

Classification II

Lecture 5

Termeh Shafie

1

The Bayes in Naive Bayes

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(\text{category} | x_1, x_2, x_3) = \frac{P(x_1, x_2, x_3 | \text{category})P(\text{category})}{P(x_1, x_2, x_3)}$$

2

The Bayes in Naive Bayes

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

for computational efficiency

$$P(\text{diabetes} | \text{HS}, \text{O}, \text{HRH}) = \frac{P(\text{HS}, \text{O}, \text{HRH} | \text{diabetes}) P(\text{diabetes})}{P(\text{HS}, \text{O}, \text{HRH})}$$

$$P(\text{diabetes}^c | \text{HS}, \text{O}, \text{HRH}) = \frac{P(\text{HS}, \text{O}, \text{HRH} | \text{diabetes}^c) P(\text{diabetes}^c)}{P(\text{HS}, \text{O}, \text{HRH})}$$

diabetes = high/low risk
HS = high sugar intake
O = Obese
HRH = high resting heart rate

3

The Bayes in Naive Bayes

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

$$P(\text{diabetes} | \text{HS}, \text{O}, \text{HRH}) \propto P(\text{HS}, \text{O}, \text{HRH} | \text{diabetes}) P(\text{diabetes})$$

$$P(\text{diabetes}^c | \text{HS}, \text{O}, \text{HRH}) \propto P(\text{HS}, \text{O}, \text{HRH} | \text{diabetes}^c) P(\text{diabetes}^c)$$

diabetes = high/low risk
HS = high sugar intake
O = Obese
HRH = high resting heart rate

4

The Naive in Naive Bayes

features are **conditionally independent** given the class label

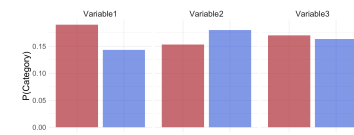
$$P(\text{HS}, \text{O}, \text{HRH}) = P(\text{HS}) \cdot P(\text{O}) \cdot P(\text{HRH})$$

diabetes = high/low risk
HS = high sugar intake
O = Obese
HRH = high resting heart rate

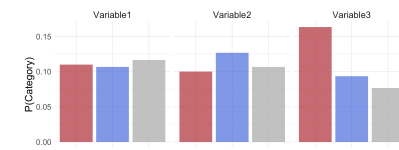
5

Naive Bayes

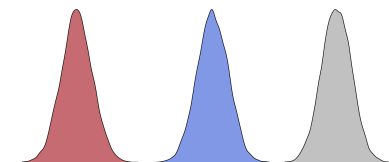
Bernoulli Naive Bayes



Categorical Naive Bayes



Gaussian Naive Bayes




6

Bernoulli Naive Bayes



	spam	dear	lunch	viagra	money
0	0.25	0.46	0.01	0.14	
1	0.32	0.05	0.53	0.67	

$$P(\text{category} | x_1, x_2, \dots, x_p) \propto \prod_{i=1}^p P(x_i | \text{category}) \cdot P(\text{category})$$

 [1,1,0,0]

