# Project Luther:
## Scraping and Modeling Major League Soccer Data

# Overview

- Introduction

- Scraping the Data

- Modeling the Data

- Challenges

- Questions

# Introduction

**Objective:**

To explain how to use the python data stack to scrape and model Major League Soccer (MLS) goalkeeper statistics in order to predict salary.

**Tools Used:**

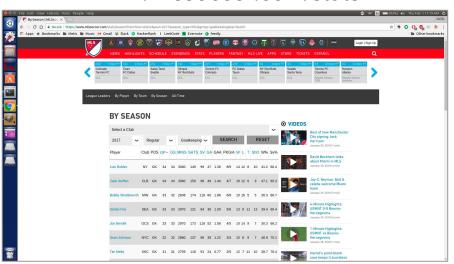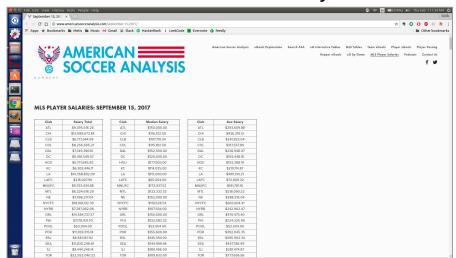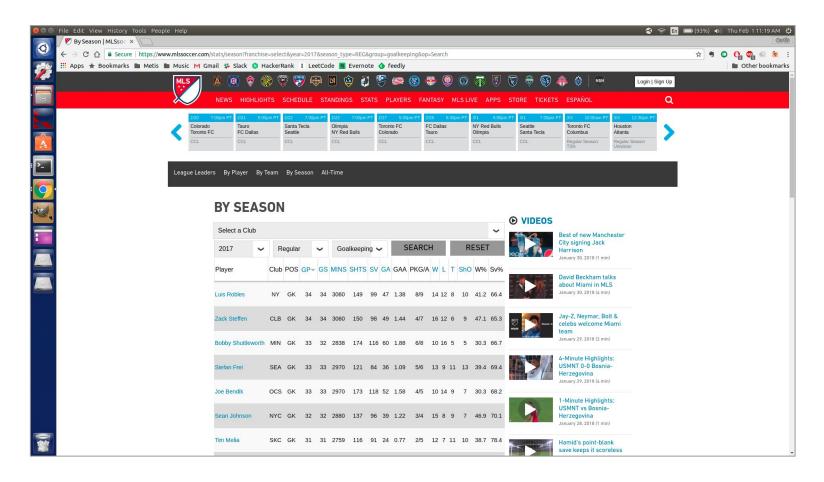# Scraping the Data

