



FIT2086 Studio 12  
Sample Exam Questions

Daniel F. Schmidt  
October 13, 2023

Contents

1 Introduction	1
2 Short Answer Questions	2
3 Maximum Likelihood Estimation	3
4 Random Variables	4
5 Confidence Intervals and p-values	5
6 Regression	6
7 Classification	7
8 Logistic Regression	8
9 Bias and Variance	9
10 Machine Learning	10
11 Appendix I: Standard Normal Distribution Table	12
12 Appendix II: Formulas	13

1 Introduction

The Studio 12 questions are examples of the type of questions you will be asked on the exam, in number and length roughly commensurate with the real exam. Please work on this questions during, and after the studio. You may ask your demonstrator for some assistance on the questions you struggle with.

2 Short Answer Questions

Please provide a short (1-2 sentences) answer to each of the following items.

A: General statement. When answering short answer questions of this form (in general), it is a good idea to use the following basic structure: your first sentence should describe what the object/item of interest is. The second and third sentences (or fourth, the 2 – 3 is a guide and not a strict requirement) should describe one or two properties of the object/item of interest. Finally, state that you can (i) identify the object of interest, and (ii) you know something about the object of interest. All the answers below follow this basic structure.

1. Bias and variance of an estimator

The bias of an estimator  $\hat{\theta}$  is given by  $E[\hat{\theta}] - \theta$ , where  $\theta$  is the population value. It is the average amount that the estimator tends to deviate from its true value. The variance of an estimator is the average squared deviation of the estimator from its average value. It measures how much we would expect the estimate to vary if we drew a new sample from the population.

2. A p-value

A p-value is used in hypothesis testing to measure evidence against the null hypothesis. A p-value is the probability of seeing a test statistic as extreme, or more extreme, than the one we have observed, just by chance, if the null hypothesis were true.

3. Classification accuracy

(1) Classification accuracy is the percentage of times our model correctly classifies an individual object. (2) Sensitivity is the proportion of classification of individuals as "successes" (or a "1") that is correct.

(1) Accuracy

As a decision tree is a supervised machine learning method, it works by sequentially splitting the data into smaller and smaller subsets by asking binary questions about the values of the attributes of each of the individuals in our data. Each leaf contains a different, simple model for that set of individuals that is used to make predictions.

5. Ridge regression

Ridge regression is a method for estimating the coefficients of a linear or logistic regression model. It works by minimizing a goodness-of-fit score (such as the sum of the squares of the residuals) plus a complexity penalty based on the sum of the coefficients. The inclusion of the complexity term helps prevent overfitting.

6. A random variable

A random variable is a variable that takes on one value from a set of values, say  $X$ , with a frequency determined by the corresponding probability distribution over  $X$ . If we generate many realizations of a random variable the observed frequencies with which it takes the different values will be close to the theoretical frequencies.

the average difference between the estimator and the population parameter

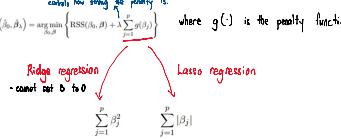
$$h_0(\hat{\theta}) = E[\hat{\theta}(Y)] - \theta$$

- If  $h_0(\hat{\theta}) < 0$ , then the estimator tends to underestimate (be smaller, on average, than) the population parameter  $\theta$ .
- If  $h_0(\hat{\theta}) > 0$ , then the estimator tends to overestimate (be greater, on average, than) the population parameter  $\theta$ .
- If  $h_0(\hat{\theta}) = 0$ , then the estimator, on average, neither overestimates or underestimates the population parameter  $\theta$ .

$$\text{Var}_\theta(\hat{\theta}) = E[(\hat{\theta}(Y) - E[\hat{\theta}(Y)])^2] = V[\hat{\theta}(Y)]$$

The variance of an estimator measures how much, on average, we expected our parameter estimate to vary under repeated sampling from our population.

$$\text{MSE}_\theta(\hat{\theta}) = b_0^2(\hat{\theta}) + \text{Var}_\theta(\hat{\theta})$$



(Week 3)

3 Maximum Likelihood Estimation

A random variable  $Y$  is said to follow an exponential distribution with a rate parameter  $\beta$ , if

$$P(Y = y | \beta) = \beta \exp(-\beta y)$$

where  $y \geq 0$  is a non-negative continuous number. Imagine we observe a sample of  $n$  non-negative real numbers  $\mathbf{y} = (y_1, \dots, y_n)$  and want to model them using an exponential distribution. (Note: remember that the data is independently and identically distributed).

1. Write down the exponential distribution likelihood function for the data  $\mathbf{y}$  (i.e., the joint probability of the data  $\mathbf{y}$  under an exponential distribution with rate parameter  $\beta$ ).