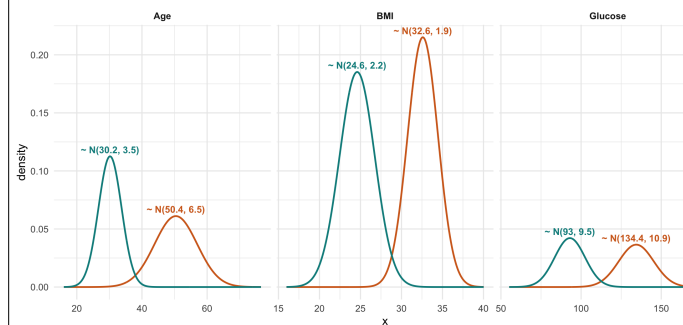
7

Gaussian Naive Bayes

predicting risk of diabetes

estimated Normal(μ , σ) distributions for low-risk and high-risk groups

prediction is based on which class makes the observed values more likely

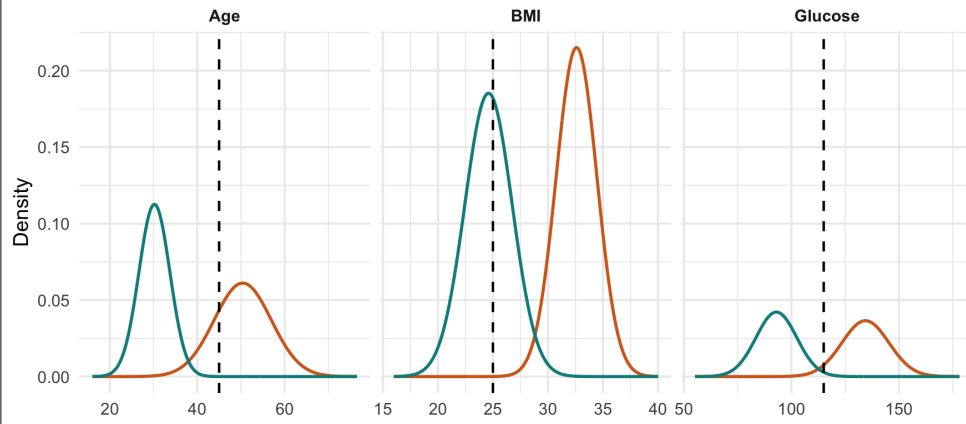


ID	BMI	Glucose	Age	Risk
1	22	82	25	low
2	24	90	30	low
3	23	88	28	low
4	26	95	35	low
5	31	125	45	high
6	33	135	50	high
7	35	142	55	high
8	30	120	42	high
9	28	110	33	low
10	34	150	60	high
⋮	⋮	⋮	⋮	⋮

8

Gaussian Naive Bayes

new observation:
Age = 45, BMI = 25, Glucose = 115

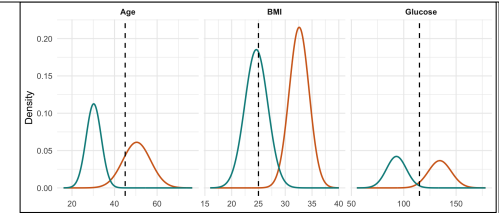


9

Gaussian Naive Bayes

class-conditional normal distribution:

$$p(x_j | y = c) = \frac{1}{\sqrt{2\pi\sigma_{c,j}^2}} \exp\left(-\frac{(x_j - \mu_{c,j})^2}{2\sigma_{c,j}^2}\right)$$



$$P(L) \times P(\text{Age} = 45 | L) \times P(\text{BMI} = 25 | L) \times P(\text{Glucose} = 115 | L)$$

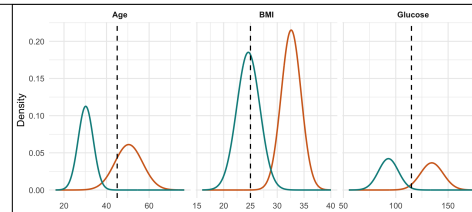
$$P(H) \times P(\text{Age} = 45 | H) \times P(\text{BMI} = 25 | H) \times P(\text{Glucose} = 115 | H)$$

10

Gaussian Naive Bayes

class-conditional normal distribution:

$$p(x_j | y = c) = \frac{1}{\sqrt{2\pi\sigma_{c,j}^2}} \exp\left(-\frac{(x_j - \mu_{c,j})^2}{2\sigma_{c,j}^2}\right)$$



Feature	Likelihood (low)	Log-likelihood (low)
Age	0.00009406119	-9.2716
BMI	0.1633826	-1.8117
Glucose	0.004344250	-5.4389

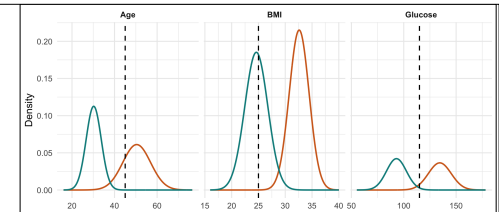
Feature	Likelihood (high)	Log-likelihood (high)
Age	0.0415668350	-3.1805
BMI	0.0002329876	-8.3645
Glucose	0.0092573944	-4.6823

