

Predicting Patient Volume at Stanford Children's Hospital

...

Conner Brown

28 August 2018

Introduction

Analytics Team at Stanford Children's Health

Work presented by Schienker shows loss of value when departments are inaccurately staffed

Goal: to help nursing managers anticipate the flux of patients in their department

Technology stack: SQL Server-Python-Tableau

Data Science Introduction

Why do we need data science?

Only 3 variables: Date, department and census count (= patient volume = ADC = ACC)

Regress on patient volume - forecast 30 days

Goal: forecast patient volume 30 days in advance and improve Mean Absolute Percentage Error to under 15% for each department

Models: ARIMA, RNN, XGBoost

Data: Type, Size, Departments

Very simple dataset

Rows 16860
Columns 3

Only 3 columns!

	EffectiveDate	DepartmentName	AverageCensusCount
6139	2016-01-01	COMPREHENSIVE CARE PGM	8
6741	2016-01-01	3 NORTH	10
7901	2016-01-01	F2-MATERNITY	21
4737	2016-01-01	PICN 1	11
4061	2016-01-01	SEQ SPEC CARE NURSERY	4

Extract relevant departments

```
raw_df['DepartmentName'].value_counts()
```

COMPREHENSIVE CARE PGM	924
LABOR & DELIVERY	924
SEQ SPEC CARE NURSERY	882
PEDI EL CAMINO	879
NICU 270	844
HEMATOLOGY/ONCOLOGY	823



Data: Statistics

Min is much
lower than 25%



	NICU	CVICU	PICU	PCU
count	844.000000	614.000000	621.000000	2560.000000
mean	28.622038	17.407166	20.713366	14.273828
std	4.240184	2.735248	2.582503	4.700779
min	16.000000	8.000000	11.000000	4.000000
25%	26.000000	16.000000	19.000000	10.000000
50%	29.000000	18.000000	21.000000	14.000000
75%	32.000000	19.000000	23.000000	16.000000
max	39.000000	21.000000	24.000000	26.000000

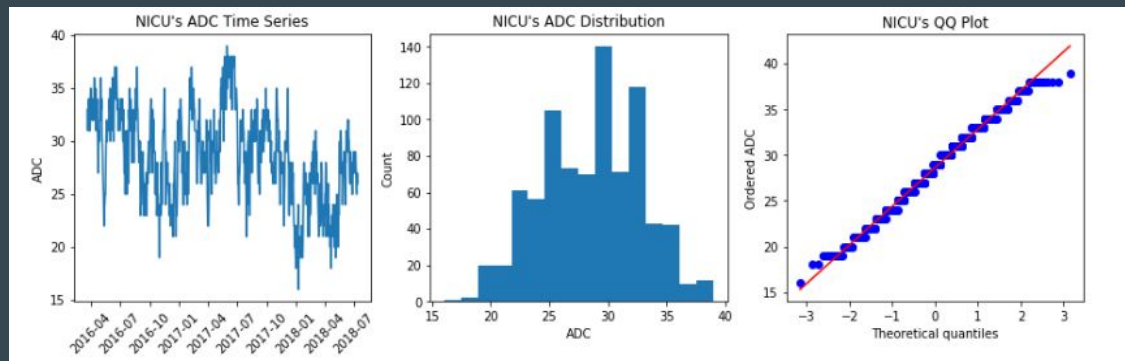


PCU count is 3x
other depts

Number of unique dates in PCU group: 858

PCU group is made
of many small depts

Data: Target Variable



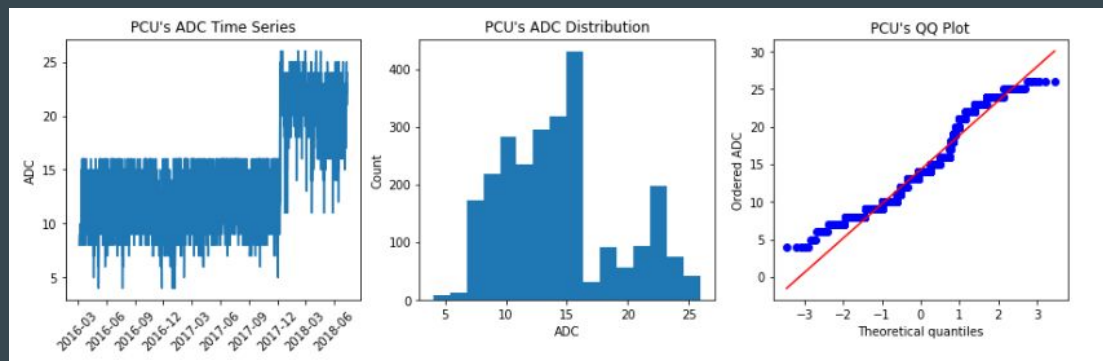
No obvious trend or frequency
in time series

