

Cost-Benefit Matrix

		Actual	
		p	n
Predicted	Y	$b(Y,p)$	$c(Y,n)$
	N	$c(N,p)$	$b(N,n)$

Cost-Benefit Matrix

		Actual	
		p	n
Predicted	Y	$b(Y,p)$	$c(Y,n)$
	N	$c(N,p)$	$b(N,n)$

$$\begin{matrix} & \text{Actual} \\ \text{Predicted } Y & \begin{pmatrix} p & n \\ 99 & -1 \\ 0 & 0 \end{pmatrix} \\ N & \end{matrix}$$

$$\text{Expected profit} = p(\mathbf{Y}, \mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N}, \mathbf{p}) \cdot b(\mathbf{N}, \mathbf{p}) + \\ p(\mathbf{N}, \mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y}, \mathbf{n}) \cdot b(\mathbf{Y}, \mathbf{n})$$

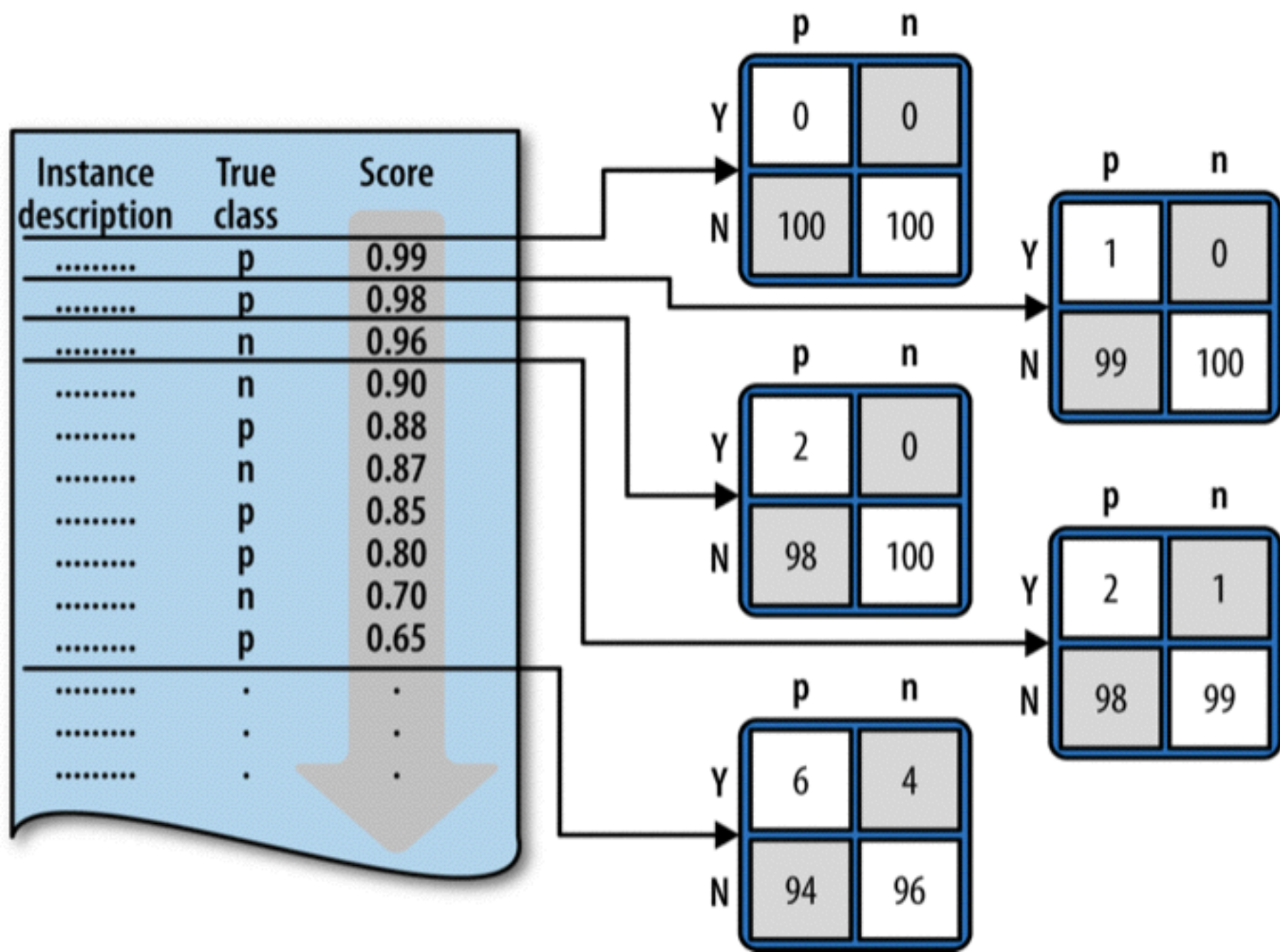
$$p(x, y) = p(y) \cdot p(x \mid y)$$

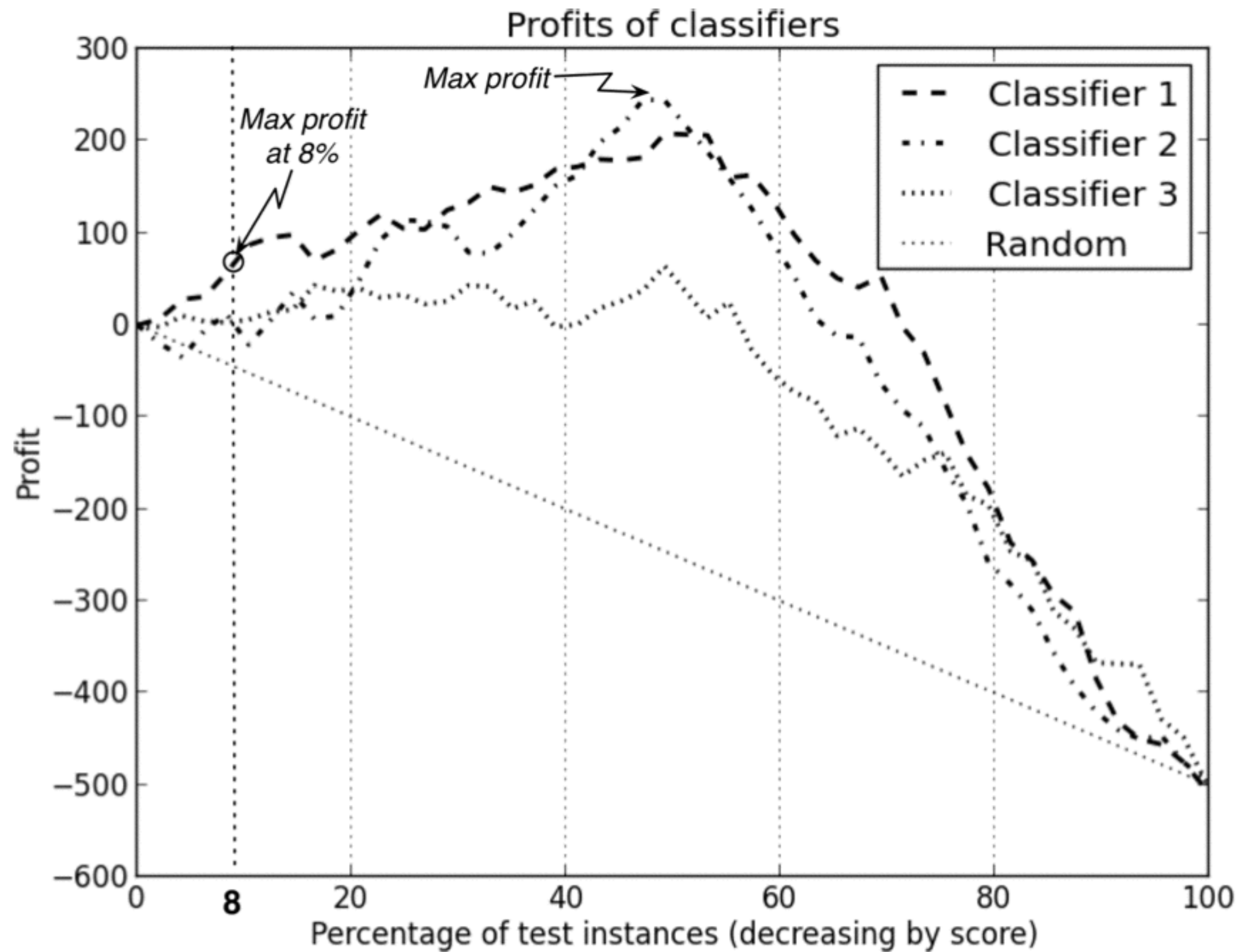
$$\text{Expected profit} = p(\mathbf{Y} \mid \mathbf{p}) \cdot p(\mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N} \mid \mathbf{p}) \cdot p(\mathbf{p}) \cdot b(\mathbf{N}, \mathbf{p}) + \\ p(\mathbf{N} \mid \mathbf{n}) \cdot p(\mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y} \mid \mathbf{n}) \cdot p(\mathbf{n}) \cdot b(\mathbf{Y}, \mathbf{n})$$