A: The data is independently and identically distributed, so the likelihood is the product of the probability for each data point

$$\begin{aligned} p(\mathbf{y} | \beta) &= \prod_{i=1}^n \beta \exp(-\beta y_i) \\ &= \beta^n \left( \prod_{i=1}^n \exp(-\beta y_i) \right) \\ &= \beta^n \exp\left(-\beta \sum_{i=1}^n y_i\right) \end{aligned}$$

where we use the fact that  $e^{-a}e^{-b} = e^{-a-b}$ .

2. Write down the negative log-likelihood function of the data  $\mathbf{y}$  under an exponential distribution with rate parameter  $\beta$ .

A: Taking negative logarithm of the above likelihood we have

$$\begin{aligned} -\log p(\mathbf{y} | \beta) &= -\log \left[ \beta^n \exp\left(-\beta \sum_{i=1}^n y_i\right) \right] \\ &= -n \log \beta + \beta \sum_{i=1}^n y_i \end{aligned}$$

where we use the facts:  $\log ab = \log a + \log b$ ,  $\log a^b = b \log a$  and  $\log e^a = a$ .

3. Derive the maximum likelihood estimator for  $\beta$ .

A: Differentiate the negative log-likelihood with respect to  $\beta$ :

$$\begin{aligned} \frac{d}{d\beta} (-\log p(\mathbf{y} | \beta)) &= -\frac{d}{d\beta} [n \log \beta] + \frac{d}{d\beta} \left\{ \beta \sum_{i=1}^n y_i \right\} \\ &= -n + \beta \sum_{i=1}^n y_i = 0 \\ &\Rightarrow \beta \sum_{i=1}^n y_i = n \\ &\Rightarrow \beta = \frac{n}{\sum_{i=1}^n y_i} \end{aligned}$$

where we use  $d \log x / dx = 1/x$ . Now set the derivative to zero and solve for  $\beta$ :

$$\begin{aligned} -n + \sum_{i=1}^n y_i &= 0 \\ \Rightarrow -n + \beta \sum_{i=1}^n y_i &= 0 \\ \Rightarrow \beta \sum_{i=1}^n y_i &= n \\ \Rightarrow \beta &= \frac{n}{\sum_{i=1}^n y_i} \end{aligned}$$

3

①

$$\begin{aligned} p(y | \mathbf{B}) &= \prod_{i=1}^n B \exp(-By_i) \\ &= B \exp(-By_1) \cdot B \exp(-By_2) \cdots B \exp(-By_n) \\ &= B^n \exp(-By_1 - By_2 - \cdots - By_n) \\ &= B^n \exp(-B \sum_{i=1}^n y_i) \end{aligned}$$

②

$$\begin{aligned} -\log p(y | \mathbf{B}) &= -\log [B^n \exp(-B \sum_{i=1}^n y_i)] \\ &= -\log B^n - \log \exp(-B \sum_{i=1}^n y_i) \\ &= -n \log B + B \sum_{i=1}^n y_i \end{aligned}$$

log ab = log a + log b  
log a^b = b log a  
log e^a = a

## ( Week 2 )

### 4 Random Variables

Suppose  $Y_1$  and  $Y_2$  are two random variables distributed as per  $Y_1 \sim \text{Po}(2)$  and  $Y_2 \sim \text{Po}(4)$ . Remember that  $\text{Po}(\lambda)$  denotes a Poisson distribution with rate parameter  $\lambda$ , which means the random variable follows the probability distribution:

$$P(Y = y) = \frac{\lambda^y e^{-\lambda}}{y!}$$

Recall that if  $Y \sim \text{Po}(\lambda)$ , then  $E[Y] = \lambda$  and  $V[Y] = \lambda$ . Let  $S = Y_1 + Y_2$  denote the sum of these two variables; then:

- What is the value of  $E[S]$ ?  
Ans:  $E[S] = E[Y_1 + Y_2] = E[Y_1] + E[Y_2] = 2 + 4 = 6$
- What is the value of  $V[S]$ ?  
Ans:  $V[S] = \frac{1}{2}V[Y_1 + Y_2] = \frac{1}{2}(V[Y_1] + V[Y_2]) = 2 + 4 = 6$  (by independence of  $Y_1, Y_2$ )
- What is the probability that  $S = 5$ ?  
Ans:  $S = Y_1 + Y_2$ , so  $S = 5$  if and only if  $Y_1 = 0$  and  $Y_2 = 5$  (as  $Y_1, Y_2$  are both non-negative integers). Therefore by independence:

$$P(S = 0) = P(Y_1 = 0)P(Y_2 = 5) = \frac{2^0 e^{-2}}{0!} \cdot \frac{4^5 e^{-4}}{5!} = e^{-2}e^{-4} = e^{-6}$$

