# Data Processing and Modeling Overview

# Removed and Modified Features

- Initial shape – 101,767 observations, 50 features

- Weight – 97% missing, removed from analysis
- Payer Code – 40% missing, removed

- Admission type/source and disposition (release) type:
  - Numerically-coded columns; curated and condensed into categories
  - Removed patients who died in the hospital

- Patient diagnosis codes:
  - Series of numeric codes XXX.xx. Converted these to categorical diagnosis types (respiratory, neurological, oncology, injury, etc)

- 24 medications – removed 7 with essentially no values.
- Converted ['No', 'Steady', 'Up', 'Down'] to ['No', 'Yes']
  - Is the patient taking this medication or not?
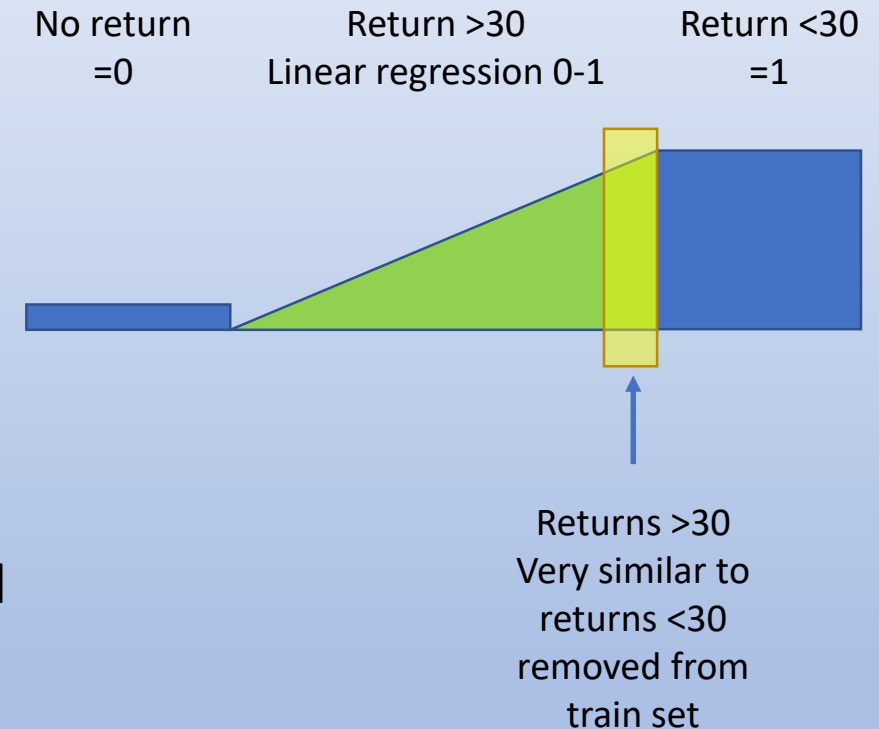
# Feature Creation

- Diabetic features:
  - Diabetic patient codes were 250.xx
  - Xx reveals additional info about the patient's diabetes
    - 250.1 – Diabetes with ketoacidosis
  - Categorized and dummified this information when available

- Diabetic medicine change:
  - Since we converted medicine codes to a simple [Yes, No], created an extra feature, "diabchange", indicating if the patient changed their dosage of *any* diabetic medication

- Primary diagnosis:
  - Specified the primary diagnosis for which patient was admitted to the hospital

# Removal of Train Set Observations

- We have a middle group (returned >30 days to the hospital)
  - These are "No" values for our model, but can offer some information to the model
  - Base AUC with everything: ~0.6662

No return =0    Return >30 Linear regression 0-1    Return <30 =1

- Removed them from the train set (TrainHL)
  - AUC: 0.6587

- Performed linear regression on TrainHL, added

back only middle observations with linear predict<X
  - X = 1.0 – 0.6672
  - **X = 0.75 – 0.6675  --  about 2% of observations removed**
  - X = 0.5 – 0.6661
  - X = 0.25 – 0.6622

