

Fairness in Machine Learning

Naive Bayes Classifier:

$$\underline{P(Y=y | X=x) = \frac{P(X=x | Y=y) P(Y=y)}{\sum_{y'} P(X=x | Y=y') P(Y=y')}} \quad \swarrow$$

Both X and Y are random variables.

Logistic Regression

$$\underline{P(Y=y | x)} = \text{Ber}(\sigma(w^T x))$$

Only Y is a random variable

Generative classifiers $\boxed{\underline{P(Y=y | x)}}$

vs.

Discriminative classifiers $\underline{P(Y=y | x)}$

$$\left[\begin{array}{ccc} X & Y & \hat{Y} \end{array} \right]$$



ignore, how to get \hat{y} from x

Measuring performance:

<u>x_i</u>	y_i	<u>\hat{y}_i</u>
x_2	y_2	\hat{y}_2
\vdots	\vdots	\vdots
x_N	y_N	\hat{y}_N

x, y, \hat{y}

$$\frac{\sum_{i=1}^N \mathbb{I}[y_i = \hat{y}_i]}{N} = \text{Accuracy}$$

$$\underline{\text{Accuracy}} = \underline{P(y = \hat{y})}$$

Accuracy is not always a good metric ^

x	y	\hat{y}
x_1	1	1
x_2	1	1
\vdots	\vdots	\vdots
x_{90}	1	1
x_{91}	0	1
x_{92}	0	1
\vdots	\vdots	\vdots
x_{100}	0	1

