

STA 141A – Fundamentals of Statistical Data Science

Department of Statistics; University of California, Davis

Instructor: Dr. Akira Horiguchi (ahoriguchi@ucdavis.edu)

Ao1 TA: Zhentao Li (ztlli@ucdavis.edu)

Ao2 TA: Zijie Tian (zjztian@ucdavis.edu)

Ao3 TA: Lingyou Pang (lyopang@ucdavis.edu)

Section 7: Classification with R

Spring 2025 (Mar 31 – Jun 05), MWF, 01:10 PM – 02:00 PM, Young 198

Based on Chapter 4 of ISL book James et al. (2021).

- For more R code examples, see R Markdown files in <https://www.statlearning.com/resources-second-edition>

Section 7: Classification.

- Logistic regression
- Alternatives to logistic regression
- Linear discriminant analysis for $p = 1$
- Idea of linear discriminant analysis for $p > 1$
- Idea of quadratic discriminant analysis for $p > 1$
- Naive Bayes
- Errors in classification
- Comparison of classification methods

1. A person arrives at the emergency room with a set of symptoms that could possibly be attributed to one of three medical conditions. Which of these medical conditions does the person have based on the symptoms given?
2. An online banking service must be able to determine whether or not a transaction being performed on the site is fraudulent, on the basis of the user's IP address, past transaction history, and so forth.
3. On the basis of DNA sequence data for a number of patients with and without a given disease, one would like to figure out which DNA mutations are disease-causing and which are not.

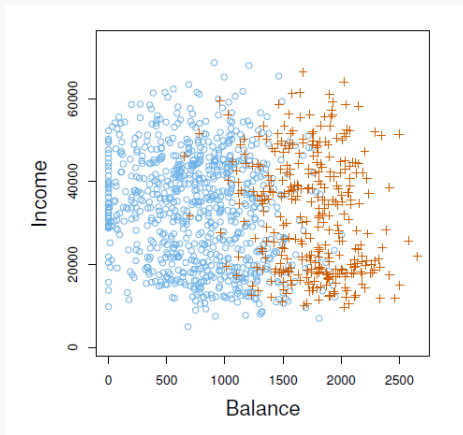


Figure 1: Image by James et al. (2021), based on the Default data set in R. The annual incomes and monthly credit card balances of a number of individuals, where the individuals who defaulted on their credit card payments are shown in orange, and those who did not are shown in blue.

Classification refers to the task of predicting *qualitative/categorical* responses

- Each response y_i is a discrete value in a predetermined category.
- Predicting a qualitative response for an observation can be referred to as *classifying* that observation, as one assigns the observation to a certain category/class.
- (In contrast, regression deals with “continuous” numeric response values.)

As in regression...

- We use training observations $(x_1, y_1), \dots, (x_n, y_n)$ to find the best estimator in the allowed class of models.
 - ▶ That is, we try to find the best estimator (among the allowed class) that fits the training data.
- We then evaluate how well the estimator performance generalizes to unseen data.
 - ▶ We can do this directly if there are also test observations $(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})$.
 - ▶ Otherwise, we can use a resampling method to estimate the estimator's generalization ability (Sec 8).

WHY NOT LINEAR REGRESSION?

- In example 1 above, a person arrives at the emergency room with a set of symptoms. We would like to treat the person based on three reasonable medical conditions: "Appendicitis", "Food poisoning", "Gastritis".
- We could assign each medical condition Y a number from 1 to 3:

$$Y = \begin{cases} 1, & \text{if "Appendicitis",} \\ 2, & \text{if "Food poisoning",} \\ 3, & \text{if "Gastritis".} \end{cases} \quad \text{or} \quad Y = \begin{cases} 1, & \text{if "Gastritis",} \\ 2, & \text{if "Appendicitis",} \\ 3, & \text{if "Food poisoning".} \end{cases}$$

Both approaches work, but imply a totally different relationship.

- Can we use linear regression for a *binary* (two levels) response? In the banking example, the two transaction categories can be coded as

$$Y = \begin{cases} 1, & \text{if "Fraudulent",} \\ 0, & \text{if "Not fraudulent".} \end{cases}$$

We have no problem with ordering here, but using linear regression here (for estimating probabilities to assign values) could lead to values smaller than zero or larger than one (unreasonable for probabilities!).

