STA 336: STATISTICAL MACHINE LEARNING

# Homework 4

Rento Saijo

February 24, 2026

## Table of Contents

# Disclosure

GPT-5.3-Codex was used to create the `YAML` portion and some `LaTeX` code to format the text/equations nicely. Page formatting code was also provided by Derin Gezgin. In the setup chunk, libraries were loaded and some helper functions were defined including but not limited to `table_latex()` and `with_family()`. See the original `RMD` file here for more details.

## Problem 1

> **Problem 1**
>
> Use the full dataset and build a model using a single predictor on student. What is
> the mathematical form of the model?

```
───────────────────────────────── Cell 1 ─────────────────────────────────
1   # Load data.
2   default_df <- ISLR2::Default
3
4   # Fit student-only logistic model.
5   model_q1 <- stats::glm(default ~ student, data = default_df, family =
    stats::binomial)
6
7   # Build coefficient table.
8   coef_q1_tbl <- tibble::tibble(
9     term     = names(stats::coef(model_q1)),
10    estimate = as.numeric(stats::coef(model_q1))
11  ) %>%
12    dplyr::mutate(estimate = round(estimate, 6))
13
14  # Print coefficient table.
15  coef_q1_tbl %>%
16    table_latex(
17      col_names = c('Term', 'Estimate'),
18      caption   = 'Logistic model coefficients with student only.'
19    )
```

Table 1: Logistic model coefficients with student only.

| Term | Estimate |
|------|----------|
| (Intercept) | -3.504128 |
| studentYes | 0.404887 |

The fitted model is

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = -3.504128 + 0.404887 \cdot \mathrm{I}\{\texttt{student} = \texttt{Yes}\},$$

where $\hat{p} = \Pr(\texttt{default} = \texttt{Yes} \mid \texttt{student})$. Equivalently,

$$\hat{p} = \frac{\exp(-3.504128 + 0.404887 \cdot \mathrm{I}\{\texttt{student} = \texttt{Yes}\})}{1 + \exp(-3.504128 + 0.404887 \cdot \mathrm{I}\{\texttt{student} = \texttt{Yes}\})}.$$

# Problem 2

PROBLEM 2

How would you interpret the coefficient of student on the model of Problem 1?

The coefficient on `studentYes` is 0.404887, so the log-odds of default are higher for students than non-students by 0.404887 in this one-predictor model. The corresponding odds ratio is

$$\exp(0.404887) = 1.4991.$$

Therefore, the model estimates that students have about $1.50\times$ the odds of default relative to non-students.

# Problem 3

> ### Problem 3
>
> What is the confusion matrix when you use the cutoff point to be 0.5? What is the
> error rate of prediction?

```
                                 ─── Cell 2 ───
20  # Build predictions with 0.5 cutoff.
21  prob_q1   <- stats::predict(model_q1, type = 'response')
22  pred_q1   <- factor(ifelse(prob_q1 > 0.5, 'Yes', 'No'), levels = c('No', 'Yes'))
23  actual_q1 <- factor(default_df$default, levels = c('No', 'Yes'))
24
25  # Build confusion matrix.
26  cm_q1     <- table(actual = actual_q1, predicted = pred_q1)
27  cm_q1_tbl <- as.data.frame.matrix(cm_q1) %>%
28    tibble::rownames_to_column(var = 'Actual')
29
30  # Compute error.
31  error_q1 <- mean(pred_q1 != actual_q1)
32  error_q1_tbl <- tibble::tibble(
33    metric = 'Prediction error rate',
34    value  = round(error_q1, 4)
35  )
36
37  # Print tables.
38  cm_q1_tbl %>%
39    table_latex(
40      col_names = c('Actual', 'Predicted No', 'Predicted Yes'),
41      caption   = 'Confusion matrix for student-only model (cutoff = 0.5).'
42    )
```

Table 2: Confusion matrix for student-only model (cutoff = 0.5).