$$\text{Expected profit} = p(\mathbf{p}) \cdot [p(\mathbf{Y} \mid \mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N} \mid \mathbf{p}) \cdot c(\mathbf{N}, \mathbf{p})] + \\ p(\mathbf{n}) \cdot [p(\mathbf{N} \mid \mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y} \mid \mathbf{n}) \cdot c(\mathbf{Y}, \mathbf{n})]$$

		Actual	
Predicted	Y	p	n
	N	99	-1
		0	0

$$\begin{aligned}
\text{expected profit} &= p(\mathbf{p}) \cdot [p(\mathbf{Y} \mid \mathbf{p}) \cdot b(\mathbf{Y}, \mathbf{p}) + p(\mathbf{N} \mid \mathbf{p}) \cdot c(\mathbf{N}, \mathbf{p})] + \\
&\quad p(\mathbf{n}) \cdot [p(\mathbf{N} \mid \mathbf{n}) \cdot b(\mathbf{N}, \mathbf{n}) + p(\mathbf{Y} \mid \mathbf{n}) \cdot c(\mathbf{Y}, \mathbf{n})] \\
&= 0.55 \cdot [0.92 \cdot b(\mathbf{Y}, \mathbf{p}) + 0.08 \cdot b(\mathbf{N}, \mathbf{p})] + \\
&\quad 0.45 \cdot [0.86 \cdot b(\mathbf{N}, \mathbf{n}) + 0.14 \cdot p(\mathbf{Y}, \mathbf{n})] \\
&= 0.55 \cdot [0.92 \cdot 99 + 0.08 \cdot 0] + \\
&\quad 0.45 \cdot [0.86 \cdot 0 + 0.14 \cdot -1] \\
&= 50.1 - 0.063 \\
&\approx \mathbf{\$50.04}
\end{aligned}$$