Normal distribution

New hospital opened in Dec. 2017

Add boolean feature
'lpch_main_open'

Keep an eye on low predicted
values

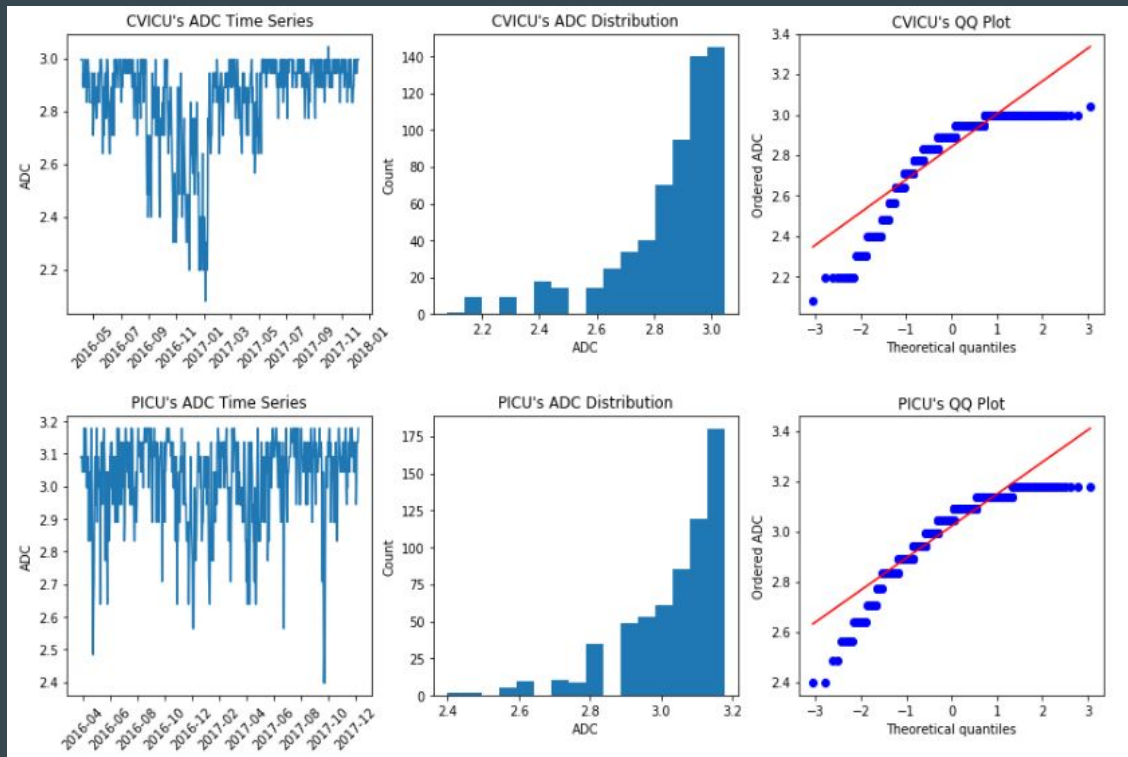


Data: Target Variable

Right skew in CVICU and
PICU departments

Non-normal QQ plot

Strong transformation does
not change normality

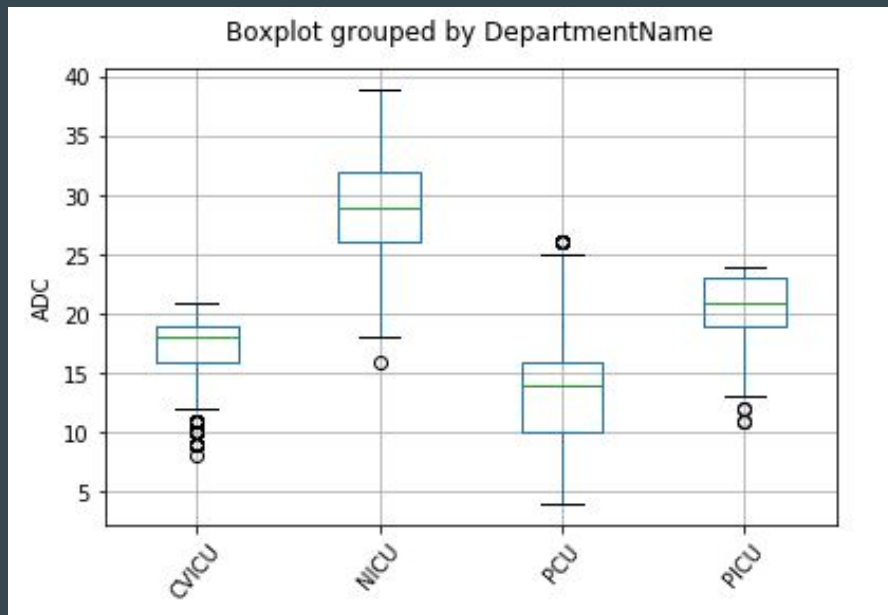


Data: Departments

Few outliers, each department seems unique

Confirmed by T-test results, must model each department separately

NICU vs CVICU	
t-statistic 57.42	pvalue 0.00e+00
NICU vs PICU	
t-statistic 41.19	pvalue 6.43e-247
NICU vs PCU	
t-statistic 78.74	pvalue 0.00e+00
