

Recall • $(x_1, y_1), \dots, (x_n, y_n)$ with $y_i \in \{0, 1\}$.

- Logistic Regression :-

Model : $P(Y=1 | \underline{x}) = \frac{e^{\beta_0 + \beta^T \underline{x}}}{1 + e^{\beta_0 + \beta^T \underline{x}}} \in (0, 1)$

- Estimate β_0, β via maximum likelihood, $\hat{\beta}_0, \hat{\beta}$.
- For a new feature \underline{x}_0 we predict:

$$\hat{y} = \begin{cases} 1 & \text{if } \frac{e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0}}{1 + e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0}} > \frac{1}{2} \\ 0 & \text{if } \frac{e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0}}{1 + e^{\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0}} < \frac{1}{2} \end{cases}$$

$\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0$

$$\begin{cases} 1 & \text{if } \hat{\beta}_0 + \hat{\beta}^T \underline{x}_0 > 0 \\ 0 & \text{if } \hat{\beta}_0 + \hat{\beta}^T \underline{x}_0 < 0 \end{cases}$$

$\hat{\beta}_0 + \hat{\beta}^T \underline{x}_0 = 0$ decision boundary

Assessing Quality of Model Fit


Given a set of data y_1, \dots, y_n and predictions $\hat{y}_1, \dots, \hat{y}_n$, how can we assess quality of the fit?


Recall $y_i \neq \hat{y}_i \in \{0, 1\}$, think of 1 = email is spam, 0 = email is okay.

$$\bullet \text{ Error rate} = \frac{1}{n} \sum_{i=1}^n \underbrace{1[y_i \neq \hat{y}_i]}_{\substack{\text{indicator} \\ = \begin{cases} 1 & y_i \neq \hat{y}_i \\ 0 & y_i = \hat{y}_i \end{cases}}} = \text{Percent misclassified}$$

- Confusion matrix

	True 1s	True 0s
Predicted 1s	10	4
Predicted 0s	2	12


 10 = predicted 10 spams that were indeed spam


 predicted 4 spam that were not spam

- sensitivity = % of true positives $\left(\frac{10}{12}\right)$
- specificity = % of true negatives $\left(\frac{12}{16}\right)$
- positive predictive value = % of predicted 1s correct $\left(\frac{10}{14}\right)$
- negative " " = % of predicted 0s correct $\left(\frac{12}{14}\right)$

- error rate = $\frac{4+2}{28}$

- False positive rate (FPR) = 1-specificity $\left(\frac{4}{16} \right)$

- True positive rate (TPR) = sensitivity $\left(\frac{10}{12} \right)$

Note

- FPR = good emails that go to spam

- TPR = spam " " " " "

