

Fairness in ML

\hat{y} ← prediction of the ML classifier.

X, Y, \hat{Y}

$$\text{Accuracy} \equiv P(Y = \hat{Y})$$

Test data

X_1	Y_1	\hat{Y}_1
X_2	Y_2	\hat{Y}_2
X_3	Y_3	\hat{Y}_3
⋮	⋮	⋮
X_{10}	Y_{10}	\hat{Y}_{10}

$$\frac{\#(Y = \hat{Y})}{10}$$

$$P(\hat{Y} = 1 | Y = 1)$$

X	Y	\hat{Y}
X_1	1	1
X_2	0	1
X_3	0	0
X_4	1	0
X_5	1	1
X_6	1	1
X_7	1	0
X_8	1	1
X_9	0	0
X_{10}	1	1

True Positive Rate

$$\text{Accuracy} = P(\hat{Y} = Y) = \frac{7}{10} = 0.7$$

$$\text{TPR} = P(\hat{Y} = 1 | \underline{Y} = 1) = \frac{5}{7}$$

recall
on +ve class

$$\text{FNR} = P(\hat{Y} = 0 | Y = 1) = \frac{2}{7}$$

$$\text{FPR} = P(\hat{Y} = 1 | Y = 0) = \frac{1}{3}$$

$$TNR = P(\hat{y}=0 | y=0) = \frac{2}{3}$$

recall on
-ve class

$$P(y=1 | \hat{y}=1) = \frac{5}{6}$$

Precision (on +ve
class)

$$P(y=0 | \hat{y}=0) = \frac{2}{4}$$

Precision on -ve
class

f-measure for +ve class

$$= \frac{2}{\frac{1}{\text{recall}_{+ve}} + \frac{1}{\text{precision}_{+ve}}}$$

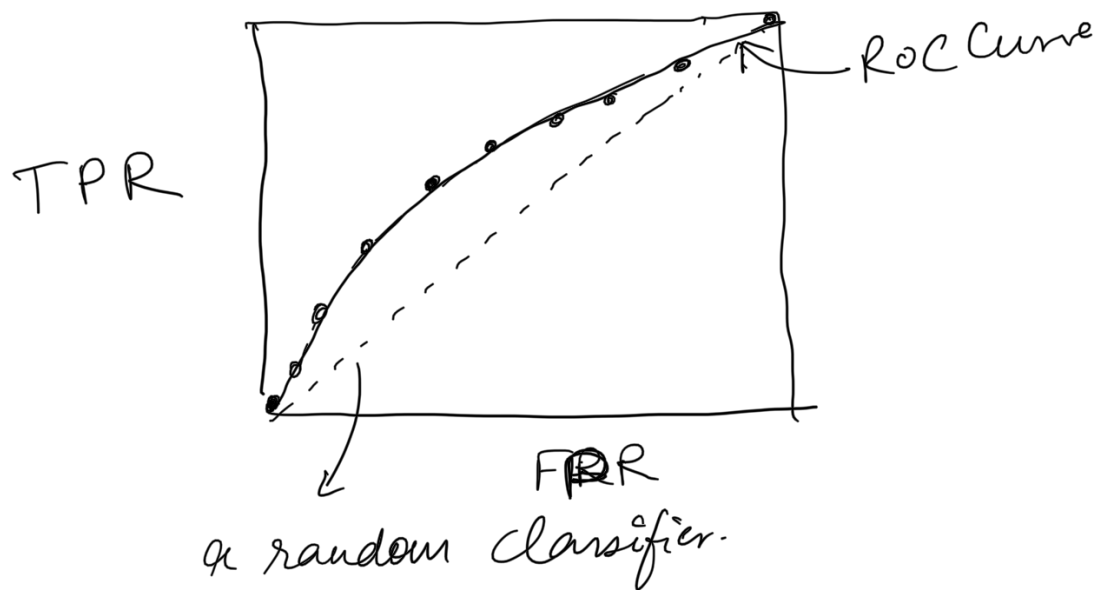
$E[y X=x]$			
X	Y	R_i	$\hat{y} \quad t=0.5$
x_1	1.	0.80	1
x_2	0.	0.52	1

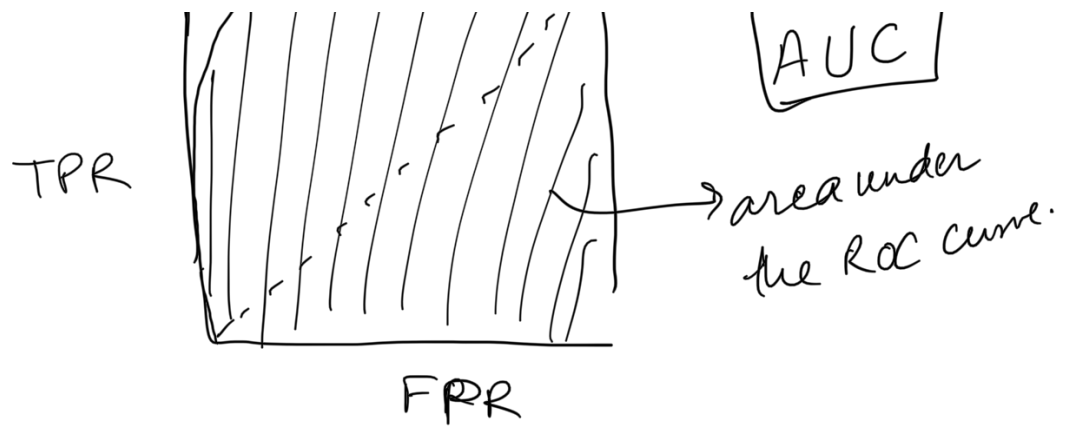
x_3	0	0.47	0
x_4	1	0.60	1
x_5	1	0.65	1
x_6	1	0.39	0
x_7	1	0.49	0
x_8	1	0.80	1
x_9	0	0.05	0
x_{10}	1	0.40	0

d

Receiver Operating Characteristic Curve (ROC)

$t=0$ ~~FPR~~ (true class) $P(\hat{Y}=1|y=0)$
 ~~TPR~~ (true class) $P(\hat{Y}=1|y=1)$





Sensitive attribute A [A-binary]

\perp \rightarrow independence

$$A \perp B \quad P(A, B) = P(A)P(B)$$

$$A \perp B | C \quad P(A, B | C) = P(A | C)P(B | C)$$

Independence

A classifier is fair if:

$$P(\hat{Y}=1 | A=a) = P(\hat{Y}=1 | A=b)$$

Friday April 23

All Gradiance quizzes will be due

on May 9th at 11.59 PM

Gradience 10, 11, 12

What threshold to use?

$$t = 0.5$$

Cost of classification

Profit