CVICU vs PICU	
t-statistic -21.84	pvalue 1.12e-89
CVICU vs PCU	
t-statistic 15.88	pvalue 1.01e-54
PICU vs PCU	
t-statistic 32.95	pvalue 4.28e-205

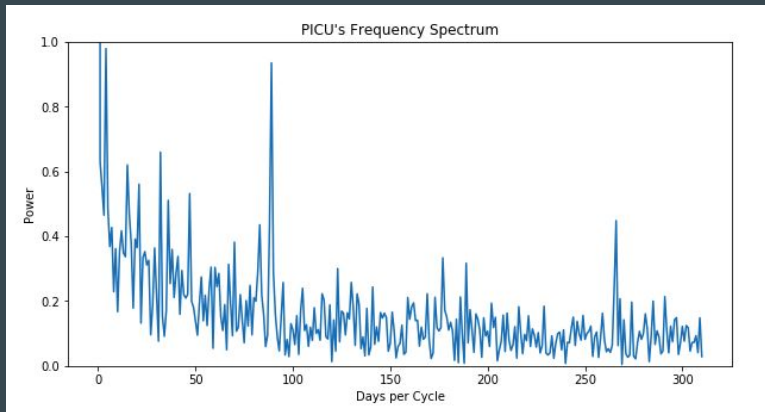


Feature Engineering

ARIMA is simplest, takes time series of values (don't even need the dates!).

RNNs are next easiest, if a row is today, then features are previous m days to today, targets are tomorrow to next n days.

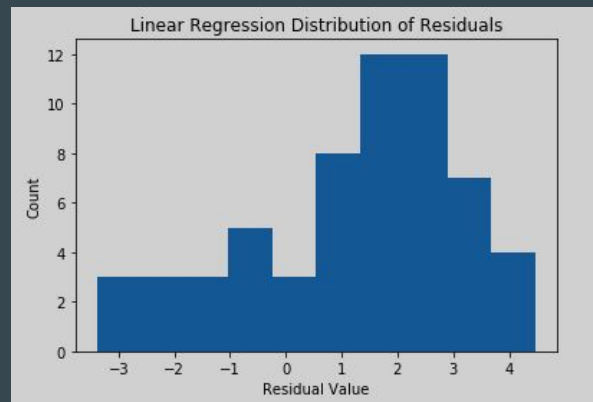
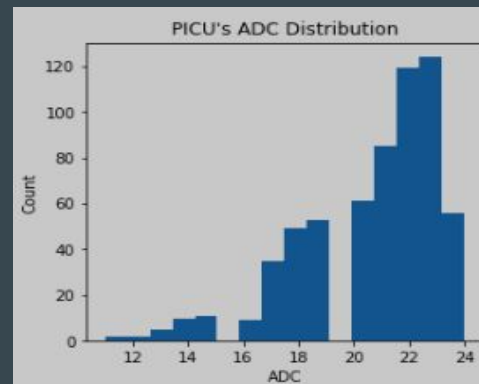
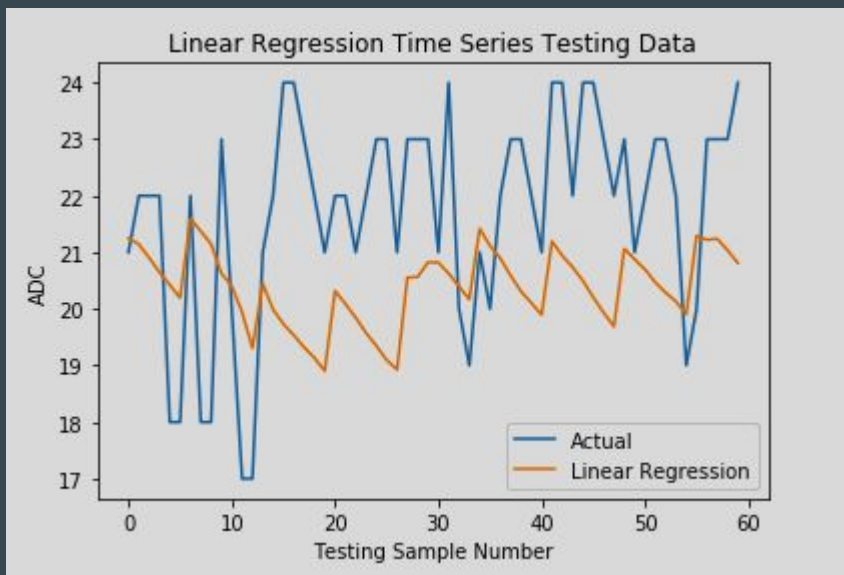
XGBoost is hardest, requires feature engineering with a priori constraints.



