

# Analysis of Career Salaries by Major, College Type and Region

## Introduction

The Wall Street Journal collected data about a student's major, type of college attended and region of college and their career earnings at starting and mid career of their careers. They also went ahead and collected data about the salary collected at various percentiles during mid career.

In this project, I conducted an EDA to determine how each of the above mentioned factors contribute to a student's prospective career earning. College applicants consider salary as an important factor before deciding which college path to follow.

## Objective

Identify how undergraduate major, college type and college region impacts an individual's career earning potential and help individuals make an informed decision.

## Collecting the Data

The data required for this project has been provided from a Kaggle [notebook](#). The data is collected and owned by the Wall Street Journal, and since it is a reputable organization, the data can be considered reliable and able to assist us in this project. The data only contains information about colleges and universities located within the United States of America making it less ideal for college applicants who plan to pursue careers in other countries. Data also does not contain college fees required for each course leaving us unable to determine which course, college and region offer the highest value for money.

## Processing the Data

For processing, analyzing and visualizing the data we need to import some key Python libraries

```
In [1]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
plt.style.use('fivethirtyeight')

First we need to import the data

In [2]: degrees_df = pd.read_csv("~/Users/aron/Documents/Portfolio Projects/College Data/degrees.csv")

In [3]: degrees_df

Out[3]:
```