11

Gaussian Naive Bayes

class-conditional normal distribution:

$$p(x_j | y = c) = \frac{1}{\sqrt{2\pi\sigma_{c,j}^2}} \exp\left(-\frac{(x_j - \mu_{c,j})^2}{2\sigma_{c,j}^2}\right)$$



0.5

$$P(L) \times P(\text{Age} = 25 | L) \times P(\text{BMI} = 25 | L) \times P(\text{Glucose} = 115 | L)$$

$$P(H) \times P(\text{Age} = 25 | H) \times P(\text{BMI} = 25 | H) \times P(\text{Glucose} = 115 | H)$$

0.5

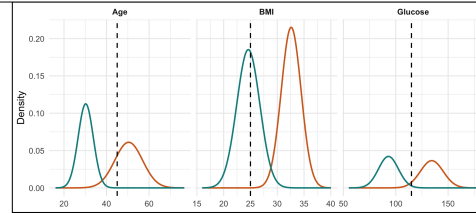
note: with log likelihood you sum up the probabilities and add the log of the prior

12

Gaussian Naive Bayes

class-conditional normal distribution:

$$p(x_j | y = c) = \frac{1}{\sqrt{2\pi\sigma_{c,j}^2}} \exp\left(-\frac{(x_j - \mu_{c,j})^2}{2\sigma_{c,j}^2}\right)$$



$$P(L | P(\text{Age} = 25, P(\text{BMI} = 25), P(\text{Glucose} = 115) \approx 0.43^*$$

$$P(H | P(\text{Age} = 25, P(\text{BMI} = 25), P(\text{Glucose} = 115) \approx 0.57^*$$

* true probabilities after normalization

13

Bayes Classifier vs. Naive Bayes

Warning: The Bayes classifier should not be confused with a Naive Bayes classifier!

- **Bayes optimal classifier** (or **Bayes classifier**) is a theoretical construct, not a practical classification method
- It is defined as the classifier that has the smallest test error rate and **assumes we know** $P(\text{category} | \text{predictors})$
- It's what we would ideally use if we knew the true data generating process

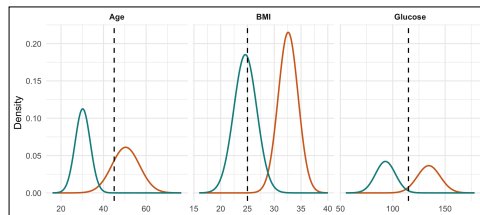
A probabilistic model-based approach to using Bayes classifier is:

1. Estimate the true distribution of test set from the training set
2. Use the Bayes optimal classifier for the estimated distribution

14

Bayes Classifier vs. Naive Bayes

Warning: The Bayes classifier should not be confused with a Naive Bayes classifier!



The Bayes optimal classifier uses the true joint distribution of Age, BMI, and Glucose to predict diabetes risk

Naive Bayes simplifies this by assuming the predictors are independent and modeling each with its own $\text{Normal}(\mu, \sigma)$ distribution.

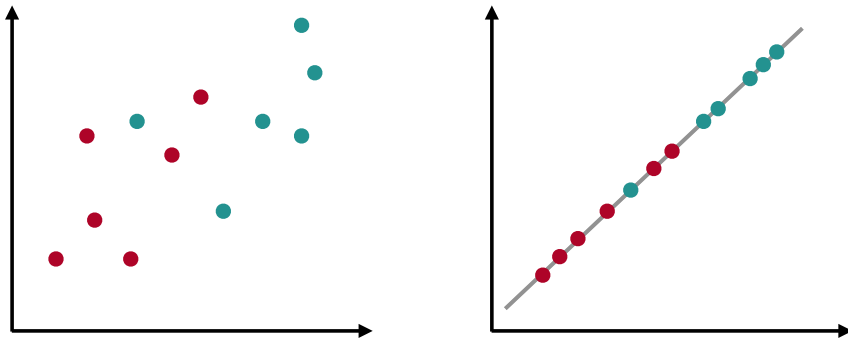
15

Linear Discriminant Analysis (LDA)

- Model the distribution of predictors in each category separately
- Use **Bayes theorem** to flip things around and obtain $P(\text{category} | \text{predictors})$
- Naive Bayes: features are **conditionally independent given the class label**
- Now: **model the joint distribution of features given the class label**
 - ▶ assume distribution of the features within each category is normally distributed
 - ▶ assume covariances of the MVN distributions are equal for both classes
 - ▶ use the Bayes optimal classifier