| Actual | Predicted No | Predicted Yes |
|--------|--------------|---------------|
| No     | 9667         | 0             |
| Yes    | 333          | 0             |

```
                                 ─── Cell 3 ───
43  error_q1_tbl %>%
44    table_latex(
45      col_names = c('Metric', 'Value'),
46      caption   = 'Error rate for student-only model (cutoff = 0.5).'
47    )
```

Table 3: Error rate for student-only model (cutoff = 0.5).

| Metric | Value |
|---|---|
| Prediction error rate | 0.0333 |

At cutoff 0.5, this model predicts every observation as `No`. The prediction error rate is $\boxed{0.0333}$ (about 3.33%).

# Problem 4

> **Problem 4**
>
> Use the full dataset and build a model using a single predictor on balance. What is
> the mathematical form of the model?

```
                                  Cell 4
48  # Fit balance-only logistic model.
49  model_q4 <- stats::glm(default ~ balance, data = default_df, family =
    stats::binomial)

50
51  # Build coefficient table.
52  coef_q4_tbl <- tibble::tibble(
53    term     = names(stats::coef(model_q4)),
54    estimate = as.numeric(stats::coef(model_q4))
55  ) %>%
56    dplyr::mutate(estimate = round(estimate, 6))

57
58  # Print coefficient table.
59  coef_q4_tbl %>%
60    table_latex(
61      col_names = c('Term', 'Estimate'),
62      caption   = 'Logistic model coefficients with balance only.'
63    )
```

Table 4: Logistic model coefficients with balance only.

| Term | Estimate |
|---|---|
| (Intercept) | -10.651331 |
| balance | 0.005499 |

The fitted model is

$$\log\left(\frac{\hat{p}}{1 - \hat{p}}\right) = -10.651331 + 0.005499\,\mathtt{balance},$$

where $\hat{p} = \Pr(\mathtt{default} = \mathtt{Yes} \mid \mathtt{balance})$. Equivalently,

$$\hat{p} = \frac{\exp(-10.651331 + 0.005499\,\mathtt{balance})}{1 + \exp(-10.651331 + 0.005499\,\mathtt{balance})}.$$

# Problem 5

How would you interpret the coefficient of balance on the model of Problem 4?

The coefficient on `balance` is 0.005499, so each one-unit increase in balance increases the log-odds of default by 0.005499. In odds terms,

$$\exp(0.005499) = 1.0055,$$

so each \$1 increase in balance multiplies the odds of default by about 1.0055.

# Problem 6

> **Problem 6**
>
> What is the confusion matrix of the model in Problem 4 when you use the cutoff
> point to be 0.5? What is the error rate of prediction?

```
                                 Cell 5
64  # Build predictions with 0.5 cutoff.
65  prob_q4   <- stats::predict(model_q4, type = 'response')
66  pred_q4   <- factor(ifelse(prob_q4 > 0.5, 'Yes', 'No'), levels = c('No', 'Yes'))
67  actual_q4 <- factor(default_df$default, levels = c('No', 'Yes'))
68
69  # Build confusion matrix.
70  cm_q4 <- table(actual = actual_q4, predicted = pred_q4)
71  cm_q4_tbl <- as.data.frame.matrix(cm_q4) %>%
72    tibble::rownames_to_column(var = 'Actual')
73
74  # Compute error.
75  error_q4 <- mean(pred_q4 != actual_q4)
76  error_q4_tbl <- tibble::tibble(
77    metric = 'Prediction error rate',
78    value  = round(error_q4, 4)
79  )
80
81  # Print tables.
82  cm_q4_tbl %>%
83    table_latex(
84      col_names = c('Actual', 'Predicted No', 'Predicted Yes'),
85      caption   = 'Confusion matrix for balance-only model (cutoff = 0.5).'
86    )
```

Table 5: Confusion matrix for balance-only model (cutoff = 0.5).

| Actual | Predicted No | Predicted Yes |
|--------|--------------|---------------|
| No     | 9625         | 42            |
| Yes    | 233          | 100           |

```
                                 Cell 6
87  error_q4_tbl %>%
88    table_latex(
89      col_names = c('Metric', 'Value'),
90      caption   = 'Error rate for balance-only model (cutoff = 0.5).'
91    )
```

Table 6: Error rate for balance-only model (cutoff = 0.5).

| Metric | Value |
|---|---|
| Prediction error rate | 0.0275 |

The balance-only model has prediction error $\boxed{0.0275}$ (about 2.75%).