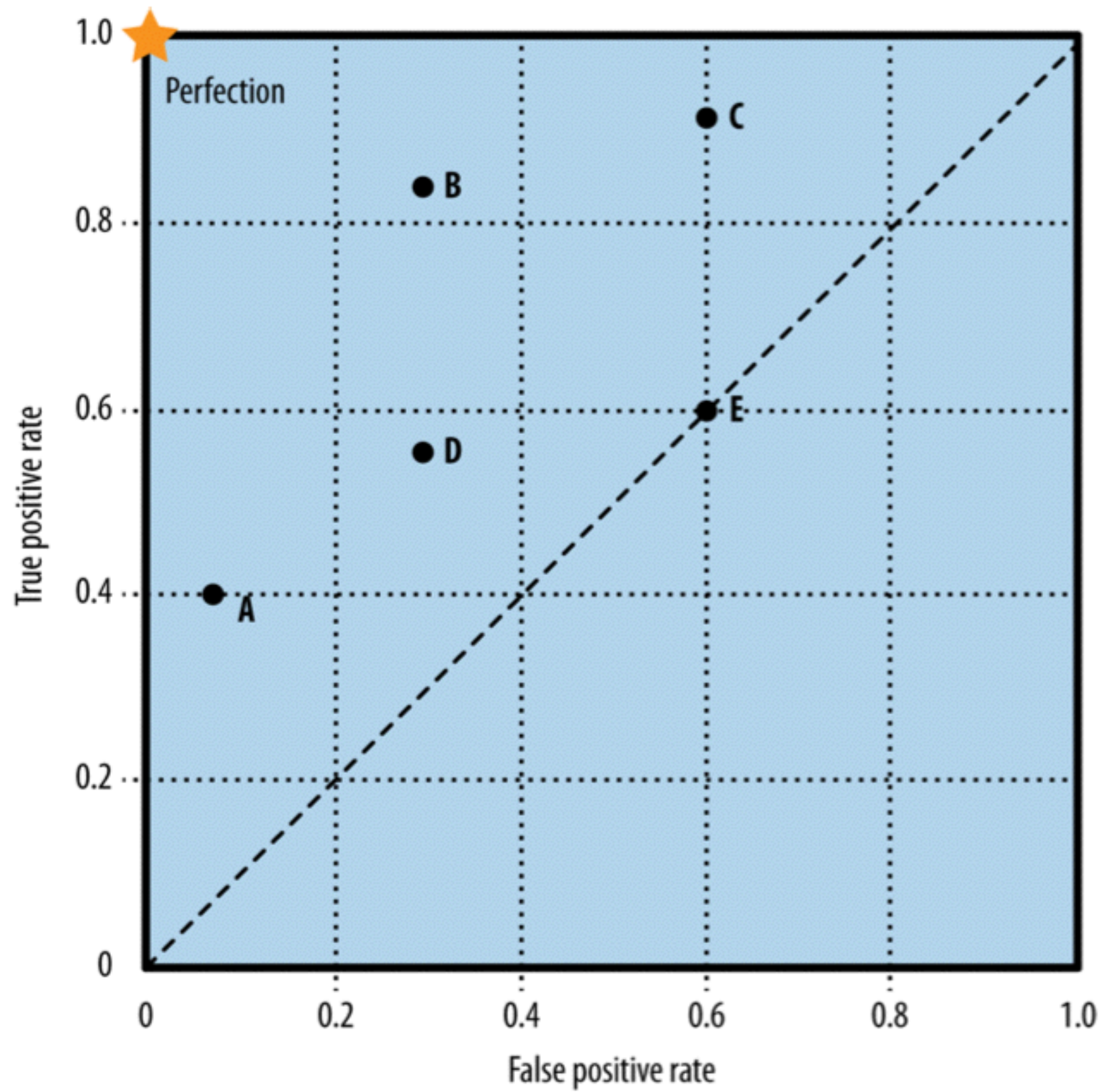




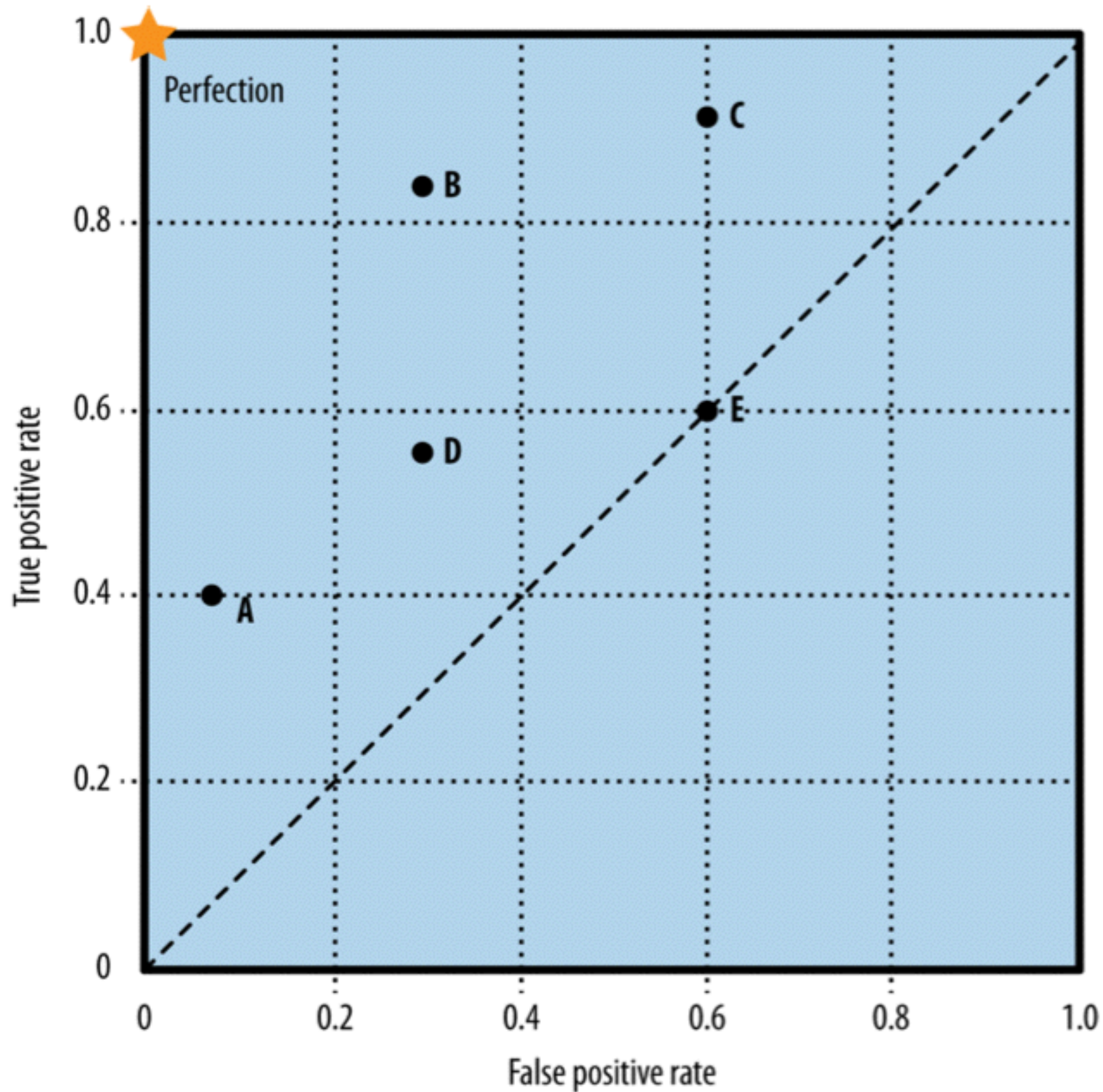
ROC Graphs & Curves

There are 2 critical conditions underlying the profit calculation

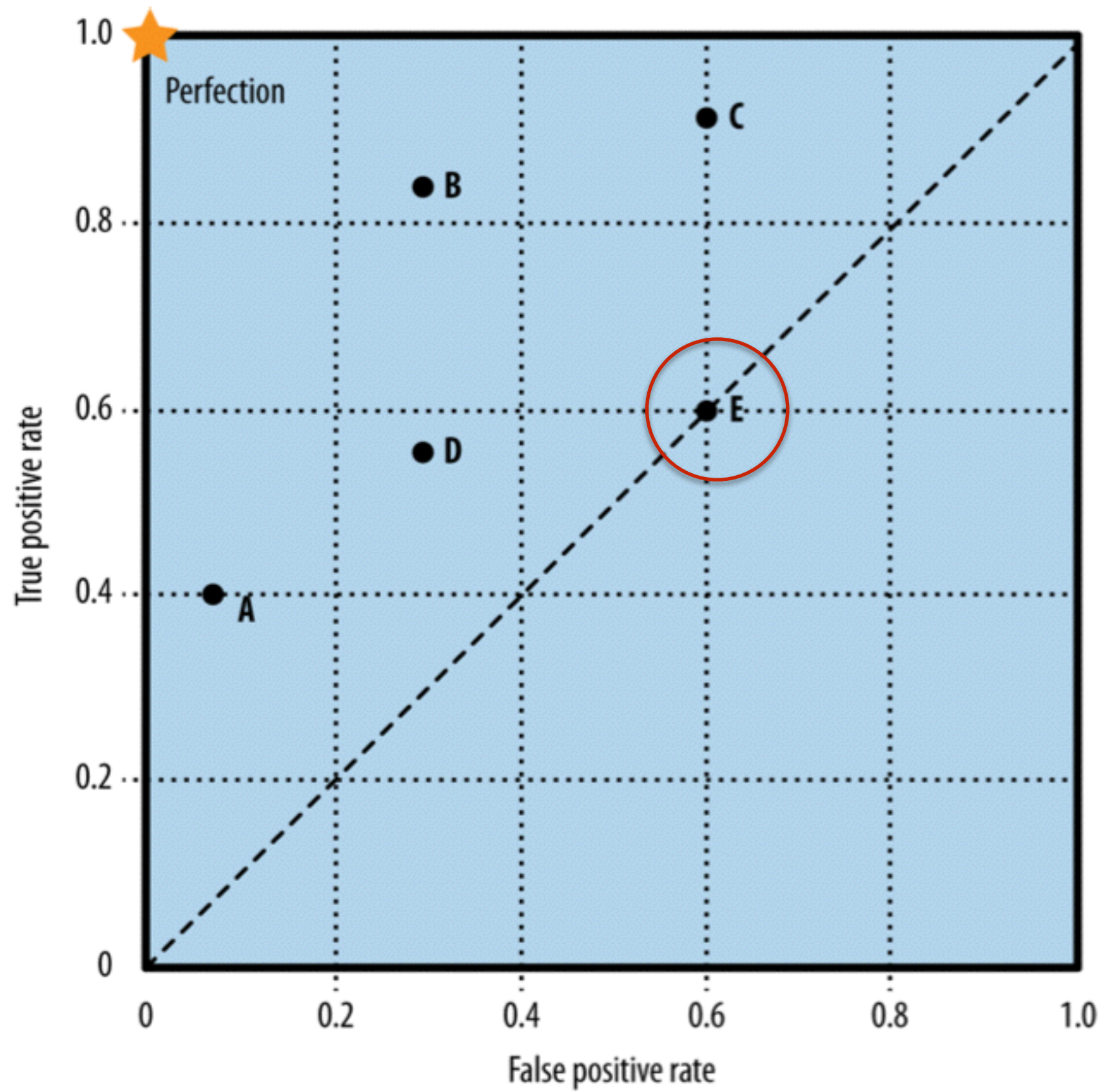
- The *class priors*; that is, the proportion of positive and negative instances in the target population, also known as the *base rate* (usually referring to the proportion of positives)
- The *costs and benefits*. The expected profit is specifically sensitive to the relative levels of costs and benefits for the different cells of the cost-benefit matrix.