PROBLEM WITH USING LINEAR REGRESSION

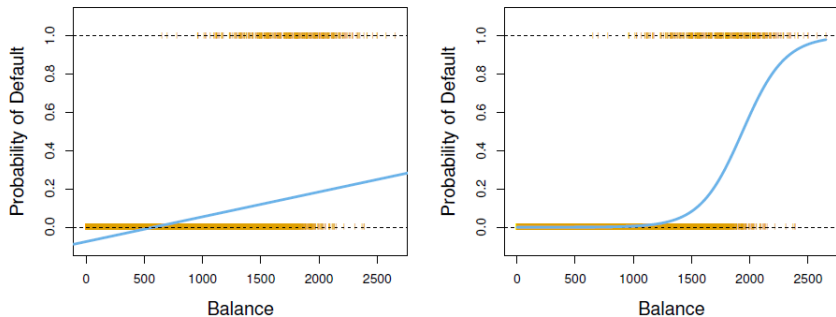


Figure 2: Image by James et al. (2021), based on the Default data set in R. Left: The estimated probability of default using linear regression, where the orange ticks indicate the values "0" for "No", and "1" for "Yes". Right: Predicted probabilities of default using logistic regression, where all probabilities lie between 0 and 1.

If not linear regression, then what can we use?

- *Logistic regression* (usable only for binary responses) is perhaps the most related, so let's start there.

CLASSIFICATION WITH R

LOGISTIC REGRESSION

IDEA: BINARY CLASSIFICATION

Classification: compute/estimate conditional probability $P(Y = k|X)$ for each class k . $P(Y=1|X)$, $P(Y=2|X)$, $P(Y=3|X)$, ..., $P(Y=K|X)$

- If there are only two classes, we only need $P(Y = 1|X)$. (Why?)

$$P(Y=0|X) + P(Y=1|X) = 1$$


$$\Rightarrow P(Y=0|X) = 1 - P(Y=1|X)$$

If we know $P(Y=1|X)$, then we also know $P(Y=0|X)$.

Suppose our two classes are coded as 0 and 1 (e.g., "no" and "yes")

- Given a value or estimate of the conditional probability

$$p(X) = P(Y = 1|X). \quad (1)$$

a default decision rule is to predict 1 (or "yes") if and only if $p(X) > 0.5$ 

- Consequence might be worse for misclassifying a "yes" than for misclassifying a "no", in which case we might use the decision rule: predict 1 (or "yes") if and only if $p(X) > 0.8$

Logistic regression models probabilities of binary responses

- Boils down to modelling probability of "yes": (1)
- Since probabilities have values in $[0, 1]$, we cannot use the linear approach

$$\underline{p(X)} = \underline{\beta_0 + \beta_1 X} \quad (2)$$

as on the left side in Figure 2, and need a function with values in $[0, 1]$.

- The *logistic function* $g(x) = \frac{e^x}{1+e^x}$ maps the real line into the interval $(0, 1)$.
- Logistic regression maps the linear form (2) into values between 0 and 1 in order to model the probabilities by

$$\underline{p(X)} = \underline{g(\beta_0 + \beta_1 X)} = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} \quad (3)$$

- (3) implies (by rearranging the terms)

$$\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}, \quad (4)$$

where this quantity is called *odds*.

- Applying the (natural) logarithm on both sides of (4) yields

$$\log \left(\frac{p(X)}{1 - p(X)} \right) = \beta_0 + \beta_1 X,$$

where the LHS is called *log odds* or *logit*. Advantage: linear in X .

- We can make predictions from estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ for β_0 and β_1 by using the estimated probability

$$\hat{p}(X) = g\left(\hat{\beta}_0 + \hat{\beta}_1 X\right) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}},$$

where we recall that g is the logistic function.

