

Classification

SYS 6018 | Spring 2025

classification.pdf

Contents

1	Classification Intro	2
1.1	Set-up	2
1.2	Binary Classification	2
2	Risk Scoring vs. Classification	2
3	Binary Classification	4
3.1	Decision Theory	4
3.2	Common Binary Loss Functions	6
3.3	Performance Metrics	7
3.4	Performance over a range of thresholds	10
3.5	More than two classes	17
3.6	Summary of Classification Evaluation	17
4	Appendix: R Code	18

1 Classification Intro

1.1 Set-up

- The outcome variable is categorical and denoted $G \in \mathcal{G}$
 - Default Credit Card Example: $\mathcal{G} = \{\text{"Yes"}, \text{"No"}\}$
 - Medical Diagnosis Example: $\mathcal{G} = \{\text{"stroke"}, \text{"heart attack"}, \text{"drug overdose"}, \text{"vertigo"}\}$
- The training data is $D = \{(X_1, G_1), (X_2, G_2), \dots, (X_n, G_n)\}$
- The optimal decision/classification is often based on the posterior probability $\Pr(G = g \mid \mathbf{X} = \mathbf{x})$

1.2 Binary Classification

- Classification is simplified when there are only 2 classes.
 - Many multi-class problems can be addressed by solving a set of binary classification problems (e.g., [one-vs-rest](#)).
- It is often convenient to transform the outcome variable to a binary $\{0, 1\}$ variable:

$$Y_i = \begin{cases} 1 & G_i = \mathcal{G}_1 \quad (\text{outcome of interest}) \\ 0 & G_i = \mathcal{G}_2 \end{cases}$$

- In the `Default` data, it would be natural to set `default=Yes` to 1 and `default=No` to 0.

2 Risk Scoring vs. Classification

Most of the models we will encounter can output a predicted probability $\hat{p}_k(x) = \widehat{\Pr}(G = k \mid X = x)$ for every class $k \in \mathcal{G}$.

Sometimes a *hard classification* needs to be made, i.e., decide on single label/class to assign the observation.

1. Hard Classification:
 - Use training data to estimate the *label* $\hat{G}(X)$
 - The loss/cost $L(G, \hat{G}(X))$ is the loss incurred by estimating G with \hat{G}
2. Risk Scoring (Soft-Classification):
 - Use training data to estimate the *probability* $\hat{p}_k(X)$
 - The loss/cost $L(G, \hat{p}(X))$ is the loss incurred by estimating G with $\hat{p}_k(X)$, where $\hat{p}(X) = [\hat{p}_1(X), \dots, \hat{p}_K(X)]$
3. Ranking:
 - Use the training data to rank the test observations according to estimated risk level.
 - The loss/cost is based on the number of outcomes of interest in the top proportion of risk.

2.0.1 Example: Recidivism Prediction

Recently the National Institute of Justice hosted a [Recidivism Forecasting Challenge](#) which challenged contestants to predict if a parolee would be arrested for another offense within the next few years. The motivation is not to determine who should be released on parole, but rather which parolees should get additional assistance/supervision.

Objective	Model Output
Classification	Predict { Yes, No } for re-offending
Scoring	Predict probability of re-offending
Ranking	Order from highest risk level to lowest

Your Turn #1 : Recidivism Prediction

1. How could you use the probability/score to make a hard classification?
2. Do you think a hard classification or probability/score is better for this scenario?
3. If there were limited resources (e.g., only N parolees could get extra assistance), which type of model output would be more useful?

3 Binary Classification

3.1 Decision Theory

- We are considering *binary* outcomes, so the outcomes $G \in \{0, 1\}$
- Let $p(x) = \Pr(G = 1 \mid X = x)$
- Loss Function: $L(\text{True Label}, \text{Estimated Label}) = L(G, \hat{G})$

Loss	Description	Name
$L(G = 0, \hat{G} = 0)$	True class is 0, Predicted class is 0	True Negative
$L(G = 1, \hat{G} = 1)$	True class is 1, Predicted class is 1	True Positive
$L(G = 0, \hat{G} = 1)$	True class is 0, Predicted class is 1	False Positive
$L(G = 1, \hat{G} = 0)$	True class is 1, Predicted class is 0	False Negative

- A model's *Expected Prediction Error (EPE)* at input X is the expected loss on new data with input X .
- The EPE (for a binary outcome) is:

$$\begin{aligned} \text{EPE}_x(g) &= E_{G|X=x} [L(G, \hat{G}(x) = g) \mid X = x] \\ &= L(1, g) \Pr(G = 1 \mid X = x) + L(0, g)(1 - \Pr(G = 1 \mid X = x)) \\ &= L(1, g)p(x) + L(0, g)(1 - p(x)) \end{aligned}$$

- Hard Decision ($\hat{G}(x) \in \{0, 1\}$): choose $\hat{G}(x) = 1$ if

$$\begin{aligned} \text{EPE}_x(1) &< \text{EPE}_x(0) \\ L(1, 1)p(x) + L(0, 1)(1 - p(x)) &< L(1, 0)p(x) + L(0, 0)(1 - p(x)) \\ p(x)(L(1, 1) - L(1, 0)) &< (1 - p(x))(L(0, 0) - L(0, 1)) \\ p(x)(L(1, 0) - L(1, 1)) &\geq (1 - p(x))(L(0, 1) - L(0, 0)) \quad (\text{multiply both sides by } -1) \\ \frac{p(x)}{1 - p(x)} &\geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)} \end{aligned}$$

Note

In most cases, there will be no loss/cost for making a correct classification. Thus it is convention to set $L(0, 0) = L(1, 1) = 0$ in these scenarios.

3.1.1 Example: Cancer Diagnosis

- Say we have a goal of estimating if a patient has cancer using medical imaging
 - Let $G = 1$ for cancer and $G = 0$ for no cancer
- Suppose we have solicited a loss function with the following values
 - $L(G = 0, \hat{G} = 0) = 0$: There is no loss for correctly diagnosis a patient without cancer.
 - $L(G = 1, \hat{G} = 1) = 0$: There is no loss (for our model) for correctly diagnosis a patient with cancer.

- $L(G = 0, \hat{G} = 1) = C_{FP}$: There is a cost of C_{FP} units if the model issues a *false positive*, estimating the patient has cancer when they don't.
- $L(G = 1, \hat{G} = 0) = C_{FN}$: There is a cost of C_{FN} units if the model issues a *false negative*, estimating the patient does not have cancer when they really do.
- In these scenarios C_{FN} is often much larger than C_{FP} ($C_{FN} \gg C_{FP}$) because the effects of not promptly treating (or further testing, etc) a patient is more severe than starting a treatment path for patients that don't actually have cancer.
- The optimal decision is to issue a positive indication for cancer if $EPE_x(1) < EPE_x(0)$. This occurs when

$$\frac{p(x)}{1-p(x)} \geq \frac{C_{FP}}{C_{FN}} \quad \text{OR} \quad p(x) \geq \frac{C_{FP}}{C_{FP} + C_{FN}} \quad \text{OR} \quad \log\left(\frac{p(x)}{1-p(x)}\right) \geq \log\left(\frac{C_{FP}}{C_{FN}}\right)$$

- The ratio of C_{FP} to C_{FN} is all that matters for the decision. Let's say that $C_{FP} = 1$ and $C_{FN} = 10$. Then if $p(x) \geq 1/11$, our model will diagnose cancer.
 - Note: $p(x) = \Pr(Y = 1 | X = x)$ is affected by the class prior $\Pr(Y = 1)$ (e.g., the portion of the population tested who have cancer), which is usually going to be small.

3.1.2 Optimal Threshold

- Recall, the optimal *hard classification* decision is to choose $\hat{G} = 1$ if:

$$\frac{p(x)}{1-p(x)} \geq \frac{L(0,1) - L(0,0)}{L(1,0) - L(1,1)}$$

- It can be convenient to use model output other than $p(x)/(1-p(x))$ to make decisions
- Some models directly output $\hat{p}(x)$
- Other models, like logistic regression, naturally work with the link function (linear part)
 - Denote $\gamma(x)$ as the *logit* of $p(x)$:

$$\gamma(x) = \log \frac{p(x)}{1-p(x)} = \log \frac{\Pr(G = 1 | X = x)}{\Pr(G = 0 | X = x)}$$

Notation

$$\gamma(x) = \log \frac{p(x)}{1-p(x)}$$

$$p(x) = \frac{e^{\gamma(x)}}{1 + e^{\gamma(x)}}$$

- Table of equivalent representations:

Score	Threshold	Threshold (simplified)
$\frac{p(x)}{1 - p(x)}$	$\frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)}$	$\frac{C_{FP}}{C_{FN}}$
$\gamma(x) = \log \frac{p(x)}{1 - p(x)}$	$\log \left(\frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)} \right)$	$\log \left(\frac{C_{FP}}{C_{FN}} \right)$
$p(x)$	$\log \left(\frac{L(0, 1) - L(0, 0)}{L(0, 1) - L(0, 0) + L(1, 0) - L(1, 1)} \right)$	$\frac{C_{FP}}{C_{FP} + C_{FN}}$

The *Threshold (simplified)* assumes $L(0, 0) = L(1, 1) = 0$

3.1.3 Using estimated values

- We will never have the actual $p(x)$ or $\gamma(x)$, so replace them with the estimated values.
- For a given threshold t and input x , the hard classification rule is $\hat{G}_t(x) = \mathbb{1}(\hat{p}(x) \geq t)$ (or $\hat{G}_t(x) = \mathbb{1}(\hat{\gamma}(x) \geq t)$ if using $\hat{\gamma}$ instead of \hat{p}).

