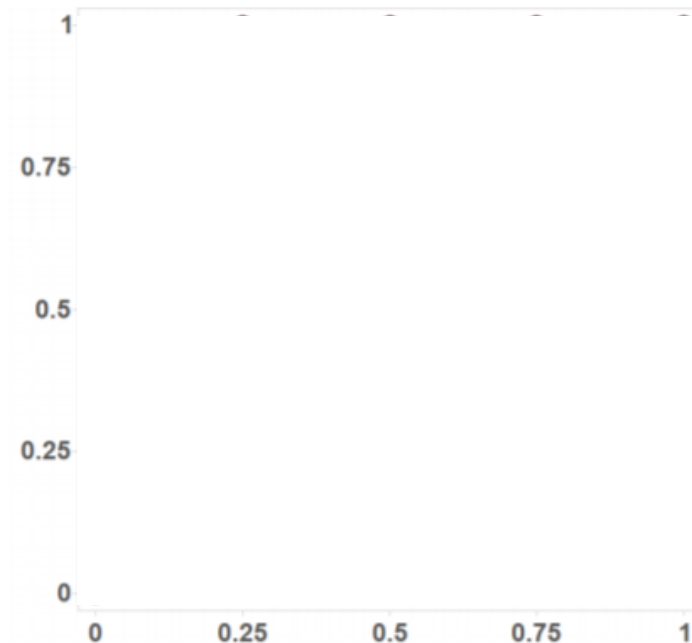
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	?	
	1	7	

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

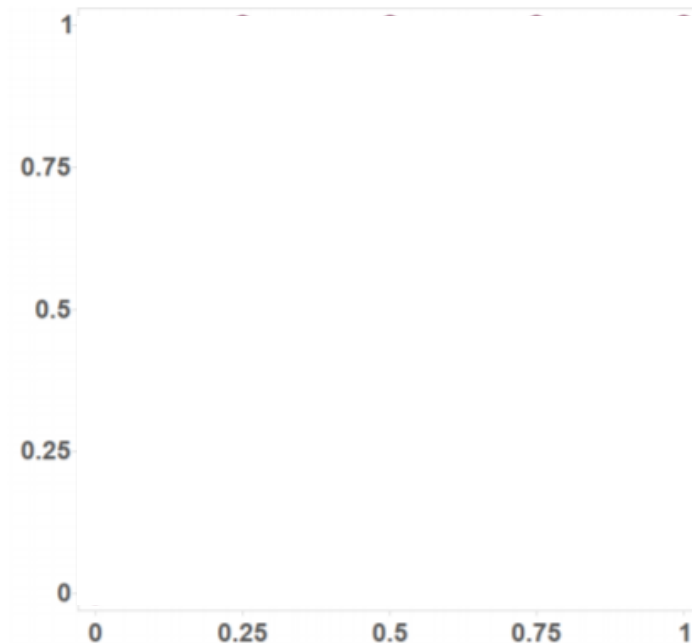
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	?	
	1	7	

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

ROC CURVE / AUC

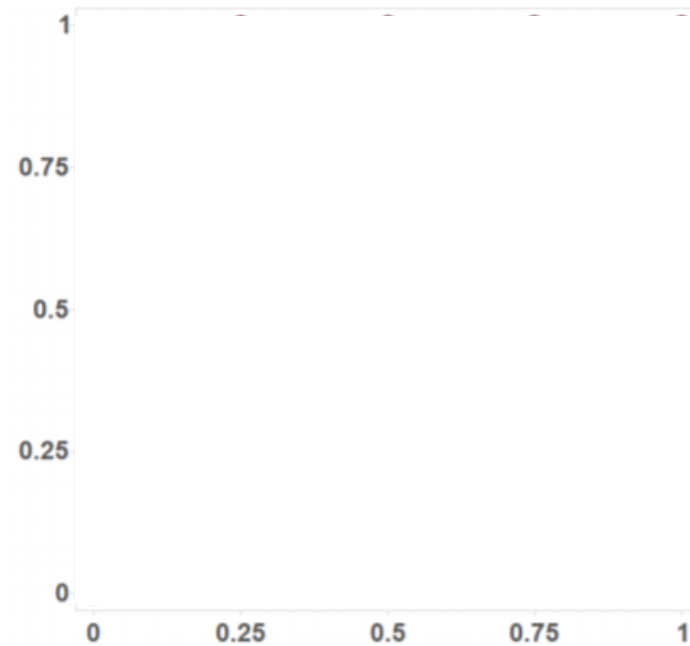
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	?	
	1	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR

CLASSIFICATION

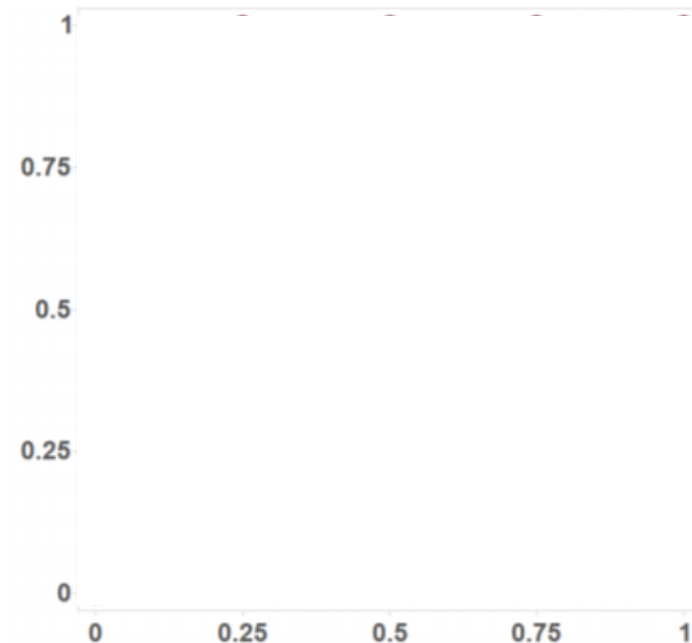
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	
	1	7	

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

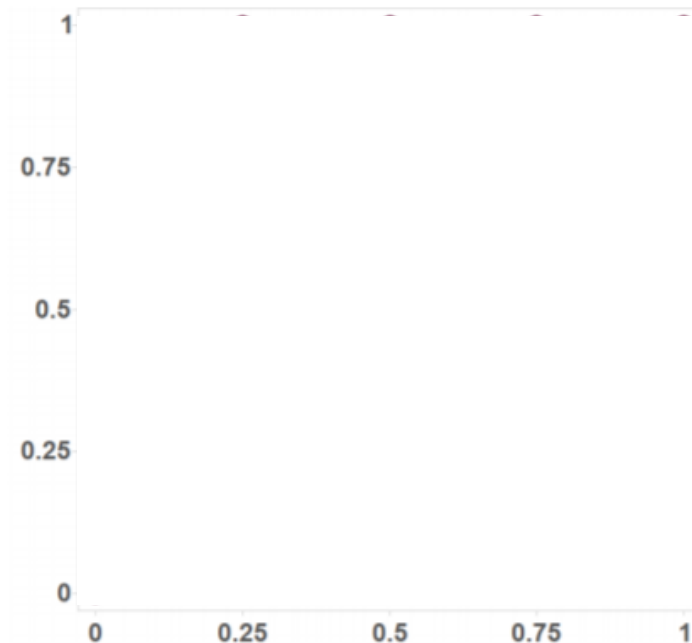
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	4
	1	7	

TPR

ROC Curve



FPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

CLASSIFICATION

ROC CURVE / AUC

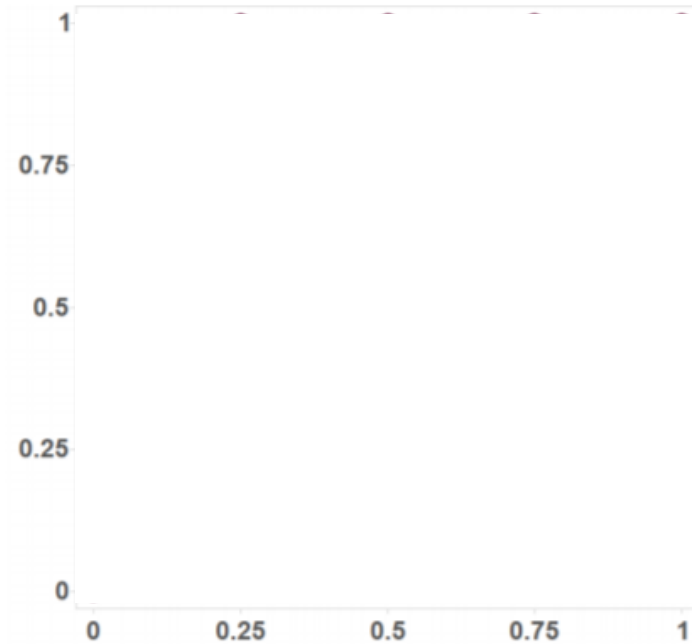
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	4
	1	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR = $3/4 \rightarrow 0.75$