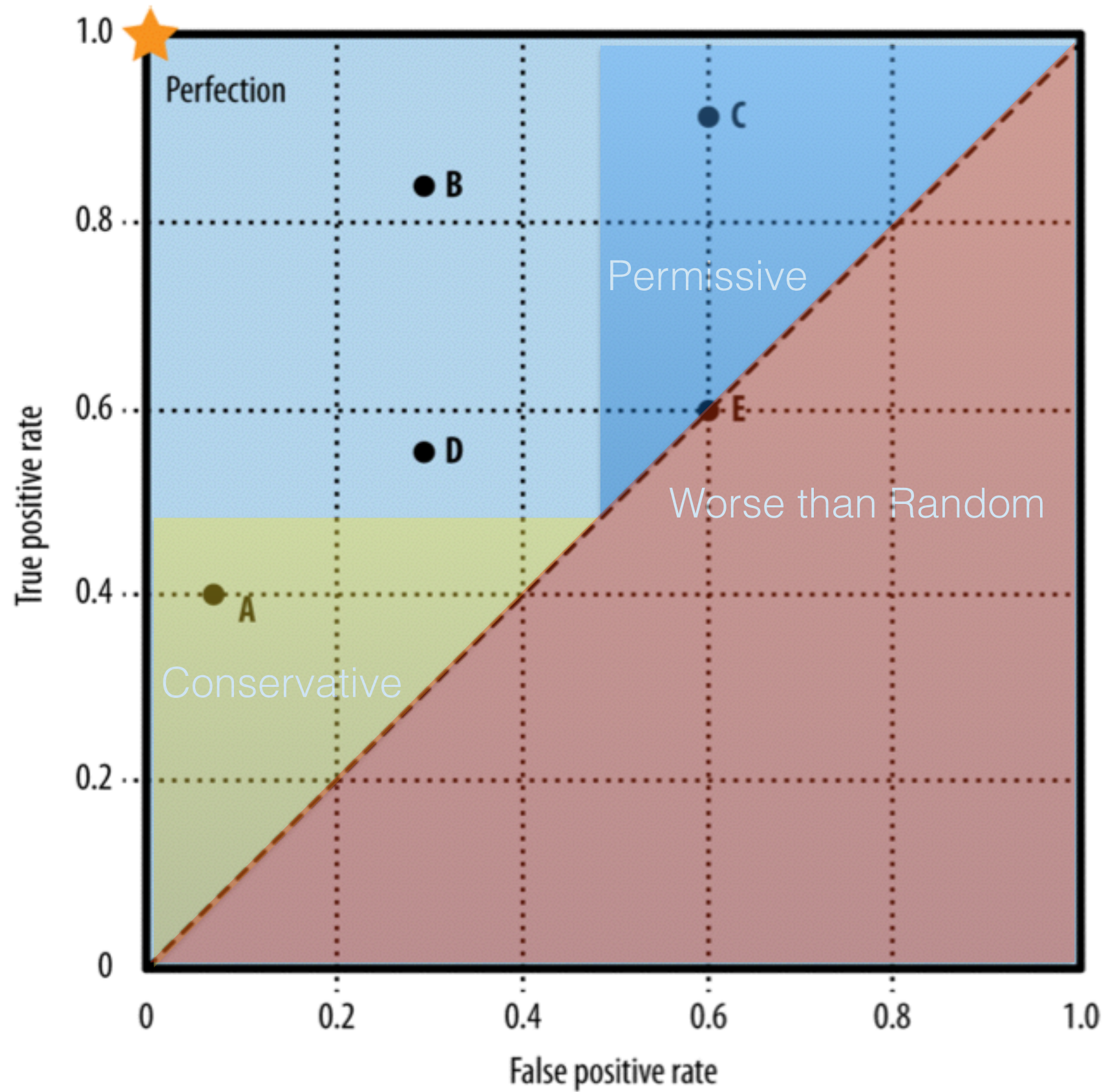


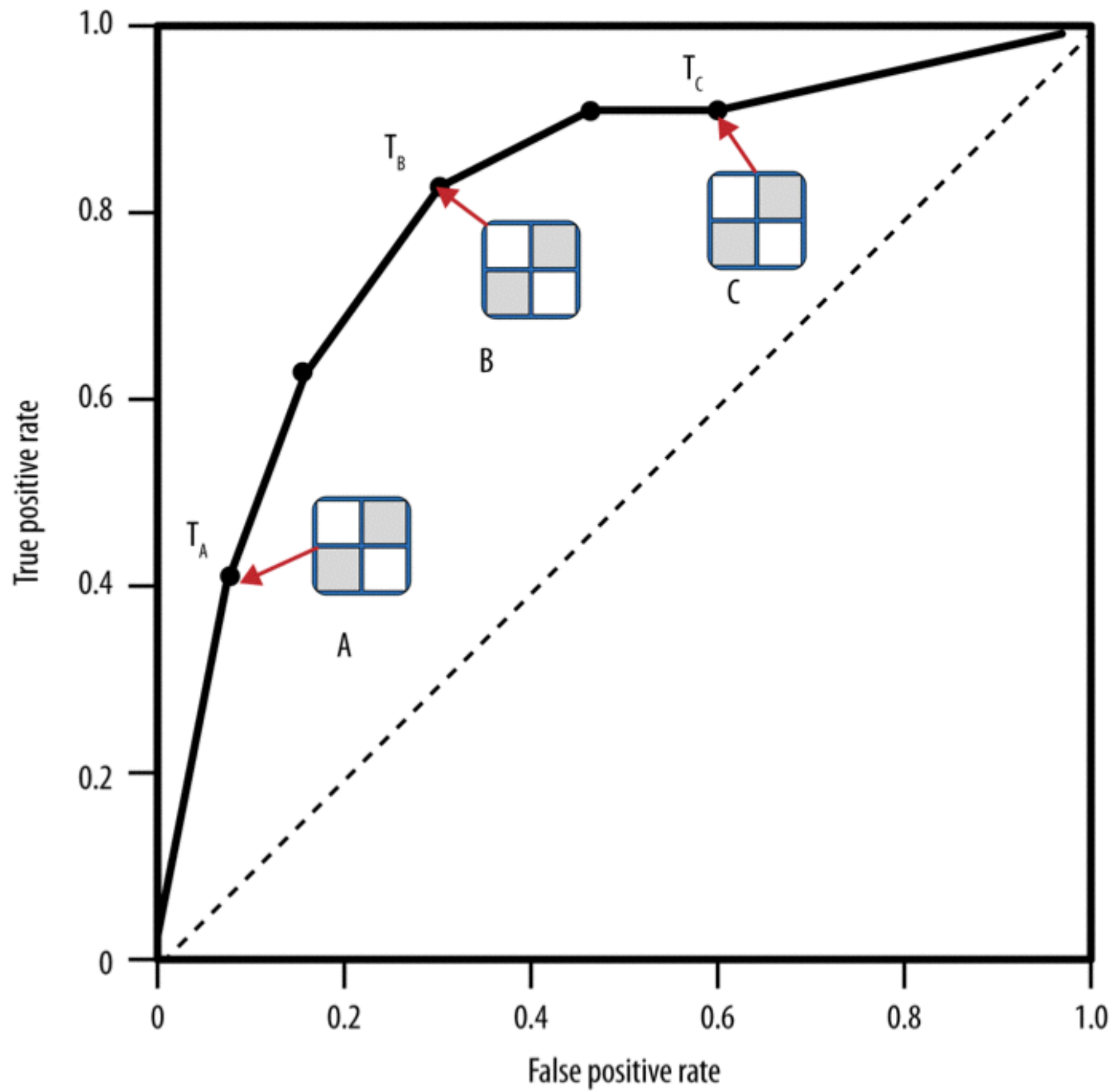
“Hit Rate”