Note

Because we have to estimate $\hat{p}(x)$ or $\hat{\gamma}(x)$, the best threshold t^* may differ from the theoretical optimal and need to be estimated. (more info about this below)

3.2 Common Binary Loss Functions

- Setting: estimate a binary outcome $G \in \{0, 1\}$ with a predicted label $\hat{G}(x)$
- **Mis-Classification Cost**

$$L(G, \hat{G}(x)) = \begin{cases} C_{FP} & G = 0, \hat{G}(x) = 1 \\ C_{FN} & G = 1, \hat{G}(x) = 0 \\ 0 & \text{otherwise} \end{cases}$$

- This requires that a *hard classification* is made.
- The theoretically optimal prediction is:

$$G^*(x) = \mathbb{1}(p(x) > C_{FP}/(C_{FP} + C_{FN}))$$

- **0-1 Loss or Misclassification Error**

$$L(G, \hat{G}(x)) = \mathbb{1}(y \neq \hat{G}(x)) = \begin{cases} 0 & G = \hat{G}(x) \\ 1 & G \neq \hat{G}(x) \end{cases}$$

- This assumes $L(0, 1) = L(1, 0)$ (i.e., false positive costs the same as a false negative)
- This requires that a *hard classification* is made.
- The theoretically optimal prediction is:

$$\begin{aligned} G^*(x) &= \mathbb{1}(p(x) > 1 - p(x)) \\ &= \mathbb{1}(p(x) > 0.50) \end{aligned}$$

3.3 Performance Metrics

3.3.1 Confusion Matrix

- Given a threshold t , we can make a *confusion matrix* to help analyze our model's performance on data
 - Data = $\{(X_i, G_i)\}_{i=1}^N$ (ideally this is hold-out/test data)
 - N_k is number of observations from class k ($N_0 + N_1 = N$)

		Model Outcome		total
		$\hat{G}_t = 1$	$\hat{G}_t = 0$	
True Outcome	$G = 1$	True Positive (TP)	False Negative (FN)	N_1
	$G = 0$	False Positive (FP)	True Negative (TN)	N_0
total		$\hat{N}_1(t)$	$\hat{N}_0(t)$	N

Table from: <https://tex.stackexchange.com/questions/20267/how-to-construct-a-confusion-matrix-in-latex>

To illustrate a confusion table in practice let's go back to the `Default` data and see how the basic logistic regression models performs.

- In order to evaluate on hold-out data, split the data into train/test (used 8000 training, 2000 testing), fit a logistic regression model on training data, and make predictions on the test data
- Note that only 3.3% of the data is default.
 - Using a threshold of $\hat{p}(x) \geq 0.10$ to make a hard classification.
 - Equivalent to $\hat{\gamma}(x) \geq \log(.10) - \log(1 - .10) = -2.1972$

```

#: train/test split
set.seed(2019)
test.ind = sample(nrow(Default), size=2000)
train.ind = -test.ind

#: fit model on training data
fit_lm = glm(y~student + balance + income, family='binomial',
            data=Default[train.ind, ])

#: Get predictions (of p(x)) on test data
p_hat = predict(fit_lm, Default[test.ind, ], type='response')

#: Make Hard classification (use .10 as cut-off)
G_hat = ifelse(p_hat >= .10, 1, 0)

#: Make Confusion Table
G_test = Default$y[test.ind] # true values

table(truth = G_test, predicted = G_hat) %>% addmargins()
#>      predicted
#> truth    0    1 Sum
#>    0   1805  128 1933

```

```
#>   1    17    50    67
#> Sum 1822 178 2000
```

3.3.2 Metrics

There are several standard evaluation metrics that can be calculated from the confusion matrix:

Metric	Definition	Estimate
Expected Cost	$\sum_{i=0}^1 \sum_{j=0}^1 L(i, j) P_X(G(X) = i, \hat{G}_t(X) = j)$	$\frac{1}{N} \sum_{i=1}^N L(G_i, \hat{G}_t(x_i))$
Mis-classification Rate	$P_{XG}(\hat{G}_t(X) \neq G(X)) =$ $P_X(\hat{G}_t(X) = 0, G(X) = 1) +$ $P_X(\hat{G}_t(X) = 1, G(X) = 0)$	$\frac{1}{N} \sum_{i=1}^N \mathbb{1}(\hat{G}_t(x_i) \neq G_i)$
False Positive Rate (FPR) {1-Specificity}	$P_X(\hat{G}_t(X) = 1 G(X) = 0)$	$\frac{1}{N_0} \sum_{i:G_i=0} \mathbb{1}(\hat{G}_t(x_i) = 1)$
True Positive Rate (TPR) {Hit Rate, Recall, Sensitivity}	$P_X(\hat{G}_t(X) = 1 G(X) = 1)$	$\frac{1}{N_1} \sum_{i:G_i=1} \mathbb{1}(\hat{G}_t(x_i) = 1)$
Precision TP/(TP + FP)	$P_X(G(X) = 1 \hat{G}_t(X) = 1)$	$\frac{1}{\hat{N}_1(t)} \sum_{i:\hat{G}_t(x_i)=1} \mathbb{1}(G_i = 1)$

N is the total number of predictions/observations, N_0 is the number of true class 0's in the data ($N_0 = \sum_{i=1}^N \mathbb{1}(y_i = 0)$), N_1 is the number of true class 1's in the data ($N_1 = \sum_{i=1}^N \mathbb{1}(y_i = 1)$), $\hat{N}_1(t)$ is the number of *predicted* class 1's using a threshold of t ($\hat{N}_1(t) = \sum_{i=1}^n \mathbb{1}(\hat{p}_i \geq t)$).

- Note: Performance estimates are best carried out on *hold-out* data!
- See [Wikipedia Page: Confusion Matrix](#) for more metrics:

Sources: [13][14][15][16][17][18][19][20] view · talk · edit

		Predicted condition			
Total population = P + N		Positive (PP)	Negative (PN)	Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) = $\frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power = $\frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate = $\frac{FN}{P} = 1 - TPR$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out = $\frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity = $\frac{TN}{N} = 1 - FPR$
Prevalence = $\frac{P}{P+N}$	Positive predictive value (PPV), precision = $\frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) = $\frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Negative likelihood ratio (LR-) = $\frac{FNR}{TNR}$	
Accuracy (ACC) = $\frac{TP+TN}{P+N}$	False discovery rate (FDR) = $\frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP (Δp) = PPV + NPV - 1	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$	
Balanced accuracy (BA) = $\frac{TPR+TNR}{2}$	F ₁ score = $\frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes-Mallows index (FM) = $\sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) = $\frac{\sqrt{TPR \times TNR \times PPV \times NPV} - \sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index = $\frac{TP}{TP + FN + FP}$	

F1 Metric

The F1 metric is the harmonic mean of precision and recall/TPR. It was first used in the context of information retrieval where there are often a massive number of irrelevant cases (negatives). In this setting, the actual number of true negatives isn't that important and the F1 metric, which doesn't consider the number of true negatives, became popular.

$$F_1 = \left(\frac{1/\text{precision} + 1/\text{recall}}{2} \right)^{-1} = \frac{2TP}{2TP + FN + FP}$$

- But even in this setting, why not use cost? You see that here false positives and false negatives show as equally important in the denominator. Wouldn't it be better to use Score = FP × C_{FP} + FN × C_{FN}?
- More reasons to avoid F1 (and Accuracy) are provided in [Common Problems With the Usage of F-Measure and Accuracy Metrics in Medical Research](#).
- Finally, what does F1 actually try to estimate? Using the definitions precision = P_X(G(X) = 1 | Ĝ_t(X) = 1) and recall = P_X(Ĝ_t(X) = 1 | G(X) = 1)

$$\begin{aligned}
 F_1 &= \left(\frac{1/\text{precision} + 1/\text{recall}}{2} \right)^{-1} \\
 &= \frac{2}{P_X(G(X) = 1 | \hat{G}_t(X) = 1)^{-1} + P_X(\hat{G}_t(X) = 1 | G(X) = 1)^{-1}} \\
 &= \frac{2}{\frac{P_X(\hat{G}_t(X)=1)}{P_X(G(X)=1, \hat{G}_t(X)=1)} + \frac{P_X(G(X)=1)}{P_X(G(X)=1, \hat{G}_t(X)=1)}} \\
 &= \frac{2P_X(G(X) = 1, \hat{G}_t(X) = 1)}{P_X(\hat{G}_t(X) = 1) + P_X(G(X) = 1)}
 \end{aligned}$$

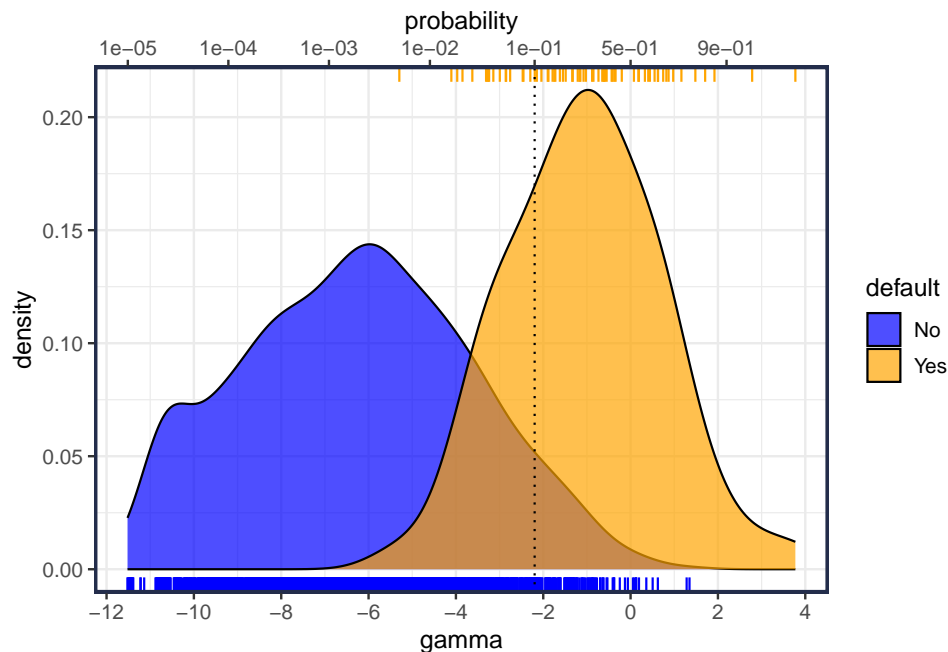
3.4 Performance over a range of thresholds

In the previous example, a hard classification was made using a threshold of $\hat{p}(x) \geq 0.10$. But performance varies as we adjust the threshold. Let's explore!