# Data Sources
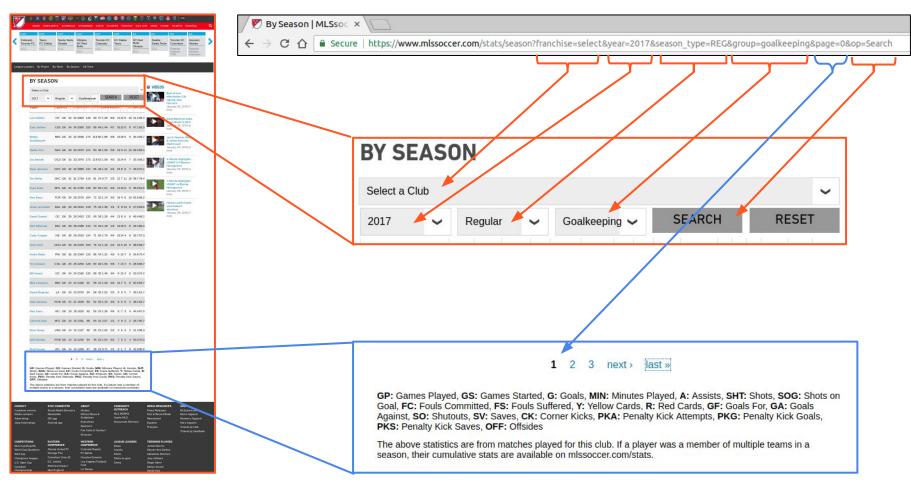
www.mlssoccer.com/stats

www.americansocceranalysis.com



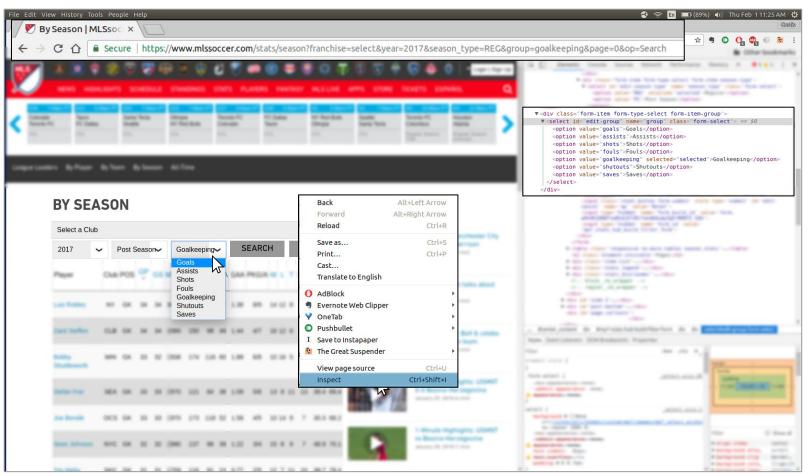Player statistics and salary information for the years 2007 through 2017.

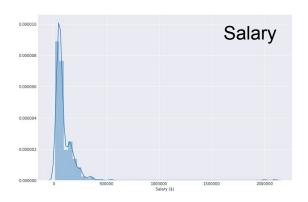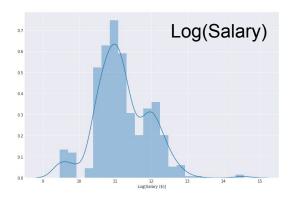# What does mlssoccer.com look like?

# Identifying URL Structure



🔒 Secure | https://www.mlssoccer.com/stats/season?franchise=select&year=2017&season_type=REG&group=goalkeeping&page=0&op=Search

## BY SEASON

| Select a Club | ⌄ |
|---|---|

| 2017 ⌄ | Regular ⌄ | Goalkeeping ⌄ | SEARCH | RESET |
|---|---|---|---|---|

**1**  2  3  next ›  last »

**GP:** Games Played, **GS:** Games Started, **G:** Goals, **MIN:** Minutes Played, **A:** Assists, **SHT:** Shots, **SOG:** Shots on Goal, **FC:** Fouls Committed, **FS:** Fouls Suffered, **Y:** Yellow Cards, **R:** Red Cards, **GF:** Goals For, **GA:** Goals Against, **SO:** Shutouts, **SV:** Saves, **CK:** Corner Kicks, **PKA:** Penalty Kick Attempts, **PKG:** Penalty Kick Goals, **PKS:** Penalty Kick Saves, **OFF:** Offsides

The above statistics are from matches played for this club. If a player was a member of multiple teams in a season, their cumulative stats are available on mlssoccer.com/stats.

# Locating Parameter Values

# Modeling the Data

# Transforming & Exploring

# Fitting a Model



Predicted vs. Actual Salary

# Evaluating the Fit



|  | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **Year** | 0.0847 | 0.008 | 10.414 | 0.000 | 0.069 | 0.101 |
| **GP** | -0.0055 | 0.084 | -0.066 | 0.948 | -0.170 | 0.159 |
| **GS** | 0.4279 | 0.100 | 4.295 | 0.000 | 0.232 | 0.624 |
| **MINS** | -0.0055 | 0.002 | -2.649 | 0.008 | -0.010 | -0.001 |
| **SHTS** | 0.0108 | 0.017 | 0.655 | 0.513 | -0.022 | 0.043 |
| **SV** | -0.0143 | 0.017 | -0.831 | 0.406 | -0.048 | 0.019 |
| **GA** | -0.0168 | 0.019 | -0.862 | 0.389 | -0.055 | 0.021 |
| **GAA** | 0.0800 | 0.043 | 1.852 | 0.065 | -0.005 | 0.165 |
| **ShO** | 0.0268 | 0.025 | 1.056 | 0.292 | -0.023 | 0.077 |
| **SvPct** | 0.0025 | 0.001 | 1.917 | 0.056 | -6.22e-05 | 0.005 |
| **W** | 0.1032 | 0.111 | 0.931 | 0.352 | -0.115 | 0.321 |
| **L** | 0.1322 | 0.111 | 1.192 | 0.234 | -0.086 | 0.350 |
| **T** | 0.1012 | 0.112 | 0.908 | 0.365 | -0.118 | 0.320 |
| **Intercept** | -159.8163 | 16.372 | -9.761 | 0.000 | -192.005 | -127.628 |