CLASSIFICATION

ROC CURVE / AUC

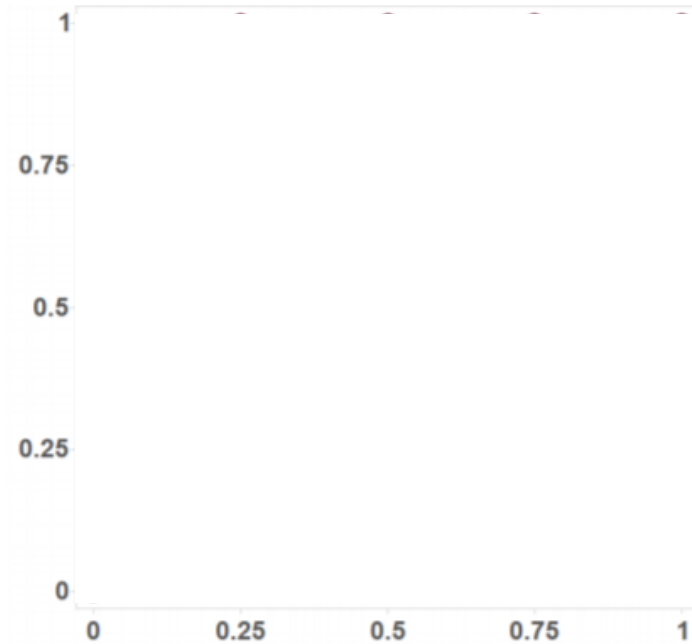
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	4
	1	7	

TPR

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



FPR = $3/4 \rightarrow 0.75$

CLASSIFICATION

ROC CURVE / AUC

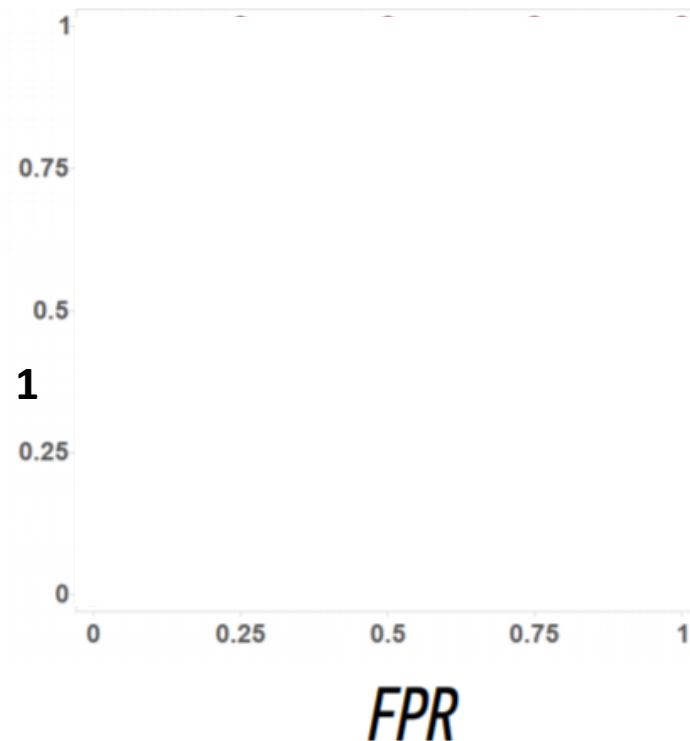
Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	4
	1	7	