I'll use $\hat{\gamma}(x)$ instead of $\hat{p}(x)$ for this illustration because it better separates the classes.

```
#: Get predictions (of gamma(x)) on test data
gamma_hat = predict(fit_lm, Default[test.ind,], type='link')
p_hat = predict(fit_lm, Default[test.ind,], type='response')
```

- The model is unable to perfectly discriminate between groups, but the *defaults* do get scored higher in general:
 - As a reference point, note that $\gamma(x) = 0 \rightarrow \Pr(Y = 1 | X = x) = 1/2$
 - $\gamma(x) = \log p(x)/(1 - p(x))$



- We can calculate performance over a range of thresholds:

```
# truth: {0,1} vector
# score: risk score with larger values correspond to label = 1.
# thres: vector of thresholds at which to calculate metric.
# Note: decision is 1 if score > thres, 0 if score <= thres.
perf_table <- function(truth, score, thres=NULL){
  if(is.null(thres)) thres = seq(min(score, max(score)), length=1000)

  x = c(-Inf, thres, Inf) %>% unique() %>% sort() # expand and clean thresholds
  tibble(truth, score) %>%
    # create groups by threshold
    mutate(
      bin = findInterval(score, x, left.open = TRUE),
      val = x[bin+1]
    ) %>%
    # counts by group/threshold
    group_by(val) %>%
    summarize(n = n(), n.1 = sum(truth), n.0 = n-n.1) %>%
    ungroup() %>%
    # add zero counts
```

```

complete(val = thres, fill = list(n=0L, n.1=0L, n.0 = 0L)) %>%
# calculate metrics
arrange(val) %>%
mutate(
  TN = cumsum(n.0), # True negatives
  FN = cumsum(n.1), # False negatives
  TP = sum(n.1) - FN, # True positives
  FP = sum(n.0) - TN, # False positives
  TPR = TP/sum(n.1), # True positive rate (TP / Positives)
  FPR = FP/sum(n.0) # False positive rate (FP / Negatives)
) %>%
# drop values outside of stated thresholds
filter(val %in% thres) %>%
# retain relevant metrics
select(-n, -n.1, -n.0, score = val)
}

thresholds = seq(0, 1, length = 1000)
perf = perf_table(truth = G_test, score = p_hat, thres = thresholds) %>%
  mutate(p_hat = score, gamma_hat = log(p_hat) - log(1-p_hat), .before=1) %>%
  select(-score)

```

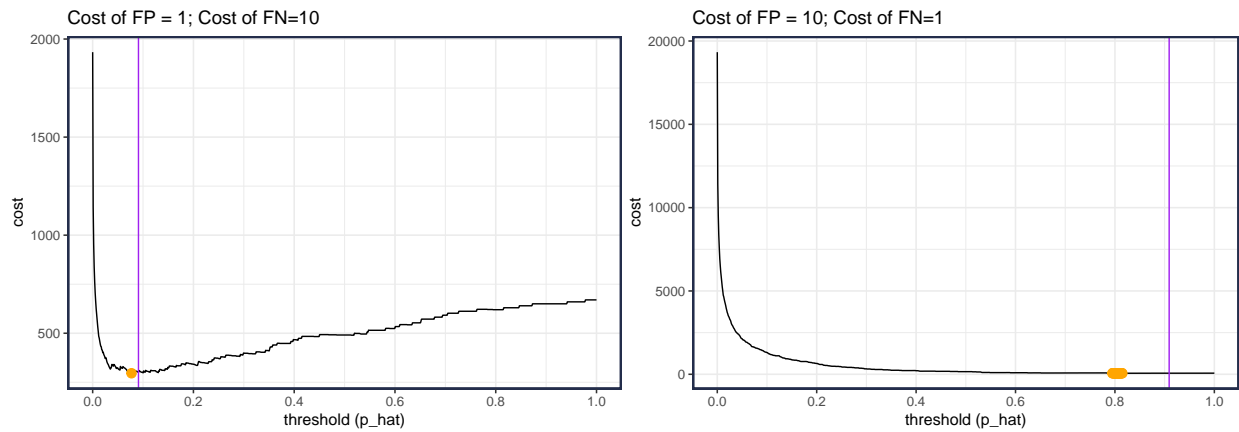
Here is a selection from the perf table

p_hat	gamma_hat	TN	FN	TP	FP	TPR	FPR
0.01	-4.594	1408	1	66	525	0.985	0.272
0.02	-3.891	1562	3	64	371	0.955	0.192
0.03	-3.475	1639	5	62	294	0.925	0.152
0.04	-3.177	1688	9	58	245	0.866	0.127
0.05	-2.943	1721	11	56	212	0.836	0.110
0.06	-2.750	1743	14	53	190	0.791	0.098
0.07	-2.586	1769	14	53	164	0.791	0.085
0.08	-2.441	1779	16	51	154	0.761	0.080
0.09	-2.313	1792	16	51	141	0.761	0.073

- Note: the perf object is *only based on the rank order* of the predictions. This means that the same results would be obtained if we used $\hat{\gamma}(x)$ or $\hat{p}(x)$ to do the ranking.
 - This is because there is a one-to-one monotone relationship between $\hat{\gamma}(x)$ and $\hat{p}(x)$.
 - The perf object grouped by both gamma_hat and p_hat so both thresholds are available. But we can switch back and forth from the relationship $\log(p/(1-p)) = \gamma$, so its easy to switch between the two depending on what is most convenient.

3.4.1 Cost Curves

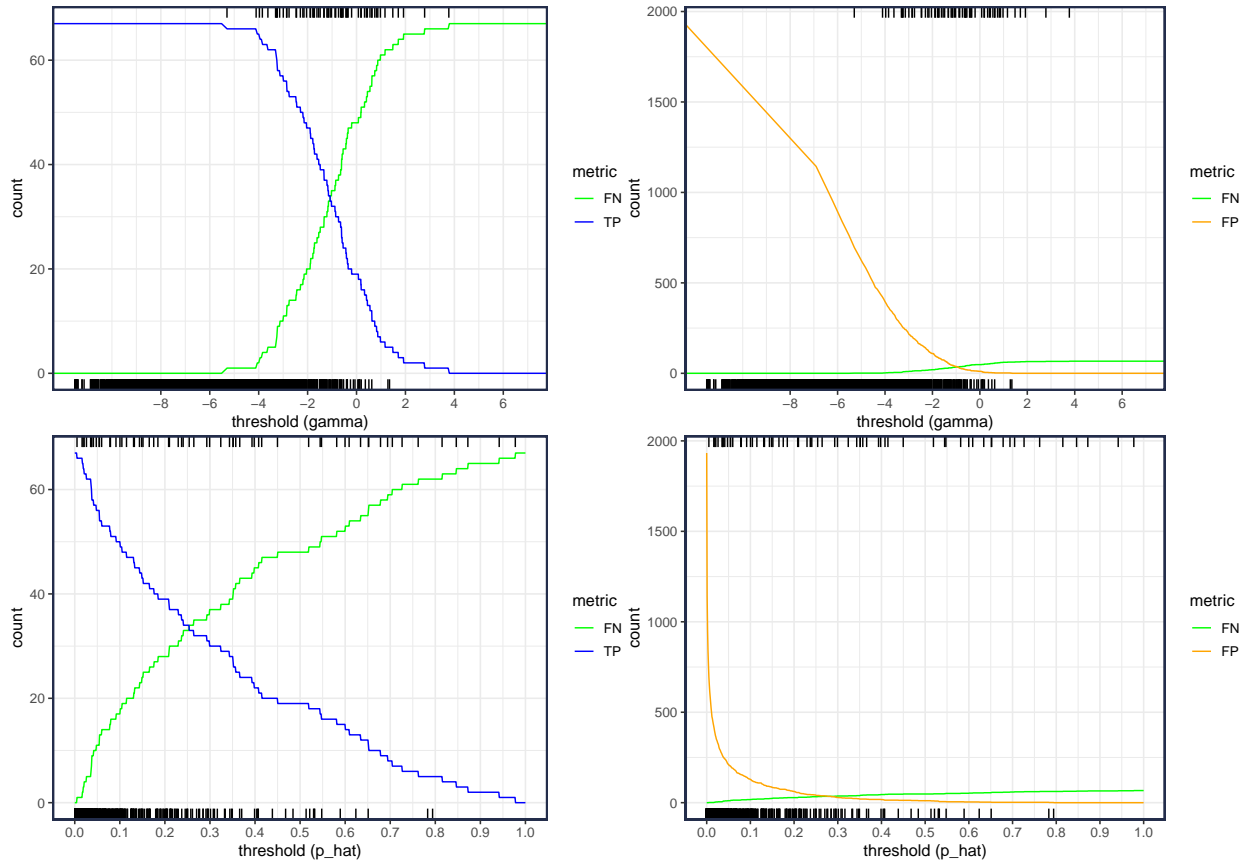
- Under the usual scenario where $L(0,0) = L(1,1) = 0$, the cost only depends on the ratio of false positive costs (C_{FP}) to false negative costs (C_{FN}).
- note: the purple is the *theoretical* optimal threshold (using $t^* = \log C_{FP}/C_{FN}$ for $\hat{\gamma}(x)$ and $C_{FP}/(C_{FP} + C_{FN})$ for $\hat{p}(x)$) and the orange point is at the optimal value for the test data.



