

# STAT 471: Midterm Exam Extra Credit

Name

Due: November 6, 2021 at 11:59pm; Time limit: 24 hours

## Contents

Instructions	1
Medicare costs at inpatient rehabilitation facilities	2
<b>1 Wrangling (40 points for correctness; 0 points for presentation)</b>	<b>2</b>
1.1 Import (8 points)	2
1.2 Clean (24 points)	2
1.3 Tidy (8 points)	4
<b>2 Exploration (25 points for correctness; 5 points for presentation)</b>	<b>4</b>
2.1 Response distribution (15 points)	4
2.2 Relationships among features (10 points)	4
<b>3 Elastic net regression (25 points for correctness; 5 points for presentation)</b>	<b>4</b>
3.1 Training and tuning (15 points)	4
3.2 Performance evaluation (10 points)	5

## Instructions

Download the Rmd file `midterm-exam-extra-credit.Rmd` from [this link](#), and place it under `stat-471-fall-2021/midterm/midterm-fall-2021/`. Set the latter directory as your working directory. Note that the link above also contains the file `irf_data_tidy.tsv`, which is the correct output of the first section of the exam. You may use this to complete the remaining sections of the exam if you are unable to complete the first section of the exam.

Use this document as a starting point for your writeup, adding your solutions after “**Solution**”. Add your R code using code chunks and add your text answers using **bold text**. Compile your writeup to PDF and submit to [Gradescope](#).

**You must complete this exam individually, but you may consult any course materials or the internet.**

We’ll need to use the following R packages and functions:

```
library(kableExtra)      # for printing tables
library(cowplot)         # for side by side plots
library(glmnetUtils)     # to run ridge and lasso
source("../functions/plot_glmnet.R") # for lasso/ridge trace plots
library(tidyverse)       # for everything else
```

# Medicare costs at inpatient rehabilitation facilities

According to [medicare.gov](https://www.medicare.gov), an inpatient rehabilitation facility (IRF) is “a hospital, or part of a hospital, that provides intensive rehabilitation to inpatients. Many patients with conditions like stroke or brain injury are transferred or admitted to an inpatient rehabilitation facility.” IRFs can incur significant medical costs, as quantified by the [Medicare spending per beneficiary \(MSPB\) measure](#). In this exam, we will be building predictive models for the MSPB based on several attributes of IRFs. We will use a dataset from [The Centers for Medicare & Medicaid Services](#), which are available for download at the URL defined below. The variables of interest are coded in these data using “measure codes”. For the purposes of this exam, we are interested in 15 of the measure codes, presented in Table 1 along with shortened variable names and descriptions. The first variable in Table 1 (MSPB) is the response and the remaining variables are features.

## 1 Wrangling (40 points for correctness; 0 points for presentation)

### 1.1 Import (8 points)

- (4 points) Import the IRF data directly from the URL below into a tibble called `irf_data_raw`. Print this tibble (no need to make a fancy table out of it).

```
# data URL
website = "https://data.cms.gov/provider-data/sites/default/files/resources"
resource_id = "6fe279310f5926cab64ad5b27b1da26a_1632923525"
filename = "Inpatient_Rehabilitation_Facility-Provider_Dec2020.csv"
url = paste(website, resource_id, filename, sep = "/")
```

- (4 points) How many IRFs are represented in these data? Note that each IRF is coded using the CMS Certification Number (CCN). How many measure codes are represented in these data?