Returns >30
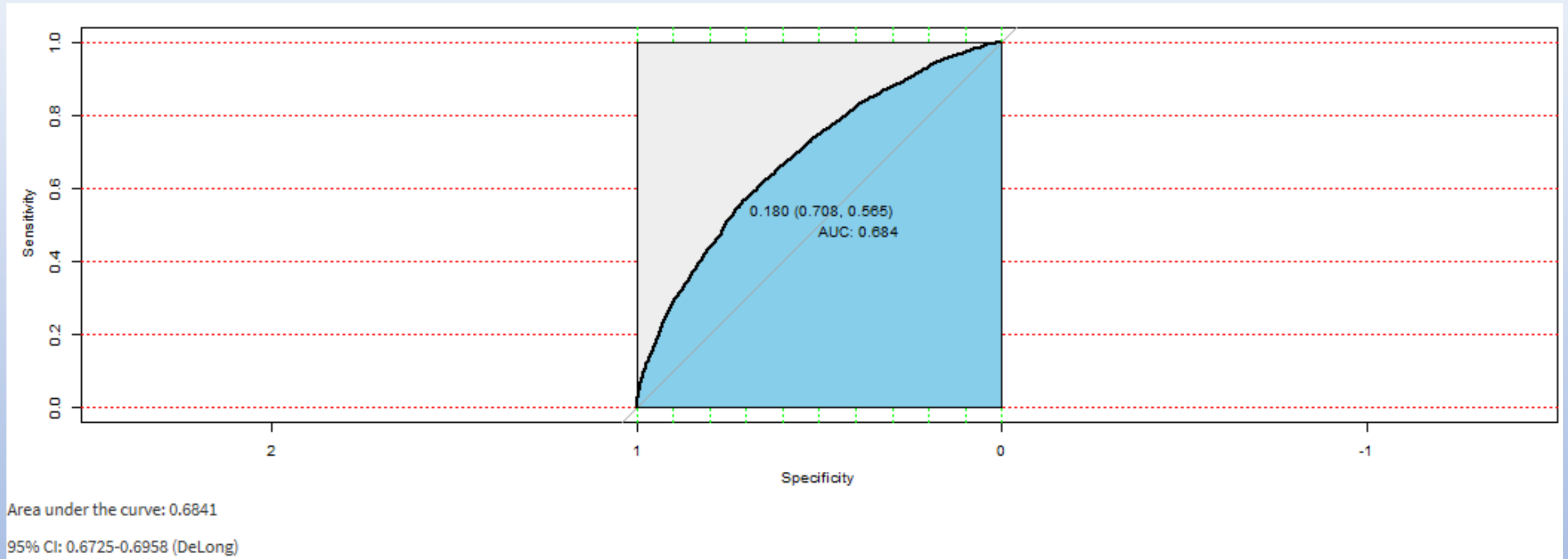Very similar to returns <30 removed from train set

# Testing and Feature Selection

- 11% of observations were a "Yes" (patient readmitted < 30 days)
- 3 models used for hyperparameter tuning:
- 5 K-fold random hyperparameter search – AUC scoring
  - Logistic regression – C=0.1, class weights = 0.2/0.8
    - **AUC = 0.668**
  - Random forest – n_trees = 1000, min_samples_split=5, min_samples_leaf=1, max_features = 'sqrt', max_depth=60, class_weight=0.2/0.8
    - **AUC = 0.671**
  - XGBoost – n_trees = 500, min_child_wt=10, max_depth=5, gamma=5
    - **AUC = 0.680**

# Model Stacking

- AUC optimization – Mapping of true positive vs false positive rate over the range of possible probability threshholds (0.0 – 1.0)
  - AUC/F1 scores are better measurements of performance for rare classifiers than accuracy
    - Max accuracy is to essentially guess "No" for almost everything

- Stacking – Grid search (i, j, k) 0:100
  - Represent percentages of LR, RF, XGB input

- Max AUC score obtained with 56% XGB, 23% RF, 21% LR
  - **AUC = 0.684**

- Used this stacked model to generate probability of return for each patient observation