- If this estimated probability is over a certain pre-defined threshold (e.g. 0.25), then we would assign $X = x$ to the class 1.
- Example: If $\hat{\beta}_0 = -9.9$ and $\hat{\beta}_1 = 0.005$, we predict the default probabilities of individuals with balance $X = \$1,000$ and $X = \$2,000$ by

$$\hat{p}(X = \underline{1,000}) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}} = \frac{e^{-9.9 + 0.005 \cdot 1,000}}{1 + e^{-9.9 + 0.005 \cdot 1,000}} \approx 0.007,$$

$$\hat{p}(X = \underline{2,000}) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X}} = \frac{e^{-9.9 + 0.005 \cdot 2,000}}{1 + e^{-9.9 + 0.005 \cdot 2,000}} \approx 0.525.$$

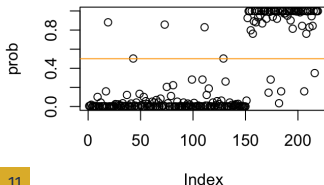
We use `glm()` for logistic regression ('glm' stands for *general linear model*).

- Must specify which variables are used, data set, and type of response.
- Must put family=binomial to specify a binary response.

```

1 library(palmerpenguins)
2 ?penguins # get help
3
4 # We work with only Adelie and Chinstrap species (we exclude Gentoo).
5 peng_binary <- na.omit(penguins[penguins$species != 'Gentoo', ])
6 logreg <- glm(species ~ bill_length_mm, data=peng_binary, family=binomial)
7 prob <- predict(logreg, type='response') # 'link' also possible
8 predicted <- ifelse(prob<.5, 'Adelie', 'Chinstrap')
9
10 plot(prob)
11 abline(a=0.5, b=0, col='orange')

```



If we have p predictors $X = (X_1, \dots, X_p)$, we have a *multiple logistic regression model*:

- The logit is

$$\log \left(\frac{p(X)}{1 - p(X)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

If $p=1$

- The probabilities are

$$p(X) = g(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}.$$

CLASSIFICATION WITH R

ALTERNATIVES TO LOGISTIC REGRESSION

Recall: logistic regression directly models $P(Y = k|X = x)$ for binary responses by using the logistic function.

Some issues:

- When there are big differences between two classes, the parameter estimates in the logistic regression model are unstable.
- If the distribution of X is approximately normal in each of the classes and the sample size is small, then other approaches are more accurate than logistic regression.

Theorem 2: Bayes' theorem

Let $\Omega \neq \emptyset$. For any events $A, B \subseteq \Omega$ with $P(B) \neq 0$ holds,

$$P(A \cap B) = P(A|B)P(B)$$

$$P(B \cap A) = P(B|A)P(A)$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}. \quad (5)$$

$$P(A|B)P(B) = P(B|A)P(A) \quad \text{divide both sides by } P(B)$$

Often $P(B)$ is unknown, but law of total probability can help compute $P(B)$.

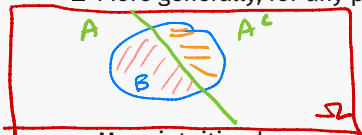
- Because the sets A and A^c partition Ω , we can write $P(B)$ as

$$P(B) = P(B \cap \Omega) = P(B \cap [A \cup A^c]) = P([B \cap A] \cup [B \cap A^c]) = P(B \cap A) + P(B \cap A^c)$$

$$P(B) = P(B|A)P(A) + P(B|A^c)P(A^c).$$

- More generally, for any partition $\{A_1, A_2, \dots\}$ of Ω , we can write $P(B)$ as

$$P(B) = \sum_{j=1}^{\infty} P(B|A_j)P(A_j).$$



- More intuition here

<https://www.youtube.com/watch?v=9wCnvr7Xw4E>

We can instead use Bayes' theorem to estimate $P(Y = k|X = x)$:

$$p_k(x) := P(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{\ell=1}^K \pi_{\ell} f_{\ell}(x)}, \quad (6)$$

where

- $\pi_k := P(Y = k)$ is the overall or *prior* probability that a randomly chosen observation comes from the k th class.
 - Here π is just a variable name — not the same as $\pi = 3.14159 \dots$!
- $f_k(X)$ is the PMF/PDF of X given that the response is from the k th class.

How to estimate these quantities π_k and f_k ?