	Undergraduate Major	Starting Median Salary	Mid-Career Median Salary	Percent change from Starting to Mid-Career Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Accounting	\$46,000.0	\$77,100.0	67.6	\$42,200.0	\$56,100.0	\$108,000.0	\$152,000.0
1	Aerospace Engineering	\$57,700.0	\$101,000.0	75.0	\$64,300.0	\$82,100.0	\$127,000.0	\$161,000.0
2	Agriculture	\$42,600.0	\$71,900.0	68.8	\$36,300.0	\$52,100.0	\$96,300.0	\$150,000.0
3	Anthropology	\$36,800.0	\$61,500.0	67.1	\$33,800.0	\$45,500.0	\$89,300.0	\$138,000.0
4	Architecture	\$41,600.0	\$76,800.0	84.6	\$50,600.0	\$62,200.0	\$97,000.0	\$136,000.0
5	Art History	\$35,800.0	\$64,900.0	81.3	\$28,800.0	\$42,200.0	\$87,400.0	\$125,000.0
6	Biology	\$38,000.0	\$64,800.0	67.0	\$36,900.0	\$47,400.0	\$94,500.0	\$135,000.0
7	Business Management	\$43,000.0	\$72,100.0	67.7	\$38,800.0	\$51,500.0	\$102,000.0	\$147,000.0
8	Chemical Engineering	\$63,200.0	\$107,000.0	69.3	\$71,900.0	\$87,300.0	\$143,000.0	\$194,000.0
9	Chemistry	\$42,600.0	\$79,900.0	87.6	\$45,300.0	\$60,700.0	\$108,000.0	\$148,000.0
10	Civil Engineering	\$53,900.0	\$90,500.0	67.9	\$63,400.0	\$75,100.0	\$115,000.0	\$148,000.0
11	Communications	\$38,100.0	\$70,000.0	83.7	\$37,500.0	\$49,700.0	\$98,800.0	\$143,000.0
12	Computer Engineering	\$61,400.0	\$105,000.0	71.0	\$66,100.0	\$84,100.0	\$135,000.0	\$162,000.0
13	Computer Science	\$55,900.0	\$95,500.0	70.8	\$56,000.0	\$74,900.0	\$122,000.0	\$154,000.0
14	Construction	\$53,700.0	\$88,900.0	65.5	\$56,300.0	\$68,100.0	\$118,000.0	\$171,000.0
15	Criminal Justice	\$35,000.0	\$56,300.0	60.9	\$32,200.0	\$41,600.0	\$80,700.0	\$107,000.0
16	Drama	\$35,900.0	\$56,900.0	58.5	\$36,700.0	\$41,300.0	\$79,100.0	\$112,000.0
17	Economics	\$50,100.0	\$98,600.0	96.8	\$50,600.0	\$70,600.0	\$145,000.0	\$210,000.0
18	Education	\$34,900.0	\$52,000.0	49.0	\$29,300.0	\$37,900.0	\$73,400.0	\$102,000.0
19	Electrical Engineering	\$60,900.0	\$103,000.0	69.1	\$69,300.0	\$83,800.0	\$130,300.0	\$168,000.0
20	English	\$38,000.0	\$64,700.0	70.3	\$33,400.0	\$46,800.0	\$93,200.0	\$133,000.0
21	Film	\$37,900.0	\$68,500.0	80.7	\$33,900.0	\$45,500.0	\$100,000.0	\$136,000.0
22	Finance	\$47,900.0	\$88,300.0	84.3	\$47,200.0	\$62,100.0	\$128,000.0	\$195,000.0
23	Forestry	\$39,100.0	\$62,600.0	60.1	\$41,000.0	\$49,300.0	\$78,200.0	\$111,000.0
24	Geography	\$41,200.0	\$65,500.0	59.0	\$40,000.0	\$50,000.0	\$90,800.0	\$132,000.0
25	Geology	\$43,500.0	\$79,500.0	82.8	\$45,000.0	\$59,600.0	\$101,000.0	\$156,000.0
26	Graphic Design	\$35,700.0	\$55,800.0	67.5	\$36,000.0	\$45,500.0	\$80,800.0	\$112,000.0
27	Health Care Administration	\$38,900.0	\$60,600.0	56.2	\$34,600.0	\$45,600.0	\$78,800.0	\$101,000.0
28	History	\$39,200.0	\$71,000.0	81.1	\$37,000.0	\$49,200.0	\$103,000.0	\$149,000.0
29	Hospitality & Tourism	\$37,000.0	\$57,500.0	52.1	\$35,500.0	\$43,600.0	\$81,900.0	\$124,000.0
30	Industrial Engineering	\$57,700.0	\$94,700.0	64.1	\$57,100.0	\$72,300.0	\$132,000.0	\$173,000.0
31	Information Technology (IT)	\$49,100.0	\$74,800.0	52.3	\$44,500.0	\$56,700.0	\$96,700.0	\$129,000.0
32	Interior Design	\$36,100.0	\$53,200.0	47.4	\$35,700.0	\$42,600.0	\$72,500.0	\$107,000.0
33	International Relations	\$40,900.0	\$80,900.0	97.8	\$38,200.0	\$56,000.0	\$111,000.0	\$157,000.0
34	Journalism	\$35,600.0	\$66,700.0	87.4	\$38,400.0	\$48,300.0	\$97,700.0	\$145,000.0
35	Management Information Systems (MIS)	\$49,200.0	\$82,300.0	67.3	\$45,300.0	\$60,500.0	\$108,000.0	\$146,000.0
36	Marketing	\$40,800.0	\$78,600.0	95.1	\$42,100.0	\$56,600.0	\$119,000.0	\$175,000.0
37	Math	\$45,400.0	\$92,400.0	103.5	\$45,200.0	\$64,200.0	\$128,000.0	\$183,000.0
38	Mechanical Engineering	\$67,900.0	\$93,600.0	61.7	\$63,700.0	\$76,200.0	\$120,000.0	\$163,000.0
39	Music	\$35,900.0	\$55,000.0	53.2	\$26,700.0	\$40,200.0	\$88,000.0	\$134,000.0
40	Nursing	\$54,200.0	\$67,000.0	23.6	\$47,600.0	\$55,400.0	\$80,900.0	\$99,300.0
41	Nutrition	\$39,900.0	\$55,300.0	38.8	\$33,900.0	\$44,500.0	\$70,500.0	\$99,200.0
42	Philosophy	\$39,900.0	\$81,200.0	103.5	\$35,500.0	\$52,800.0	\$127,000.0	\$168,000.0
43	Physician Assistant	\$74,300.0	\$91,700.0	23.4	\$66,400.0	\$75,200.0	\$108,000.0	\$124,000.0
44	Physics	\$50,300.0	\$97,300.0	93.4	\$56,000.0	\$74,200.0	\$132,000.0	\$178,000.0
45	Political Science	\$40,800.0	\$78,200.0	91.7	\$41,200.0	\$55,300.0	\$114,000.0	\$168,000.0
46	Psychology	\$35,900.0	\$60,400.0	68.2	\$31,600.0	\$42,100.0	\$87,500.0	\$127,000.0
47	Religion	\$34,100.0	\$52,000.0	52.5	\$29,700.0	\$36,500.0	\$70,900.0	\$96,400.0
48	Sociology	\$36,500.0	\$58,200.0	59.5	\$30,700.0	\$40,400.0	\$81,200.0	\$118,000.0
49	Spanish	\$34,000.0	\$53,100.0	56.2	\$31,000.0	\$40,000.0	\$76,800.0	\$96,400.0