# AUC – Visualization of true/false positive rates at different probability threshholds
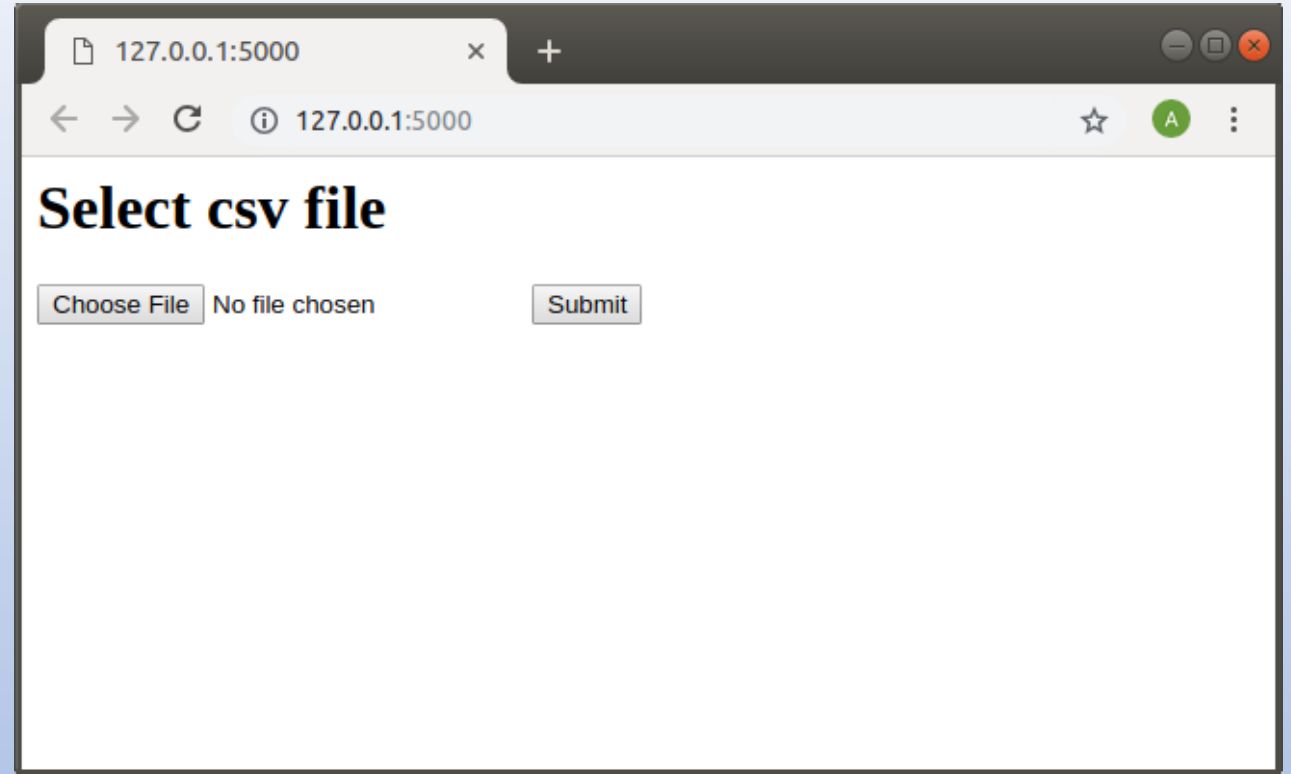


- Decision about who to target for a possible readmission reduction plan will depend on various factors, including:
  Resource availability, balance of true/false positives, possible downstream benefits

# Flask Pipeline

# Pipeline

- Input: a csv file with patient information

- Output: the same csv file with the predictions for each patient and probability of patient being readmitted according to stacked model

# Flask code

```python
12    app = Flask(__name__)
13
14    def transform(text_file_contents):
15        return text_file_contents.replace("=", ",")
16
17
18    @app.route('/')
19    def form():
20        return """
21            <html>
22                <body>
23                    <h1>Select csv file</h1>
24
25                    <form action="/transform" method="post"
26                    enctype="multipart/form-data">
27                        <input type="file" name="data_file" />
28                        <input type="submit" />
29                    </form>
30                </body>
31            </html>
32        """
```

```python
33    @app.route('/transform', methods=["POST"])
34    def transform_view():
35        f = request.files['data_file']
36        if not f:
37            return "No file"
38
39        stream = io.StringIO(f.stream.read().decode("UTF8"), newline=None)
40
41        csv_input = csv.reader(stream)
42        #print("file contents: ", file_contents)
43        #print(type(file_contents))
44
45        # for row in csv_input:
46        #     print(row)
47
48        stream.seek(0)
49        diabetes01 = transform(stream.read())
50        # print(pd.read_csv(StringIO(diabetes01)))
51
52        result = conduct_model(clean_data(pd.read_csv(StringIO(diabetes01))))
53
54        s = io.StringIO()
55
56        result.to_csv(s)
57
58        response = make_response(s.getvalue())
59        response.headers["Content-Disposition"] = "attachment; filename=result.csv"
60        return response
61
62    if __name__ == "__main__":
63        app.debug = True
64        app.run(host='0.0.0.0', port=5000, debug=True)
65
```

# Questions?