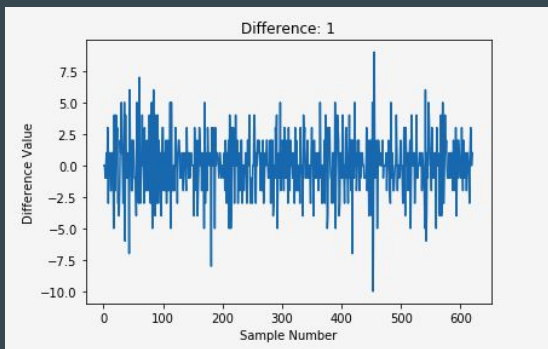
```
# sequential split
def seq_split(df, window_size=7):
    X = df.copy()
    df_temp = df.copy()
    for i in range(window_size):
        X = pd.concat([df_temp.shift(i+1), X], axis=1)
    return X.dropna(axis=0)
```

Modeling: Baseline (Linear Regression)

Feature engineering - date values,
new hospital boolean, rolling
statistics



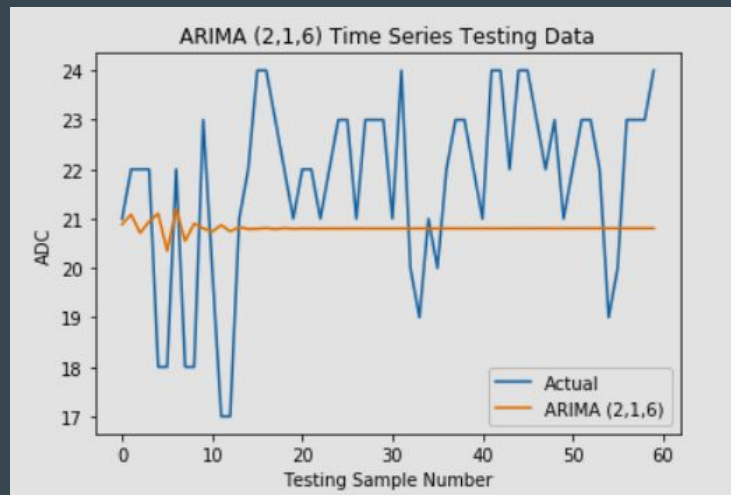
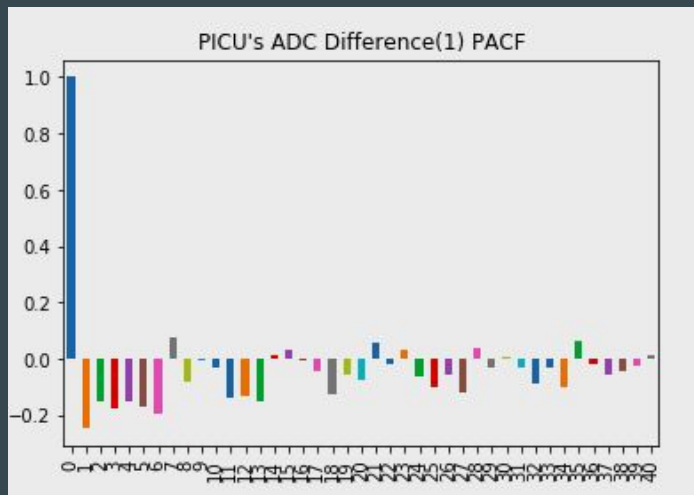
Modeling: ARIMA



