# Overview

**Database Creation:**

- downloading college data CSV, web scraping using Python packages
- establishing relational tables with constraints
- utilizing regular expression functions for data cleaning
- joins/subqueries to create appropriate tables

**Analysis (SQL, Python, R):**

- general theme of our analysis: looking at mid-career salary (10+ years experience)
- exploratory data analysis: correlation analysis, boxplot state comparison
- NWC specific, general Oregon/Washington analysis; school size differentiation
- fitting a basic regression model with mid-career salary as the target variable

# Data Collection

## College Data (CSV file):

- public Kaggle dataset from 2020 with easily downloadable format



---

## Oregon/Washington Schools (Webscraping):

- Data comes from PayScale's College Salary Report (2021)

```
salary_html_oregon = requests.get('https://www.payscale.com/college-salary-report/best-schools-by-state/bachelors/orego
df_oregon = pd.read_html(salary_html_oregon)[0]
cols = ['Rank by Mid-Career Pay', 'School Name', 'School Type', 'Early Career Pay($)', 'Mid-Career Pay($)', '% High Mea
df_oregon.columns = cols
df_oregon.head(3) # looking at the data will need to use some regex selection in SQL
```
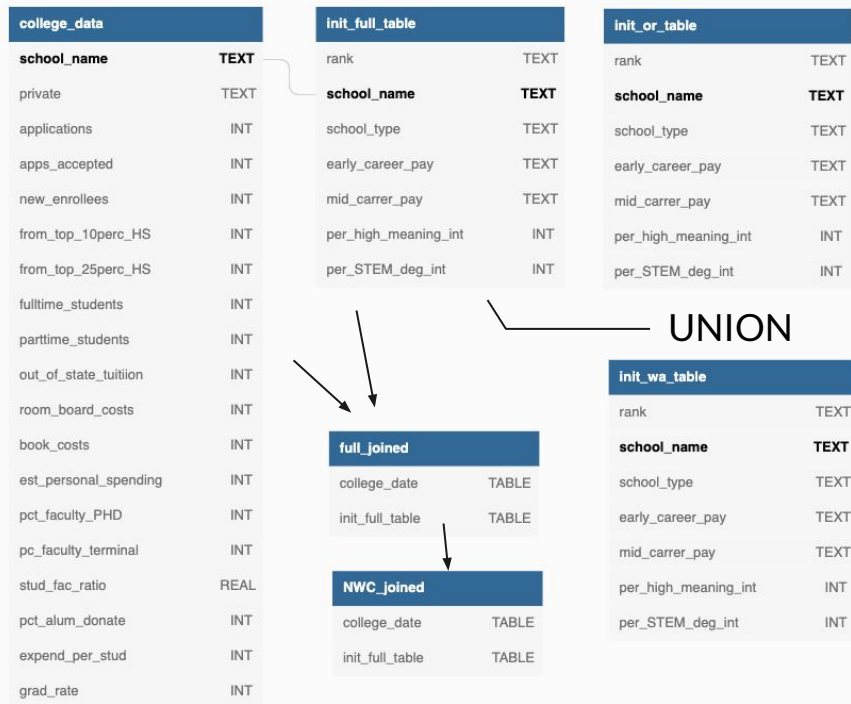
# Regular Expressions