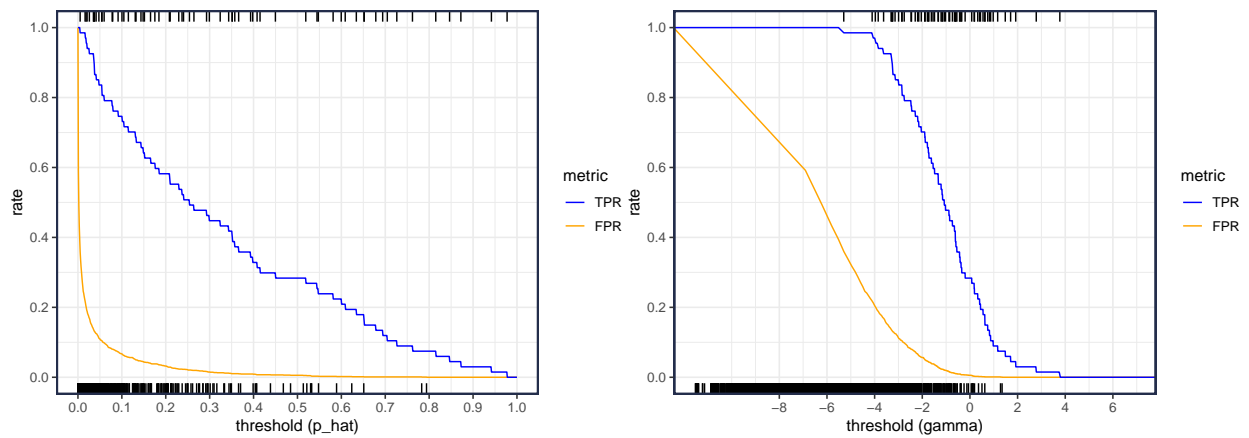
Optimal Threshold

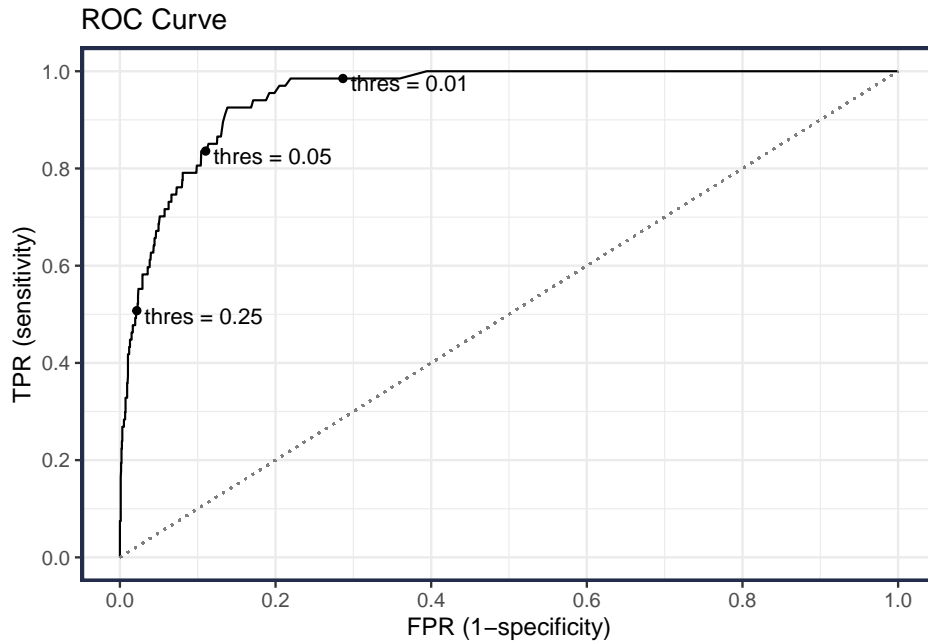
- The *theoretically* optimal threshold is based on the *true* $\gamma(x) = \log \frac{p(x)}{1-p(x)}$ (for a given cost ratio of FP to FN)
- The observed optimal threshold will differ when the model's estimate $\hat{\gamma}(x) \neq \gamma(x)$
 - Hopefully, they are close and it won't make much difference which one you use. But I'd take the estimated threshold if I had sufficient data.
- Note that the estimated values depend on the prior class probabilities. If you suspect these may differ in the future, then you should adjust the threshold.

3.4.2 General Performance as function of threshold (select metrics)



3.4.3 ROC Curve (Receiver Operating Characteristic)



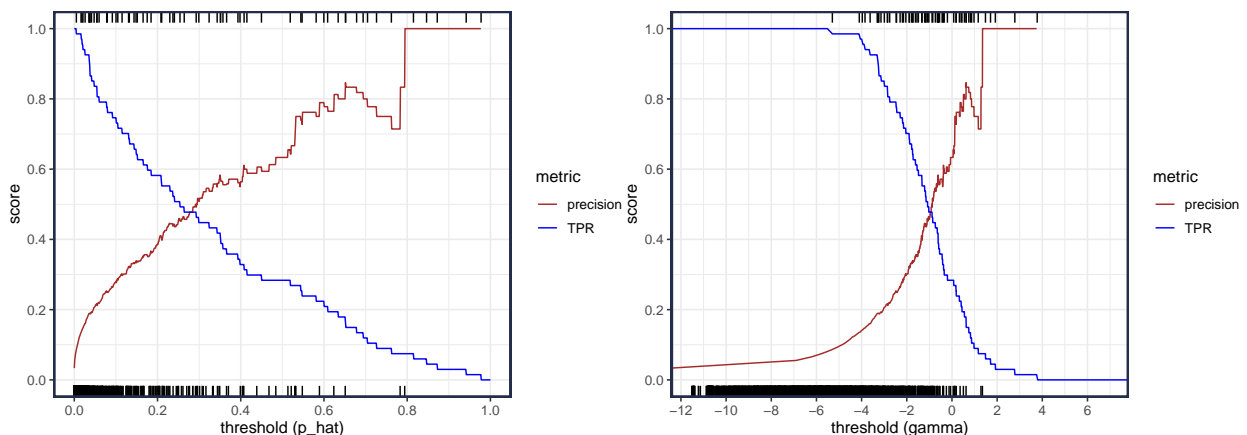


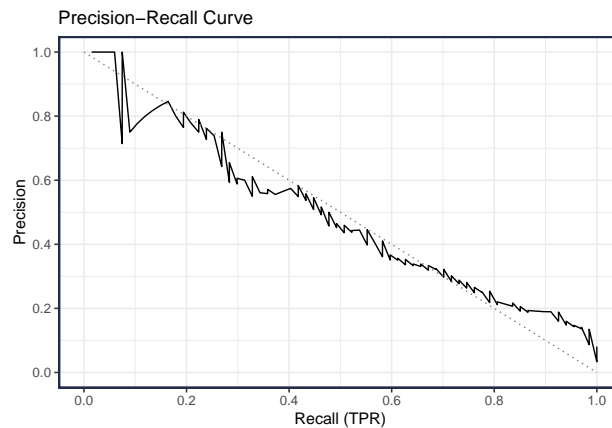
AUROC

- The *area under the ROC curve* (AUROC or AUC) is a popular performance metric
- I don't think it is a great way to compare classifiers for several reasons
 - The main reason is that in a real application you can almost always come up with an estimated cost/loss for the different decisions
 - To say it another way, comparisons should be made at a single point on the curve; the entire FPR region should not factor into the comparison.
- The AUROC is equal to the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.
 - AUROC is proportional to the [Mann-Whitney U statistic](#)

3.4.4 Precision Recall Curves

- Popular for information retrieval/ranking
- The *precision* metric is not monotonic wrt threshold, hence the sawteeth pattern.





3.4.5 R Code

Once we have the FP, TP, TN, FN values for a set of thresholds (like what is in the `perf` object), then we have everything we need to calculate any metric (e.g., gain, lift, F1, ...).

- But I will mention the `yardstick` R package which offers some functionality you may find convenient
- List of the [metrics included in the yardstick package](#)
- The `roc_curve()` function will generate the hard classification metrics for a range of thresholds.
 - *1-specificity* is the False Positive Rate (FPR) and *sensitivity* is the True Positive Rate (TPR).

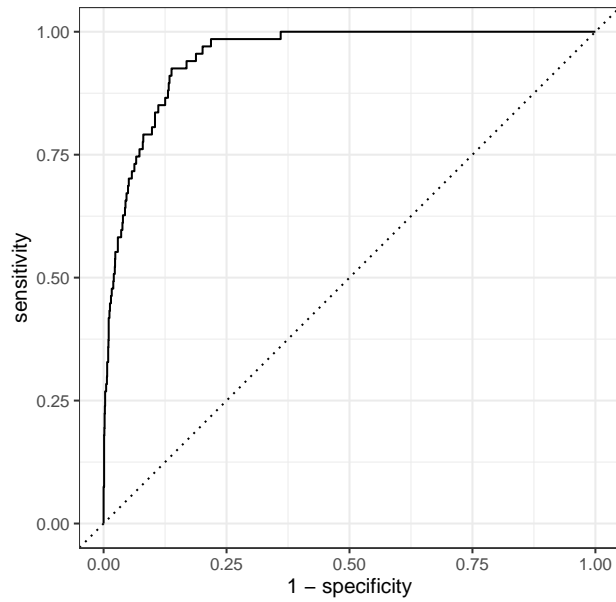
```
library(yardstick) # for evaluation functions

#: The yardstick package requires the categorical outcome be
#: a factor with 1st level the outcome of interest.
test_perf =
  tibble(
    truth = Default$y[test.ind] %>% factor(levels=c(1,0)),
    gamma_hat
  )

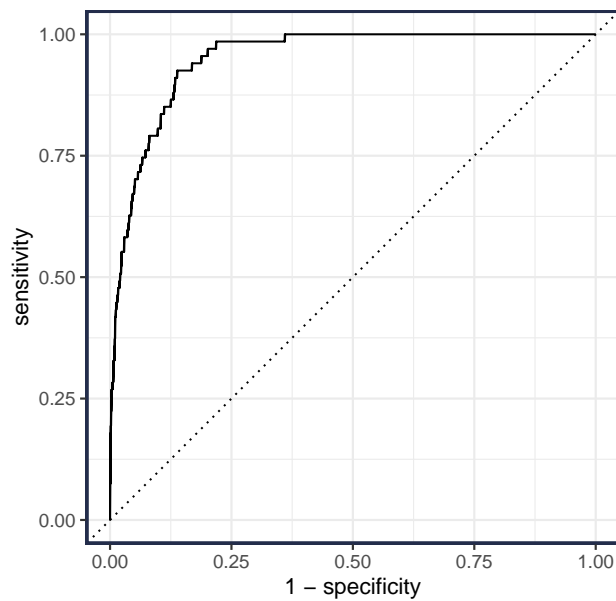
#: ROC plots
ROC =
  test_perf %>%
  yardstick::roc_curve(truth, gamma_hat)

ROC
#> # A tibble: 2,002 x 3
#>   .threshold specificity sensitivity
#>   <dbl>         <dbl>         <dbl>
#> 1    -Inf             0             1
#> 2    -11.5            0             1
#> 3    -11.5    0.000517         1
#> 4    -11.5    0.00103         1
#> 5    -11.5    0.00155         1
#> 6    -11.5    0.00207         1
#> # i 1,996 more rows

autoplot(ROC) # autoplot() method
```



```
# manual ROC plot (same as autoplot())
ROC %>%
  ggplot(aes(1-specificity, sensitivity)) + geom_line() +
  geom_abline(lty=3) +
  coord_equal()
```