### 1.2 Clean (24 points)

- (4 points) Create a new tibble called `irf_data_1` from `irf_data_raw` that contains only the columns CMS Certification Number (CCN), Facility Name, City, State, Measure Code, Score, renamed to `facility_id`, `facility_name`, `city`, `state`, `measure_code`, `score`, respectively. Print this tibble (no need to create a fancy table).
- (4 points) The column `measure_code` contains the measure codes. Create a new tibble called `irf_data_2` from `irf_data_1` that contains only the rows representing measure codes in Table 1. Print this tibble (no need to create a fancy table). [Hint: Use the tibble `variables` underlying Table 1.]
- (4 points) Table 8 of the [data dictionary](#) contains all of the measure codes in the original data (listed as “Provider Variables”). Why might we want to use only the subset of measure codes in Table 1 for our predictive models?
- (4 points) The column `score` contains the value of each measure, with missing values indicated using the string `Not Available`. What fraction of IRFs have any missing scores? [Hint: Use the `any` function within `summarise`; see `?any`.]
- (4 points) Create a new tibble called `irf_data_3` from `irf_data_2` that contains only those IRFs for which none of the scores are missing. Print this tibble (no need to create a fancy table). How many IRFs remain? [Hint: First create a list of IRFs for which none of the scores are missing using code from the previous bullet point.]
- (4 points) Convert the `score` column from `character` to `numeric`. Replace the `measure_code` column with a column named `feature`, which is a factor with `levels` equal to the first column of Table 1 and `labels` equal to the second column of Table 1. [Hint: You can create a factor variable out of a character variable with given `levels` and `labels` using the `factor()` function; check out `?factor`.] Store the resulting tibble in `irf_data_4`. Print this tibble (no need to create a fancy table).

Table 1: Features and response variable of interest in the inpatient rehabilitation facility data.

Measure code	Variable name	Variable description
I_020_01_MSPB_SCORE	mspb	Medicare Spending Per Beneficiary (MSPB)
I_006_01_SIR	uti_rate	Catheter-associated urinary tract infection rate
I_008_02_OBS_RATE	abilities_goals	Percentage of patients whose functional abilities were assessed and functional goals were included in their treatment plan
I_009_03_ADJ_CHG_SFCR_SCORE	self_care_ability_improved	Patients' ability to care for themselves changed between facility admission and discharge
I_010_03_ADJ_CHG_MOBL_SCORE	mobility_improved	Patients' ability to move around changed between facility admission and discharge
I_011_03_OBS_RATE	self_care_expectation_met	Percentage of patients who achieve or exceed a self-care ability expected for their condition at discharge
I_012_03_OBS_RATE	mobility_expectation_met	Percentage of patients who achieve or exceed the level of movement expected for their condition at discharge
I_013_01_OBS_RATE	fall_rate	Percentage of IRF patients who experience one or more falls with major injury during their IRF stay
I_015_01_SIR	cdi_rate	Clostridium difficile infection rate
I_016_01_OBS_RATE	flu_vax_rate	Influenza Vaccination Coverage Among Healthcare Personnel
I_017_01_PPR_PD_RSRR	readmission_rate_after	Rate of potentially preventable hospital readmissions 30 days after discharge from an IRF
I_018_01_PPR_WI_RSRR	readmission_rate_during	Rate of potentially preventable hospital readmissions during the IRF stay
I_019_02_DTC_RS_RATE	rate_return	Rate of successful return to home and community from an IRF
I_021_01_OBS_RATE	medications_reviewed	Percentage of patients whose medications were reviewed and who received follow-up care when medication issues were identified
I_022_01_ADJ_RATE	pressure_ulcers	Percentage of patients with pressure ulcers/injuries that are new or worsened

### 1.3 Tidy (8 points)

- (4 points) Tidy the tibble `irf_data_4`, storing the result in a new tibble called `irf_data_tidy`. Print this tibble (no need to create a fancy table).
- (4 points) What does each row in `irf_data_tidy` represent? How many rows are there?

## 2 Exploration (25 points for correctness; 5 points for presentation)

The tidied dataset `irf_data_tidy` is provided for you [here](#). You may use this to complete the remainder of the exam if you are unable to complete the tidying yourself.

### 2.1 Response distribution (15 points)

- (5 points) Print the mean, standard deviation, minimum, and maximum of `mspb` in a nice table.
- (5 points) Produce a histogram of `mspb`, with a vertical dashed line at the mean. Comment on the shape of this distribution.
- (5 points) Produce a (nice) table of the IRFs with the top 5 `mspb` values; for each IRF print its `facility_name`, `city`, `state`, and `mspb`. What cities (or city) do the two most costly IRFs come from?