| | | | | | | |
|---|---|---|---|---|---|---|
| **Dep. Variable:** | Log_Salary | **R-squared:** | 0.590 | **Omnibus:** | 36.689 | **Durbin-Watson:** | 1.923 |
| **Model:** | OLS | **Adj. R-squared:** | 0.576 | **Prob(Omnibus):** | 0.000 | **Jarque-Bera (JB):** | 71.917 |
| **Method:** | Least Squares | **F-statistic:** | 43.23 | **Skew:** | 0.526 | **Prob(JB):** | 2.42e-16 |
| **Date:** | Thu, 01 Feb 2018 | **Prob (F-statistic):** | 1.94e-67 | **Kurtosis:** | 4.776 | **Cond. No.** | 1.58e+06 |
| **Time:** | 21:13:45 | **Log-Likelihood:** | -268.82 | | | | |
| **No. Observations:** | 405 | **AIC:** | 565.6 | | | | |
| **Df Residuals:** | 391 | **BIC:** | 621.7 | | | | |
| **Df Model:** | 13 | | | | | | |
| **Covariance Type:** | nonrobust | | | | | | |

# Generalizing the Fit



| Regularized Model Scores on Test Data | | |
|---|---|---|
| RidgeCV | LassoCV | ElasticNetCV |
| alpha = 0.1 | alpha ~ $5.9e^{-4}$ | alpha ~ $3.7e^{-4}$ (L1 = 0.975) |
| $R^2$ ~ 0.538, $R^2$-adj ~ 0.497 | $R^2$ ~ 0.544, $R^2$-adj ~ 0.505 | $R^2$ ~ 0.546, $R^2$-adj ~ 0.506 |

# Challenges & Improvements

**Issues:**

- Salary is often fixed for several years at a time due to contract structure.
- Extreme outliers skew the distribution of salaries positively and the features we have cannot explain this variance.
- Many players bring intangible value to their teams that isn't captured in their stats.
- Data integrity issues (errors in the MLS database)

**Possible Solutions:**

- Add additional features to the dataset:
  - Contract signing year and any performance incentives.
  - Complete player history, not just their time in the MLS (i.e., international experience, time in other leagues)
- Data validation with other sources.

# Questions?

# Backup

# Querying Example

```python
# Define base URL
base_url = 'http://www.mlssoccer.com/stats/season'

# Define query parameters
params = {'franchise': select,
          'group': 'goalkeeping',
          'season_type': 'REG',
          'year': 2017,
          'page': 0,
          'op': 'Search'})

# Send request
response = requests.get(url, params)
response.text
```

'<!DOCTYPE html>\n<!--[if IEMobile 7]><html class="iem7"  lang="en" dir="ltr"><![endif]-->\n<!--[if lte IE 6]><html class="lt-ie9 lt-ie8 lt-ie7"  lang="en" dir="ltr"><![endif]-->\n<!--[if (IE 7)&(!IEMobile)]><html class="lt-ie9 lt-ie8"  lang="en" dir="ltr"><![endif]-->\n<!--[if IE 8]><html class="lt-ie9"  lang="en" dir="ltr"><![endif]-->\n<!--[if IE 9]><html class="eq-ie9"  lang="en" dir="ltr"><![endif]-->\n<!--[if (gte IE 9)|(gt IEMobile 7)]><!-->\n<html class="no-ie"  lang="en" dir="ltr"\n  xmlns:og="http://ogp.me/ns#"\n  xmlns:article="http://ogp.me/ns/article#"\n  xmlns:book="http://ogp.me/ns/book#"\n  xmlns:profile="http://ogp.me/ns/profile#"\n  xmlns:video="http://ogp.me/ns/video#"><!--<![endif]-->\n<head profile="http://www.w3.org/1999/xhtml/vocab">\n\n  <meta charset="utf-8" />\n<link rel="apple-touch-icon" sizes="144x144" href="/sites/league/themes/league/img/apple-touch-icon-144x144-precomposed.png" />\n<link rel="apple-touch-icon" sizes="152x152" href="/sites/league/themes/league/img/apple-touch-icon-152x152-precomposed.png" />\n<link rel="apple-touch-icon" sizes="120x120" href="/sites/league/themes/league/img/apple-touch-icon-120x120-precomposed.png" />\n<link rel="shortcut icon" href="https://league-mp7static.mlsdigital.net/favicon.ico?MbTlG7NMwJYaZIfo0mOlfyfDvF9eMba2" type="image/vnd.microsoft.icon" />\n<link rel="apple-touch-icon" sizes="72x72" href="/sites/league/themes/league/img/apple-touch-icon-72x72-precomposed.png" />\n<link rel="apple-touch-icon" sizes="60x60" href="/sites/league/themes/league/img/apple-touch-icon-60x60-precomposed.png" />\n<link rel="apple-touch-icon" sizes="57x57" href="/sites/league/themes/league/img/apple-touch-icon.png" />\n<link rel="apple-touch-icon" sizes="76x76" href="/sites/league/themes/league/img/apple-touch-icon-76x76-precomposed.png" />\n<link rel="apple-touch-icon" sizes="114x114" href="/sites/league/themes/league/img/apple-touch-icon-114x114-precomposed.png" />\n<meta name="apple-itunes-app" content="app-id=397303467" />\n<script src="//cdns.gigya.com/js/gigyaGAIntegration.js"></script>\n<script src="//cdns.gigya.com/js/socialize.js?apiKey=3_qXcJkloa6NFF9zexvt85l9soAHM8lMBWhxcXyhpo3eqangPp8bQONNH8vunw-rTE&amp;lang=en">{lang: "en"}</script>\n<script>var _sf_startpt=