Similarly, we need to import the remaining 2 datasets

```
In [4]: college_type_df = pd.read_csv("~/Users/aron/Documents/Portfolio Projects/College Data/college_type.csv")

In [5]: region_df = pd.read_csv("~/Users/aron/Documents/Portfolio Projects/College Data/region.csv")

In order to perform analysis on the data, the salaries have to be converted into 'float' datatype

In [6]: degrees_df[degrees_df.columns[1:3]] = degrees_df[degrees_df.columns[1:3]].replace({'$',''}, regex=True).astype(float)

In [7]: degrees_df[degrees_df.columns[4:]] = degrees_df[degrees_df.columns[4:]].replace({'$',''}, regex=True).astype(float)

To see if the data type conversion was successful, we check using a preview of the data

In [8]: degrees_df.head()

Out[8]:
```

	Undergraduate Major	Starting Median Salary	Mid-Career Median Salary	Percent change from Starting to Mid-Career Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Accounting	46000.0	77100.0	67.6	42200.0	56100.0	108000.0	152000.0
1	Aerospace Engineering	57700.0	101000.0	75.0	64300.0	82100.0	127000.0	161000.0
2	Agriculture	42600.0	71900.0	68.8	36300.0	52100.0	96300.0	150000.0
3	Anthropology	36800.0	61500.0	67.1	33800.0	45500.0	89300.0	138000.0
4	Architecture	41600.0	76800.0	84.6	50600.0	62200.0	97000.0	136000.0

```
In [9]: college_type_df[college_type_df.columns[2:]] = college_type_df[college_type_df.columns[2:]].replace({'$',''}, regex=True).astype(float)

In [10]: college_type_df.head()

Out[10]:
```

	School Name	Region	Starting Median Salary	Mid-Career Median Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Massachusetts Institute of Technology (MIT)	Engineering	72200.0	126000.0	76800.0	99200.0	168000.0	220000.0
1	California Institute of Technology (CIT)	Engineering	75500.0	123000.0	NaN	104000.0	161000.0	NaN
2	Harvey Mudd College	Engineering	71800.0	122000.0	NaN	96000.0	160000.0	NaN
3	Polytechnic University of New York, Brooklyn	Engineering	62400.0	114000.0	66800.0	94300.0	143000.0	190000.0
4	Cooper Union	Engineering	62200.0	114000.0	NaN	80200.0	142000.0	NaN

```
In [11]: region_df[region_df.columns[2:]] = region_df[region_df.columns[2:]].replace({'$',''}, regex=True).astype(float)

In [12]: region_df.head()

Out[12]:
```

	School Name	Region	Starting Median Salary	Mid-Career Median Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Stanford University	California	70400.0	129000.0	68400.0	93100.0	184000.0	257000.0
1	California Institute of Technology (CIT)	California	75500.0	123000.0	NaN	104000.0	161000.0	NaN
2	Harvey Mudd College	California	71800.0	122000.0	NaN	96000.0	160000.0	NaN
3	University of California, Berkeley	California	59900.0	112000.0	59500.0	81000.0	149000.0	201000.0
4	Occidental College	California	51900.0	105000.0	NaN	54800.0	157000.0	NaN

The College Type dataset and Region dataset both contain some missing values. These values must be removed in order to do analysis. If possible, more research would be done to obtain those values but since that is not feasible, the rows containing missing values will be deleted.

```
In [13]: college_type_df.dropna(axis=0, how='any', thresh=None, subset=None, inplace=False)

Out[13]:
```

	School Name	Region	Starting Median Salary	Mid-Career Median Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Massachusetts Institute of Technology (MIT)	Engineering	72200.0	126000.0	76800.0	99200.0	168000.0	220000.0
3	Polytechnic University of New York, Brooklyn	Engineering	62400.0	114000.0	66800.0	94300.0	143000.0	190000.0
5	Worcester Polytechnic Institute (WPI)	Engineering	61000.0	114000.0	80000.0	91200.0	137000.0	180000.0
6	Carnegie Mellon University (CMU)	Engineering	61800.0	111000.0	63300.0	80100.0	150000.0	209000.0
7	Rensselaer Polytechnic Institute (RPI)	Engineering	61300.0	110000.0	71600.0	85500.0	140000.0	182000.0
...	...	...	...	...	...	...	...	...
264	Austin Peay State University	State	37700.0	59200.0	32200.0	40500.0	73900.0	96200.0
265	Pittsburg State University	State	40400.0	58200.0	25600.0	48000.0	84800.0	117000.0
266	Southern Utah University	State	41900.0	56500.0	30700.0	39700.0	78400.0	116000.0
267	Montana State University - Billings	State	37900.0	50600.0	22600.0	31800.0	78500.0	98900.0
268	Black Hills State University	State	35300.0	43900.0	27000.0	32200.0	60900.0	87600.0