### 2.2 Relationships among features (10 points)

- (5 points) Some of the features appear to be related to each other, e.g. `self_care_ability_improved` and `self_care_expectation_met`; `mobility_improved` and `mobility_expectation_met`; `readmission_rate_during` and `readmission_rate_after`. Create scatter plots of these three pairs of features, displaying all three side by side in a single figure, and comment on the degree to which each pair appears to be correlated.
- (5 points) Discuss how the findings in the above bullet point might impact the fitted coefficients of ridge and lasso regressions.

## 3 Elastic net regression (25 points for correctness; 5 points for presentation)

Next, let's train penalized regression models to predict the case-fatality ratio based on the available features. We use the following train-test split:

```
num_train_samples = 300
set.seed(1) # seed set for reproducibility (DO NOT CHANGE)

# uncomment the three commented lines below to create train-test split
# train_samples = sample(1:nrow(irf_data_tidy), num_train_samples)
# irf_train = irf_data_tidy %>% filter(row_number() %in% train_samples)
# irf_test = irf_data_tidy %>% filter(!(row_number() %in% train_samples))
```

### 3.1 Training and tuning (15 points)

- (3 points) Using `irf_train`, fit a 10-fold cross-validated elastic net regression of `mspb` on the 14 other features in Table 1, where the cross-validation is over `alpha` as well as `lambda`.

```
set.seed(127) # set seed for reproducibility (DO NOT CHANGE)
# TODO: Fit cross-validated elastic net regression
```

- (3 points) Produce a plot of the CV error (minimized across `lambda` for each `alpha`) versus `alpha`. Which value of `alpha` performs best? What does this suggest about the number of features actually impacting the response?
- (3 points) Produce a regular CV plot (showing CV error as a function of `lambda`) for the `alpha` selected in the previous bullet point. Based on this plot, how many features have nonzero coefficients in the model selected by the one-standard error rule?
- (3 points) Produce the trace plot for the `glmnet` model from the previous bullet point. [Hint: To make the plot easier to read, change the y axis scale by appending `+ scale_y_continuous(limits = c(-0.02, 0.02))` after your `plot_glmnet` call.] What is the first feature that comes in with a positive coefficient? What is the first feature that comes in with a negative coefficient? Do the signs of these coefficients make sense to you? Why or why not?
- (3 points) Produce a nice table with the features selected by the above `glmnet` model and their standardized coefficients, ordered by decreasing magnitude. According to these coefficients, if an IRF increases its `rate_return` by 10 while keeping other things equal, how would this change the expected `mspb`?

### 3.2 Performance evaluation (10 points)

Let's compare the performance of elastic net models for different `alpha`. While the fit objects for each `alpha` can be extracted from the original fit object from Section 3.1, we haven't learned how to do this. We will instead separately re-train models for different `alpha`.

- (5 points) Train three cross-validated elastic net models—one for `alpha = 0`, one for `alpha = 0.5`, and one for `alpha = 1`—this time cross-validating over `lambda` only.

```
set.seed(127) # set seed for reproducibility for alpha = 0 (DO NOT CHANGE)
# TODO: Fit elastic net for alpha = 0

set.seed(127) # set seed for reproducibility for alpha = 0.5 (DO NOT CHANGE)
# TODO: Fit elastic net for alpha = 0.5

set.seed(127) # set seed for reproducibility for alpha = 1 (DO NOT CHANGE)
# TODO: Fit elastic net for alpha = 1
```

- (5 points) Compute the RMSE prediction error for each of these three models, with `lambda` chosen based on the one-standard-error rule. Print these in a nice table, together with the corresponding `alpha` values (use `digits = 4` in your call to `kable`). Comment on the ordering among the prediction errors, and the extent to which it is consistent with the CV error versus `alpha` plot from Section 3.1.