16

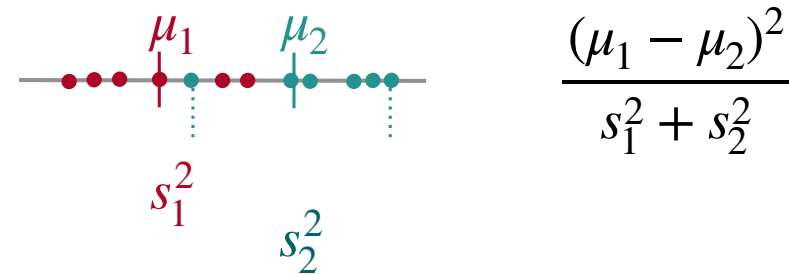
Linear Discriminant Analysis (LDA)



17

Linear Discriminant Analysis (LDA)

- 1) maximize the distance between the means
- 2) minimize the variation (s^2) within each category



18

LDA with One Predictor

- $f_k(x)$ is normal with following density in one dimension:

$$f_k(x) = \frac{1}{\sqrt{2\pi\sigma_k}} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)$$

where μ_k and σ_k^2 are mean and variance of k th class and assume variances are equal

- Plug this into Bayes theorem

$$P_k(x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)} = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}$$

- The Bayes classifier assigns an observation to where the above is the largest which is equivalent to the largest discriminant score: $\delta_k(x) = x \frac{\mu_k}{\sigma^2} + \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$
- This is **the linear discriminant classifier**



19

LDA with Multiple Predictors

- Each class k has a multivariate normal distribution:

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right)$$

where p = number of predictors

$\boldsymbol{\mu}_k$ = mean vector of class k

Σ = common covariance matrix (**same for all classes**)

- Plug this into Bayes theorem

$$P_k(\mathbf{x}) = \frac{\pi_k f_k(\mathbf{x})}{\sum_{l=1}^K \pi_l f_l(\mathbf{x})}$$

- The Bayes classifier assigns an observation to the class with the largest discriminant score

$$\delta_k(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k + \frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log(\pi_k).$$

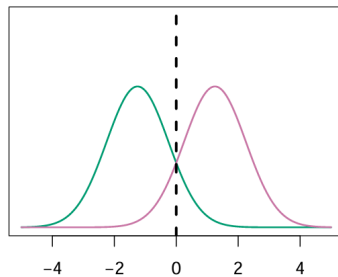
- This is **the linear discriminant classifier** (linear in \mathbf{x}).



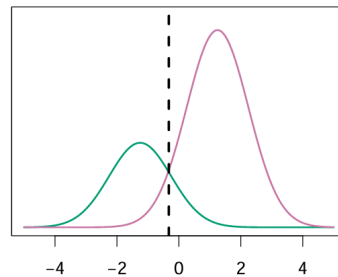
20

Linear Discriminant Analysis (LDA)

$$\pi_1 = .5, \pi_2 = .5$$



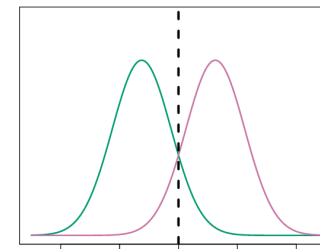
$$\pi_1 = .3, \pi_2 = .7$$



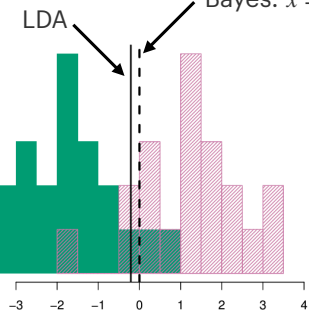
- the dashed line represents the Bayes decision boundary (Bayes Classifier)
- we classify a new point to which density is highest
- when priors are different, take them into account and compare $\pi_k f_k(x)$
- on the right, we favor the pink class the decision boundary has shifted to the left