The yardstick package also has an AUC function `roc_auc()`:

```
#: Area under ROC (AUROC)
test_perf %>%
  roc_auc(truth, gamma_hat)
#> # A tibble: 1 x 3
#>   .metric .estimator .estimate
#>   <chr>   <chr>       <dbl>
#> 1 roc_auc binary       0.951
```


3.5 More than two classes

Nothing really changes when there are more than two classes; we still want to choose the (hard) label that has **minimum** EPE:

$$\begin{aligned} \text{EPE}_x(g) &= \mathbb{E}_{G|X=x} [L(G, \hat{G}(x) = g) \mid X = x] \\ &= \sum_k L(G = k, \hat{G} = g) \Pr(G = k \mid X = x) \end{aligned}$$

$$\begin{aligned} G^*(x) &= \arg \min_g \widehat{\text{EPE}}_x(g) \\ &= \arg \min_g \sum_k L(G = k, \hat{G} = g) \widehat{\Pr}(G = k \mid X = x) \end{aligned}$$

3.6 Summary of Classification Evaluation

Decision Theory

A prediction is not a decision!

- Ask yourself: do I really need to make a hard classification? Or are risk scores/probabilities better for end user (who is in charge of *making the decisions*)?
- Use cost! The other metrics are probably not going to give you what you really want.
 - Resist the pressure to use AUROC, Accuracy, F1.
 - If you don't know the exact cost(FP)/cost(FN) ratio, then report performance for a reasonable range of values.
- For Binary Classification Problems, the optimal decision is to choose $\hat{G}(x) = 1$ if

$$\begin{aligned} \frac{p(x)}{1 - p(x)} &\geq \frac{L(0, 1) - L(0, 0)}{L(1, 0) - L(1, 1)} \\ &= \frac{\text{FP} - \text{TN}}{\text{FN} - \text{TP}} \end{aligned}$$

- Consider the connection to Decision Theory, make the decision that maximizes *expected utility*. The losses define the utility.
- In practice, we need to use an *estimated* $\hat{p}(x)$ or $\hat{\gamma}(x)$ and *estimated* threshold.
- Model parameters are usually estimated with a different metric than what's used for evaluation.
 - E.g., Estimate logistic regression parameters by minimizing Log-loss (i.e., maximum likelihood)
 - E.g., Hinge Loss for Support Vector Machines (SVM)
 - But Total Cost, MAE, F1, AUROC are used for evaluation (and tuning parameter estimation).
 - Reason: its difficult to estimate model parameters with such loss functions (e.g., non-differentiable, non-unique, etc.)

4 Appendix: R Code

Set-up

```

#: Load Required Packages
library(ISLR)
library(glmnet)
library(yardstick)
library(tidyverse)

#-----#
#: Default Data
# From the ISLR package
# The outcome variable is `default`
#-----#
library(ISLR)
data(Default, package="ISLR") # load the Default Data

#: Create binary column (y)
Default = Default |> mutate(y = ifelse(default == "Yes", 1L, 0L))
#: Summary Stats (Notice only 333 (3.3%) have defaulted)
summary(Default)
#> default student balance income y
#> No :9667 No :7056 Min. : 0 Min. : 772 Min. :0.0000
#> Yes: 333 Yes:2944 1st Qu.: 482 1st Qu.:21340 1st Qu.:0.0000
#> Median : 824 Median :34553 Median :0.0000
#> Mean : 835 Mean :33517 Mean :0.0333
#> 3rd Qu.:1166 3rd Qu.:43808 3rd Qu.:0.0000
#> Max. :2654 Max. :73554 Max. :1.0000

#: Train/Test Split
library(rsample)
set.seed(2024)
Default_split =
  initial_split(
    Default,
    prop = 1 - 500/nrow(Default), # n_test = 500
    strata = y # stratify on y
  )

```

Penalized Logistic Regression

First need to convert data to numeric model matrix

```

#: using the recipe/bake method
library(tidymodels)
rec_fit = recipe(y~student + balance + income,
  data = training(Default_split)) |>
  step_dummy(all_nominal(), one_hot = TRUE) |>
  prep()

#: training data
X.train = rec_fit |>
  bake(
    all_predictors(),
    new_data = training(Default_split),
    composition="matrix"
  )

```

```

Y.train = training(Default_split)$y

#: testing data
X.test = rec_fit |>
  bake(
    all_predictors(),
    new_data = testing(Default_split),
    composition="matrix"
  )
Y.test = testing(Default_split)$y

```

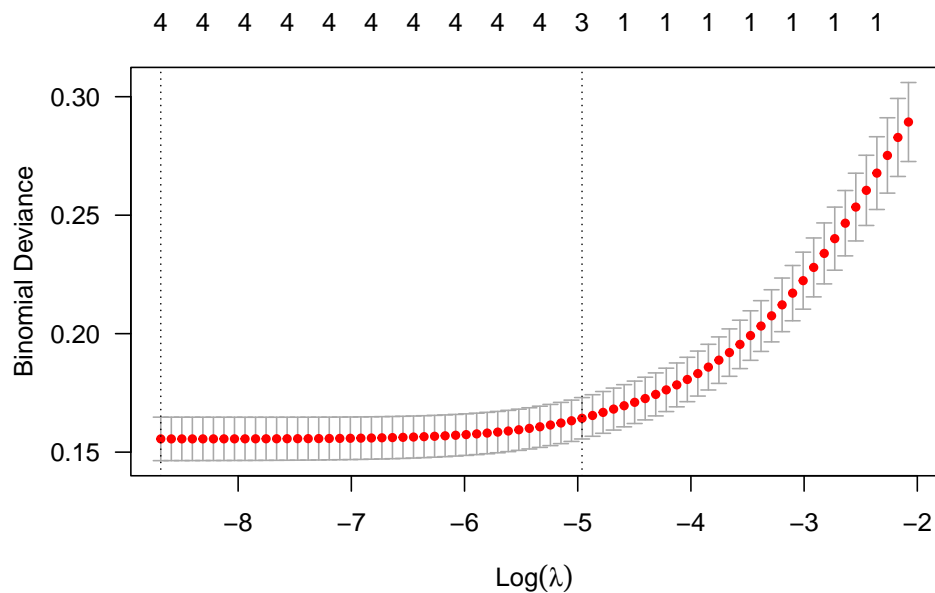
Now fit the elastic net

```

#: Elastic net with alpha = .5. Use CV to select lambda.
library(glmnet)
set.seed(2020) # seed controls the cross-validation splits
fit_enet =
  cv.glmnet(X.train, Y.train,
            alpha=.5,
            family="binomial")

#: CV performance plot
plot(fit_enet, las=1)

```



Performance Metrics and Curves

```

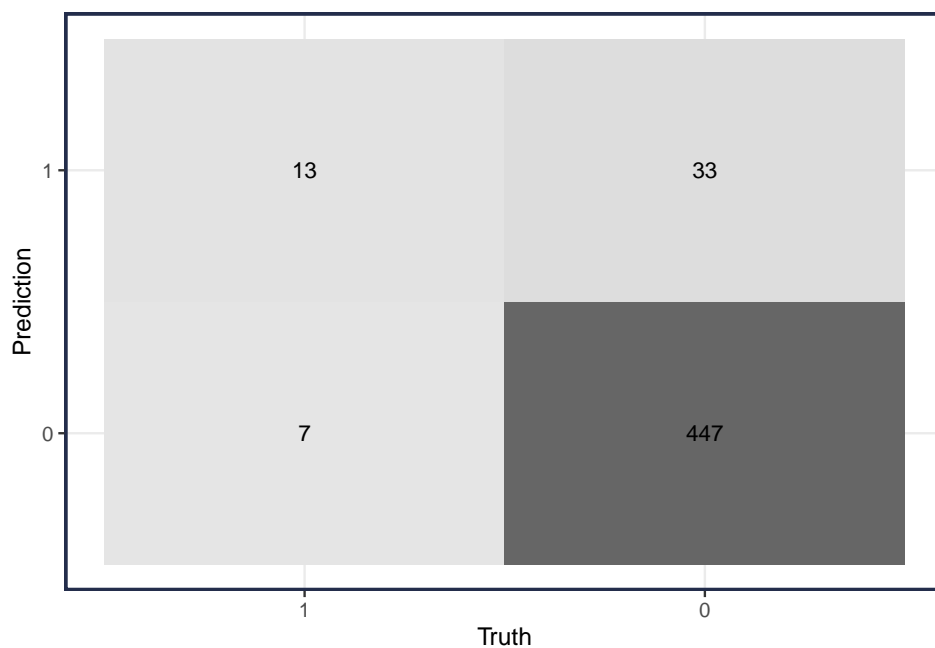
perf_data =
  testing(Default_split) |>
  mutate(
    # p_hat(x)
    p_hat = predict(fit_enet, X.test, s = "lambda.min", type = "response")[,1],
    # gamma_hat(x) = log(p) - log(1-p)
    gamma_hat = predict(fit_enet, X.test, s = "lambda.min", type = "link")[,1],
    # Make Hard classification (use .10 as cut-off)
    G_hat = ifelse(p_hat >= .10, 1, 0)
  ) |>
  as_tibble()

