

# Checkpoint 2 Exercises

## Solution Sheet

Zak Varty

---

This checkpoint question sheet is for you to assess your own progress through Ethics Part 1 and to identify any areas that need further clarifications. It is recommended that you make a written attempt at all questions before solutions are made available to you.

This activity is expected to take approximately **90 minutes** of effort. The available marks for each question are indicated in square brackets, with a total of **26 marks** available.

---

### Confusion at the bank

A bank has developed two new methods of assessing eligibility for a business loan. The head of operations at the bank wishes to compare two new methods to their current approach. They are also concerned about the fairness of the methods, having heard anecdotal evidence that female applicants were being disadvantaged by their current approach.

You have been hired as a consultant to advise on the assessment of the accuracy and fairness of the bank's loan eligibility methods. To assist you, the bank has provided historical loan outcomes from the last year, along with the predictions made by the current and new eligibility methods. These are supplied in `loan_outcomes.csv`.

- a) Construct the confusion matrix for the benchmark method and the two new methods. [3]

#### Solution (a)

**Solution:** We can do this by hand or write a function to print the confusion matrices for any such experiment.

```

print_confusion_matrices <- function(df){

  # Construct confusion matrix for each test
  CM0 <- data.frame(
    test_pos = c(sum(df$curr_test & df$repaid) , sum(df$curr_test & !df$repaid)),
    test_neg = c(sum(!df$curr_test & df$repaid), sum(!df$curr_test & !df$repaid))
  )

  CM1 <- data.frame(
    test_pos = c(sum(df$new_test_1 & df$repaid) , sum(df$new_test_1 & !df$repaid)),
    test_neg = c(sum(!df$new_test_1 & df$repaid), sum(!df$new_test_1 & !df$repaid))
  )

  CM2 <- data.frame(
    test_pos = c(sum(df$new_test_2 & df$repaid) , sum(df$new_test_2 & !df$repaid)),
    test_neg = c(sum(!df$new_test_2 & df$repaid), sum(!df$new_test_2 & !df$repaid))
  )

  # assign meaningful rownames
  rownames(CM0) <- c("repaid", "defaulted")
  rownames(CM1) <- c("repaid", "defaulted")
  rownames(CM2) <- c("repaid", "defaulted")

  # Print tables to console
  cat("Current Test \n")
  print(CM0)
  cat("\n")

  cat("New Test 1 \n")
  print(CM1)
  cat("\n")

  cat("New Test 2 \n")
  print(CM2)
  cat("\n")

  # Quietly return tables
  invisible(list(CM0, CM1, CM2))
}

```

Applying this to our dataset, we find:

```
loan_outcomes <- read.csv("loan_outcomes.csv")
print_confusion_matrices(df = loan_outcomes)
```

Current Test

	test_pos	test_neg
repaid	607	179
defaulted	21	163

New Test 1

	test_pos	test_neg
repaid	713	73
defaulted	21	163

New Test 2

	test_pos	test_neg
repaid	705	81
defaulted	4	180

**Mark Scheme:** One mark per correct two-way table.

- b) Calculate the true positive rate and false positive rate for each method. Use these values to recommend which test the bank should use. You should explain the reasoning behind this recommendation. [5]

### Solution (b)

#### **Solution:**

Recall the definitions of True Positive Rate and False Positive rate.

$$TPR = \frac{TP}{TP + FN} = \frac{\#(\text{Accepted \& Repaid})}{\#(\text{Repaid})}$$

$$FPR = \frac{FP}{FP + TN} = \frac{\#(\text{Accepted \& Defaulted})}{\#(\text{Defaulted})}$$

