

MACHINE LEARNING WITH PYTHON

NAÏVE BAYES

Themistoklis Diamantopoulos

Categorical Problem

- Decide whether the traffic is going to be high based on the weather and the day

Weather	Day	HighTraffic
Hot	Work	No
Cold	Vacation	No
Hot	Vacation	Yes
Hot	Work	Yes
Hot	Work	Yes
Cold	Vacation	No
Cold	Vacation	Yes

Naïve Bayes

- Independent features
- Bayes Theorem

$$P(c | x) = \frac{P(x_1 | c) \cdot P(x_2 | c) \cdot \dots \cdot P(x_n | c) \cdot P(c)}{P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_n)}$$

Weather	Day	HighTraffic
Hot	Work	No
Cold	Vacation	No
Hot	Vacation	Yes
Hot	Work	Yes
Hot	Work	Yes
Cold	Vacation	No
Cold	Vacation	Yes

$$P(\text{Yes} | \text{Hot}, \text{Vacation}) = \frac{P(\text{Hot} | \text{Yes}) \cdot P(\text{Vacation} | \text{Yes}) \cdot P(\text{Yes})}{P(\text{Hot}) \cdot P(\text{Vacation})} = \frac{3/4 \cdot 2/4 \cdot 4/7}{4/7 \cdot 4/7} = 21/32 = 0.65625$$



$$P(\text{No} | \text{Hot}, \text{Vacation}) = \frac{P(\text{Hot} | \text{No}) \cdot P(\text{Vacation} | \text{No}) \cdot P(\text{No})}{P(\text{Hot}) \cdot P(\text{Vacation})} = \frac{1/3 \cdot 2/3 \cdot 3/7}{3/7 \cdot 3/7} = 14/27 = 0.51852$$

When the weather is Hot and the day is Vacation, the traffic is High (prob: 0.56)

Classification Evaluation

- Confusion Matrix

Predicted \ Actual Class		Actual Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN
		$P = TP + FN$	$N = FP + TN$

- Evaluation Metrics

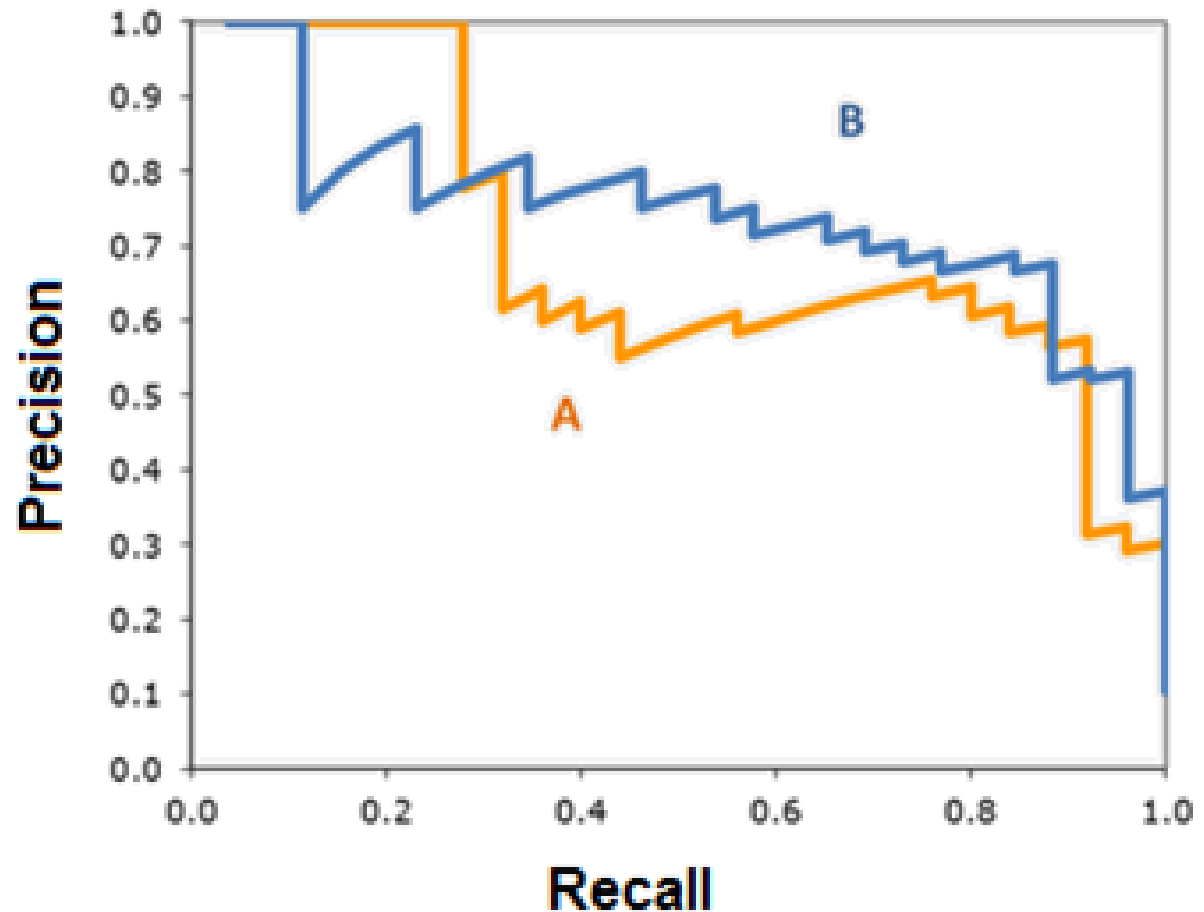
- Accuracy = $(TP + TN) / (P + N)$
- Precision = $TP / (TP + FP)$
- Recall = $TP / (TP + FN)$

Precision and Recall

- Tradeoff between Precision & Recall

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

$$Recall = \frac{TP}{TP + FN}$$



ROC Curve

- True Positive Rate (also known as sensitivity or recall)

$$TPR = \frac{TP}{TP + FN}$$

- False Positive Rate (also known as specificity)

$$FPR = \frac{FP}{FP + TN}$$

- AUC (Area Under the Curve)