- What is the value of  $E[Y_1^2]$ ?  
Ans:  $E[Y_1^2] = E[Y_1](E[Y_1]) = 2 \times 4 = 8$  (by independence of  $Y_1, Y_2$ )

- What is the value of  $E[Y_1^2]?$   
Ans: We can use the relationship:

$$\underline{E[Y_1^2] = E[Y_1^2] - E[Y_1]^2}$$

Then we have

$$\underline{E[Y_1^2] = E[Y_1^2] - E[Y_1]^2}$$

so recalling that  $E[Y_1] = 2$  and  $V[Y_1] = 2$  we have  $E[Y_1^2] = 4 + 2 = 6$ .

### Expectation & Variance rules

- Given a R.V.  $X_1, X_2$ ,
- $E[X_1 + X_2] = E[X_1] + E[X_2]$
- $E[X_1 X_2] = E[X_1] \cdot E[X_2]$   $\Rightarrow$  i.i.d
- $E[bX + c] = b \cdot E[X] + c$

$$V[cX] = c^2 V[X]$$

$$\left. \begin{array}{l} V[X_1 + X_2] = V[X_1] + V[X_2] \\ V[X_1 - X_2] = V[X_1] + V[X_2] \end{array} \right\} i.i.d.$$

$$V[X] = E[X^2] - E[X]^2$$

4

### 5 Confidence Intervals and p-values

#### Week 4 | 5

Consider a drug targeting obesity being considered for introduction to the market by the Therapeutic Goods Administration (TGA). The drug has been demonstrated to substantially reduce BMI, but the TGA are concerned about possible side-effects. They have measured cholesterol levels (in millilitres per 1 mmol/L) on a cohort of 7 individuals who have been administered our drug. The measured values are:

$$y = (5.2, 5.05, 5.35, 5.03, 5.43, 5.36).$$

The population standard deviation for cholesterol levels is unknown, so we assume the general distribution is approximately normal, and that the population standard deviation of cholesterol levels for individuals in our sample is the same as the population standard deviation of cholesterol levels for the general population.

- Using our sample, estimate the population mean cholesterol level of people being administered the drug. Calculate a 95% confidence interval for the population mean cholesterol level. Summarise your result.

Ans: The sample mean of sample is

$$\bar{y} = \frac{1}{7} (5 + 5.2 + 5.05 + 5.35 + 5.03 + 5.43 + 5.36) \approx 5.2$$

where we rounded 5.292 down to 5.2. The error is

$$w_p = \frac{0.6}{\sqrt{7}} \approx 0.227$$

The formula for the 95% confidence interval for a mean with known variance is

$$CI = (\bar{y} - 1.96 w_p, \bar{y} + 1.96 w_p)$$

so we have

$$CI = (5.2 - 1.96 \times 0.227, 5.2 + 1.96 \times 0.227) = (4.75, 5.64)$$

(Summary of CI) The estimated mean cholesterol level in our sample of size  $n = 7$  of individuals being prescribed our drug of interest is  $5.2 \pm 0.227$ . We are 95% confident that the population mean cholesterol level of people being administered our drug is between 4.75 and 5.64.

- The population mean cholesterol level in the general population is known to be  $1.767 \pm 0.227$ . The TGA wants to know two things: (i) is the population mean cholesterol level in people being given the drug different from the general population, and (ii) is it lower than in the general population. Specify appropriate null and alternative hypotheses and perform a hypothesis test to provide evidence against each null hypothesis. What is your conclusion regarding these two questions?

Ans: First, always state the hypothesis you are testing clearly. This helps to make sure you are doing the right thing.

- To answer the first part of the question we are testing the null hypothesis  $H_0: \mu = 4.8$  vs  $H_A: \mu \neq 4.8$ . From the Lecture notes regarding testing the population mean of a normal population with known variance, we must first calculate the sample mean for our sample to have above ( $\bar{y} = 5.2 \pm 0.227$ ). Then we calculate the z-score