```

```

# : Make Confusion Table (yardstick)
library(yardstick)
# Note: the yardstick package functions, like conf_mat(), requires that hard
#       classifications have *factor* inputs (instead of *character*)
cm = perf_data |>
  mutate(
    y = factor(y, levels=c("1", "0")),
    G_hat = factor(G_hat, levels=c("1", "0"))
  ) %>%
  conf_mat(truth = y, estimate = G_hat)
cm
#>      Truth
#> Prediction 1  0
#>      1   13 33
#>      0    7 447
autoplot(cm, type = "heatmap")

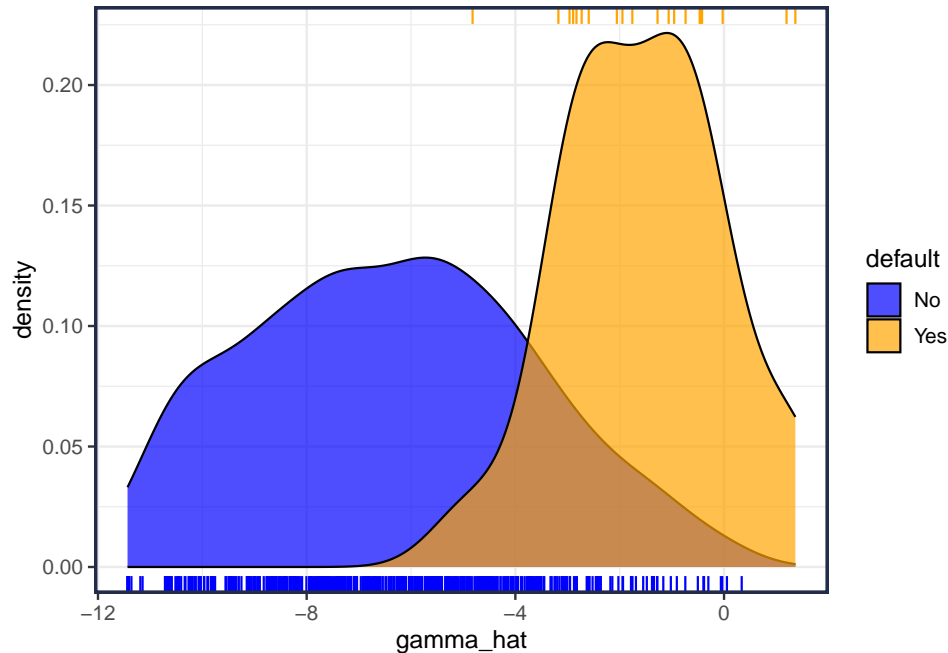
```



```

# : Visualize Performance by score
perf_data |>
  ggplot(aes(gamma_hat, fill=default)) +
  geom_density(alpha=.70) + # add kernel density estimates
  geom_rug(data=. %>% filter(default == 'Yes'), # add rug to top
    aes(color=default), sides='t') +
  geom_rug(data=. %>% filter(default == 'No'), # add rug to bottom
    aes(color=default), sides='b') +
  scale_fill_manual(values=c(Yes="orange", No="blue")) + # modify fill colors
  scale_color_manual(values=c(Yes="orange", No="blue"), guide="none") # modify colors

```



```

#: Get performance data (by threshold)
# This table has one row for every threshold. The columns are the elements
# of the confusion table plus FPR, TPR
perf_table =
  perf_data |>
  #- group_by() + summarize() in case of ties
  group_by(gamma_hat, p_hat) |>
  summarize(
    n = n(),
    n.1 = sum(y),
    n.0 = n-sum(y)
  ) |>
  ungroup() |>
  #- calculate metrics
  arrange(gamma_hat) %>%
  mutate(
    FN = cumsum(n.1), # false negatives
    TN = cumsum(n.0), # true negatives
    TP = sum(n.1) - FN, # true positives
    FP = sum(n.0) - TN, # false positives
    N = cumsum(n), # number of cases predicted to be 1
    TPR = TP/sum(n.1), FPR = FP/sum(n.0)
  ) %>%
  #- only keep relevant metrics
  select(-n, -n.1, -n.0, gamma_hat, p_hat)

## Note: gamma = log(p_hat) - log(1-p_hat) = log(p_hat / (1-p_hat))

```

```

perf_table |>
  head(10)
#> # A tibble: 10 x 9
#>   gamma_hat    p_hat    FN    TN    TP    FP    N    TPR    FPR
#>   <dbl>    <dbl> <int> <int> <int> <int> <int> <dbl> <dbl>
#> 1   -11.4 0.0000108     0     1    20   479     1     1 0.998
#> 2   -11.4 0.0000108     0     2    20   478     2     1 0.996
#> 3   -11.4 0.0000108     0     3    20   477     3     1 0.994

```

```
#> 4      -11.4 0.0000111      0      4      20      476      4      1 0.992
#> 5      -11.4 0.0000111      0      5      20      475      5      1 0.990
#> 6      -11.4 0.0000111      0      6      20      474      6      1 0.988
#> # i 4 more rows
```

```
## Alternatively, using yardstick::roc_curve()
```

```
perf_data |>
  mutate(
    y = factor(y, levels = c(1,0)),
  ) |>
  roc_curve(y, p_hat)
#> # A tibble: 502 x 3
#>   .threshold specificity sensitivity
#>   <dbl>         <dbl>         <dbl>
#> 1 -Inf             0             1
#> 2  0.0000108       0             1
#> 3  0.0000108       0.00208       1
#> 4  0.0000108       0.00417       1
#> 5  0.0000111       0.00625       1
#> 6  0.0000111       0.00833       1
#> # i 496 more rows
```

```
#: Make performance curves
```

```
col_lines = c(TP = "blue", FP="orange", FN="green", TN="brown",
              TPR = "blue", FPR="orange", FNR="green", TNR="brown")
```

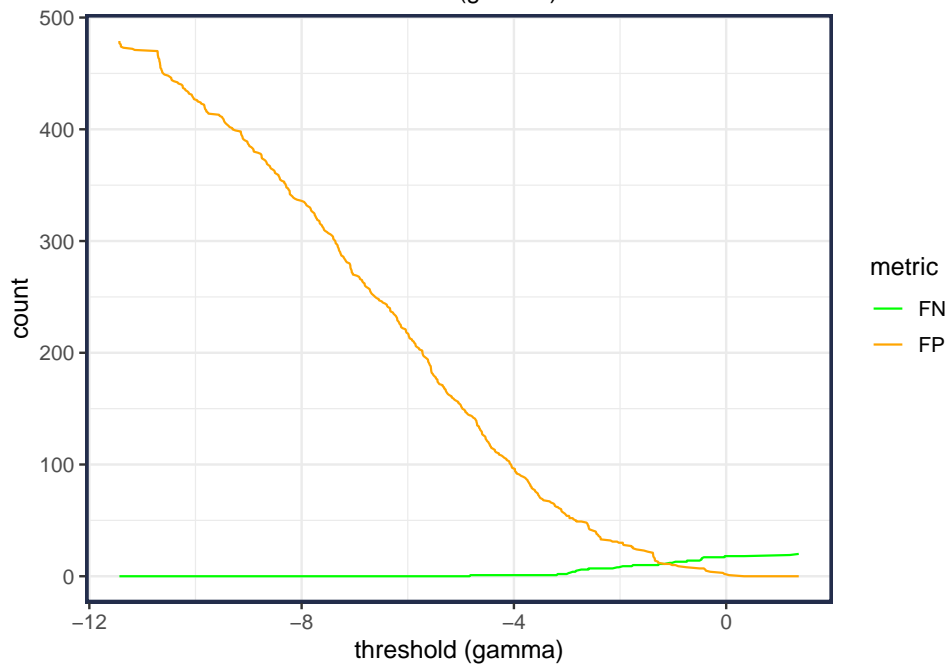
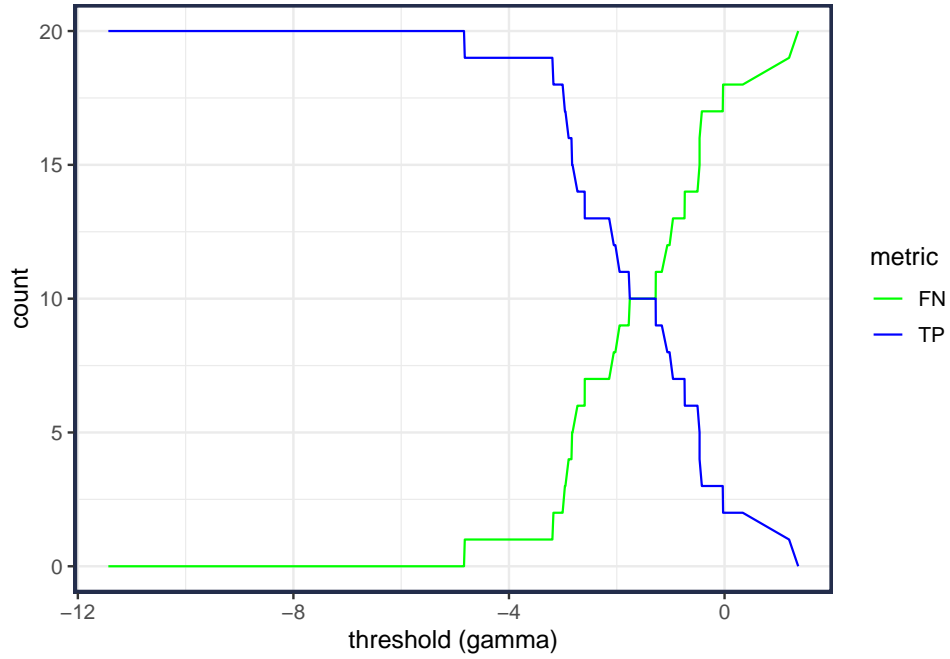
```
#: Make performance curves
```

