

$$\hat{p} \rightarrow \hat{y} \quad \hat{y} = \mathbb{1}_{\hat{p} \geq p_{th}}$$

confusion table

	$\hat{y}$	0	1	
$y$	0	TN	FP	$N$
	1	FN	TP	$P$
		$P_N$	$P_P$	$n$

these three quantities are static i.e. not a function of your model

$$\text{Sensitivity} = \text{Recall} := TP / P \in [0, 1]$$

$$\text{Specificity} := TN / N \in [0, 1]$$

tradeoff

$$\text{False Positive Rate (FPR)} := FP / N = 1 - \text{Specificity}$$

Consider a "dumb" probability estimation model  $g_{pr}(xvec)$  that produces random guesses from  $U(0,1)$ .  $\hat{y} = \mathbb{1}_{g_{pr}(\tilde{x}) \geq p_{th}}$

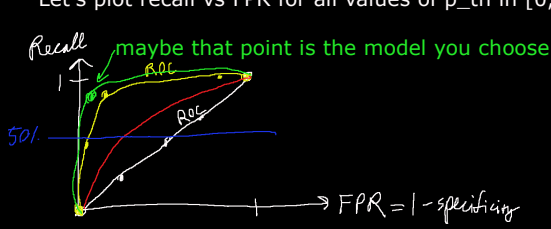
What does the confusion table for this dumb model look like?

	$\hat{y}$	0	1	
$y$	0	$TN = p_{th} N$	$FP = (1 - p_{th}) N$	$N$
	1	$FN = p_{th} P$	$TP = (1 - p_{th}) P$	$P$
		$P_N = p_{th} n$	$P_P = (1 - p_{th}) n$	$n$

$$\text{Recall} = \frac{TP}{P} = \frac{(1 - p_{th}) P}{P} = 1 - p_{th}$$

$$\text{FPR} = \frac{FP}{N} = \frac{(1 - p_{th}) N}{N} = 1 - p_{th}$$

Let's plot recall vs FPR for all values of  $p_{th}$  in  $[0, 1]$ :



In the dumb model,  $y$  is independent of  $\hat{p}$ . In a smarter model, then  $g_{pr}$  is connected to  $y$  thus you can have high recall with high specificity. If

$g_{pr} = \hat{y} = 1$  always  $\Rightarrow$  recall = 1 but specificity = 0 and

$g_{pr} = \hat{y} = 0$  always  $\Rightarrow$  specificity = 1 but recall = 0.

Tracing out the performance (recall and specificity) for all  $p_{th}$  in  $[0, 1]$  gives you the yellow line which is named the "receiver-operator curve" (ROC). The ROC can be used to compare probability estimation models by calculating the area under the curve (AUC) also called the "c-statistic":

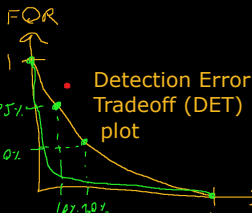
$$AUC = \int_0^1 ROC dFPR \in [0, 1]$$

An  $AUC > 1/2$  has predictive power. AUC's closer and closer to 1 indicate better and better models. AUC is kind of a scoring rule (e.g. Brier).

AUC is a metric that gauges the overall fit of a probability estimation model. It is not a metric that gauges performance of a classification model (i.e. one of the point on the ROC curve).

My opinion is that ROC (and AUC) is silly in a prediction context since when you use your model to predict the future, you don't know what the  $y$ 's are nor  $P$  and  $N$ . What's critical is your prediction errors (in the columns: FDR vs FOR).

	$\hat{y}$	0	1	
$y$	0	TN	FP	
	1	FN	TP	
		$P_N$	$P_P$	



$$FOR = \frac{FN}{PN} = \frac{p_{th} P}{p_{th} n} = p_y, \text{const}$$

$$FDR = \frac{FP}{PP} = \frac{(1 - p_{th}) N}{(1 - p_{th}) n} = 1 - p_y, \text{const}$$

$$\hat{y} = 1 \text{ always} \Rightarrow FDR = 1 - p_y$$

$$\hat{y} = 0 \text{ always} \Rightarrow FOR = p_y$$

Thus, there's no nice "reference line"

The DET is traced out by varying  $p_{th}$  in  $[0, 1]$ .