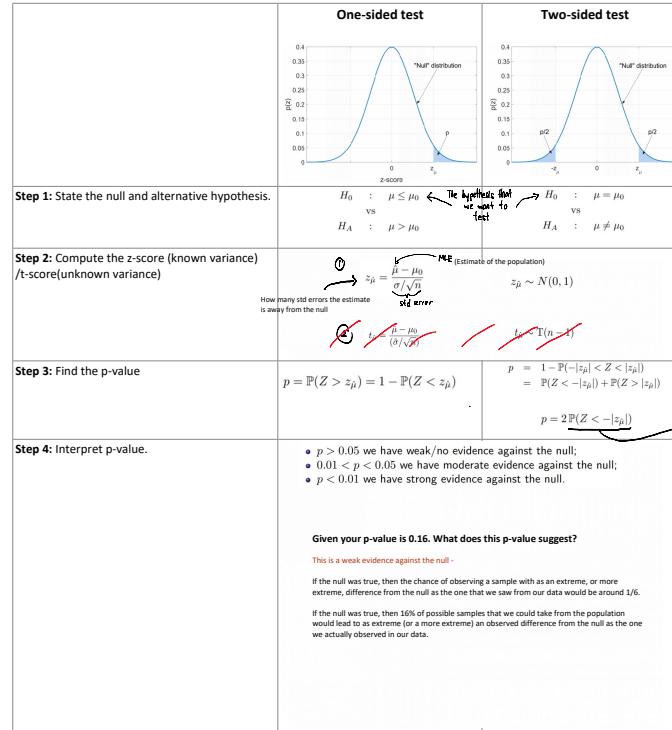
$$z_p = \frac{\bar{y} - \mu_0}{\sigma/\sqrt{n}} = \frac{5.2 - 4.8}{0.2/\sqrt{7}} = 1.763$$

where we can treat our population standard deviation as known and equal to  $\sigma = 0.6$  (as per the assumptions made in the question), and our sample size is  $n = 7$ . Then, using this z-score we do a one-sided test, to see our p-value is

$$p = 2F(Z < -|z_p|)$$

which is this, look at the Standard Normal Distribution in the Appendix. The first column is the absolute value of the z-score, the second column is the cumulative probability up to that z-score. Move down the rows to  $|z| = 1.763$  (which is the closest entry to our z-score), and we see that  $F(Z < -1.767) \approx 0.061$ . Our p-value is twice this, as we are doing a one-tailed test, so we have  $p \approx 0.072$ . This suggests there is some weak evidence against the null that the population mean cholesterol level of people using the drug is the same as the population mean cholesterol level in the general population.

Overall, looking at both tests, we see there is some evidence to suggest that we can reject the null that the drug does not affect the mean cholesterol level of individuals compared individuals in the general population, but it is not very conclusive. A larger study is probably required.



5

## Week 6

### 6 Regression

1. Please explain the intuition behind the principle of least squares that is used to fit a linear model with predictor  $\mathbf{x} = (x_1, \dots, x_n)$  to target variable  $y = (y_1, \dots, y_n)$ , and write down the least-squares objective function. Hint: we want to minimize the sum of squared errors (residuals) between the model predictions and the data values  $y$ . The idea is to find the straight line that most closely fits the data we have observed, which we hope will also be true if future data coming from the same source. The least-squares objective function is:

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

with  $\beta_1$  the coefficient relating the predictor to the target, and  $\beta_0$  the intercept.

2. If one of our predictors in a regression or logistic regression model is categorical, how can we handle it?

At first our predictor is a categorical variable with  $K$  categories, we can handle this by creating  $K - 1$  new dummy variables (predictors). The new variable number  $k$  will take on a "1" if an individual is in category  $k + 1$  and a 0 otherwise. These are called indicator variables as they indicate which category an individual is in.

3. Imagine we model a person's blood pressure in mmHg (BP) using a linear regression. Two predictors are fitted as part of the model: (i) the person's age in years (AGE), and (ii) the amount of alcohol they consume on average per week (ALCOTH). The model arrived at:

$$E[BP] = 51 + 1.4 \cdot AGE + 0.6 \cdot ALCOTH$$

(a) From this model, how does a person's blood pressure change as their age and alcohol consumption vary?

A: i. For each additional year a person lives, their expected blood pressure will increase by 1.4mmHg.

ii. For each additional standard drink a person consumes on average per week, their expected blood pressure will increase by 0.6mmHg.

(b) If a person is 33 years old, and drinks on average 2.5 standard drinks per week, what is their expected blood pressure?

A: The predicted expected blood pressure for such an individual is  $51 + 1.4 \times 33 + 0.6 \times 2.5 = 98.7\text{mmHg}$ .

4. The  $R^2$  statistic is defined as

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$$

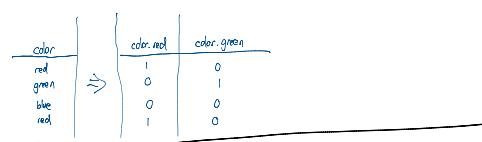
where RSS is the residual sum-of-squares and TSS is the total sum-of-squares.

(a) What does an  $R^2$  of 1 mean?

A: An  $R^2$  of 1 means that RSS = 0 and indicates that the model perfectly fits the data.