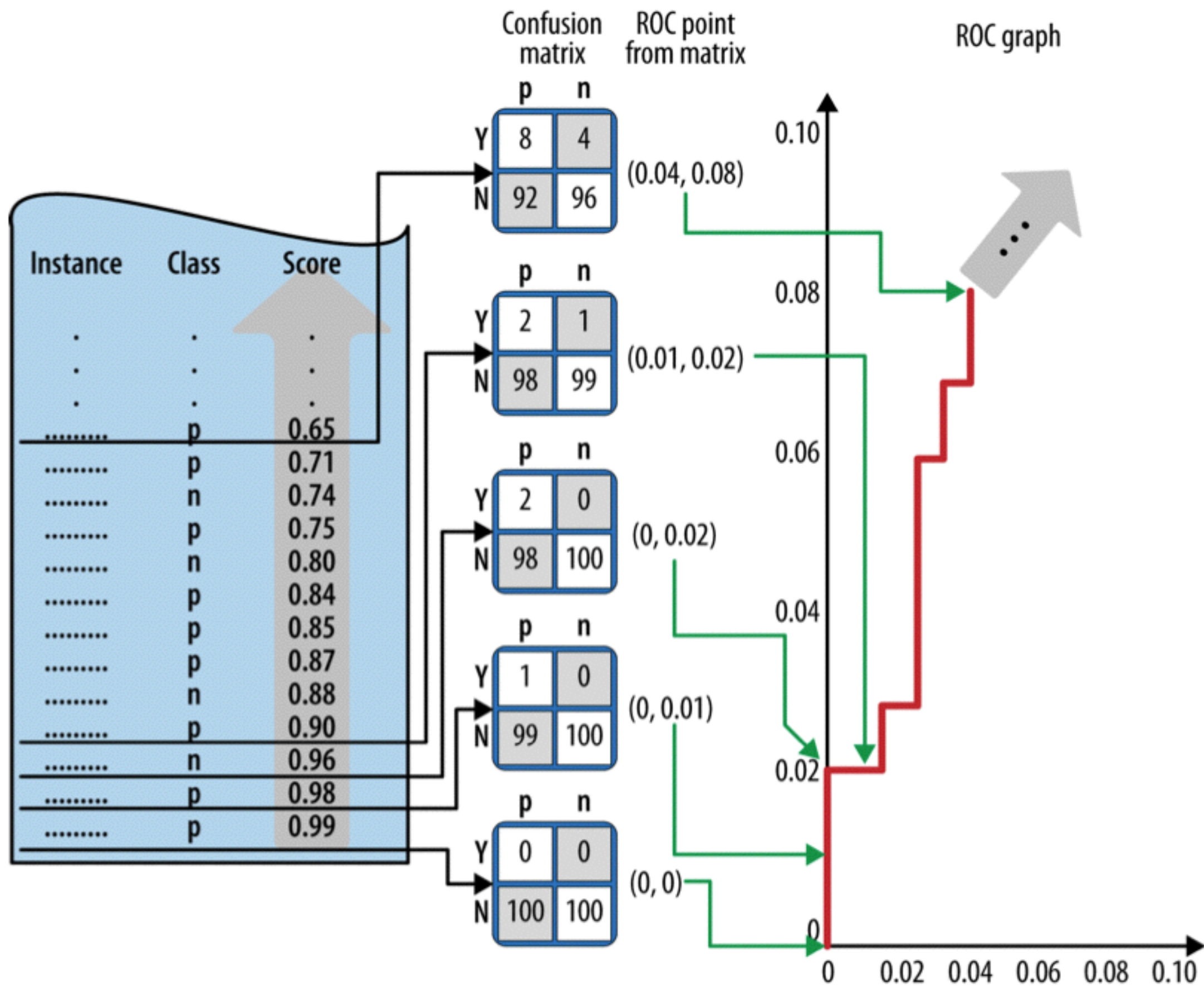


“False Alarm Rate”

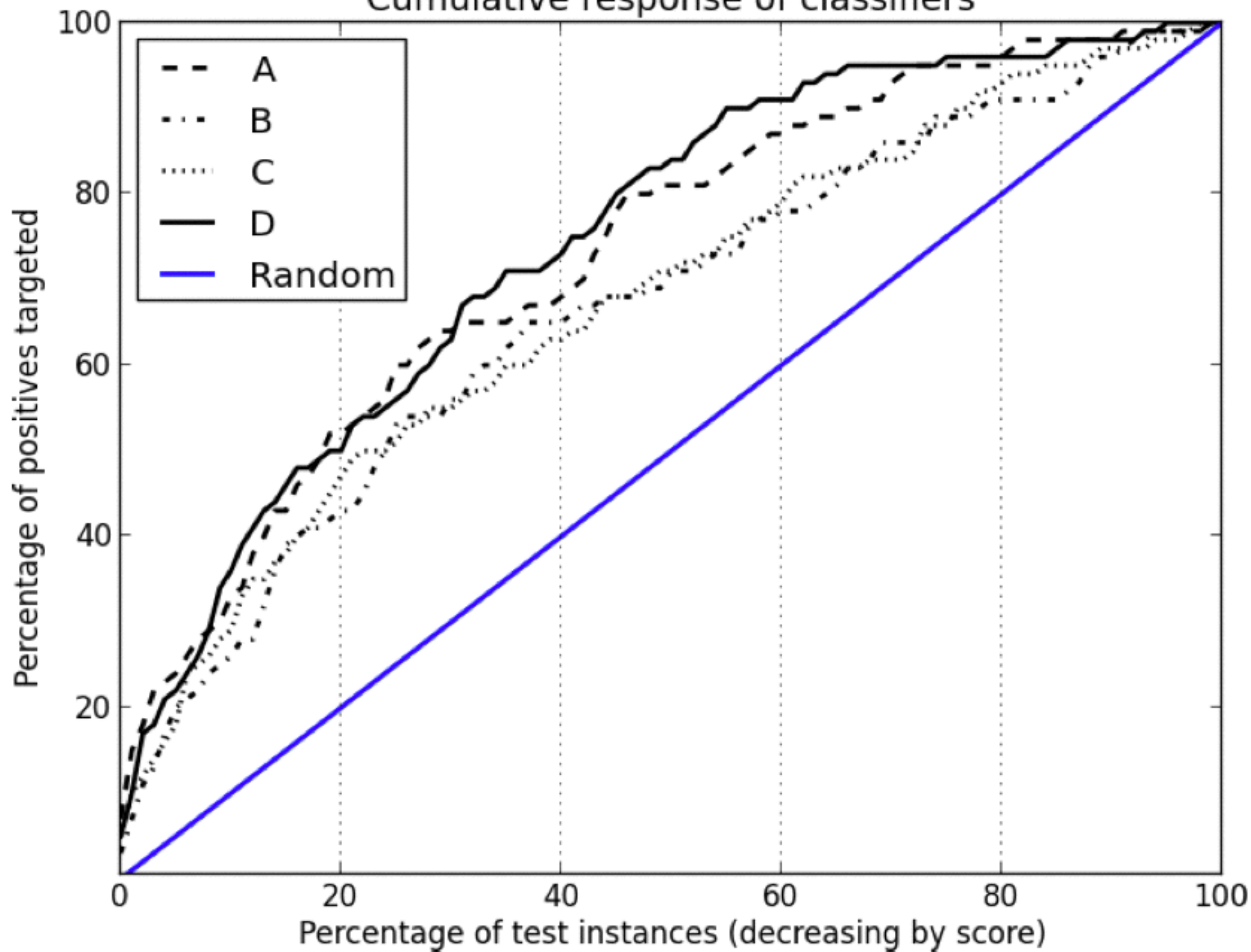








Cumulative response of classifiers



Model	Accuracy
Classification tree	95%
Logistic regression	93%
<i>k</i> -Nearest Neighbor	100%
Naive Bayes	76%