Difference for stationarity

PACF for parametric grid search

(2,1,6) AIC of 2671

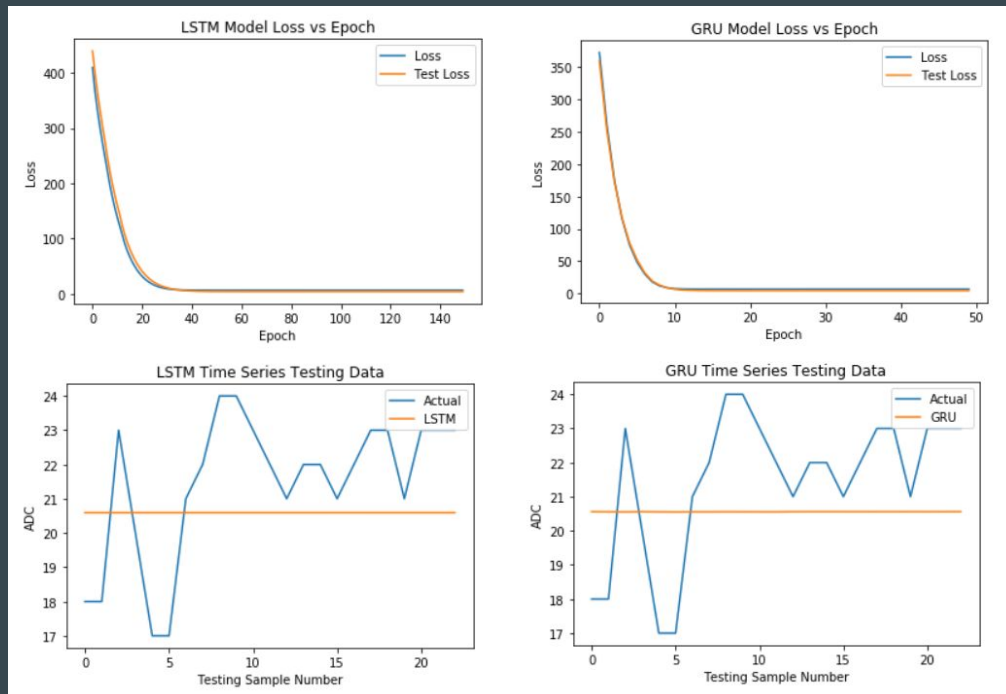


Modeling: RNNs

GRU (Gated Recurrent Unit)

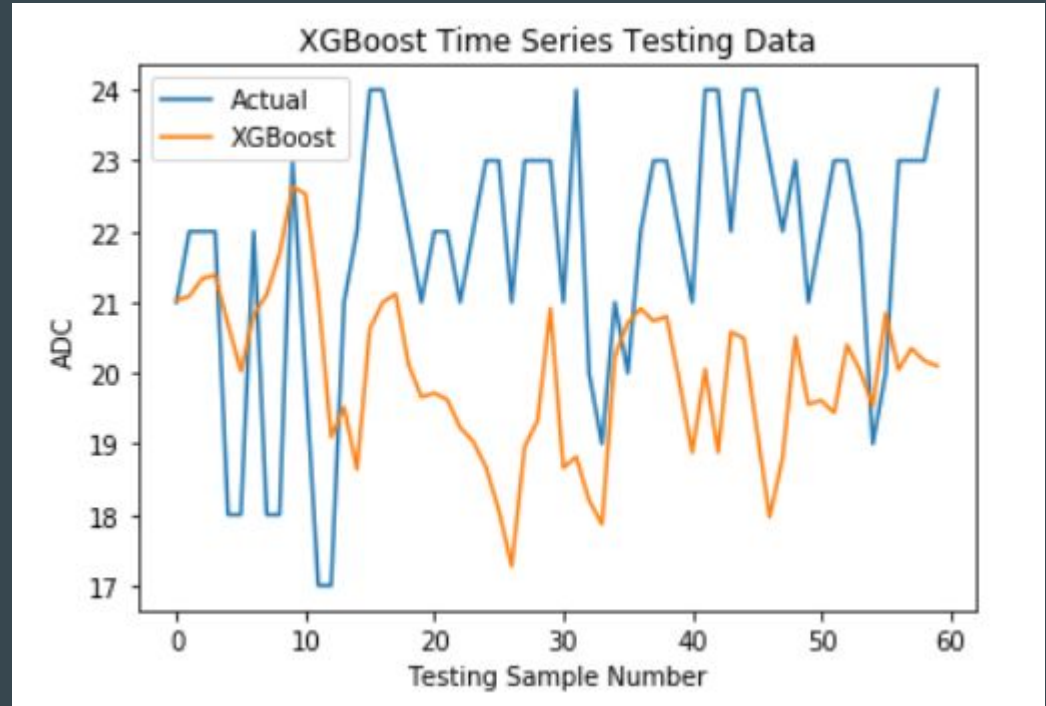
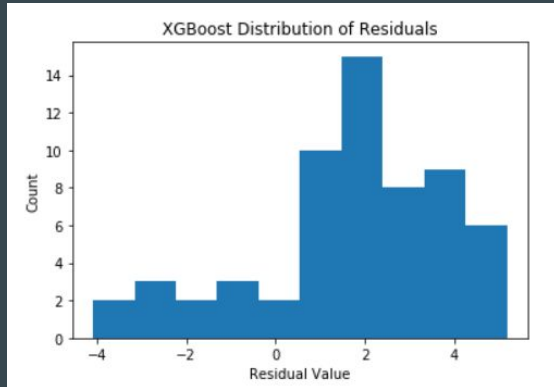
LSTM (Long Short Term Memory)