21

Linear Discriminant Analysis (LDA)



$$\text{Bayes: } x = \frac{\mu_1 + \mu_2}{2}$$

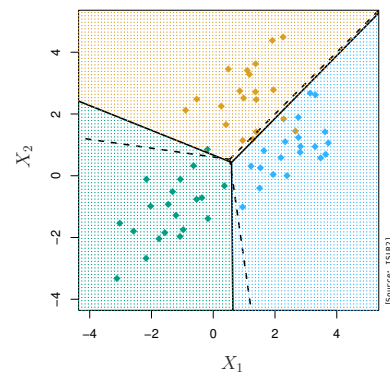
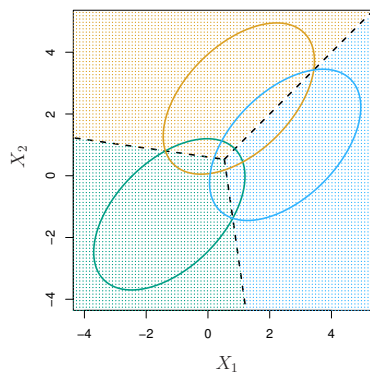


$$\mu_1 = -1.5, \mu_2 = 1.5, \sigma_1^2 = \sigma_2^2 = 1, \pi_1 = \pi_2 = 0.5$$

- typically we don't know these parameters
- in that case, we estimate them and plug them into the rule

22

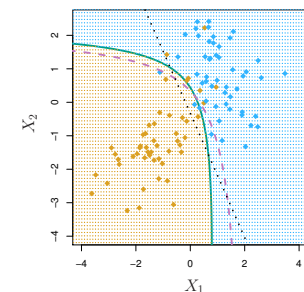
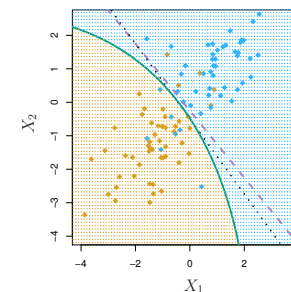
LDA with Three Classes



23

Quadratic Discriminant Analysis (QDA)

does not assume a common covariance across classes for these MVNs



- LDA uses a linear decision rule (less flexible)
 - constrains such that it uses same covariance matrix for each class
- QDA uses a quadratic decision rule (more flexible)
 - allows each class k to have a different covariance matrix

24

KNN vs. Logistic Regression vs. LDA vs. QDA

- **KNN**: good when complex boundaries and n is sufficiently large
- **Logistic regression** and **LDA**: good when linear boundaries or p is big relative to n
 - LDA extends better to multi-class problems
 - LDA is more stable during estimation
 - Logistic regression is more robust to outliers
- **QDA**: good when quadratic (or moderately complex) boundaries and n is moderately big

25

Measuring Classification Performance

- Measure of classification performance is
error rate = fraction of points that are classified incorrectly

- The **training error rate** is

$$\text{training error} = \frac{1}{n} \sum_{i=1}^n I(\hat{y}_i \neq y_i)$$

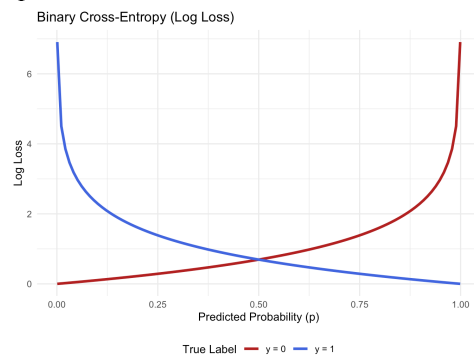
- The (expected) **test error rate** is given by $\mathbf{E} \left(I(\hat{Y}_0 \neq Y_0) \right)$
- We have to construct \hat{f} to **minimize the test error rate**
 \implies we need a **loss function** $L(\hat{y}, y)$ for penalizing errors in $\hat{y} = \hat{f}(x)$ when truth is y
- Strongly contingent on application

26

Binary Cross Entropy

$$l(\beta_0, \beta_1) = -\frac{1}{n} \sum_{i=1}^n y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i))$$

Loss function



27

Assessing Model Performance

- Did it make the correct prediction?
 - accuracy
 - sensitivity
 - specificity
- How well does to perform in distinguishing classes correctly?