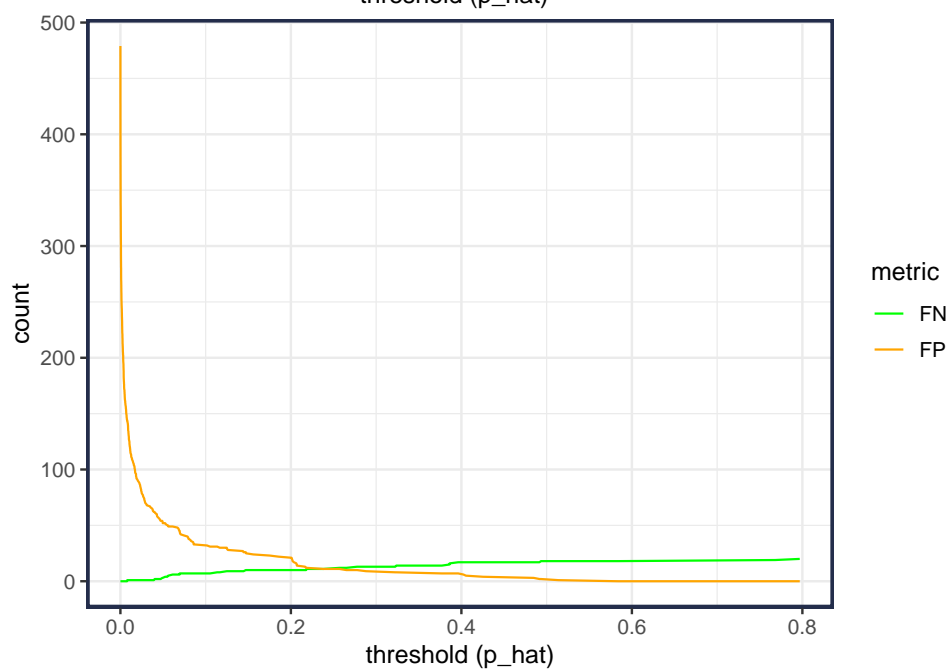
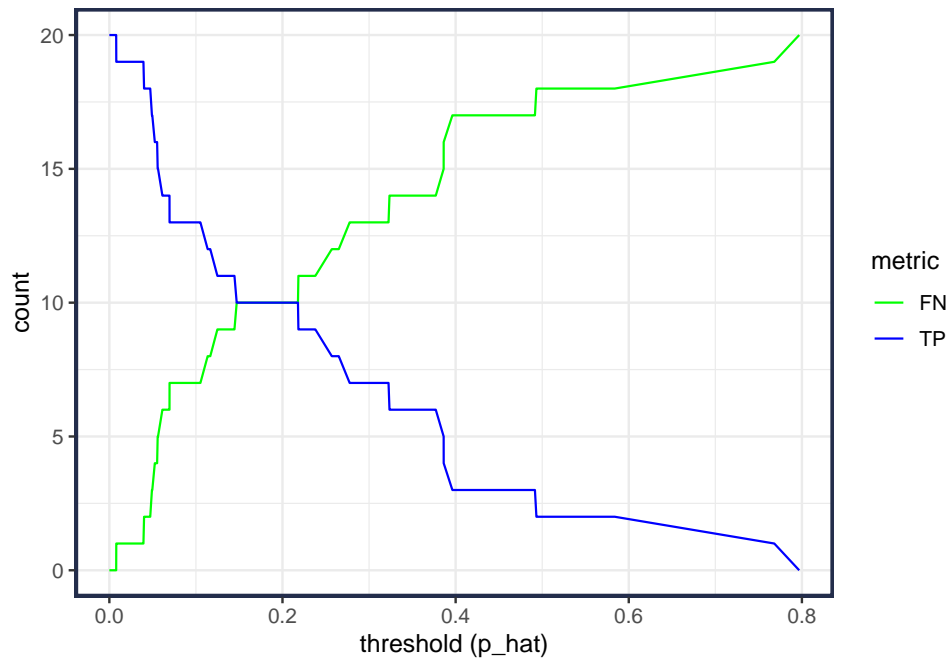
```
perf_table %>%
  select(threshold = gamma_hat, FN, TP) %>%
  gather(metric, n, -threshold) %>%
  ggplot(aes(threshold, n, color=metric)) + geom_line() +
  labs(x= "threshold (gamma)", y="count") +
  scale_color_manual(values=col_lines)
```

```
perf_table %>%
  select(threshold = gamma_hat, FN, FP) %>%
  gather(metric, n, -threshold) %>%
  ggplot(aes(threshold, n, color=metric)) + geom_line() +
  labs(x= "threshold (gamma)", y="count") +
  scale_color_manual(values=col_lines)
```

```
perf_table %>%
  select(threshold = p_hat, FN, TP) %>%
  gather(metric, n, -threshold) %>%
  ggplot(aes(threshold, n, color=metric)) + geom_line() +
  labs(x= "threshold (p_hat)", y="count") +
  scale_color_manual(values=col_lines)
```

```
perf_table %>%
  select(threshold = p_hat, FN, FP) %>%
  gather(metric, n, -threshold) %>%
  ggplot(aes(threshold, n, color=metric)) + geom_line() +
  labs(x= "threshold (p_hat)", y="count") +
  scale_color_manual(values=col_lines)
```





```

#: Make Cost curves
perf_table %>%
  mutate(cost = 1*FP + 10*FN) %>% # use 1:10 costs
  ggplot(aes(p_hat, cost)) + geom_line() +
  geom_point(data=. %>% filter(cost==min(cost)), size=3, color='orange') + # # optimal from test data
  geom_vline(xintercept = 1/11, color='purple') + # theoretical optimal
  ggtitle('Cost of FP = 1; Cost of FN=10') +
  labs(x="threshold (p_hat)")

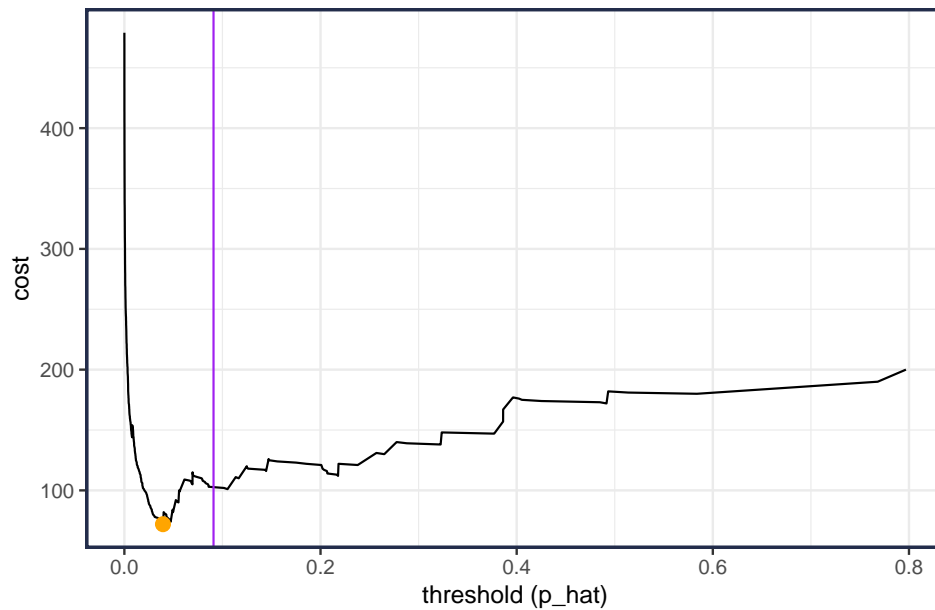
perf_table %>%
  mutate(cost = 10*FP + 1*FN) %>% # use 10:1 costs
  ggplot(aes(p_hat, cost)) + geom_line() +
  geom_point(data=. %>% filter(cost==min(cost)), size=3, color='orange') + # optimal from test data

```

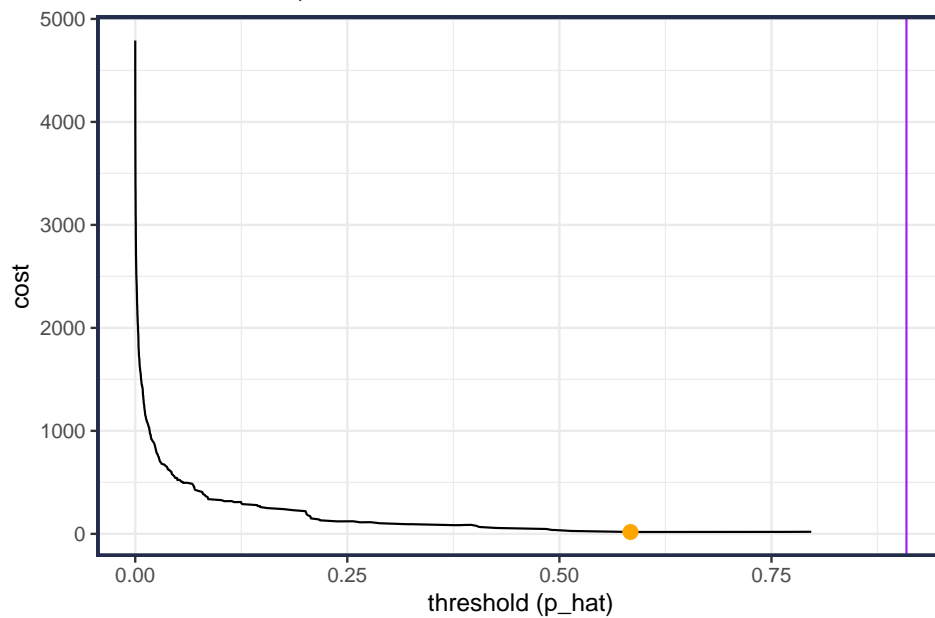


```
geom_vline(xintercept = 10/11, color='purple') + # theoretical optimal
ggtitle('Cost of FP = 10; Cost of FN=1') + labs(x="threshold (p_hat)")
```