	Width	Epochs	Batch Size
GRU	64	50	8
LSTM	32	150	8



Modeling: XGBoost

Train on random splitting of data
- more versatile, assumes patient volume does not depend on a short term sequence, but rather the properties of a given day



Evaluation

ARIMA is most accurate

RNNs take the longest to train

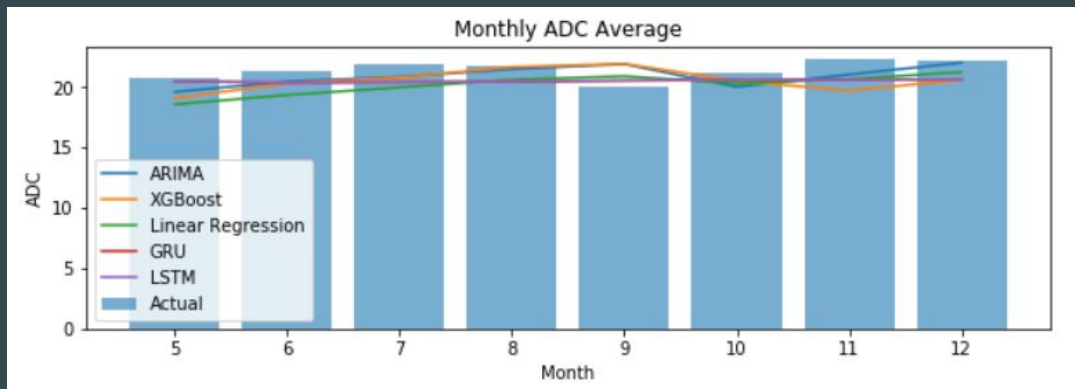
XGBoost and linear regression
show adaptability

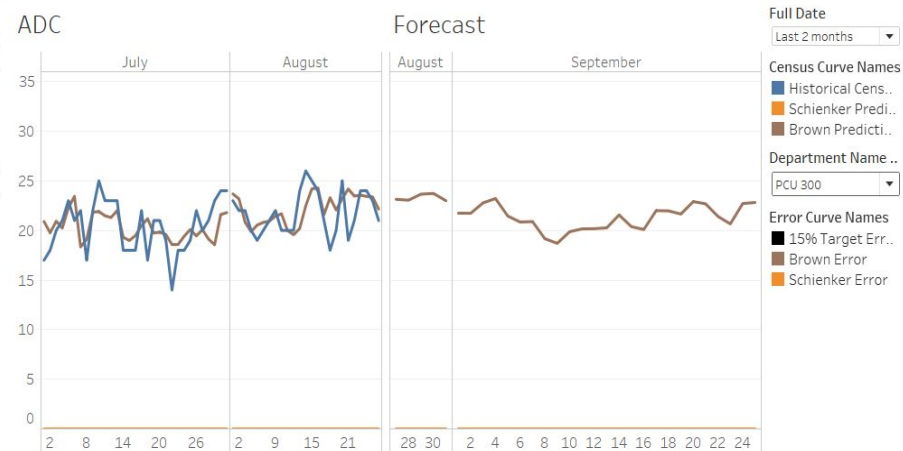
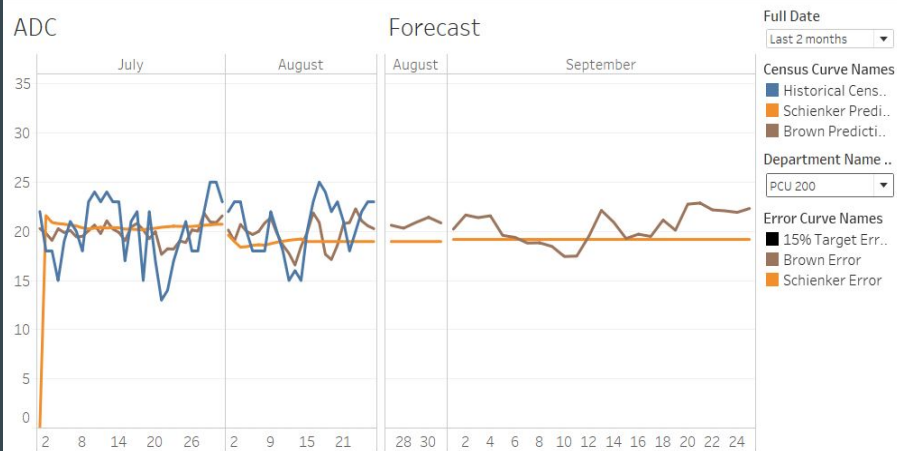
ARIMA only has short term (~10
days) adaptability