- π_k — the proportion of observed elements in the k th class
 - E.g. if there are 3, 2, 5 elements in the classes 1, 2, 3, respectively, then the estimated probabilities are $\hat{\pi}_1 = \frac{3}{10}$, $\hat{\pi}_2 = \frac{2}{10}$, $\hat{\pi}_3 = \frac{5}{10}$.
- Estimating f_k is more challenging — approaches will be discussed in the next few slides.

CLASSIFICATION WITH R

LINEAR DISCRIMINANT ANALYSIS FOR $p = 1$

This approach estimates $f_k(x)$ by making some simplifying assumptions:

- f_k is a normal/Gaussian PDF, i.e. for all x holds

$$\underline{f_k(x)} = \frac{1}{\sqrt{2\pi}\underline{\sigma_k}} \exp \left\{ -\frac{1}{2\sigma_k} (x - \underline{\mu_k})^2 \right\}. \quad (7)$$

- Same variance parameter across all K classes: $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_K^2 = \sigma^2$.

K classes $\Rightarrow \mu_1, \dots, \mu_K$ & σ^2

With these assumptions, we plug in this PDF (7) into (6) to get

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu_k)^2\right\}}{\sum_{\ell=1}^K \pi_\ell \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu_\ell)^2\right\}}. \quad (8)$$

- Recall: *Bayes classifier* assigns obsn x to class $\arg \max_{k \in \{1, 2, \dots, K\}} p_k(x)$.
 $\log(LS) = \log(\pi_k) - \frac{1}{2\sigma^2}(x - \mu_k)^2$
 $-\log\left[\sum_{\ell=1}^K \pi_\ell \exp\left\{-\frac{1}{2\sigma^2}(x - \mu_\ell)^2\right\}\right]$

- This is equivalent to assigning x to class for which *discriminant function*

$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k), \quad k \in \{1, 2\}, \quad (9)$$

is largest. (Why? Take the log of (8).)

- If $K = 2$, classifier assigns x to class 1 if $\delta_1(x) > \delta_2(x)$, to class 2 otherwise.
- The *Bayes decision boundary* is the point x for which $\delta_1(x) = \delta_2(x)$.
 - What does this inequality $\delta_1(x) > \delta_2(x)$ and boundary simplify to if $\pi_1 = \pi_2$?

EXAMPLE FOR DECISION BOUNDARY

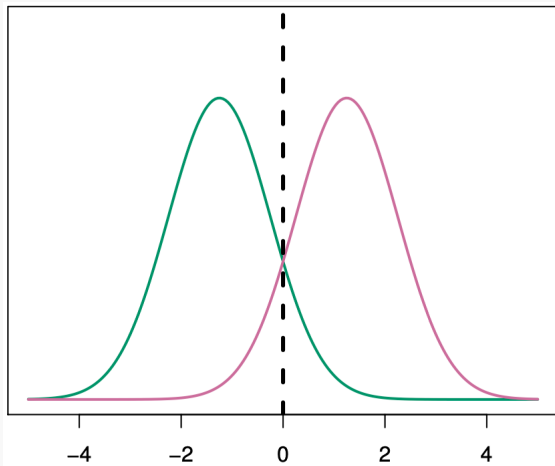


Figure 3: Image by James et al. (2021). Two pdfs of normal distributions with means $\mu_1 = -1.25$ and $\mu_2 = 1.25$, respectively, and variance $\sigma^2 = 1$. The dashed vertical line represents the Bayes decision boundary, so we assign the observation to class 1 if x is left of the line, and to class 2 otherwise.

In the plot above, we can calculate Bayes classifier because we know values for all parameters $\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \sigma^2$.

- In practice, we must estimate these parameters to apply Bayes classifier.
- *Linear discriminant analysis* (LDA) method plugs the estimates

$$\hat{\pi}_k = \frac{n_k}{n}, \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i, \quad \hat{\sigma}^2 = \frac{1}{n-K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2,$$

into (9) to get the discriminant function

$$\hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k). \quad (10)$$

- $\hat{\pi}_k$ — proportion of all training observations from k th class.
- $\hat{\mu}_k$ — average of all training observations from k th class.
- $\hat{\sigma}^2$ — weighted average of sample variances for each class.

This discriminant function $\hat{\delta}_k(x)$ is linear in x , hence the name LDA.