Functions: `regexp_split_to_array(text, re), regexp_match(str, re)`

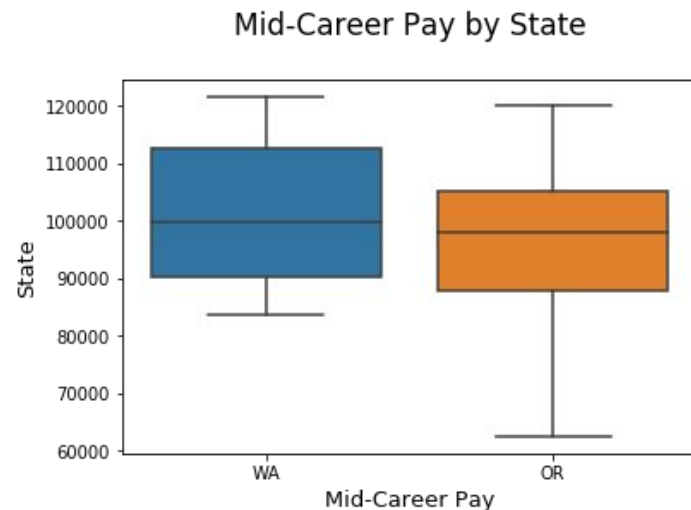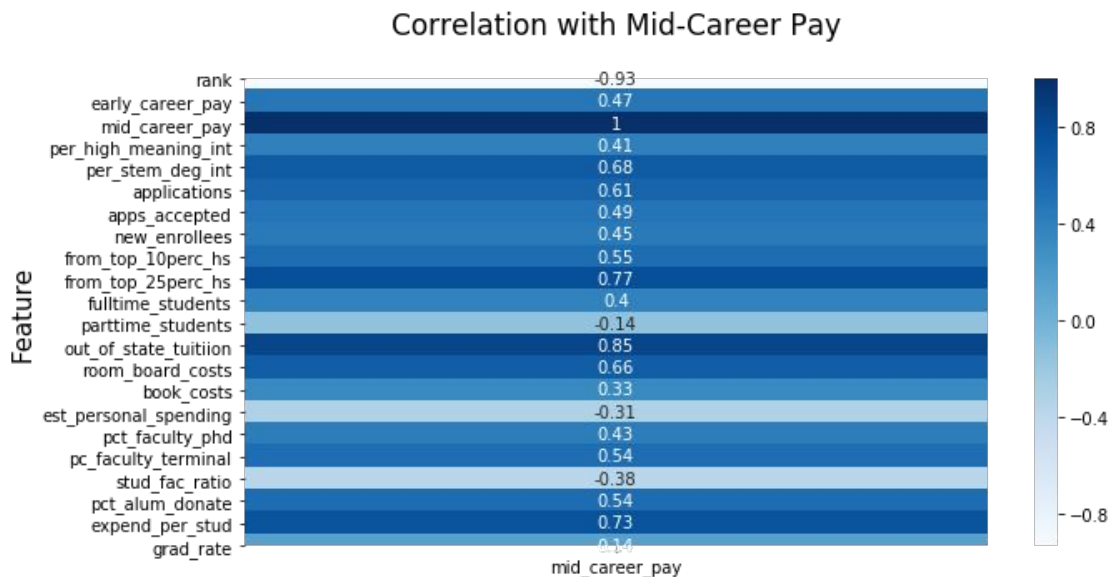| rank<br>text | * school_name<br>text | school_type<br>text | early_career_pay<br>text | mid_carrer_pay<br>text | per_high_meaning<br>text | per_stem_deg<br>text |
|---|---|---|---|---|---|---|
| Rank:1 | School Name:Reed College | School Type:Liber | Early Career Pay:$61,600 | Mid-Career Pay:$120,100 | % High Meaning:40% | % STEM Degrees:31% |
| Rank:2 | School Name:Willamette Univ | School Type:Liber | Early Career Pay:$55,600 | Mid-Career Pay:$116,100 | % High Meaning:53% | % STEM Degrees:13% |
| Rank:3 | School Name:University of Po | School Type:Priva | Early Career Pay:$64,000 | Mid-Career Pay:$112,300 | % High Meaning:51% | % STEM Degrees:31% |

| rank<br>integer | school_name<br>text | school_type<br>text | early_career_pay<br>integer | mid_career_pay<br>integer | per_high_meaning_int<br>integer | per_stem_deg_int<br>integer |
|---|---|---|---|---|---|---|
| 8 | Lewis and Clark College | Liberal Arts School, Priva | 53900 | 104500 | 57 | 11 |
| 1 | Reed College | Liberal Arts School, Priva | 61600 | 120100 | 40 | 31 |
| 2 | Willamette University | Liberal Arts School, Priva | 55600 | 116100 | 53 | 13 |

