

# DS 102 Discussion 1

Wednesday, September 2, 2020

1. **ROC Curves.** In lecture we defined and discussed ROC curves, or “receiver operating characteristic” curves. ROC curves plot the true positive rate (TPR) and false positive rates (FPR) for a binary classifier at different decision thresholds. Recall that the TPR and FPR are defined as:

$$\text{TPR} = \frac{\# \text{ true positives}}{\# \text{ positives}}, \quad \text{FPR} = \frac{\# \text{ false positives}}{\# \text{ negatives}},$$

where “true positives” are examples where the model made a positive decision and the label was positive, and “positives” are examples where the label was positive.

In this exercise, we will consider the ROC curve on an example dataset. Let  $Y$  be the label,  $X_1, X_2$  be features, and consider the model function  $f(X_1, X_2) = 3X_1 + 2X_2 + 1$ .

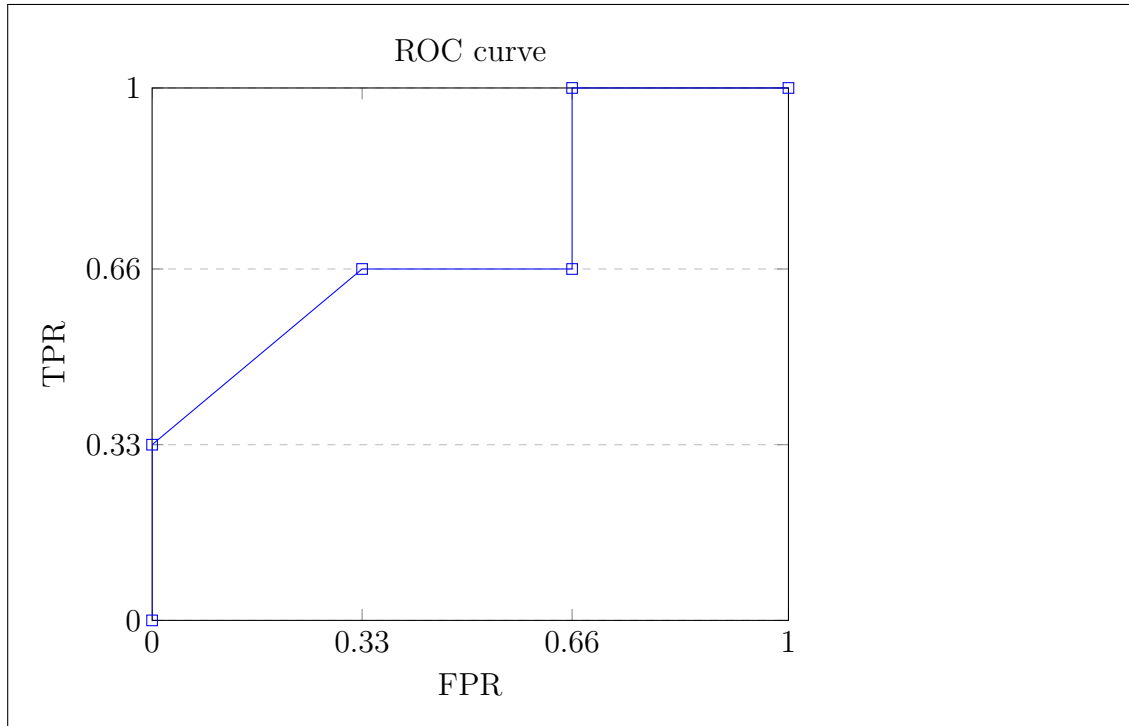
Table 1: Example dataset

$Y$	$f(X_1, X_2)$	$X_1$	$X_2$
0	-1	-1	0.5
1	-0.5	-1	0.75
0	0	-1	1
1	1	0.2	-0.3
1	0.25	-0.25	0
0	0.25	-0.05	-0.3

- (a) Plot the ROC curve for the model  $f(X_1, X_2)$  with respect to the label  $Y$ .

**Solution:** At a given decision threshold  $\alpha$ , if  $f(X_1, X_2) > \alpha$ , then the decision is a positive classification, and if  $f(X_1, X_2) \leq \alpha$ , then the decision is a negative classification. Since the model  $f(X_1, X_2)$  only takes five different values on this dataset, only five different decision thresholds lead to different true and false positive rates.

- $\alpha < -1$ : TPR = 1, FPR = 1
- $-1 \leq \alpha < -0.5$ : TPR = 1, FPR =  $\frac{2}{3}$
- $-0.5 \leq \alpha < 0$ : TPR =  $\frac{2}{3}$ , FPR =  $\frac{2}{3}$
- $0 \leq \alpha < 0.25$ : TPR =  $\frac{2}{3}$ , FPR =  $\frac{1}{3}$
- $0.25 \leq \alpha < 1$ : TPR =  $\frac{1}{3}$ , FPR = 0
- $1 \leq \alpha$ : TPR = 0, FPR = 0



- (b) Suppose that we can choose two decision thresholds  $\alpha_1$  and  $\alpha_2$ , and for each data example, we flip a coin to decide which decision threshold to use for that example. Choose  $\alpha_1$  and  $\alpha_2$ , and probabilities for using  $\alpha_1$  and  $\alpha_2$ , such that in expectation, the true positive rate is  $\frac{1}{3}$  and the false positive rate is  $\frac{1}{3}$ .

**Solution:** If we choose a decision threshold  $\alpha_1 = -0.5$ , the TPR is  $\frac{2}{3}$  and the FPR is  $\frac{2}{3}$ . If we choose a decision threshold of  $\alpha_2 = 1.0$ , the TPR is 0 and the FPR is 0. If for each data point, we choose the decision threshold  $\alpha = \alpha_1$  with probability  $\frac{1}{2}$  and  $\alpha = \alpha_2$  with probability  $\frac{1}{2}$ , then the expected TPR is  $\frac{1}{2} * \frac{2}{3} + \frac{1}{2} * 0 = \frac{1}{3}$ . Likewise, the expected FPR is  $\frac{1}{2} * \frac{2}{3} + \frac{1}{2} * 0 = \frac{1}{3}$ .