# Parsing and Compiling

```python
# Parse HTML response
soup = BeautifulSoup(response.text,'lxml')
stats_table = soup.find('table')

# Extract salary data
stat_header = []
stat_data = []
for row in stat_table.findAll('tr'):
    row_data = []

    # Get row type and check if header or data row
    row_type = row.findChild().name
    if row_type == 'th':
        # Extract header
        for h in stat_table.findAll('th'):
            stat_header.append(h.text)
    else:
        # Extract data
        for data in row.findAll('td'):
            row_data.append(data.text)
        stat_data.append(row_data)

# Compile stat data into dataframe
stat_df = pd.DataFrame(stat_data, columns=stat_header)
```

| Player | Club | POS | GP | GS | MINS | SHTS | SV | GA | GAA | W | L | T | ShO | Wpct | SvPct | Year | Season | PKG | PKA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bouna Coundoul | COL | GK | 30.0 | 30.0 | 2668.0 | 158.0 | 120.0 | 32.0 | 1.07 | 9.0 | 12.0 | 8.0 | 9.0 | 30.0 | 75.9 | 2007 | REG | 2.0 | 2.0 |
| Kevin Hartman | KC | GK | 30.0 | 30.0 | 2700.0 | 159.0 | 110.0 | 45.0 | 1.50 | 11.0 | 12.0 | 7.0 | 5.0 | 36.7 | 69.2 | 2007 | REG | 6.0 | 7.0 |
| Matt Reis | NE | GK | 30.0 | 30.0 | 2700.0 | 169.0 | 120.0 | 43.0 | 1.43 | 14.0 | 8.0 | 8.0 | 10.0 | 46.7 | 71.0 | 2007 | REG | 3.0 | 3.0 |
| Joe Cannon | LA | GK | 29.0 | 29.0 | 2610.0 | 171.0 | 119.0 | 46.0 | 1.59 | 9.0 | 13.0 | 7.0 | 5.0 | 31.0 | 69.6 | 2007 | REG | 4.0 | 5.0 |
| Troy Perkins | DC | GK | 29.0 | 29.0 | 2610.0 | 155.0 | 117.0 | 32.0 | 1.10 | 16.0 | 6.0 | 7.0 | 8.0 | 55.2 | 75.5 | 2007 | REG | 1.0 | 3.0 |
| Kendall McIntosh | POR | GK | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2017 | REG | 0.0 | 0.0 |
| Josh Saunders | ORL | GK | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2017 | REG | 0.0 | 0.0 |
| Eric Kronberg | MTL | GK | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2017 | REG | 0.0 | 0.0 |
| Bryan Meredith | SEA | GK | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2017 | REG | 0.0 | 0.0 |
| Ryan Meara | NY | GK | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 2017 | REG | 0.0 | 0.0 |

# Linking Datasets