Model	Accuracy (%)	AUC
Classification Tree	91.8 \pm 0.0	0.614 \pm 0.014
Logistic Regression	93.0 \pm 0.1	0.574 \pm 0.023
<i>k</i> -Nearest Neighbor	93.0 \pm 0.0	0.537 \pm 0.015
Naive Bayes	76.5 \pm 0.6	0.632 \pm 0.019

	p	n
Y	127 (3%)	848 (18%)
N	200 (4%)	3518 (75%)

Here is the *k*-Nearest Neighbors confusion matrix on the same test data:

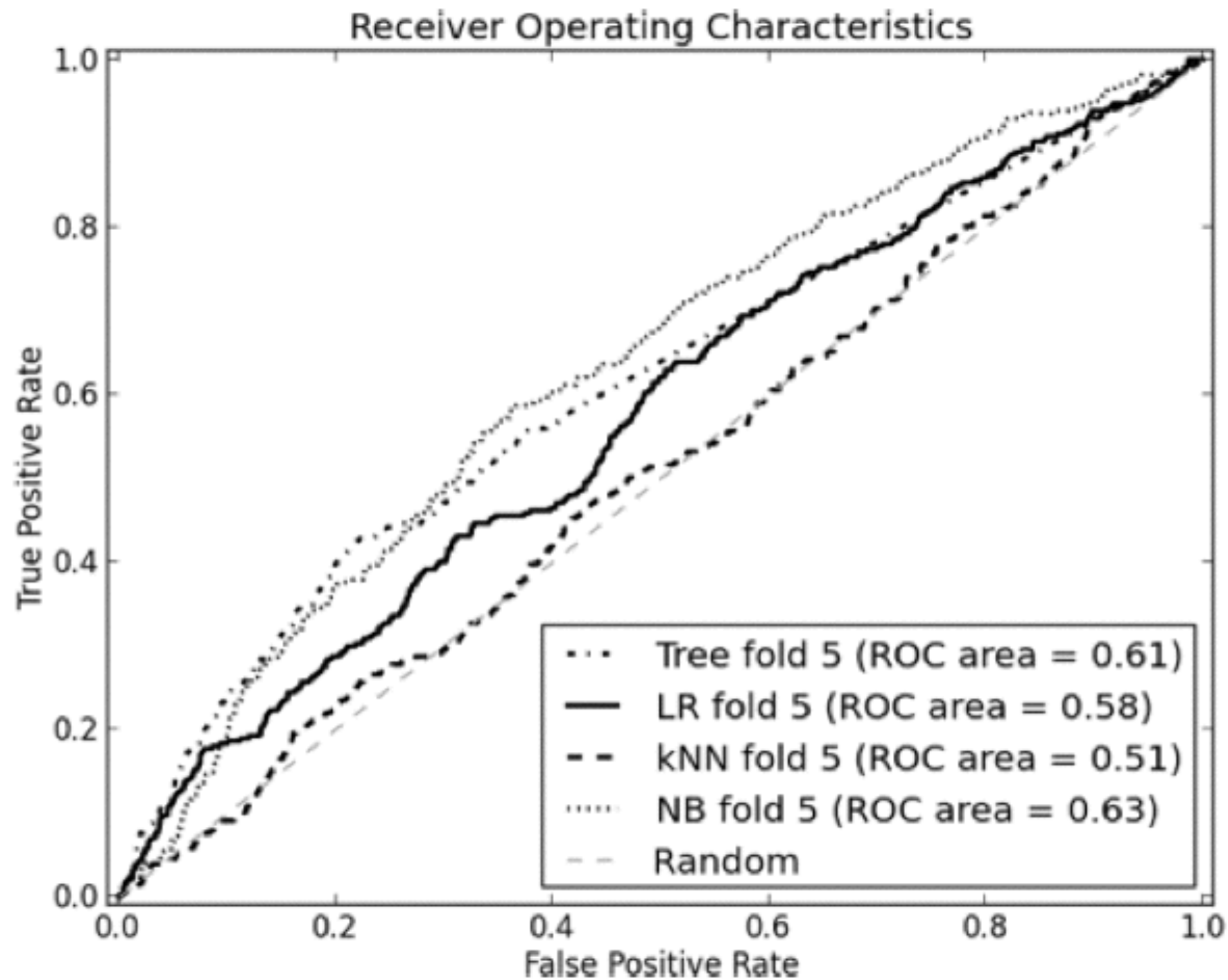
Model	Accuracy
Classification tree	95%
Logistic regression	93%
<i>k</i> -Nearest Neighbor	100%
Naive Bayes	76%

Model	Accuracy (%)	AUC
Classification Tree	91.8 ± 0.0	0.614 ± 0.014
Logistic Regression	93.0 ± 0.1	0.574 ± 0.023
<i>k</i> -Nearest Neighbor	93.0 ± 0.0	0.537 ± 0.015
Naive Bayes	76.5 ± 0.6	0.632 ± 0.019

	p	n
Y	127 (3%)	848 (18%)
N	200 (4%)	3518 (75%)

Here is the *k*-Nearest Neighbors confusion matrix on the same test data:

	p	n
Y	3 (0%)	15 (0%)
N	324 (7%)	4351 (93%)



ROC curves of the classifiers on one fold of cross-validation for the churn problem