To see this calculation broken down explicitly:

$$\begin{aligned}
 E_{\alpha}[\text{TPR}] &= \frac{E_{\alpha}[\sum_{i=1}^n \mathbf{1}(f(X_1^i, X_2^i) > \alpha)]}{\sum_{i=1}^n \mathbf{1}(Y^i > 0)} \\
 &= \frac{\frac{1}{2} \sum_{i=1}^n \mathbf{1}(f(X_1^i, X_2^i) > \alpha_1) + \frac{1}{2} \sum_{i=1}^n \mathbf{1}(f(X_1^i, X_2^i) > \alpha_2)}{\sum_{i=1}^n \mathbf{1}(Y^i > 0)} \\
 &= \frac{\frac{1}{2} * 2}{3} + \frac{\frac{1}{2} * 0}{3} \\
 &= \frac{1}{3}
 \end{aligned}$$

Note that by choosing probabilities for using each decision threshold  $\alpha_1$  and  $\alpha_2$ , the expected TPR and FPR are a **convex combination** of the TPRs and

FPRs of  $\alpha_1$  and  $\alpha_2$  individually.

- (c) Is it possible to choose two decision thresholds  $\alpha_1, \alpha_2$  and probabilities of using each decision threshold such that the expected true positive rate is  $\frac{1}{3}$ , and the expected false positive rate is  $\frac{2}{3}$ ?

**Solution:** No, this is not possible for this dataset. No single threshold or combination of two thresholds can achieve a false positive rate of  $\frac{2}{3}$  without also increasing the true positive rate up to  $\frac{2}{3}$  (try it!). Specifically, any (TPR, FPR) pair that is *not* a convex combination of the (TPR, FPR) pairs for single thresholds (plotted on the ROC curve) cannot be achieved by choosing two decision thresholds probabilistically.

2. **Hypothesis Testing.** As discussed in lecture, one can imagine different metrics for quantifying how “good” a decision is. For example, we would like our decisions to have both high true positive rate and low false positive rate. Our goal as statisticians is to develop reasonable strategies for doing well on both metrics. In other words, how should we pick a point on the ROC curve? Once we pick a point, how do we achieve it?

The Neyman-Pearson Lemma offers one solution. To be concrete, we focus on the case of hypothesis testing. We call the probability of a false positive under null hypothesis  $H_0$  the *significance level*  $\alpha$  of a test, and we call the probability of a true positive under the alternative hypothesis  $H_1$  the *power* of a test.

The Neyman-Pearson formulation prescribes the following point on the ROC curve: fix a significance level you are willing to tolerate, then pick the point that maximizes power. The Neyman-Pearson Lemma prescribes how to achieve this point:

**Lemma (Neyman & Pearson, 1933)** *Suppose  $\theta_1 < \theta_0$ . For any significance level  $\alpha \in [0, 1]$ , the following likelihood-ratio test maximizes power among all tests with level at most  $\alpha$ :*

$$\delta(x) = \begin{cases} \text{Reject Null} & : \frac{f_{\theta_0}(x)}{f_{\theta_1}(x)} \leq \eta \\ \text{Accept Null} & : \frac{f_{\theta_0}(x)}{f_{\theta_1}(x)} > \eta \end{cases}$$

where  $f_{\theta_0}, f_{\theta_1}$  are the likelihoods under the null and alternative distributions, respectively, and  $\eta$  is the real value such that  $\Pr(\delta(X) = 1 \mid H_0) = \alpha$ .

**Example.** Suppose that you have a sample from a distribution with probability density function  $f_{\theta}(x) = \theta x^{\theta-1}$  where  $0 < x < 1$ . You would like to design a test to discern between the null hypothesis that  $\theta = 4$ , and the alternative hypothesis that  $\theta = 3$ .

- (a) Derive the most powerful test for this problem such that the significance level is less than  $\alpha$ .

**Solution:** Leveraging the Neyman-Pearson Lemma, we design a likelihood-ratio test. The likelihood ratio has the form:

$$\frac{f_{\theta_0}(X)}{f_{\theta_1}(X)} = \frac{4x^3}{3x^2} = \frac{4x}{3}.$$

Now we need to solve for  $\eta$  such that the significance level is  $\alpha$ , or

$$Pr(x \leq \frac{3}{4}\eta \mid H_0) = \alpha.$$

That is, we need

$$\int_0^{0.75\eta} f_{\theta_0}(x)dx = \int_0^{0.75\eta} 4x^3 dx = \alpha.$$

Solving for this gives

$$0.75^4 \eta^4 = \alpha$$

which yields  $\eta = \frac{4}{3}\alpha^{0.25}$ .

(b) What is the power of the test,  $Pr(\delta(X) = 1 \mid H_1)$ ?

**Solution:** Pattern matching from above, we need to calculate  $Pr(x \leq \alpha^{0.25} \mid H_1)$ . More explicitly,

$$\begin{aligned} Pr\left(\frac{f_{\theta_0}(X)}{f_{\theta_1}(X)} \leq \eta \mid H_1\right) &= Pr\left(\frac{4x}{3} \leq \frac{4}{3}\alpha^{0.25} \mid H_1\right) \\ &= Pr(x \leq \alpha^{0.25} \mid H_1) \\ &= \int_0^{\alpha^{0.25}} 3x^2 dx \\ &= \alpha^{\frac{3}{4}} \end{aligned}$$