# Problem 7

PROBLEM 7

Split the data into halves and use one half as the training and the other half as the test. Fit a model on the training data using `student`, `balance`, and `income` as predictors.

```
Cell 7
92  # Create train/test split.
93  set.seed(20060527)
94  n_default      <- nrow(default_df)
95  train_idx      <- sample(seq_len(n_default), n_default / 2)
96  default_train <- default_df[train_idx, ]
97  default_test  <- default_df[-train_idx, ]
98
99  # Fit training model.
100 model_q7 <- stats::glm(
101   default ~ student + balance + income,
102   data  = default_train,
103   family = stats::binomial
104 )
105
106 # Build split table.
107 split_q7_tbl <- tibble::tibble(
108   data_set     = c('Training', 'Test'),
109   observations = c(nrow(default_train), nrow(default_test))
110 )
111
112 # Build coefficient table.
113 coef_q7_tbl <- as.data.frame(stats::coef(summary(model_q7))) %>%
114   tibble::rownames_to_column(var = 'term') %>%
115   dplyr::rename(
116     estimate  = Estimate,
117     std_error = `Std. Error`,
118     z_value   = `z value`,
119     p_value   = `Pr(>|z|)`
120   ) %>%
121   dplyr::mutate(dplyr::across(-term, ~ round(.x, 6)))
122
123 # Create probabilities for later questions.
124 prob_train_q7 <- stats::predict(model_q7, newdata = default_train, type =
    'response')
125 prob_test_q7  <- stats::predict(model_q7, newdata = default_test, type =
    'response')
126
127 # Print tables.
128 split_q7_tbl %>%
129   table_latex(
130     col_names = c('Data set', 'Observations'),
131     caption   = 'Train/test split sizes.'
```

```
132    )
```

Table 7: Train/test split sizes.

| Data set | Observations |
|----------|--------------|
| Training | 5000 |
| Test | 5000 |

──────── Cell 8 ────────
```
133  coef_q7_tbl %>%
134    table_latex(
135      col_names = c('Term', 'Estimate', 'Std. Error', 'z value', 'Pr(>|z|)'),
136      caption   = 'Logistic regression coefficient summary.'
137    )
```

Table 8: Logistic regression coefficient summary.

| Term | Estimate | Std. Error | z value | $\Pr(>|z|)$ |
|------|----------|-----------|---------|-------------|
| (Intercept) | -11.362494 | 0.729852 | -15.568209 | 0.000000 |
| studentYes | -0.439823 | 0.345068 | -1.274598 | 0.202452 |
| balance | 0.005885 | 0.000340 | 17.319003 | 0.000000 |
| income | 0.000008 | 0.000012 | 0.693137 | 0.488224 |

# Problem 8

> **PROBLEM 8**
>
> What is the training error of the model in Problem 7?

```
                              ─── Cell 9 ───
138  # Build training predictions.
139  pred_train_q7   <- factor(ifelse(prob_train_q7 > 0.5, 'Yes', 'No'), levels =
     c('No', 'Yes'))
140  actual_train_q7 <- factor(default_train$default, levels = c('No', 'Yes'))
141
142  # Compute training error.
143  train_error_q7  <- mean(pred_train_q7 != actual_train_q7)
144  train_error_tbl <- tibble::tibble(
145    metric = 'Training error rate',
146    value  = round(train_error_q7, 4)
147  )
148
149  # Print training error table.
150  train_error_tbl %>%
151    table_latex(
152      col_names = c('Metric', 'Value'),
153      caption   = 'Training error with cutoff = 0.5.'
154    )
```

Table 9: Training error with cutoff = 0.5.

| Metric | Value |
|---|---|
| Training error rate | 0.0264 |

The training error is $\boxed{0.0264}$ (about 2.64%).