Goal

try to minimize FPR and simultaneously maximize TPR

Note

As a baseline comparison, compare against random guessing

Ex Dataset w/ 100 positives + 1000 negatives

Guess 1 with 90% probability

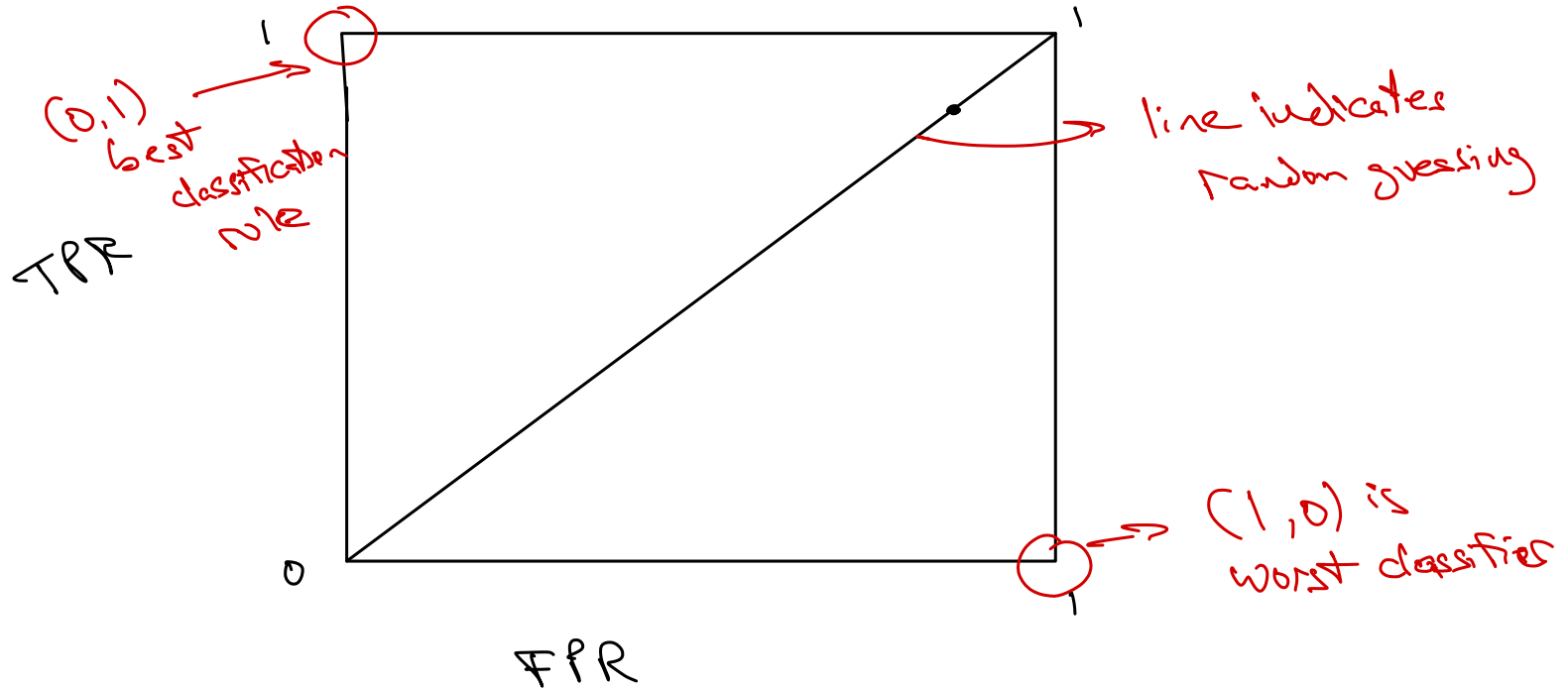
	true 1s	true 0s
pred 1s	90	900
pred 0s	10	100

$$\bullet \text{ FPR} = \frac{900}{1000} = 90\%$$

$$\bullet \text{ TPR} = \frac{90}{100} = 90\%$$

DEF A receiver operating characteristic (ROC) graph

Plots the TPR (y-axis) against FPR (x-axis)



For any given model + classification/decision rule,
we get a single point in ROC space.