Again, we can write a function to calculate these metrics for any dataset with the given structure.

```

print_TPRs_and_FPRs <- function(df){
  # TPR = TP / P = (predicted to repay and did / did repay)
  TPR0 <- sum(df$curr_test & df$repaid) / sum(df$repaid)
  TPR1 <- sum(df$new_test_1 & df$repaid) / sum(df$repaid)
  TPR2 <- sum(df$new_test_2 & df$repaid) / sum(df$repaid)

  # FPR = FP / N = (predicted to repay and didn't / didn't repay)
  FPR0 <- sum(df$curr_test & !df$repaid) / sum(!df$repaid)
  FPR1 <- sum(df$new_test_1 & !df$repaid) / sum(!df$repaid)
  FPR2 <- sum(df$new_test_2 & !df$repaid) / sum(!df$repaid)

  TPRs <- c(TPR0, TPR1, TPR2)
  FPRs <- c(FPR0, FPR1, FPR2)
  distances <- sqrt((TPRs - 1)^2 + (FPRs - 0)^2)

  out <- data.frame(test = 0:2, TPR = TPRs, FPR = FPRs, distance = distances)

  # Print TPR, FPR and interpretation to console
  print(out)
  cat(
    "\n Preferred test is",
    which.min(distances) - 1,
    "(assuming equal treatment of type 1 and 2 errors).\"
  )

  # Quietly return data frame
  invisible(out)
}

```

And apply this to our particular dataset.

```
print_TPRs_and_FPRs(df = loan_outcomes)
```

	test	TPR	FPR	distance
1	0	0.7722646	0.11413043	0.2547335
2	1	0.9071247	0.11413043	0.1471448
3	2	0.8969466	0.02173913	0.1053214

Preferred test is 2 (assuming equal treatment of type 1 and 2 errors).

**Mark Scheme:**

- Each pair of correct TPR and FPR (half mark for each). [1] [1] [1]
- Recommending the test that is closest to (0,1) on a ROC plot (the test with the shorted 'distance' value in the R output) [1]
- Justifying as having the greatest predictive power (or words to that effect) [1]

*Note: Selecting a different test is okay if justified sufficiently. e.g. bank is risk-adverse and so finds a false-positive defaulting on their loan more damaging than a false-negative where no loan is given.*

- c) Disaggregate the data by sex and calculate the TPR and FPR of each test by sex. In a few sentences, interpret these results for the head of operations. You should comment on whether the performance of the current test appears to differ by sex. You should also identify the test with the greatest predictive power for males and for non-males. [5]

#### Solution (c)

**Solution:** We can now make use of the functions we previously wrote to subset our data by sex and then recalculate the TPRs and FPRs.

```
print_TPRs_and_FPRs_by_sex <- function(df){

  df_m <- df[df$sex == 1,]
  df_f <- df[df$sex == 0,]

  cat('TEST PERFORMANCE ON MALES: \n')
  print_TPRs_and_FPRs(df_m)
  cat("\n ----- \n")

  cat('TEST PERFORMANCE ON NON-MALES: \n')
  print_TPRs_and_FPRs(df_f)
  cat("\n ----- \n")

}
```

```
print_TPRs_and_FPRs_by_sex(df = loan_outcomes)
```

TEST PERFORMANCE ON MALES:

	test	TPR	FPR	distance
1	0	0.8056426	0.090909091	0.21456759
2	1	0.9153605	0.083333333	0.11877832
3	2	0.9843260	0.007575758	0.01740878

Preferred test is 2 (assuming equal treatment of type 1 and 2 errors).

-----  
TEST PERFORMANCE ON NON-MALES:


	test		TPR	FPR	distance
1	0	0.7494647	0.17307692	0.3045055	
2	1	0.9014989	0.19230769	0.2160664	
3	2	0.8372591	0.05769231	0.1726644	

Preferred test is 2 (assuming equal treatment of type 1 and 2 errors).

-----  
**Mark Scheme:**

- Sex specific TPRs correct [2]
- Sex specific FPRs correct [2]
- Correct interpretation given in context of loan applications [1]

d) Write a short, contextualised description of error parity, equalised odds and equalised opportunity for the head of operations. [6]

 **Solution (d)**

**Solution:**

**Error parity:** Males and non-males should be mis-classified with the same probability/proportion. Mis-classification occurs here when someone would have repaid a loan they were not given or else failed to repay a loan that they were given.