(b) What is an advantage of using the  $R^2$  statistic over simply using the residual sum-of-squares (RSS) to assess model fit?

A: The RSS depends on the unit of measurement of the target  $y$ . In contrast, the  $R^2$  statistic is a unit-free measure of goodness-of-fit and is easier to interpret.



**Write down the linear regression equation**

$$E[y] = \beta_0 + \beta_1 x_1 + \dots$$

↗ start  
 ↗ start  
 ↗ intercept  
 ↗ variable name

	No Breast Cancer ( $C = 0$ )	Breast Cancer ( $C = 1$ )	
Non-dense Breasts ( $D = 0$ )	0.15	0.10	
Dense Breasts ( $D = 1$ )	0.20	0.15	

$\boxed{0.35 \Rightarrow P(C=1)}$

Table 1: Population joint probabilities of having dense breasts ( $D$ ), and breast cancer by age 60 ( $C$ ).

### 7 Classification (week 7)

Breast cancer is one of the leading causes of death in Western populations. It is believed that non-dense breasts, which is defined as the sum of all fat tissue in a woman's breast, is strongly associated with the risk of developing breast cancer. We define a woman's breasts to be "dense" if they contain over 70% of non-fat tissue. Table 1 shows the joint probabilities of having dense breasts and contracting breast cancer by age 60.

1. What is the probability of contracting breast cancer by age 60 given that a woman does not have dense breasts?

A: Use the conditional probability formula  $P(C = 1 | D = 0) = P(C = 1, D = 0) / P(D = 0)$ :

$$P(C = 1 | D = 0) = \frac{0.15}{0.15 + 0.3} = 0.34$$

2. What is the probability of contracting breast cancer by age 60 given that a woman does have dense breasts?

A: Use the conditional probability formula  $P(C = 1 | D = 1) = P(C = 1, D = 1) / P(D = 1)$ :

$$P(C = 1 | D = 1) = \frac{0.15}{0.2 + 0.15} = 0.7333$$

3. What is the odds of contracting breast cancer given we have dense breasts?

A: The odds is the ratio of probabilities:

$$\frac{P(C = 1 | D = 1)}{P(C = 0 | D = 1)} = \frac{0.7333}{1 - 0.7333} \approx 2.75$$

4. Do you think that having dense breasts is a good predictor of contracting breast cancer by age 60, and why/why not?

A: Yes, it is a good predictor as you are almost twice the probability of contracting breast cancer by age 60 (0.7333/0.4 = 1.83) if you have dense breasts than if you don't.

$$\frac{P(C = 1 | D = 1)}{P(C = 0 | D = 0)} = \frac{0.7333}{0.34} \approx 1.83.$$

8

### 8 Logistic Regression (Week 7)

Imagine that we have built a logistic regression model to predict the probability of heart disease,  $H$ , given a person's age (AGE) in years and their cholesterol level (CHOL) in mg/dL.

$$\log(\text{Odds}(H)) = -7 + 0.0082 \cdot \text{CHOL} + 0.1 \cdot \text{AGE}$$

From this model, what is the effect of age of heart disease?

A: For every year a person lives, the log-odds of heart-disease increase by 0.1.

2. A person is 70 years old and has 260 mg/dL of cholesterol, what is the probability that they will have heart disease?

A: Plugging the numbers into the above formula to get the log-odds:

$$\log(\text{Odds}(H)) = -7 + 0.0082 \cdot 260 + 0.1 \cdot 70 \approx 2.132$$

Then we can use the logistic transformation to get the probability:

$$P(H = 1) = \frac{1}{1 + e^{-2.132}} = 0.894$$

3. Using this probability, if you were asked to predict if they have heart disease, what would you say?

A: The probability of having heart disease is 0.894 > 0.5, so we would say that we would predict this individual to have heart disease.

4. If my amount of cholesterol was non-linearly related to the log-odds of having heart disease, what could you do to try and improve the model?

A: There are two different possible answers:

- We could use some non-linear transformations of CHOL, for example, adding the square or cube of CHOL (or even more polynomial transformations), or the log of CHOL, etc into the logistic regression model as new predictors.