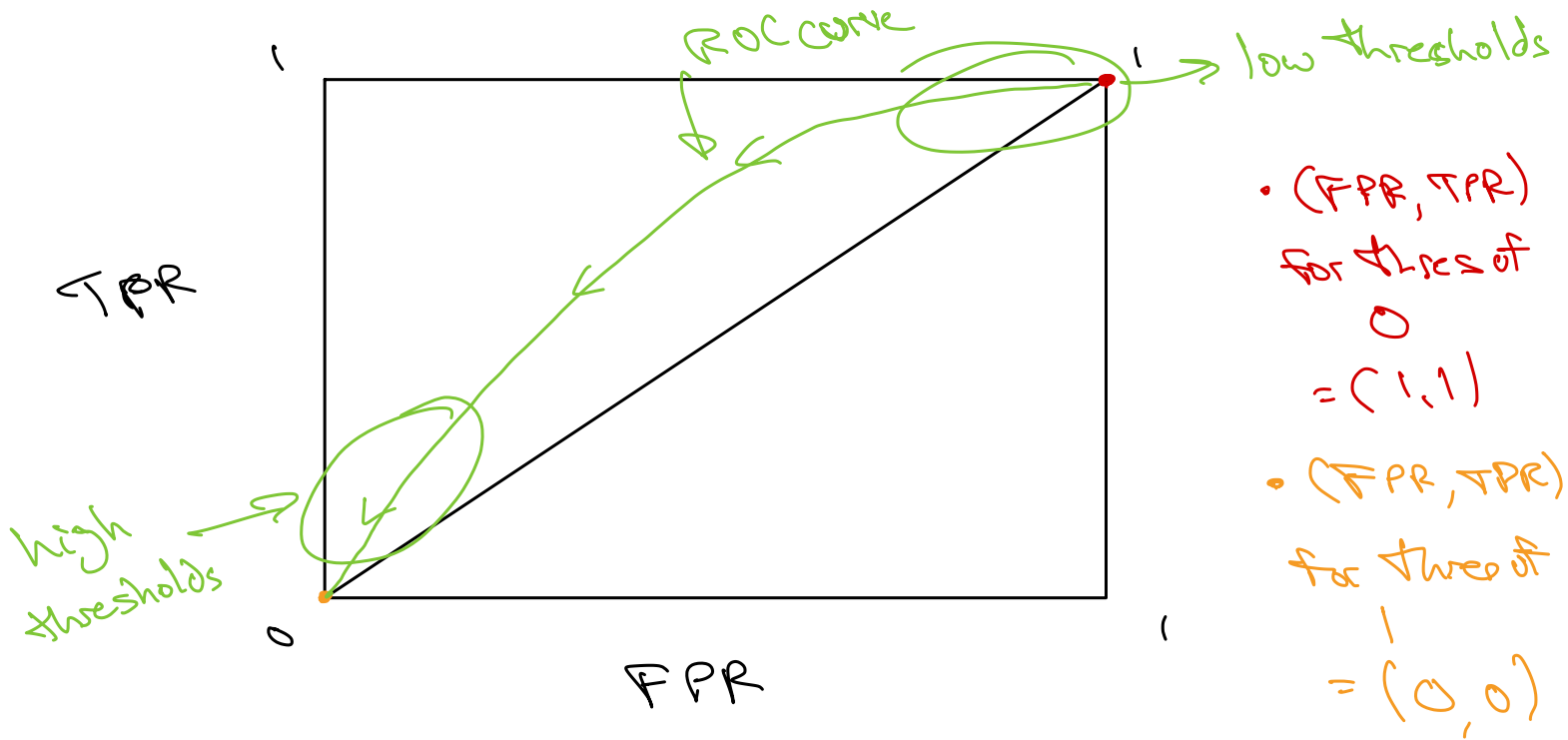
Note

Usual classification rule is $\hat{y} = 1$ if $\hat{p} > \frac{1}{2}$.

What if we vary $\frac{1}{2}$?

Ex: $\hat{y} = 1$ if $\hat{p} > \underline{0.75}$

As we vary the classification threshold from 0 to 1,
we sweep out a set of points in ROC space
that is called a ROC curve.



Note One numeric summary of quality of a model fit is the area under curve (AUC) $\in [0, 1]$

where goodness of fit diag.

In regression we used R^2 as a goodness-of-fit measure,
in logistic regression models the equivalent idea is
deviance :

$$G^2 = 2 \sum_{i=1}^n \left[y_i \log \left(\frac{y_i}{\hat{p}(x_i)} \right) + (1 - y_i) \log \left(\frac{1 - y_i}{1 - \hat{p}(x_i)} \right) \right]$$

The i th deviance residual is

$$\text{dev}_i = \pm \left[-2 \left[y_i \log \hat{p}(x_i) + (1 - y_i) \log (1 - \hat{p}(x_i)) \right] \right]^{1/2}$$

where \pm is $+$ if $y_i \geq \hat{p}(x_i)$

Think

$$y_i = 1$$

$$\bullet \quad p \approx 1 \quad \Rightarrow \text{der} \approx 0$$

$$\bullet \quad p^2 \approx 0 \quad \Rightarrow \text{der} \approx 0$$