231 rows x 8 columns

```
In [14]: region_df.dropna(axis=0, how='any', thresh=None, subset=None, inplace=False)

Out[14]:
```

	School Name	Region	Starting Median Salary	Mid-Career Median Salary	Mid-Career 10th Percentile Salary	Mid-Career 25th Percentile Salary	Mid-Career 75th Percentile Salary	Mid-Career 90th Percentile Salary
0	Stanford University	California	70400.0	129000.0	68400.0	93100.0	184000.0	257000.0
3	University of California, Berkeley	California	59900.0	112000.0	59500.0	81000.0	149000.0	201000.0
5	Cal Poly San Luis Obispo	California	57200.0	101000.0	55000.0	74700.0	133000.0	170000.0
6	University of California at Los Angeles (UCLA)	California	52600.0	101000.0	51300.0	72500.0	139000.0	193000.0
7	University of California, San Diego (UCSD)	California	51100.0	101000.0	51700.0	75400.0	131000.0	177000.0
...	...	...	...	...	...	...	...	...
315	State University of New York (SUNY) at Potsdam	Northeastern	38000.0	70300.0	35100.0	51200.0	100000.0	179000.0
316	Niagara University	Northeastern	36900.0	69700.0	44000.0	57000.0	92000.0	128000.0
317	State University of New York (SUNY) at Fredonia	Northeastern	37800.0	66200.0	32900.0	44200.0	93300.0	181000.0
318	University of Southern Maine	Northeastern	39400.0	63600.0	40400.0	47900.0	85700.0	117000.0
319	Mercy College	Northeastern	43700.0	62600.0	35600.0	47300.0	99000.0	134000.0

273 rows x 8 columns

## Analyzing the Data

### Analyzing Salaries by Undergraduate Major

First, we are gonna determine which undergraduate degree offers the the highest starting median salary

```
In [15]: high_major_start_sal = degrees_df.loc[degrees_df["Starting Median Salary"].idxmax()]

In [16]: high_major_start_sal

Out[16]:
```

	Physician Assistant
Starting Median Salary	74390.0
Mid-Career Median Salary	91700.0
Percent change from Starting to Mid-Career Salary	23.4
Mid-Career 10th Percentile Salary	66490.0
Mid-Career 25th Percentile Salary	71990.0
Mid-Career 75th Percentile Salary	108690.0
Mid-Career 90th Percentile Salary	124990.0
Name:	43, dtype: object

```
In [17]: degrees_sorted_start = degrees_df.sort_values('Starting Median Salary', ascending=False).head(10)

In [18]: start_sal = degrees_sorted_start['Starting Median Salary'].head(10)

In [19]: start_sal = start_sal.iloc[:1]

In [20]: undergrad_major_start = degrees_sorted_start['Undergraduate Major'].head(10)

In [21]: undergrad_major_start = undergrad_major_start.iloc[:1]

In [22]: plt.barh(undergrad_major_start, start_sal)
plt.title('Top 10 Jobs by Starting Median Salary')
plt.xlabel('Salary in USD')
plt.show()
```



Next, we will find the highest paying career at mid-career

```
In [23]: high_major_mid_sal = degrees_df.loc[degrees_df["Mid-Career Median Salary"].idxmax()]

In [24]: high_major_mid_sal

Out[24]:
```

	Chemical Engineering
Starting Median Salary	63200.0
Mid-Career Median Salary	107690.0
Percent change from Starting to Mid-Career Salary	69.3
Mid-Career 10th Percentile Salary	71990.0
Mid-Career 25th Percentile Salary	87390.0
Mid-Career 75th Percentile Salary	143900.0
Mid-Career 90th Percentile Salary	194990.0
Name:	8, dtype: object

```
In [25]: degrees_sorted_mid = degrees_df.sort_values('Mid-Career Median Salary', ascending=False)

In [26]: mid_sal = degrees_sorted_mid['Mid-Career Median Salary'].head(10)

In [27]: mid_sal

Out[27]:
```

8	107690.0
12	105990.0
19	105990.0
1	101990.0
17	98690.0
44	97200.0
13	95590.0
39	94790.0
38	93690.0
37	92490.0

Name: Mid-Care