- We could potentially switch from using a logistic regression to a decision tree or random forest (i.e., a model that allows non-linearities naturally)

$$\text{conditional probability rule}$$

$$P(C = 1 | D = 0) = \frac{P(C = 1, D = 0)}{P(D = 0)}$$

$$\begin{aligned}
 \log \left( \frac{\theta}{1-\theta} \right) &\approx 1.132 \\
 \frac{\theta}{1-\theta} &= e^{1.132} \\
 \theta &= e^{1.132} - 1 e^{1.132} \\
 \theta + 1 e^{1.132} &= e^{1.132} \\
 \theta (1 + e^{1.132}) &= e^{1.132} \\
 \theta &= \frac{e^{1.132}}{1 + e^{1.132}} \approx 0.894
 \end{aligned}$$

9

### ( Week 3 & 8)

**9 Bias and Variance**

Imagine you have  $n$  random variables  $Y_1, \dots, Y_n$ , with mean  $E[Y_i] = \mu$  and variance  $V[Y_i] = \sigma^2$ . Consider the sample mean

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

Answer the following questions:

1. Prove that the sample mean  $\bar{Y}$  is an unbiased estimator of the population mean.

A: The bias of the sample mean is

$$\begin{aligned} \text{bias} &= E[\bar{Y}] - \mu \\ \text{so we have} \\ E[\bar{Y}] &= E\left[\frac{1}{n} \sum_{i=1}^n Y_i\right] \\ &= \frac{1}{n} E\left[\sum_{i=1}^n Y_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n E[Y_i] \\ &= \frac{1}{n} \sum_{i=1}^n \mu \\ &= \mu \end{aligned}$$

Plugging this into the above formula for bias yields a bias of zero proving the sample mean is unbiased.

2. Prove that the sample mean  $\bar{Y}$  has variance  $V[\bar{Y}] = \sigma^2/n$ .

A: The variance of the sample mean is

$$\begin{aligned} V[\bar{Y}] &= V\left[\frac{1}{n} \sum_{i=1}^n Y_i\right] \\ &= \frac{1}{n^2} V\left[\sum_{i=1}^n Y_i\right] \\ &= \frac{1}{n^2} V[\sum_{i=1}^n Y_i] \\ &= \frac{1}{n^2} \sum_{i=1}^n V[Y_i] \\ &= \frac{1}{n^2} n \sigma^2 \\ &= \frac{\sigma^2}{n} \end{aligned}$$

3. Prove that the sample mean is a consistent estimator of the population mean.

A: An estimator  $\hat{\theta}$  is a consistent estimator of  $\theta$  if

$$E[(\hat{\theta} - \theta)^2] \rightarrow 0 \text{ as } n \rightarrow \infty$$

Remember the mean-squared-error is given by the sum of the bias-squared and variance:

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[(\hat{\theta} - \mu)^2] + V[\hat{\theta}] \\ &= \frac{\sigma^2}{n} \end{aligned}$$

It is clear that as  $n \rightarrow \infty$ ,  $\sigma^2/n \rightarrow 0$ .

If an estimator  $\hat{\theta}$  satisfies  
 $b_n(\hat{\theta}) \rightarrow 0$ ,  
 $V_n(\hat{\theta}) \rightarrow 0$ ,  
as the sample size  $n \rightarrow \infty$ , then we can conclude that the estimator is consistent.

10

### 10 Machine Learning (Week 9)

```
> ctree(diabetes ~ ., data = diab, control = ctree::ctree_control(max_n = 10))
node), split, n, deviance, yval
denotes terminal node
1) root 442 262100 157.00 *
 2) SBM < 26.95 171 366600 96.31 *
 3) SBM < 4.69015 224 191500 119.00 *
 4) SBM < 26.95 171 366600 96.31 *
 5) SBM < 4.69015 224 191500 119.00 *
 6) SBM < 27.75 118 475100 162.70 *
 7) SBM < 27.75 118 475100 162.70 *
 14) BP < 99.5 77 305400 208.60 *
 28) BP < 99.5 33 173800 178.20 *
 29) BP < 99.5 43 173800 178.20 *
 15) BM < 32.75 31 66120 268.90 *
```

Figure 1: R output describing a decision tree learned using cross-validation for the diabetes progression dataset.

1. We have collected data on  $n = 442$  diabetic people. Figure 1 shows the R output after using the `ctree` package to learn a decision tree to predict their degree of diabetes progression (a non-negative integer) using three predictors in the dataset. The predictors used were as follows:  $\text{BM}$  is body-mass index ( $\text{kg}/\text{m}^2$ ),  $\text{BP}$  is blood pressure ( $\text{mmHg}$ ) and  $\text{SBM}$  is serum-bilirubin measurement (in milligrams).

- (a) How many "leaf" nodes does the tree have?

A: The leaf nodes are terminal nodes, and are stored – in this case, there are 6 leaf nodes.

- (b) If  $\text{BM} = 23$ ,  $\text{BP} = 20$ ,  $\text{SBM} = 5.5$  what is the degree of diabetes progression predicted by this tree?