**Equalised odds:** The probability of correctly accepting or rejecting the loan application should be the same, regardless of sex. This requires **both** of the following:

$$\begin{aligned}\Pr(\text{ Accepted } | \text{ repaid \& male } ) &= \Pr(\text{ Accepted } | \text{ repaid \& not male } ) \\ \Pr(\text{ Rejected } | \text{ defaulted \& male } ) &= \Pr(\text{ Rejected } | \text{ defaulted \& not male } ).\end{aligned}$$

**Equalised opportunity:** This is a weaker version of equalised odds. The probability of correctly accepting loan application should be the same for those who will repay the loan, regardless of sex. This requires only that:

$$\Pr(\text{ Accepted } | \text{ repaid \& male } ) = \Pr(\text{ Accepted } | \text{ repaid \& not male } )$$

**Mark Scheme:**

Two marks per metric. First mark for correct definition and second mark for each correctly put into context.

- e) Do tests 0, 1 and 2 appear to satisfy error parity by sex? Briefly describe how you might formally assess this. [3]

💡 Solution (e)

**Solution:**

We can again repurpose our existing function to calculate sex-specific error rates.

```
print_error_rates_by_sex <- function(df){
  df_m <- df[df$is_male == 1,]
  df_f <- df[df$is_male == 0,]

  test_0_male_errors <- mean(df_m$curr_test != df_m$repaid)
  test_0_female_errors <- mean(df_f$curr_test != df_f$repaid)

  test_1_male_errors <- mean(df_m$new_test_1 != df_m$repaid)
  test_1_female_errors <- mean(df_f$new_test_1 != df_f$repaid)

  test_2_male_errors <- mean(df_m$new_test_2 != df_m$repaid)
  test_2_female_errors <- mean(df_f$new_test_2 != df_f$repaid)

  test = c(0,1,2)
  male_error_rate <- c(test_0_male_errors, test_1_male_errors, test_2_male_errors)
  female_error_rate <- c(test_0_female_errors, test_1_female_errors, test_2_female_errors)
  relative_error_rate <- female_error_rate / male_error_rate

  print(data.frame(test, male_error_rate, female_error_rate, relative_error_rate))
}
```

```
print_error_rates_by_sex(df = loan_outcomes)
```

	test	male_error_rate	female_error_rate	relative_error_rate
1	0	0.16407982	0.2427746	1.479613
2	1	0.08425721	0.1078998	1.280600
3	2	0.01330377	0.1522158	11.441554

- Tests 0 and 1 appear to satisfy error parity since their relative error rate is approximately 1. [1]
- Test 2 has a much higher error rate for people who are not male. [1]
- Formally, we should use a hypothesis test to account for sampling variability in the proportion of misclassified individuals. [1]

- f) Describe the practical and statistical issues in assessing fairness through the use of fairness metrics on this data. [4]

💡 Solution (f)

**Solution:**

Practically, it is **not possible to satisfy all fairness metrics at once** and so a decision has to be made as to what type of fairness the bank is concerned with. Additionally, some types of fairness **can only be assessed retrospectively**, because they rely on knowing the true outcome of granting the loan.

Statistically, we want to ensure that the population error rates are equal but must base this on the sample proportions. When samples are small our **power to detect differences may be limited**. If **sampled individuals are not representative of the population** (e.g. only see outcomes for those who were granted loans) then the bias it introduces must be accounted for.

*Note: Four points given here as examples, only three valid points are required for full marks. These need not be the points listed here but should include at least one practical and one statistical issue.*

**Mark Scheme:**

- One valid statistical issue [1]
- One valid practical issue [1]
- One additional issue, practical or statistical [1]
- Clear writing and presentation of ideas [1]