$$\text{accuracy} = 0.9$$

class imbalance

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

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

x	y	\hat{y}
x_1	1	1
x_2	1	0
x_3	0	0
x_4	0	0
x_5	1	1
x_6	1	0
x_7	0	1
x_8	0	0

$$Y \in \{0, 1\}$$

give a loan - 1
not give a loan - 0

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

x_9	1	1
x_{10}	0	0

$$P(\hat{y}=1 | y=1) = \frac{3}{5} \quad \text{True Positive Rate}$$

recall +ve class

$$P(\hat{y}=0 | y=0) = \frac{4}{5} \quad \text{True negative rate}$$

recall -ve class

$$P(\hat{y}=0 | y=1) = \frac{2}{5} \quad \text{False negative rate}$$

$$P(\hat{y}=1 | y=0) = \frac{1}{5} \quad \text{False positive rate}$$

$$P(y=1 | \hat{y}=1) = \frac{3}{4} \quad \text{Precision +ve class}$$

logistic regression

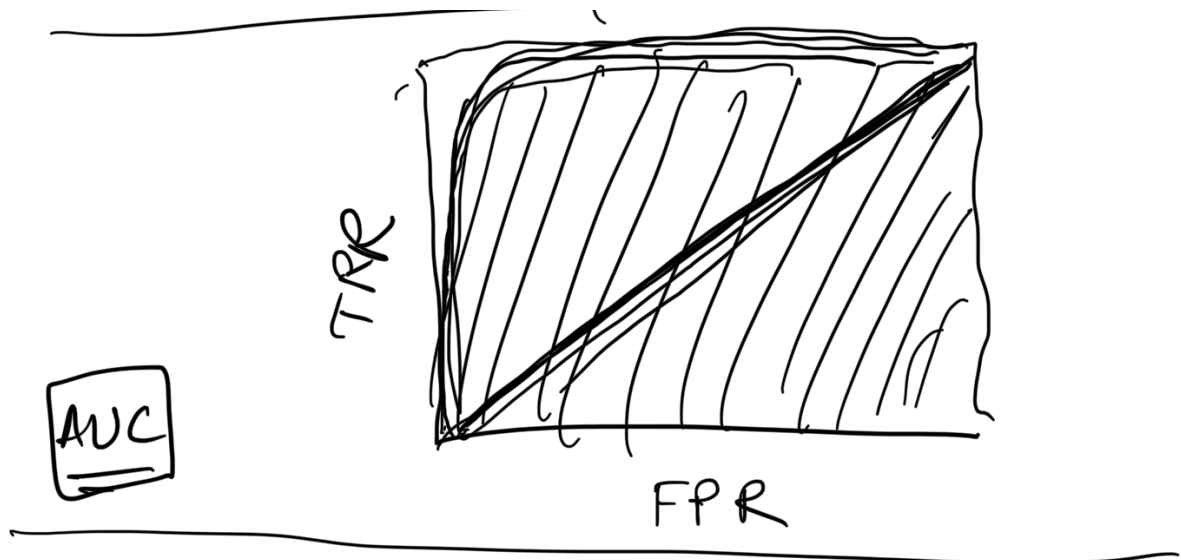
$$r = \sigma(w^T x)$$

$$\sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$

$$\text{If } r \geq \underline{0.5} \quad y = 1$$

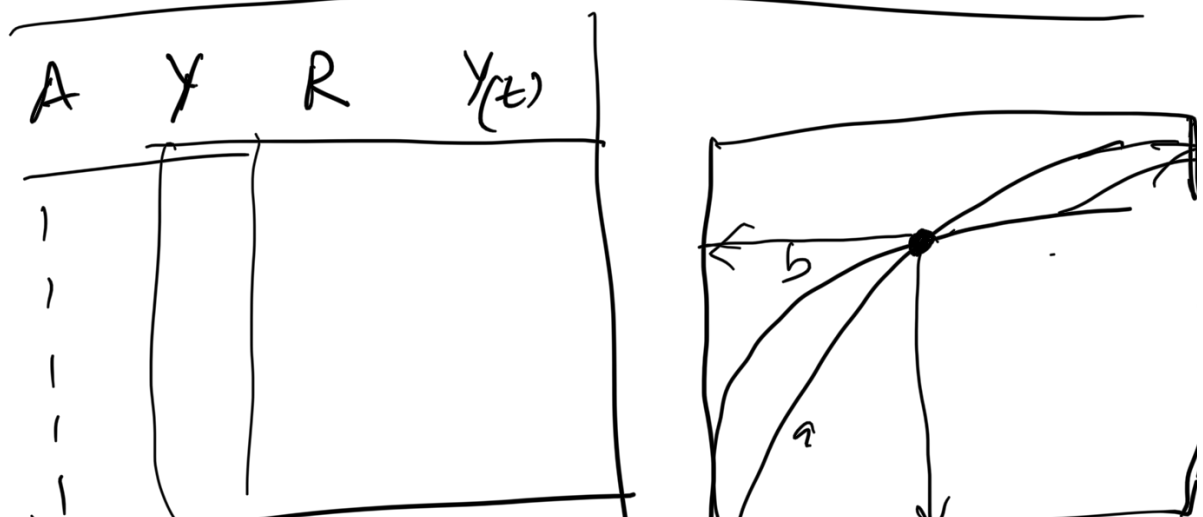
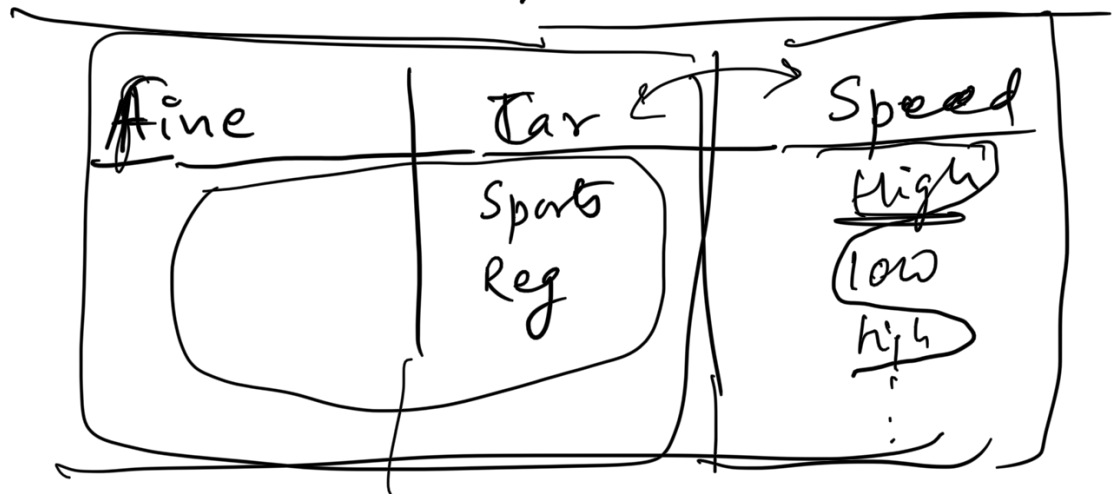
$$r < 0.5 \quad y = 0$$

... 1 ...



Sensitive Attribute **A**

$$A \cap X = \emptyset$$



0	
0	
0	
0	
0	
0	

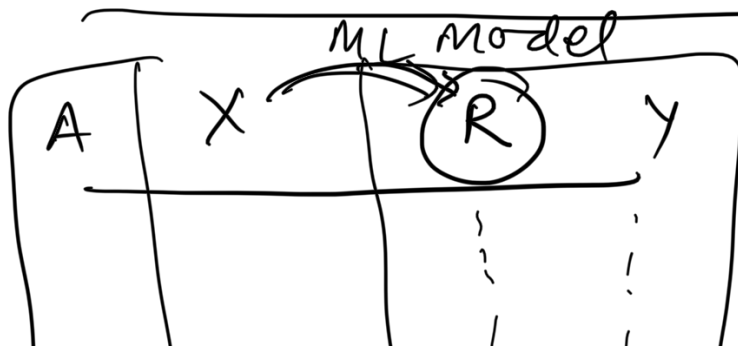
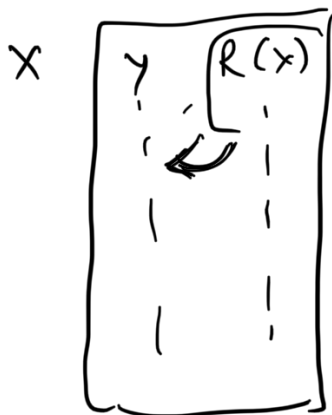
Calibration of Scores

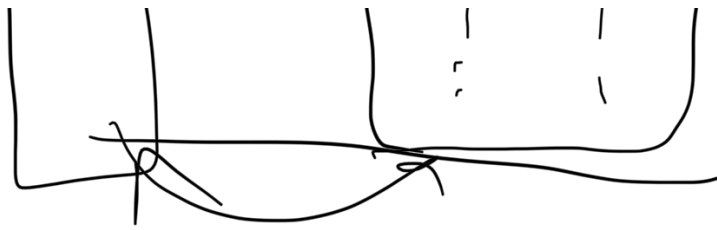
$R(x)$ \rightarrow score

$$P(Y=1 | R=r) = r$$

$w^T x$

Platt Scaling





Strategy to convert $R \rightarrow \hat{y}$

~~tp~~ tp \rightarrow # of applicants given a loan who repaid back

fp \rightarrow # applicants given a loan who did not pay back

Profit -

$$x tp - 6 x fp$$

$$x (tp - 6 fp)$$