EXAMPLE FOR DECISION BOUNDARY

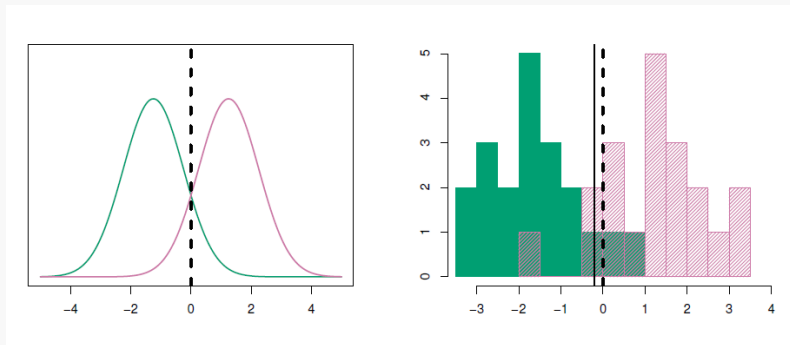


Figure 4: Image by James et al. (2021). Left: Two pdfs of normal distributions with means $\mu_1 = -1.25$ and $\mu_2 = 1.25$, respectively, and variance $\sigma^2 = 1$. The dashed vertical line represents the Bayes decision boundary, so we assign the observation to class 1 if $x < 0$ and class 2 otherwise. Right: 20 observations were drawn from each of the two classes, and are shown as histograms. The Bayes decision boundary is shown as a dashed vertical line, and the solid vertical line represents the LDA decision boundary estimated from the training data.

CALCULATION EXAMPLE

For $K = 2$ classes (class "0" and "1"), we have the five data points

$$(9, 1), (8, 0), (6, 0), (7, 1), (4, 0)$$

and want to calculate the LDA discrimination function (10) for $k = 0$ and $k = 1$.

Recall (stringent) assumptions: f_k is normal/Gaussian, and variance is same across all K classes.

- Next we introduce extensions that loosen these assumptions at the cost of increased computation and/or “too much” flexibility.

CLASSIFICATION WITH R

IDEA OF LINEAR DISCRIMINANT ANALYSIS FOR $p > 1$

- Linear discriminant analysis can also be conducted for $p > 1$ predictors.
- For $p > 1$, it is also assumed that $X = (X_1, \dots, X_p)$ is normally distributed, but since X is a vector, it is drawn from a *multivariate normal distribution* with a certain vector of means and covariance matrix.
- Discriminant function can be derived in a manner similar to as in $p = 1$.

PLOT OF MULTIVARIATE NORMAL DISTRIBUTION

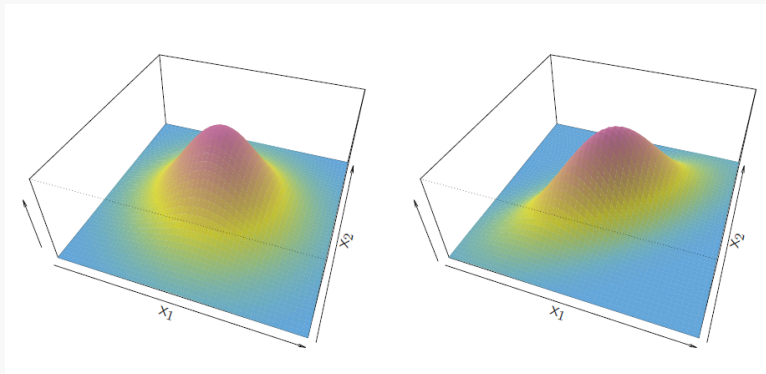


Figure 5: Image by James et al. (2021). Left: A two-dimensional normal distribution with $p = 2$ uncorrelated predictors. Right: A two-dimensional normal distribution with $p = 2$ predictors having a correlation of 0.7.

CLASSIFICATION WITH R

IDEA OF QUADRATIC DISCRIMINANT ANALYSIS FOR $p > 1$

Each class k is now allowed to have its own “covariance matrix” Σ_k .