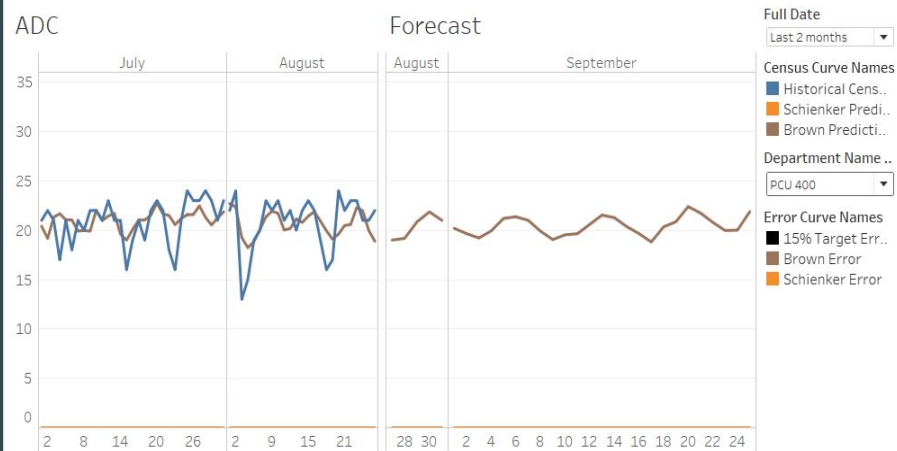
XGBoost has most adaptability
and error within constraints

	Target (Under 15% MAPE)	3.209
ARIMA	0.956 +/- 0.565	
XGBoost	1.308 +/- 0.765	
LinReg	1.425 +/- 0.535	
GRU	1.024 +/- 0.544	
LSTM	1.023 +/- 0.537	

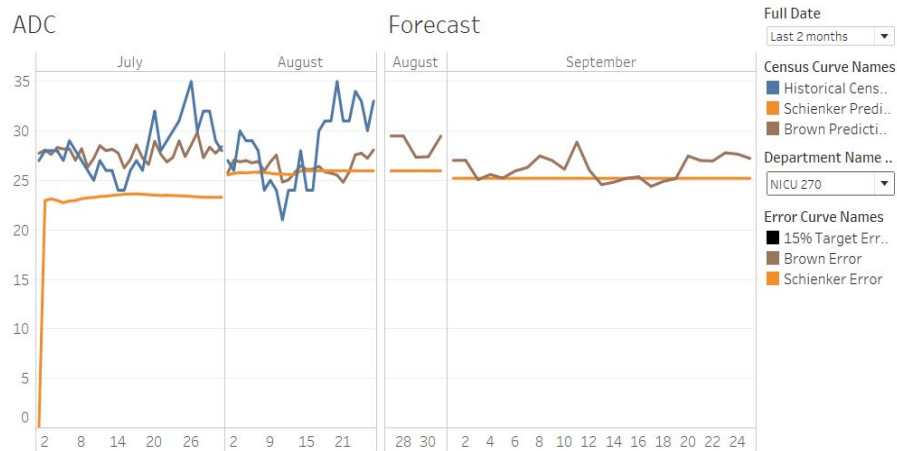
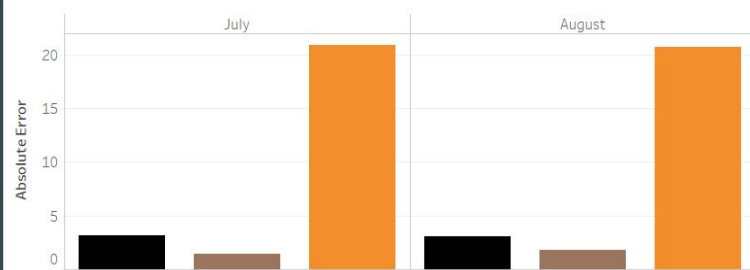
	Average Training Time in seconds
ARIMA	2.793283
XGB	0.115379
LinReg	0.001510
GRU	9.562463
LSTM	23.730989