# Problem 9

> **Problem 9**
>
> What is the test error of the model in Problem 7?

```
155  # Build test predictions.
156  pred_test_q7    <- factor(ifelse(prob_test_q7 > 0.5, 'Yes', 'No'), levels =
     c('No', 'Yes'))
157  actual_test_q7 <- factor(default_test$default, levels = c('No', 'Yes'))
158
159  # Compute test error.
160  test_error_q7   <- mean(pred_test_q7 != actual_test_q7)
161  test_error_tbl <- tibble::tibble(
162    metric = 'Test error rate',
163    value  = round(test_error_q7, 4)
164  )
165
166  # Print test error table.
167  test_error_tbl %>%
168    table_latex(
169      col_names = c('Metric', 'Value'),
170      caption   = 'Test error with cutoff = 0.5.'
171    )
```

Table 10: Test error with cutoff = 0.5.

| Metric | Value |
|---|---|
| Test error rate | 0.0266 |

The test error is 0.0266 (about 2.66%).

# Problem 10

> ### PROBLEM 10
>
> Draw the ROC curve for model in Problem 7 and find the best cutoff point. What does the best cutoff point mean?

---- Cell 11 ----

```
172  # Build ROC.
173  roc_q7 <- Epi::ROC(
174    form = default ~ student + balance + income,
175    data = default_train,
176    plot = 'ROC',
177    MX   = TRUE
178  )
```
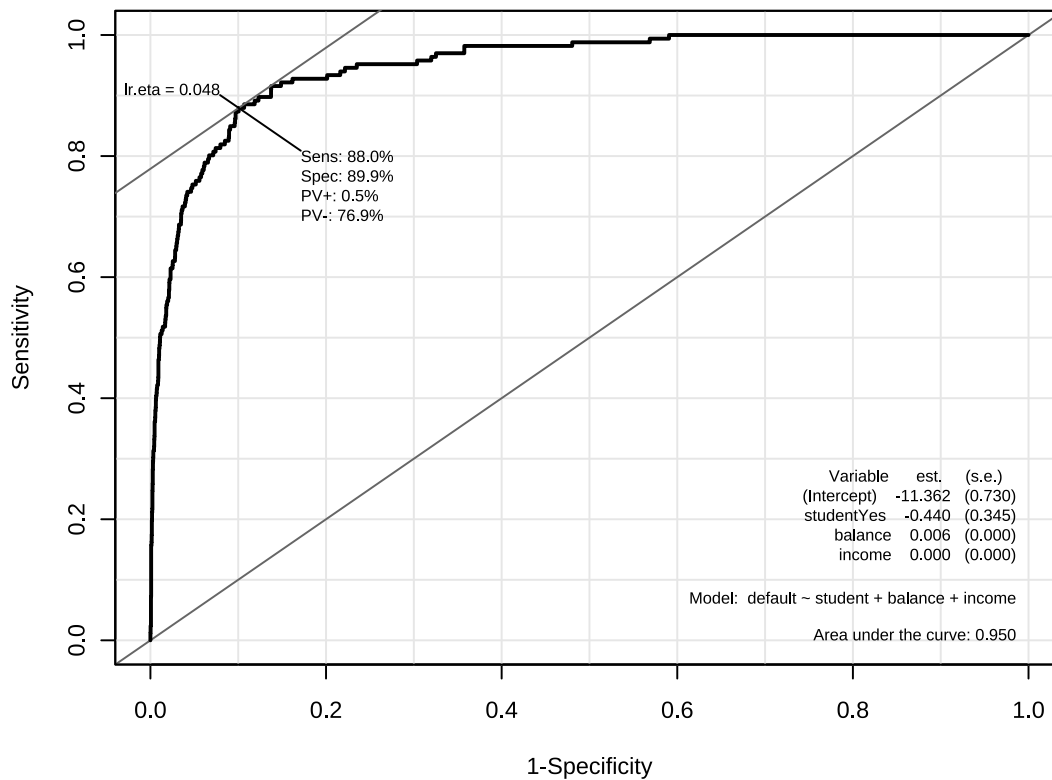


Figure 1: ROC curve on training data for the multivariable model.