$TPR = \frac{4}{4} \rightarrow 1$

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

ROC Curve



CLASSIFICATION

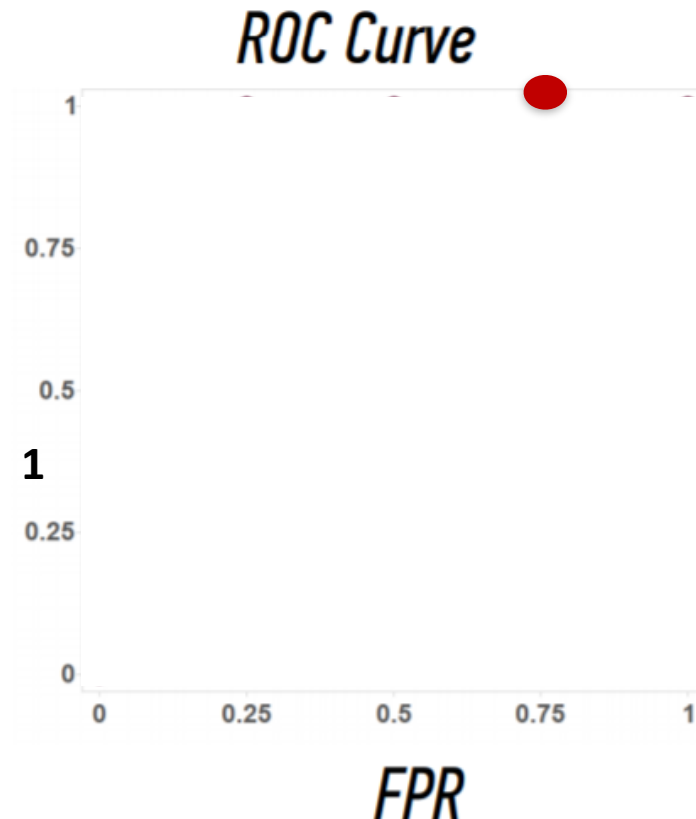
ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