- For $p = 1$, same as saying each class has its own variance parameter σ_k^2 .
- Why use QDA over LDA, or vice-versa? *Bias-variance trade-off* (recall from Section 3, slide 25/32).
- For p predictors, number of parameters to estimate is
 - ▶ quadratic in p for QDA,
 - ▶ linear in p for LDA.
- LDA is much less flexible, and so has much lower variance.
- But if covariance matrices wildly differ between classes, then LDA can suffer from high bias.
- Roughly speaking,
 - ▶ use LDA if n is relatively small, which makes it crucial to reduce variance,
 - ▶ use QDA if n is very large, or if assumption of a common covariance matrix for the K classes is clearly wrong.

PLOT OF MULTIVARIATE NORMAL DISTRIBUTION

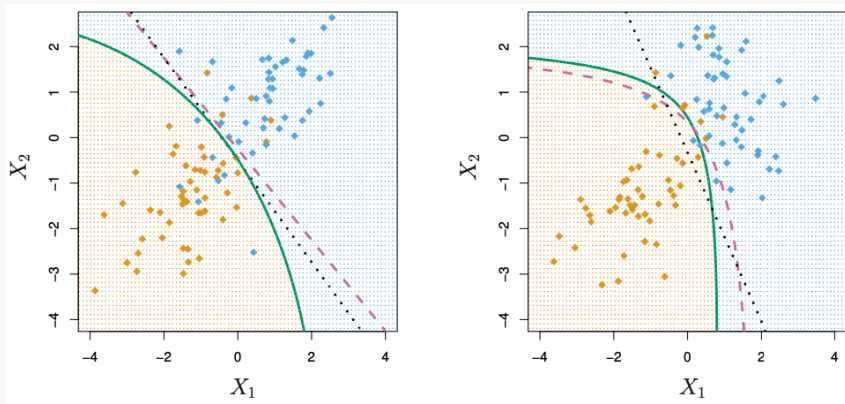


Figure 6: Image by James et al. (2021). Left: The Bayes (purple dashed), LDA (black dotted), and QDA (green solid) decision boundaries for a two-class problem with $\Sigma_1 = \Sigma_2$. The shading indicates the QDA decision rule. Since the Bayes decision boundary is linear, it is more accurately approximated by LDA than by QDA. Right: here $\Sigma_1 \neq \Sigma_2$. Since the Bayes decision boundary is non-linear, it is more accurately approximated by QDA than by LDA.

CLASSIFICATION WITH R

NAIVE BAYES

The *naive Bayes classifier* relies on Bayes' theorem.

- LDA/QDA make distribution assumptions about $f_k(x)$.
- Instead, the naive Bayes classifier assumes for each class $k = 1, \dots, K$:

Within the k th class, the p predictors are independent.

- Mathematically, this means that for each class k :

$$f_k(x) = f_{k,1}(x_1) \times f_{k,2}(x_2) \times \dots \times f_{k,p}(x_p) \quad (11)$$

where f_{kj} is the PDF/PMF of the j th predictor for obsns in the k th class.

- This assumption, although often unrealistic, produces decent results, especially when n is too small to effectively estimate the joint distribution $f_k(x)$ of predictors x within each class k .
- Without strong simplifying assumptions, estimating a joint distribution typically requires a huge amount of data.
- Bias-variance tradeoff: introduce some bias to reduce variance.
- Plug (11) into (6) to get the posterior probability

$$p_k(x) = P(Y = k | X = x) = \frac{f_{k,1}(x_1) \cdot f_{k,2}(x_2) \cdot f_{k,p}(x_p) \pi_k}{\sum_{\ell=1}^K f_{\ell,1}(x_1) \cdot f_{\ell,2}(x_2) \cdot f_{\ell,p}(x_p) \pi_{\ell}} \quad (12)$$

To estimate the one-dimensional density functions $f_{k,j}$ for all classes $k = 1, \dots, K$ and all predictors $j = 1, \dots, p$, there are several options:

- If X_j is quantitative, we can assume $X_j|Y = k \sim \mathcal{N}(\mu_{k,j}, \sigma_{k,j}^2)$ (as in LDA) or use a nonparametric approach.
- If X_j is qualitative, we could count the proportion of training observations for the j th predictor corresponding to each class k .
 - E.g. suppose we want to predict whether a student studies more than 10 hours per week based on their major (so $p = 1$). We survey 100 people:

	Math major	Art major	Poli Sci major
Study > 10 hr /wk	20	25	5
Study \leq 10 hr /wk	15	25	10

We use proportions (i.e., divide each cell by column sum) to estimate the “true” PMFs $f_{k,j}$ (each column is a PMF):

	Math major	Art major	Poli Sci major
Study > 10 hr /wk	0.6	0.55	0.35
Study \leq 10 hr /wk	0.4	0.45	0.65

CLASSIFICATION WITH R

ERRORS IN CLASSIFICATION

CONFUSION MATRIX

In classification, observations can be assigned to the wrong class.