---- Cell 12 ----

```
179  # Extract ROC grid.
180  roc_q7_res <- roc_q7$res
181
182  # Find best cutoff with MX criterion.
```

```r
183  best_row_q7 <- roc_q7_res %>%
184     dplyr::mutate(mx_value = sens + spec) %>%
185     dplyr::filter(mx_value == max(mx_value)) %>%
186     dplyr::arrange(dplyr::desc(lr.eta)) %>%
187     dplyr::slice(1)
188  best_cutoff_q7 <- best_row_q7$lr.eta[[1]]
189
190  # Build best-cutoff table.
191  best_cutoff_tbl <- tibble::tibble(
192    metric = c(
193       'Best cutoff',
194       'Sensitivity',
195       'Specificity',
196       'False positive rate',
197       'MX value (sensitivity + specificity)',
198       'AUC'
199    ),
200    value  = round(c(
201       best_cutoff_q7,
202       best_row_q7$sens[[1]],
203       best_row_q7$spec[[1]],
204       1 - best_row_q7$spec[[1]],
205       best_row_q7$mx_value[[1]],
206       roc_q7$AUC
207    ), 4)
208  )
209
210  # Print best-cutoff table.
211  best_cutoff_tbl %>%
212    table_latex(
213       col_names = c('Metric', 'Value'),
214       caption   = 'Best cutoff chosen from training ROC.'
215    )
```

Table 11: Best cutoff chosen from training ROC.

| Metric | Value |
|---|---|
| Best cutoff | 0.0477 |
| Sensitivity | 0.8795 |
| Specificity | 0.8995 |
| False positive rate | 0.1005 |
| MX value (sensitivity + specificity) | 1.7790 |
| AUC | 0.9498 |

Using (`Epi::ROC` with `MX = TRUE`), the best cutoff is $\boxed{0.0477}$. This is the threshold that maximizes sensitivity + specificity on the training data.

# Problem 11

> **PROBLEM 11**
>
> Use the best cutoff point on the test data. What is the test error now?

```
─── Cell 13 ───
216  # Build test predictions with best cutoff.
217  pred_test_best_q7 <- factor(ifelse(prob_test_q7 > best_cutoff_q7, 'Yes', 'No'),
     levels = c('No', 'Yes'))
218
219  # Compute test error with best cutoff.
220  test_error_best_q7 <- mean(pred_test_best_q7 != actual_test_q7)
221  test_error_best_tbl <- tibble::tibble(
222    metric = 'Test error rate (best cutoff)',
223    value = round(test_error_best_q7, 4)
224  )
225
226  # Print test error table.
227  test_error_best_tbl %>%
228    table_latex(
229      col_names = c('Metric', 'Value'),
230      caption = 'Test error using best cutoff from training ROC.'
231    )
```

Table 12: Test error using best cutoff from training ROC.

| Metric | Value |
|---|---|
| Test error rate (best cutoff) | 0.1016 |

Using the best cutoff from Problem 10, the test error is $\boxed{0.1016}$ (about 10.16%).