	Predicted: NO	Predicted: YES	
Actual: NO	1	3	4
Actual: YES	0	4	4
	1	7	

TPR
 $= 4/4 \rightarrow 1$

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

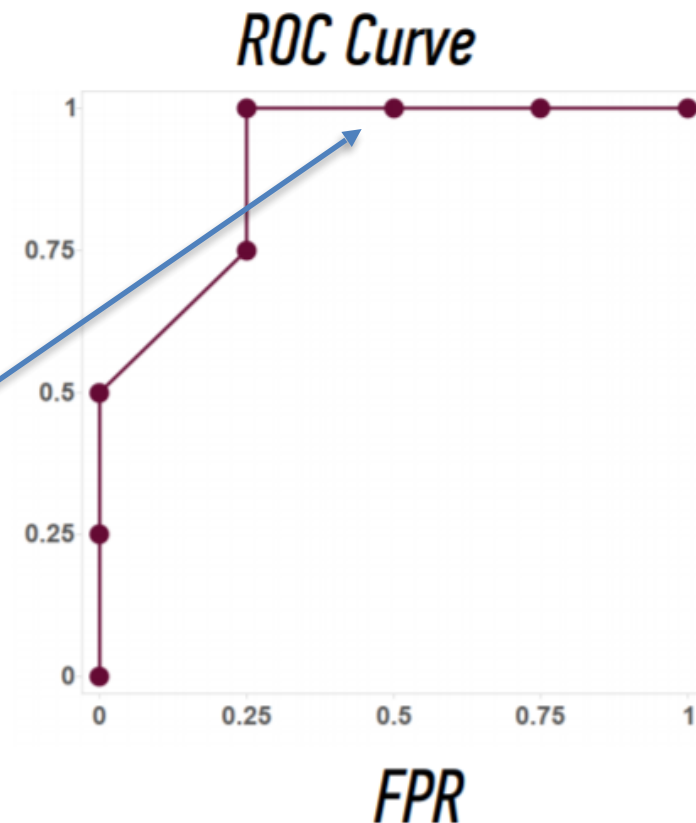


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

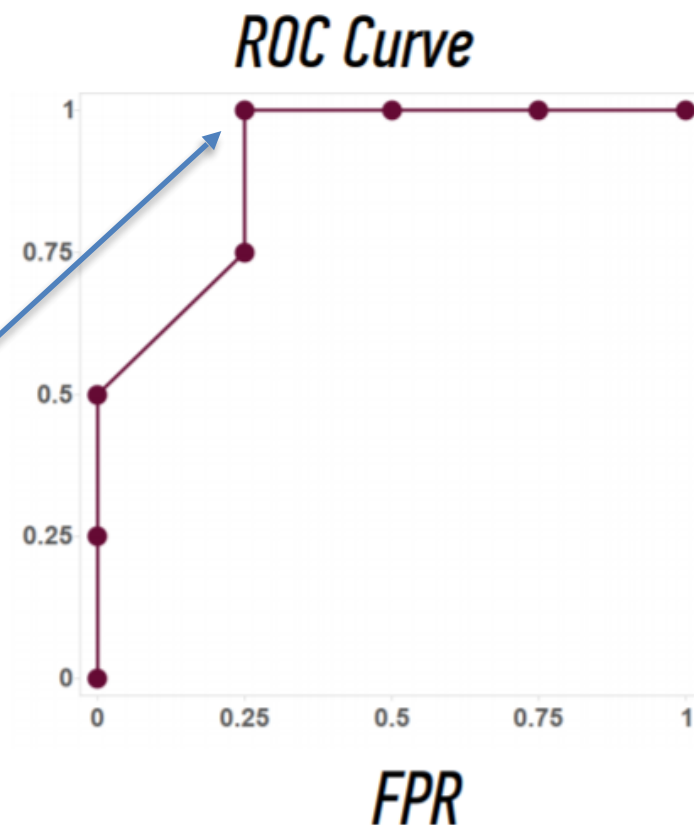


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

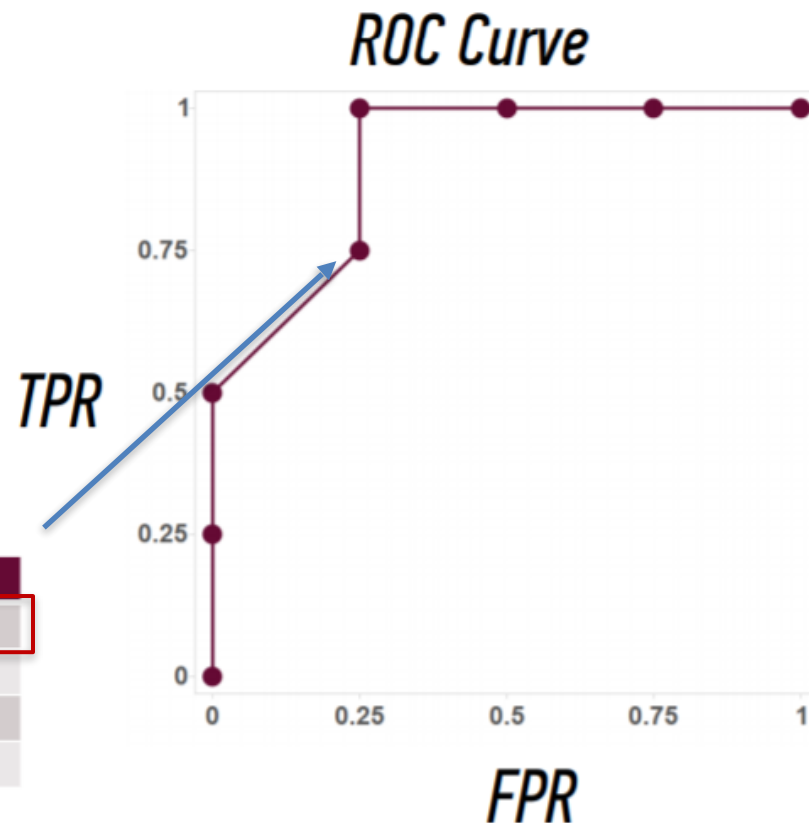


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

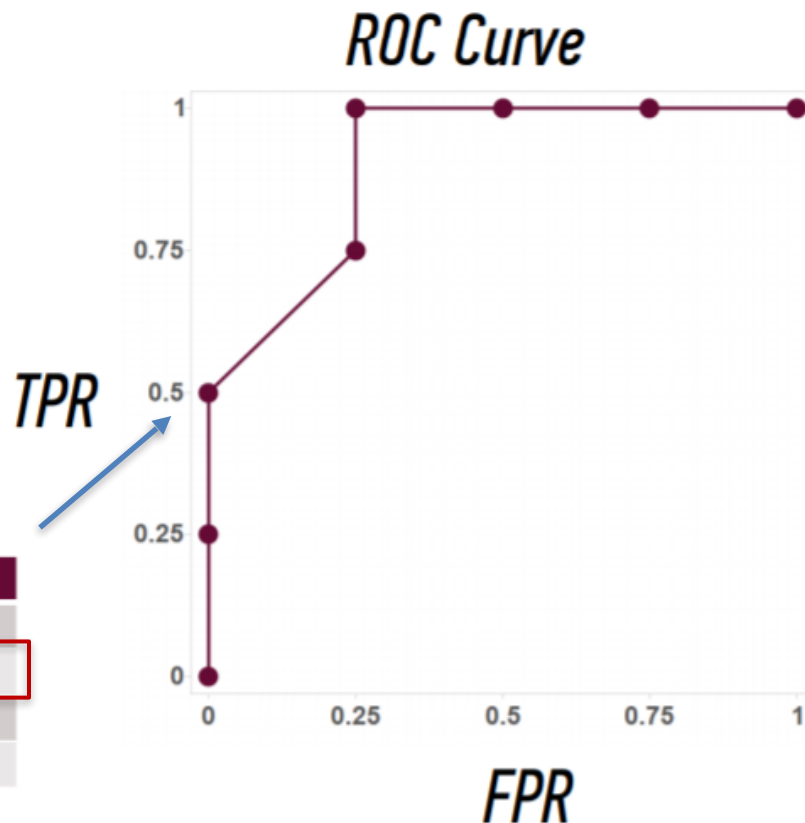


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

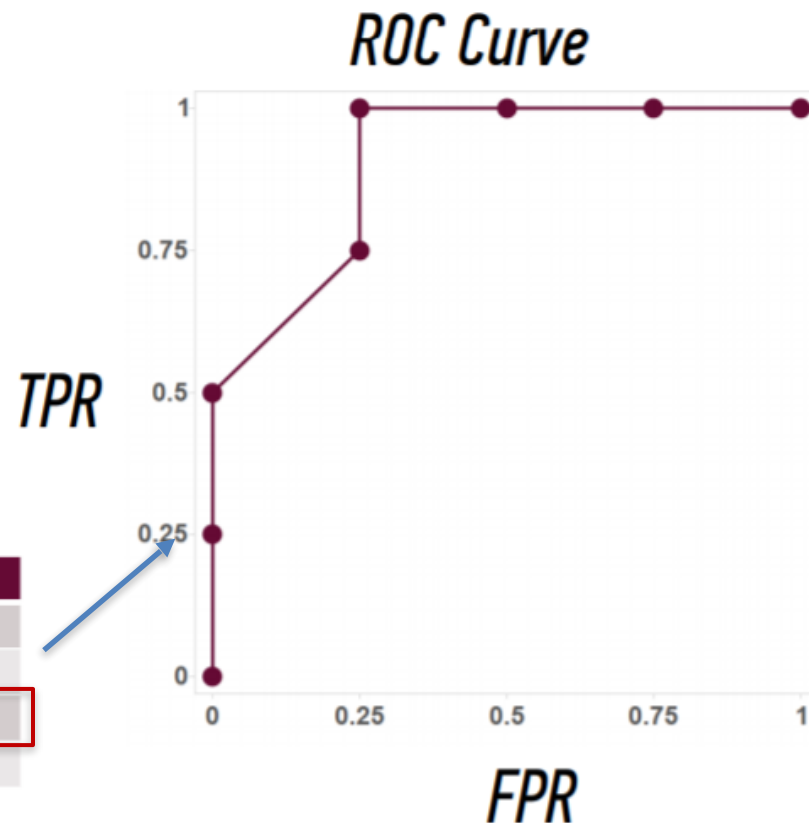


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0

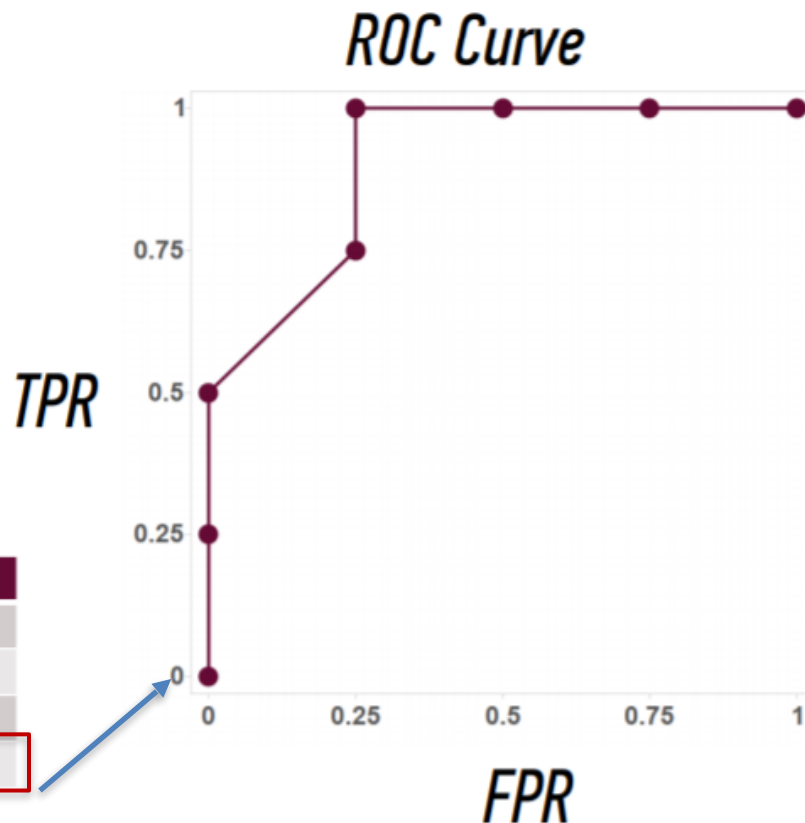


CLASSIFICATION

ROC CURVE / AUC

Email Number	Score	True Label
5	0.99	Spam
8	0.82	Spam
2	0.60	Spam
1	0.60	Ham
7	0.48	Spam
3	0.22	Ham
4	0.10	Ham
6	0.02	Ham

Cutoff	TPR (y)	FPR (x)	Cutoff	TPR (y)	FPR (x)
0	1	1	0.50	0.75	0.25
0.05	1	0.75	0.65	0.5	0
0.15	1	0.5	0.85	0.25	0
0.25	1	0.25	1	0	0



CONFUSION VS ROC?

Discuss in your groups:

- ☐ What information do you take away from each of these evaluation techniques?