- In binary classification one can make two mistakes: *false positives* and *false negatives*.
- Examples of not default vs default: cancer vs no cancer, spam vs not spam.
- A *confusion matrix* displays both error types.

		<i>True class</i>		
		– or Null	+ or Non-null	Total
<i>Predicted class</i>	– or Null	True Neg. (TN)	False Neg. (FN)	N*
	+ or Non-null	False Pos. (FP)	True Pos. (TP)	P*
	Total	N	P	

Name	Definition	Synonyms
False Pos. rate	FP/N	Type I error, 1–Specificity
True Pos. rate	TP/P	1–Type II error, power, sensitivity, recall
Pos. Pred. value	TP/P*	Precision, 1–false discovery proportion
Neg. Pred. value	TN/N*	

Figure 7: Tables by James et al. (2021). A confusion matrix compares the LDA predictions to the true default statuses for the 10,000 training observations in the Default data set, using a modified threshold value that predicts default for any individuals whose posterior default probability exceeds 20 %.

CONFUSION MATRIX

```
1 # Using peng_binary and predicted from previous slide
2
3 pb_species <- factor(peng_binary$species, levels=c('Adelie', 'Chinstrap'))
4 table(pb_species, predicted)
```

```
> table(pb_species, predicted)
      predicted
pb_species  Adelie Chinstrap
Adelie      141      5
Chinstrap    6     62
```

```
1 # pb_species line is not necessary, but what happens if we instead did:
2 table(peng_binary$species, predicted)
```

Recall: Bayes classifier assigns observation to class for which the posterior probability is greatest.

- For binary responses, assign x to default class if

$$P(\text{default} = \text{Yes} | X = x) > 0.5.$$

- Weights both types of mistakes (FN and FP) the same.
- But sometimes care more about lowering false negatives. E.g., a credit card company trying to detect a fraudulent charge.
- Can change the threshold from 0.5 to e.g., 0.2.
- What happens as threshold decreases w.r.t. TP rate and FP rate?

The *ROC curve* simultaneously displays both types of errors for all thresholds.

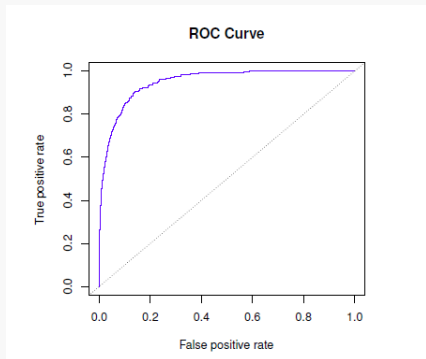


Figure 8: Image by James et al. (2021). An ROC curve for LDA classifier on Default data. Dotted line represents “no information” classifier, i.e., one that doesn’t use predictors.

- Overall performance of a classifier, summarized over all possible thresholds, is given by the *area under the ROC curve (AUC)*.
- The larger the AUC, the better the classifier.

CLASSIFICATION WITH R

COMPARISON OF CLASSIFICATION METHODS

Analytical (or mathematical) comparison:

- Classifiers with a linear decision boundary are special cases of naive Bayes. Hence, LDA is a special case of naive Bayes. (This is not obvious.)
- No method uniformly dominates others: The appropriate model depends on the predictor's distribution in each class as well as n and p .
- K -nearest neighbors (KNN) is a nonparametric approach. Hence, one can expect that it dominates naive Bayes and LDA when the true decision boundary is highly non-linear. However, KNN requires many observations relative to the number of predictors to perform well.

For an *empirical* (or data-based) comparison, see Section 4.5.2.