Once again, your classification model will be one point on the DET curve (just like it's one point on the ROC curve) but I can visually see the tradeoff of the two errors that are critical to prediction.

Bias-Variance Tradeoff in Regression Modeling,  $y \in \mathbb{R}$

$$y = f(\tilde{x}) + \delta, \text{ and } xvec \text{ is a constant}$$

Let's assume (I)  $\delta$  is a realization from  $\Delta$  a r.v. which is mean-independent from  $xvec$  and has expectation 0 i.e.

$$E[\Delta | \tilde{x}] = 0 \Rightarrow E[\Delta] = 0$$

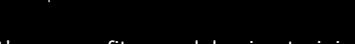
Makes sense since we defined  $f$  to be a function that extracts all useful information from  $x$  regarding  $y$ . This means  $y$  is also a realization from r.v.  $Y$ :

$$Y = f(\tilde{x}) + \Delta \Rightarrow E[Y | X = \tilde{x}] = E[f(\tilde{x}) + \Delta | X = \tilde{x}] = f(\tilde{x}) + E[\Delta | X = \tilde{x}] = f(\tilde{x})$$

In textbooks,  $f$  is called the "conditional expectation function" (CEF).

Let's assume (II) homoskedasticity i.e. constant variance.

$$\sigma^2 = \text{Var}[\Delta | X = \tilde{x}] = E[\Delta^2 | X = \tilde{x}] - E[\Delta | X = \tilde{x}]^2 = E[\Delta^2 | X = \tilde{x}]$$



Let's say we fit a model using training data  $D$  to get model  $g$ :

$$y = g + e = g + (f - g) + \delta \Rightarrow e = f - g + \delta$$

$$\stackrel{I, II}{\Rightarrow} Y = g + (f - g) + \Delta \Rightarrow E = f - g + \Delta$$

Let's say we now predict for the next observation  $x_*$ .

$$Y_* = g(x_*) + (f(x_*) - g(x_*)) + \Delta_* \Rightarrow E_* = f - g + \Delta_*$$

$$\begin{aligned} \text{Bias}(\tilde{x}_*) &:= E[Y_* - g(\tilde{x}_*) | \tilde{x}_*] = E[E_* | \tilde{x}_*] = E[f - g + \Delta_* | \tilde{x}_*] \\ &= f - g + E[\Delta_* | \tilde{x}_*] \stackrel{I}{=} f(\tilde{x}_*) - g(\tilde{x}_*) \end{aligned}$$

$$\text{MSE}(\tilde{x}_*) = E[(Y_* - g(\tilde{x}_*))^2 | \tilde{x}_*] = E[(Y_* - f(\tilde{x}_*))^2 | \tilde{x}_*] = E[\Delta_*^2 | \tilde{x}_*] = \sigma^2$$

If our model is "perfect" i.e. no misspecification nor estimation error, then the MSE is  $\sigma^2$ , the irreducible squared error due to ignorance. If our model is not perfect,  $g(\tilde{x}_*) \neq f(\tilde{x}_*)$  then we show  $\text{MSE} \geq \sigma^2$

$$\begin{aligned} &= E[Y_*^2 - 2Y_* g(\tilde{x}_*) + g(\tilde{x}_*)^2 | \tilde{x}_*] = E[(f + \Delta_*)^2] - 2g(\tilde{x}_*) E[f + \Delta_*] + g(\tilde{x}_*)^2 \\ &= f^2 + 2f E[\Delta_*] + E[\Delta_*^2] - 2g f + g^2 = \sigma^2 + f^2 - 2g f + g^2 \\ &= \sigma^2 + (f - g)^2 = \sigma^2 + \text{Bias}(\tilde{x}_*)^2 > \sigma^2 \end{aligned}$$

This calculation above assumed a fixed dataset  $D$  with fixed observations  $\{<xvec_1, y_1>, \dots, <xvec_n, y_n>\}$  and one fixed  $xstar$ . Now, let's consider random dataset  $D$ 's. In these random datasets,  $Y_1, \dots, Y_n$  are rvs (but  $xvec_1, \dots, xvec_n$  are still constants) due to  $\Delta_1, \dots, \Delta_n$  being rvs (i.e.  $\delta_1, \dots, \delta_n$  are realizations creating the random  $y_1, \dots, y_n$ ).