- ☐ What decisions can be made from each tool?

CLASSIFICATION

CODING

NAÏVE BAYES THEOREM


*Suppose we have a dataset with features x_1, \dots, x_n and a class label C .
What can we say about classification using Bayes' theorem?*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Bayes' theorem can help us to determine the probability of a record belonging to a class, given the data we observe.


NAÏVE BAYES THEOREM

*This term is the **prior probability of C**. It represents the probability of a record belonging to class C before the data is taken into account.*

$$P(\text{class } C | \{x_i\}) = \frac{P(\{x_i\} | \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$



NAÏVE BAYES THEOREM

*This term is the **likelihood function**. It represents the joint probability of observing features $\{x_i\}$ given that that record belongs to class C .*

$$P(\text{class } C | \{x_i\}) = \frac{P(\{x_i\} | \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$



NAÏVE BAYES THEOREM

This term is the normalization constant. It doesn't depend on C , and is generally ignored.

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


NAÏVE BAYES THEOREM

*This term is the **posterior probability** of C . It represents the probability of a record belonging to class C after the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


*The idea of Bayesian inference, then, is to **update** our beliefs about the distribution of C using the data (“evidence”) at our disposal.*

NAÏVE BAYES THEOREM

Q: What piece of the puzzle we've seen so far looks like it could intractably difficult in practice?

A: Estimating the full likelihood function.

$$P(\{x_i\}|C) = P(\{x_1, x_2, \dots, x_n\}|C)$$

Observing this exactly would require us to have enough data for every possible combination of features to make a reasonable estimate.

NAÏVE BAYES THEOREM

Q: So what can we do about it?

A: Make a simplifying assumption. In particular, we assume that the features x_i are conditionally independent from each other:

$$P(\{x_i\}|C) = P(\{x_1, x_2, \dots, x_n\}|C) \approx P(x_1|C) * P(x_2|C) * \dots * P(x_n|C)$$

This “naïve” assumption simplifies the likelihood function to make it tractable.