# Relational Database

# Exploratory Data Analysis - Visualization



Correlation with Mid-Career Pay

| Feature | mid_career_pay |
|---|---|
| rank | -0.93 |
| early_career_pay | 0.47 |
| mid_career_pay | 1 |
| per_high_meaning_int | 0.41 |
| per_stem_deg_int | 0.68 |
| applications | 0.61 |
| apps_accepted | 0.49 |
| new_enrollees | 0.45 |
| from_top_10perc_hs | 0.55 |
| from_top_25perc_hs | 0.77 |
| fulltime_students | 0.4 |
| parttime_students | -0.14 |
| out_of_state_tuitiion | 0.85 |
| room_board_costs | 0.66 |
| book_costs | 0.33 |
| est_personal_spending | -0.31 |
| pct_faculty_phd | 0.43 |
| pc_faculty_terminal | 0.54 |
| stud_fac_ratio | -0.38 |
| pct_alum_donate | 0.54 |
| expend_per_stud | 0.73 |
| grad_rate | 0.14 |



Mid-Career Pay by State

# Analysis Findings

7 of the 9 NWC schools (excluding George Fox, Whitworth) have a higher than average median mid-career salary for their respective state

Three out of the top ten schools come from the NWC – Whitman, Willamette, Puget Sound.

Willamette ranks 2nd in Oregon for median mid-career salary, 2nd in the NWC (Whitman)

We are 95% confident that the true difference in mean between the median salaries of Washington and Oregon schools (WA - OR) is between $-2220.00 and $13906.19. This suggests that there is not a significant difference between the two as $0 is contained in the interval.

[ -2220.00 , 13906.19] dollars ($)

# Small, Medium, Large Schools

We were able to group the colleges based on size (small: x<3000, medium: 3000>x>9999, large: x > 10000 students) and run certain aggregates across them.

We were able to learn that while large universities get more of the top 25 and top 10 percent high school students, they still tend to have the lowest average graduation rate.

While small universities tend to have a higher average acceptance rate, they tend to have a higher average of mid-career pay (in OR and WA)

We were able to group based on the two states (OR and WA) and find that while the average percentage of PHD professors is higher in Oregon, both the graduation rate and the mid career pay is on average lower than that of Washington.

# Analysis on WA and OR schools

We were able to use the partitioning method on the table on the size of the school to see the top 10 highest mid career pay. Through this analysis we were able to learn that out of the 10 schools in this list, 7 of the schools were small schools, 2 of the schools were large, and 1 of the schools was medium sized. It comes to show that regardless of the size of the university, small schools tend to have a higher average mid career pay than other school sizes.

Although, this it is important to note that this is only representative of universities in Washington and Oregon, and not statistically significant throughout all U.S. schools.

```sql
SELECT
    j.school_name,
    j.school_size,
    ROUND(AVG(j.mid_career_pay) OVER (PARTITION BY j.school_size ORDER BY j.rank DESC) , 0) AS "mid_pay"
FROM full_joined AS j
ORDER BY "mid_pay" DESC
LIMIT 10;
```

| | school_name | school_size | mid_pay |
|---|---|---|---|
| 1 | Washington State Unive | large | 110100 |
| 2 | Seattle University | small | 109293 |
| 3 | Reed College | small | 109293 |
| 4 | University of Oregon | large | 107400 |
| 5 | Willamette University | small | 107383 |
| 6 | Whitman College | small | 106591 |
| 7 | University of Portland | small | 106591 |
| 8 | Gonzaga University | small | 104500 |
| 9 | Western Washington Un | medium | 103950 |
| 10 | University of Puget Sour | small | 102875 |

# Fitting a Lasso Regression Predictor

**Target variable:** mid-career pay
**Features:** all other features, w/o school name, extended school type

**Testing instances**

> NWC schools

**Categorical encoding:**
state: 'OR' = 1, 'WA' = 0
private: 'Yes' = 1, 'No' = 0

**Training instances**

> All OR/WA schools without NWC schools included

**Model**: Lasso(alpha = X)

**Predictions**

> Output of NWC school predictions, assess

```
In [332]:  lasso = Lasso(alpha=0.25)
           lasso.fit(X_train, y_train)
           y_pred = lasso.predict(X)
           print(f'R2 Score: {r2_score(y, y_pred)}')
           print(f'RMSE: {np.sqrt(MSE(y, y_pred))}')
```

**R^2 Score:** approx. 0.72

**RMSE:** approx. $5200

# Thank you!