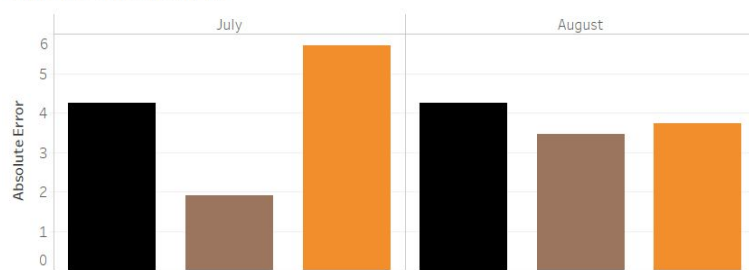


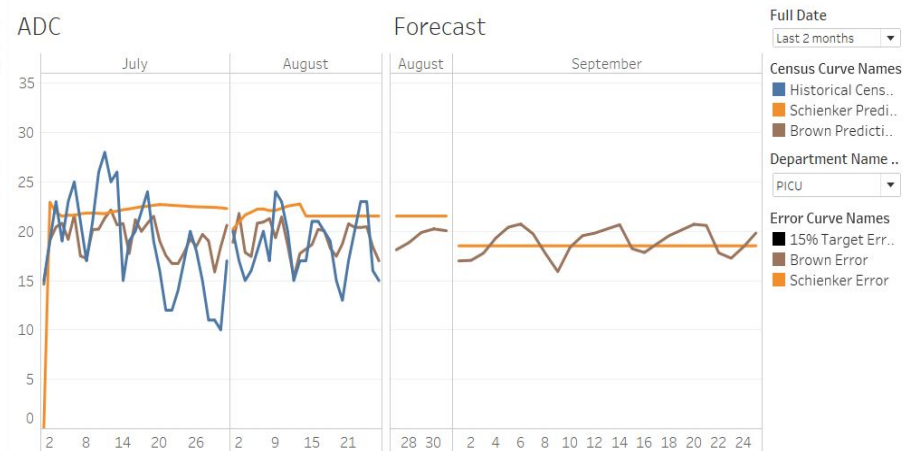
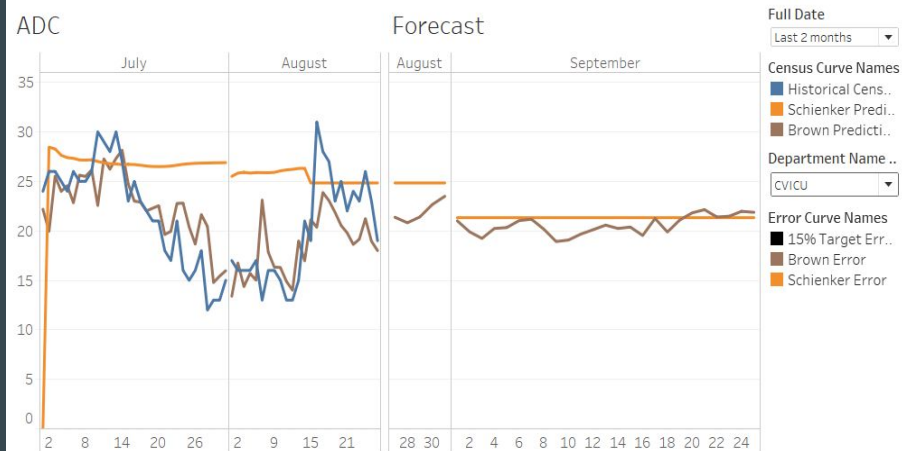


ADC Historical Error

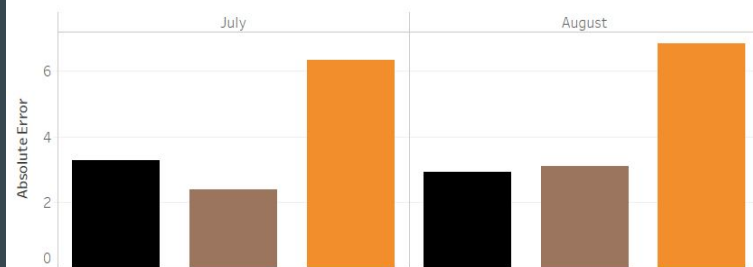


ADC Historical Error

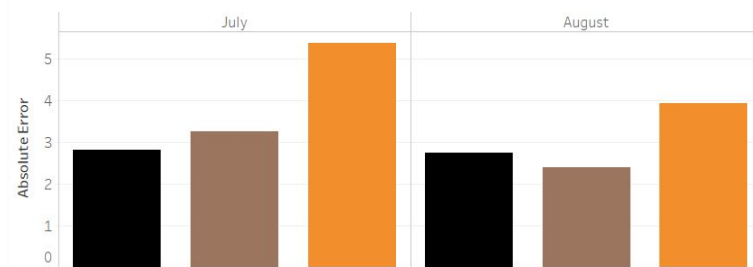




ADC Historical Error



ADC Historical Error



Deployment

Model.py

Job.py

Schedule.py

Run Schedule.bat

```
def schienker_predictions(df):  
def xgb_features(df):  
def xgb_train(df, xgb):  
def open_connection():  
def df_to_sql(df):
```

```
def query_to_df():  
def iter_job(i=0):  
def job(j=0):
```

```
from Job import job  
import schedule  
import time  
  
if __name__ == '__main__':  
    schedule.every().day.at("4:00").do(job)  
    while True:  
        schedule.run_pending()  
        time.sleep(3600)
```

```
@echo off  
py Schedule.py  
pause
```

Thanks!

Brendan Watkins, Ganga Palakkatil

Dave Samuel