A: To find the degree of diabetes progression we simply need to traverse the tree for the predictors we have. First, we note that we have  $\text{SBM} > 4.69015$ , so we move to Node #3. Then, our  $\text{BM} < 27.75$  so we move to Node #14. Finally, our  $\text{BP} < 99.5$ , so we move to Node #28. The degree of diabetes is the last number before the "star", so we see that  $E[\text{DIABETES}] = 162.7$ .

- (c) What combination of predictors leads to the greatest degree of diabetes progression?

A: To answer this question, we must find the leaf node that has the largest degree of diabetes progression, and then work back down the tree to figure out what combination of predictor values we need to arrive at the leaf node. In this case, the leaf node with the largest degree of diabetes is Node #15, which has a degree of diabetes progression in the tree. To arrive at this node, we need to go from the root to Node #3 ( $\text{SBM} > 4.69015$ ), then to Node #7 ( $\text{BM} < 27.75$ ), then to Node #15 ( $\text{BM} > 32.75$ ). So to summarize, we see that

- $\text{SBM} > 4.69015$
- $\text{BM} > 32.75$ .

2. Discuss one advantage that a decision tree has in comparison to a linear regression model.

A: A decision tree is a non-linear supervised learning algorithm. In this respect, one of the advantages it has over linear regression is that it can capture non-linear relationships between the predictors and the targets if they exist, while a linear regression model is restricted to learning only linear relationships.

3. The k-nearest neighbors method is a commonly used machine learning algorithm.

(a) Figure 2 shows a scatter plot of a training sample of women, both with and without breast cancer. The  $x$  and  $y$  axis are the predictors volume of dense tissue in the woman's breast and body mass index, respectively. The question mark shows a new individual we have obtained data on, but for whom we do not know breast cancer status.

not know breast status. Would a k-nearest neighbour algorithm, using standard Euclidean distance and  $k = 3$  nearest neighbors, predict them to have breast cancer or not? Please justify your answer.

A: We would predict the woman to have breast cancer, as 3 of her 4 nearest neighbors all have breast cancer.

(b) Looking at the configuration of the data points in Figure 3, do you think that a logistic regression model would be appropriate for separating women with and without breast cancer on the basis of the volume of dense tissue in a woman's breast and her body mass index? If so, why do you think so, and if not, why do you not think so?

A: From the configuration of the data points it does not appear that a logistic regression is appropriate. The logistic regression straight line will not do a particularly good job of separating the women with breast cancer from those without, while it does appear that even a reasonable simple non-linear curve would separate them quite well.

11

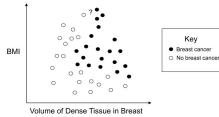


Figure 2: Scatter-plot of body mass index (BMI) against percent dense-tissue in a woman's breast for women with and without breast cancer.

not know breast status. Would a k-nearest neighbour algorithm, using standard Euclidean distance and  $k = 3$  nearest neighbors, predict them to have breast cancer or not? Please justify your answer.

A: We would predict the woman to have breast cancer, as 3 of her 4 nearest neighbors all have breast cancer.

(b) Looking at the configuration of the data points in Figure 3, do you think that a logistic regression model would be appropriate for separating women with and without breast cancer on the basis of the volume of dense tissue in a woman's breast and her body mass index? If so, why do you think so, and if not, why do you not think so?

A: From the configuration of the data points it does not appear that a logistic regression is appropriate. The logistic regression straight line will not do a particularly good job of separating the women with breast cancer from those without, while it does appear that even a reasonable simple non-linear curve would separate them quite well.

12

## 11 Appendix I: Standard Normal Distribution Table

$ z $	$\Pr(Z < - z )$	$\Pr(Z <  z )$	$ z $	$\Pr(Z < - z )$	$\Pr(Z <  z )$
0.000	0.500000	0.500000	2.047	0.020333	0.979667
0.003	0.629243	0.370757	2.140	0.016196	0.983804
0.016	0.626202	0.373796	2.233	0.012789	0.987211
0.029	0.620000	0.399904	2.326	0.010920	0.989980
0.032	0.354912	0.645688	2.419	0.007790	0.992210
0.045	0.520924	0.479767	2.512	0.006009	0.993991
0.058	0.288775	0.711625	2.605	0.004958	0.995462
0.061	0.237471	0.742529	2.698	0.004091	0.996559
0.074	0.228382	0.771618	2.791	0.003030	0.997370
0.087	0.201257	0.798763	2.884	0.001965	0.998055
0.090	0.170125	0.828757	2.977	0.001457	0.998443
1.023	0.153925	0.849007	3.070	0.001071	0.998629
1.116	0.132151	0.897449	3.163	0.000781	0.999219
1.209	0.113273	0.886727	3.256	0.000565	0.999435
1.302	0.096100	0.903597	3.349	0.000406	0.999594
1.395	0.081455	0.93845	3.442	0.000289	0.999711
1.488	0.068325	0.931674	3.535	0.000234	0.999796
1.581	0.056894	0.943106	3.628	0.000143	0.999857
1.674	0.047024	0.932976	3.721	0.000099	0.999901
1.767	0.038577	0.940123	3.814	0.000068	0.999932
1.860	0.031110	0.968590	3.907	0.000047	0.999933
1.953	0.025381	0.974619	> 4.000	< 0.000032	> 0.999968

Table 2: Cumulative Distribution Function for the Standard Normal Distribution  $Z \sim N(0, 1)$

13

## 12 Appendix II: Formulas

Formulas are collated on this page, some of which may be useful in answering this exam.