## Problem 12

> **PROBLEM 12**
>
> For a student, with 0 balance, and 20000 income, what is the predicted probability of default? Use the model on training data.

```
Cell 14
232  # Build student profile.
233  new_student_q12 <- data.frame(
234    student = factor('Yes', levels = levels(default_train$student)),
235    balance = 0,
236    income  = 20000
237  )
238
239  # Compute predicted probability.
240  p_q12 <- as.numeric(stats::predict(model_q7, newdata = new_student_q12, type =
     'response'))
241  q12_tbl <- tibble::tibble(
242    metric = 'Predicted probability',
243    value  = signif(p_q12, 6)
244  )
245
246  # Print probability table.
247  q12_tbl %>%
248    table_latex(
249      col_names = c('Metric', 'Value'),
250      caption   = 'Predicted probability for student profile.'
251    )
```

Table 13: Predicted probability for student profile.

| Metric | Value |
|---|---|
| Predicted probability | 0.0000088 |

The predicted probability is $\boxed{0.00000882}$.

# Problem 13

> ### Problem 13
>
> For a non-student, with 0 balance, and 20000 income, what is the predicted probability of default? Use the model on training data.

```
Cell 15
252  # Build non-student profile.
253  new_nonstudent_q13 <- data.frame(
254    student = factor('No', levels = levels(default_train$student)),
255    balance = 0,
256    income  = 20000
257  )
258
259  # Compute predicted probability.
260  p_q13 <- as.numeric(stats::predict(model_q7, newdata = new_nonstudent_q13, type =
     'response'))
261  q13_tbl <- tibble::tibble(
262    metric = 'Predicted probability',
263    value  = signif(p_q13, 6)
264  )
265
266  # Print probability table.
267  q13_tbl %>%
268    table_latex(
269      col_names = c('Metric', 'Value'),
270      caption   = 'Predicted probability for non-student profile.'
271    )
```

Table 14: Predicted probability for non-student profile.

| Metric | Value |
|---|---|
| Predicted probability | 0.0000137 |

The predicted probability is $\boxed{0.00001369}$.

# Problem 14

---
**Problem 14**

What is the difference between the predicted values of Problems 12 and 13? Does this difference come from the coefficient of the student variable? How are they related?
---

Cell 16

```
272  # Compute difference and student effect.
273  diff_q14 <- p_q12 - p_q13
274  beta_student_q14 <- as.numeric(stats::coef(model_q7)['studentYes'])
275  odds_ratio_student_q14 <- exp(beta_student_q14)
276
277  # Build relation table.
278  q14_tbl <- tibble::tibble(
279    metric = c(
280      'p(student = Yes) - p(student = No)',
281      'student coefficient (log-odds scale)',
282      'student odds ratio'
283    ),
284    value = c(diff_q14, beta_student_q14, odds_ratio_student_q14)
285  ) %>%
286    dplyr::mutate(value = signif(value, 6))
287
288  # Print relation table.
289  q14_tbl %>%
290    table_latex(
291      col_names = c('Metric', 'Value'),
292      caption   = 'Probability difference and student coefficient relation.'
293    )
```

Table 15: Probability difference and student coefficient relation.

| Metric | Value |
|---|---|
| p(student = Yes) - p(student = No) | -0.0000049 |
| student coefficient (log-odds scale) | -0.4398230 |
| student odds ratio | 0.6441500 |

The difference is $\boxed{p_{12} - p_{13} = -4.87 \times 10^{-6}}$, so the student profile has a slightly lower predicted default probability here. This difference is directly related to the coefficient. Let

$$r = \frac{p_{12}/(1-p_{12})}{p_{13}/(1-p_{13})} = e^{\hat{\beta}_{\text{studentYes}}} = 0.6442.$$

Solving this for $p_{12}$ in terms of $p_{13}$,

$$p_{12} = \frac{r\,p_{13}}{1 - p_{13} + r\,p_{13}},$$

so

$$p_{12} - p_{13} = \frac{p_{13}(1 - p_{13})(r - 1)}{1 - p_{13} + r\, p_{13}}.$$

This shows the probability difference comes from the student coefficient via $r = e^{\hat{\beta}_{\texttt{studentYes}}}$. Plugging in $r = 0.6442$ and $p_{13} = 1.368903 \times 10^{-5}$ gives

$$p_{12} - p_{13} \approx -4.87 \times 10^{-6},$$

which matches the observed difference.