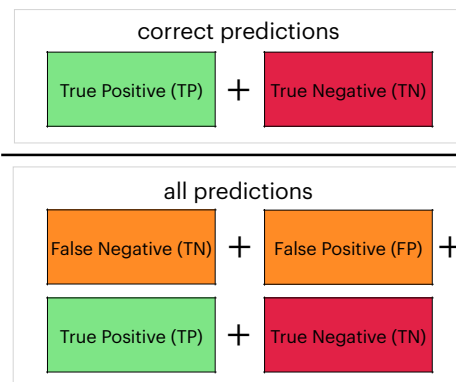
28

Confusion Matrix

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (TN)
	Negative	False Positive (FP)	True Negative (TN)

29

Confusion Matrix: Accuracy

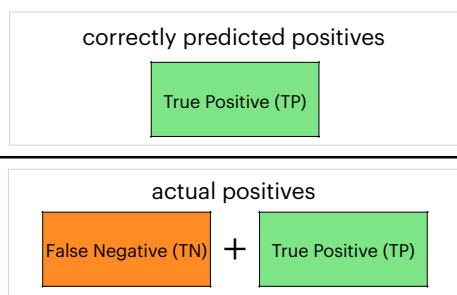


How often is the model correct?

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (TN)
	Negative	False Positive (FP)	True Negative (TN)

30

Confusion Matrix: Sensitivity/Recall

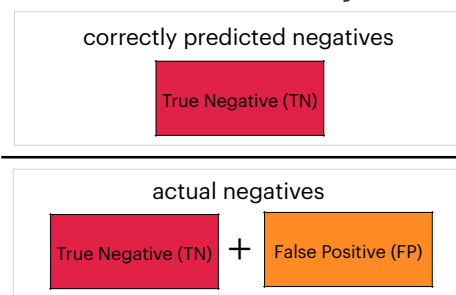


How often is the model correct for Positive Cases?

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (TN)
	Negative	False Positive (FP)	True Negative (TN)

31

Confusion Matrix: Specificity

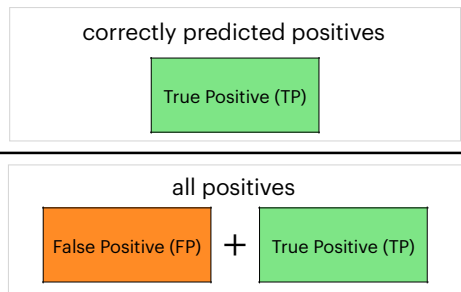


How often is the model correct for Negative Cases?

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (TN)
	Negative	False Positive (FP)	True Negative (TN)

32

Confusion Matrix: Precision



How many of the predicted
Positives are correct?

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

33

Confusion Matrix: F1 Score

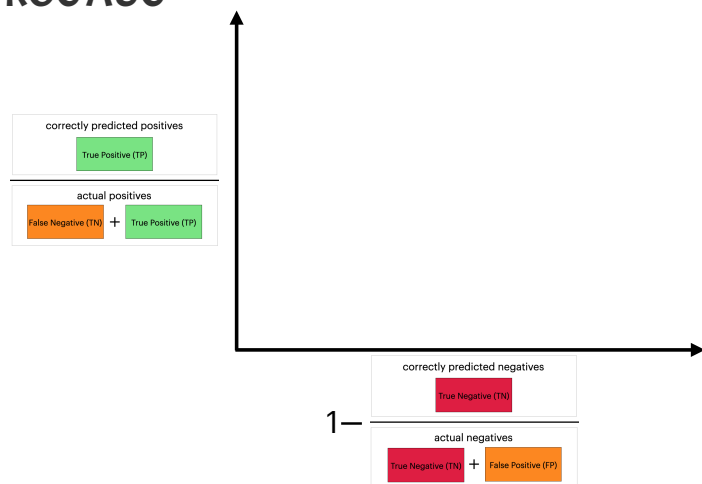
$$2 \times \frac{1}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

Combination of
Precision (how often predicted
positives ARE positive) and
Recall (how often we correctly
predict actual positives)

		Predicted	
		Positive	Negative
Actual	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