CLASSIFICATION

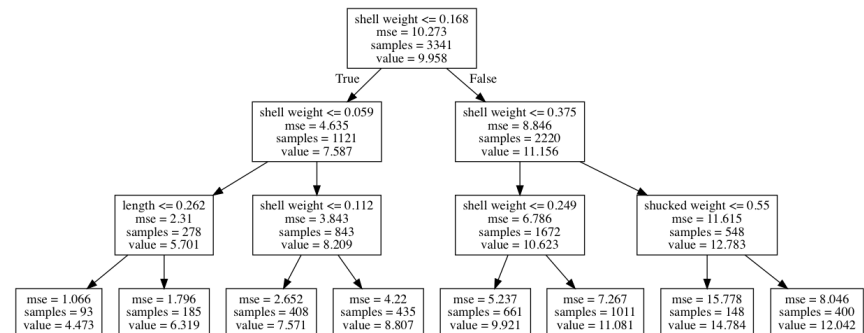
CODING

CLASSIFICATION

DECISION TREES

1. Find the purest split
(using gini index)
2. Find the next purest split
3. Continue until max depth is reached

Tuning: max_depth,
min_sample_leaf



$$\text{Gini Index} = 1 - \sum_j p_j^2$$

CLASSIFICATION

DECISION TREES

Advantages

1. Easy to interpret and make for straightforward visualizations.
2. The internal workings are capable of being observed and thus make it possible to reproduce work.
3. Can handle both numerical and categorical data.
4. Perform well on large datasets
5. Are extremely fast

Disadvantages

1. Purest split at each step might lead to local maximum not global maximum
2. Prone to overfitting, leads to high variance from sample to sample

CLASSIFICATION

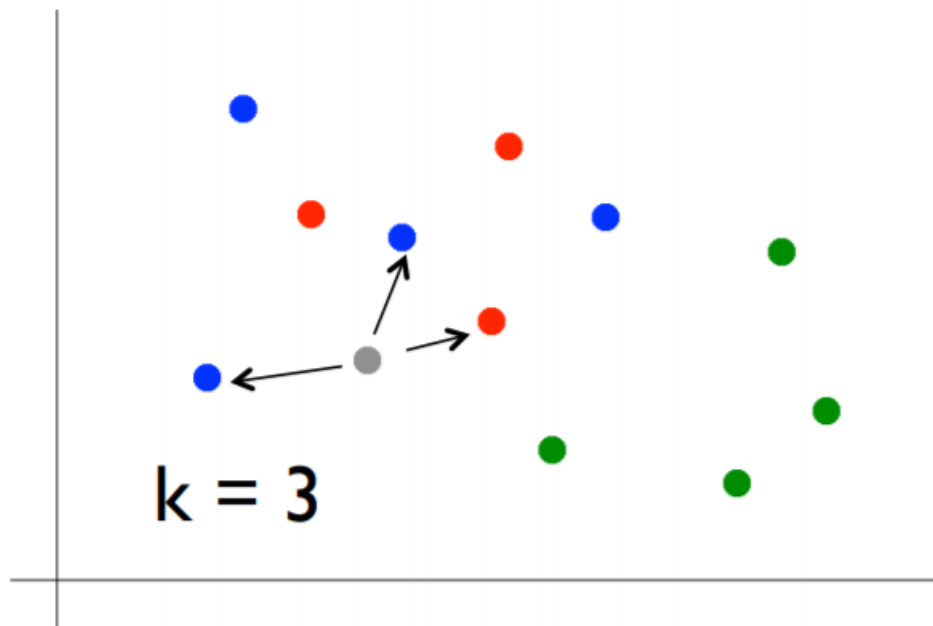
CODING

CLASSIFICATION

KNN

1. Pick a value for k
2. Find colors of k nearest neighbors
3. Assign the most common color to the gray dot

Tuning: `k_neighbors`



CLASSIFICATION

KNN

Advantages

1. Can learn complex topics by local approximation (simple methods)
2. No assumptions or cost of learning

Disadvantages

1. Does not handle categorical values well
2. Can't be interpreted
3. Computationally expensive

CLASSIFICATION

CODING

CLASSIFICATION

Q&A