### Probability and Random Variables

Expectation of a RV:  $E[X] = \sum x \cdot \Pr(X = x)$

Marginal probability formula:  $\Pr(Y = y) = \sum_{x \in X} \Pr(Y = y, X = x)$

Conditional probability formula:  $\Pr(Y = y | X = x) = \frac{\Pr(Y = y, X = x)}{\Pr(X = x)}$

### Differentiation

$$\frac{d}{dx} \{e^{f(x)}\} = e^{\frac{d}{dx} f(x)}$$

$$\frac{d}{dx} \{x^k\} = kx^{k-1}$$

$$\frac{d}{dx} \{\log x\} = \frac{1}{x}$$

$$\text{Chain rule: } \frac{d}{dx} \{f(g(x))\} = \frac{d}{dg(x)} \{f(g(x))\} \cdot \frac{d}{dx} \{g(x)\}$$

### Confidence Interval and Hypothesis Test for Mean with Known Variance

Let  $\bar{\mu}$  be the sample mean of a sample of size  $n$  with population variance  $\sigma^2$ . Then a  $100(1 - \alpha)\%$  confidence interval for  $\mu$  is

$$\left( \bar{\mu} - z_{\alpha/2} \sqrt{\frac{\sigma^2}{n}}, \bar{\mu} + z_{\alpha/2} \sqrt{\frac{\sigma^2}{n}} \right)$$

where  $z_{\alpha/2}$  is the  $(100(1 - \alpha)/2)$ -percentile of the standard normal distribution  $N(0, 1)$ . To test the null hypothesis  $H_0: \mu = \mu_0$ , calculate

$$p = \begin{cases} \frac{2 \Pr(Z < -|z_{\alpha/2}|)}{1 - \Pr(Z < z_{\alpha/2})} & \text{if } H_0: \mu = \mu_0 \text{ vs } H_A: \mu \neq \mu_0 \\ 1 - \Pr(Z < z_{\alpha/2}) & \text{if } H_0: \mu \leq \mu_0 \text{ vs } H_A: \mu > \mu_0 \\ \Pr(Z < z_{\alpha/2}) & \text{if } H_0: \mu \geq \mu_0 \text{ vs } H_A: \mu < \mu_0 \end{cases}$$

where  $Z \sim N(0, 1)$ , and

$$z_{\alpha/2} = \frac{\bar{\mu} - \mu_0}{\sqrt{\sigma^2/n}}$$

### Confidence Interval and Hypothesis Test for Difference of Means with Known Variances

Let  $\bar{\mu}_1, \bar{\mu}_2$  be the sample means from two samples of size  $n_1, n_2$  and  $\sigma_1^2, \sigma_2^2$  be the known population variances of the two samples. The  $100(1 - \alpha)\%$  confidence interval for  $\mu_1 - \mu_2$  is

$$\left( \bar{\mu}_1 - \bar{\mu}_2 - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{\mu}_1 - \bar{\mu}_2 + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

where  $z_{\alpha/2}$  is the  $(100(1 - \alpha)/2)$ -percentile of the standard normal distribution  $N(0, 1)$ . To test the null hypothesis  $H_0: \mu_1 = \mu_2$ , calculate

$$p = \begin{cases} \frac{2 \Pr(Z < -|z_{\alpha/2}|)}{1 - \Pr(Z < z_{\alpha/2})} & \text{if } H_0: \mu_1 = \mu_2 \text{ vs } H_A: \mu_1 \neq \mu_2 \\ 1 - \Pr(Z < z_{\alpha/2}) & \text{if } H_0: \mu_1 \leq \mu_2 \text{ vs } H_A: \mu_1 > \mu_2 \\ \Pr(Z < z_{\alpha/2}) & \text{if } H_0: \mu_1 \geq \mu_2 \text{ vs } H_A: \mu_1 < \mu_2 \end{cases}$$

where  $Z \sim N(0, 1)$ , and

$$z_{\alpha/2} = \frac{\bar{\mu}_1 - \bar{\mu}_2}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

14