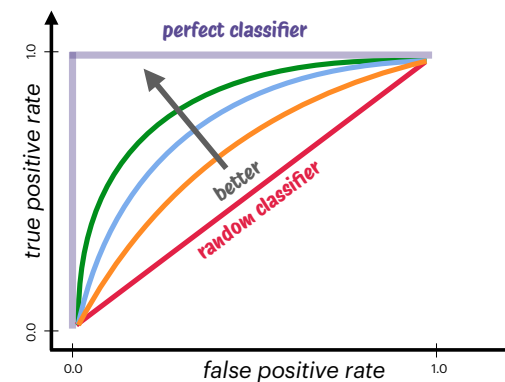
34

ROC AUC



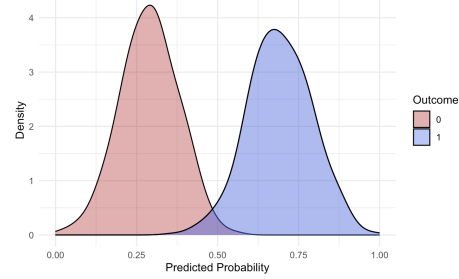
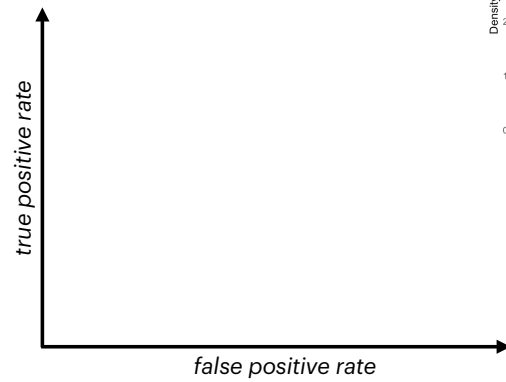
35

ROC AUC



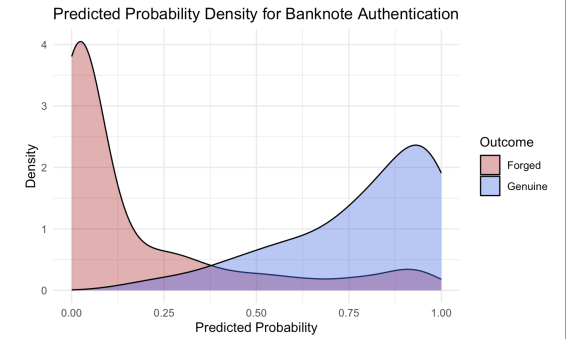
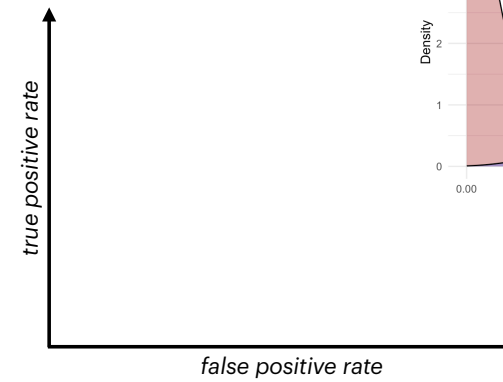
36

ROC AUC



37

ROC AUC



more realistic...

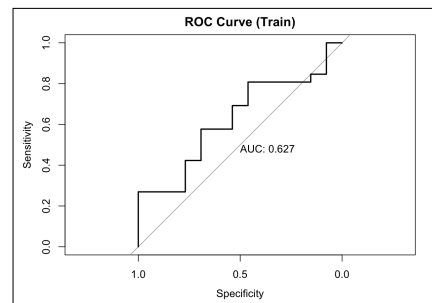
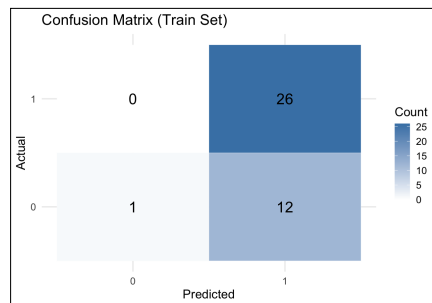
38

This Week's Practical

Classification and Model Evaluation



spam or ham text message



39