Cost of FP = 1; Cost of FN=10



Cost of FP = 10; Cost of FN=1



```
#: Make ROC curve
perf_table %>%
  ggplot(aes(FPR, TPR)) +
  geom_path() +
  labs(x='FPR (1-specificity)', y='TPR (sensitivity)') +
  geom_segment(x=0, xend=1, y=0, yend=1, lty=3, color='grey50') +
  scale_x_continuous(breaks = seq(0, 1, by=.20)) +
  scale_y_continuous(breaks = seq(0, 1, by=.20)) +
  ggtitle("ROC Curve")
```

```
## Using yardstick package
library(yardstick) # for evaluation functions
# Notes:
# - for ROC curve and AUROC, it doesn't matter if the estimates/predictions
#   are p_hat or gamma_hat
# - for

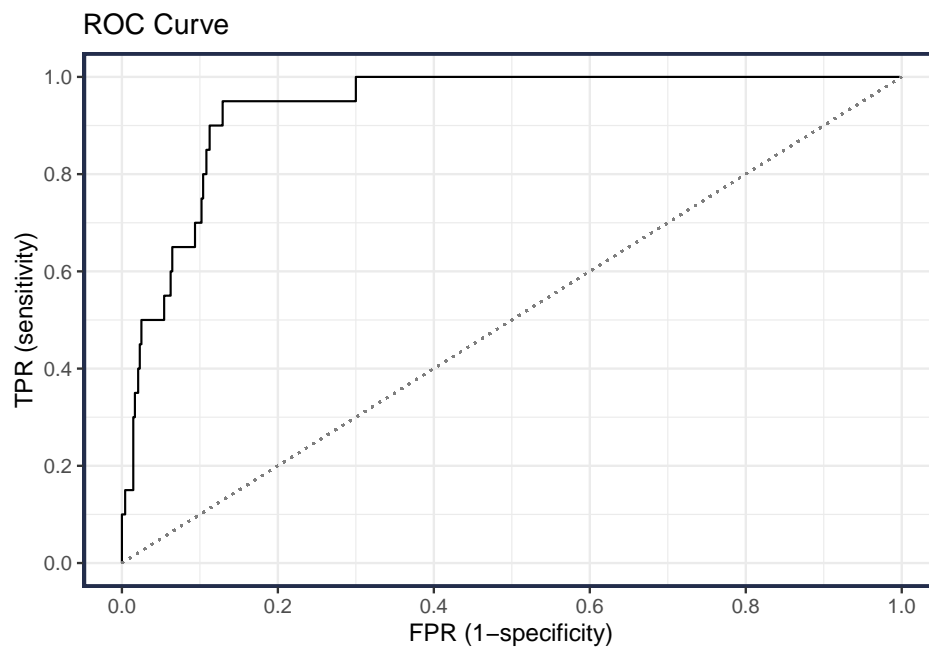
#: ROC plots
ROC =
  perf_data %>%
  mutate(y = factor(y, levels=c(1,0))) %>%
  yardstick::roc_curve(y, gamma_hat)

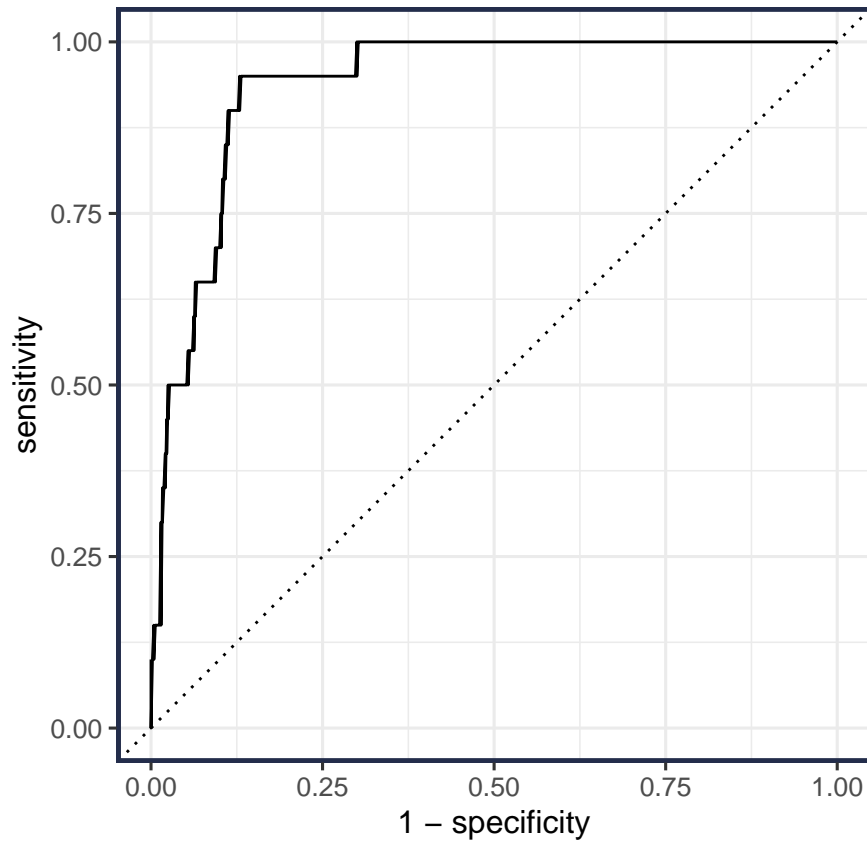
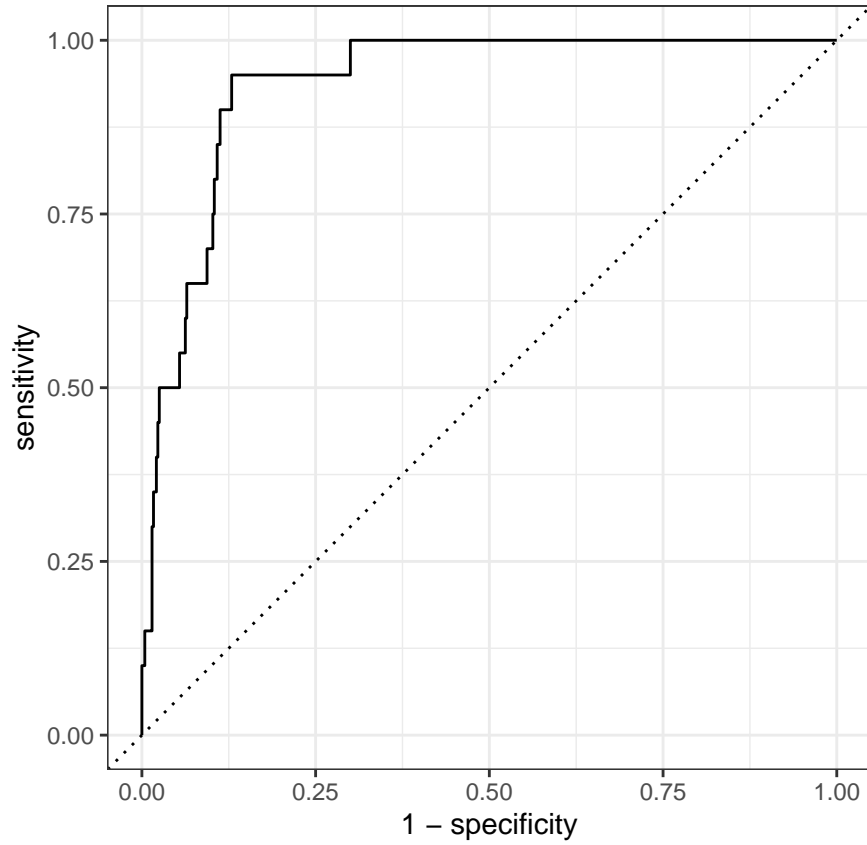
autoplot(ROC) # autoplot() method

ROC |> # same as autoplot()
  ggplot(aes(1-specificity, sensitivity)) + geom_line() +
  geom_abline(lty=3) +
  coord_equal()

#: Area under ROC (AUROC)
perf_data |>
  mutate(y = factor(y, levels=c(1,0))) |>
  roc_auc(y, gamma_hat)
#> # A tibble: 1 x 3
#>   .metric .estimator .estimate
#>   <chr>   <chr>       <dbl>
#> 1 roc_auc binary         0.937

yardstick::roc_auc_vec(factor(perf_data$y, 1:0), perf_data$gamma_hat)
#> [1] 0.9368
```





```
#: Log Loss Metric
perf_data |>
  mutate(y = factor(y, levels=c(1,0))) |>
  mn_log_loss(y, p_hat)
#> # A tibble: 1 x 3
#>   .metric      .estimator .estimate
#>   <chr>        <chr>        <dbl>
#> 1 mn_log_loss binary          0.104
```

```
#: Precision-Recall
perf_table %>%
  mutate(threshold = p_hat, precision = TP/(TP + FP)) %>%
  select(threshold, TPR, precision) %>%
  gather(metric, n, -threshold) %>%
  ggplot(aes(threshold, n, color=metric)) + geom_line() +
  scale_x_continuous(breaks = seq(0, 1, by=.20)) +
  scale_y_continuous(breaks = seq(0, 1, by=.20)) +
  scale_color_manual(values = c(TPR="blue", precision="brown")) +
  labs(x="thredhold (p_hat)", y="score")
#> Warning: Removed 1 row containing missing values or values outside the scale range
#> (`geom_line()`).
```

```
perf_table %>%
  mutate(threshold=p_hat, precision = TP/(TP + FP)) %>%
  ggplot(aes(TPR, precision)) + geom_line() +
  scale_x_continuous(breaks = seq(0, 1, by=.20)) +
  scale_y_continuous(breaks = seq(0, 1, by=.20)) +
  labs(x='Recall (TPR)', y='Precision', # (TP/(TP+FP))
       title="Precision-Recall Curve")
#> Warning: Removed 1 row containing missing values or values outside the scale range
#> (`geom_line()`).
```

