

# Job Salary Prediction

## Data Robot

Dhanashree Mane

Sanya Purwar

Arjoo Gangwal



# What are we doing?

- Salary is a **big deciding factor** for a majority of people
- We are trying to build a model that would **predict the average salary** based on
  - Job title
  - Location
  - Years of experience
  - Company Size
  - Company Sector

The perfect  
**Job**

- 
- ✓ good salary
  - ✓ benefits
  - ✓ job security
  - ✓ close to home



Sanya has 2 years of experience in Software Industry

Offer 1:

Location: San Francisco

Company-Size: Small-size

Company-Sector: Software/IT

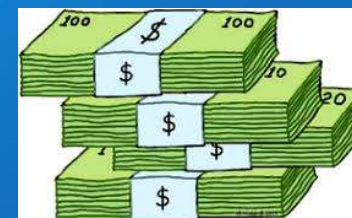


Offer 2:

Location: Washington

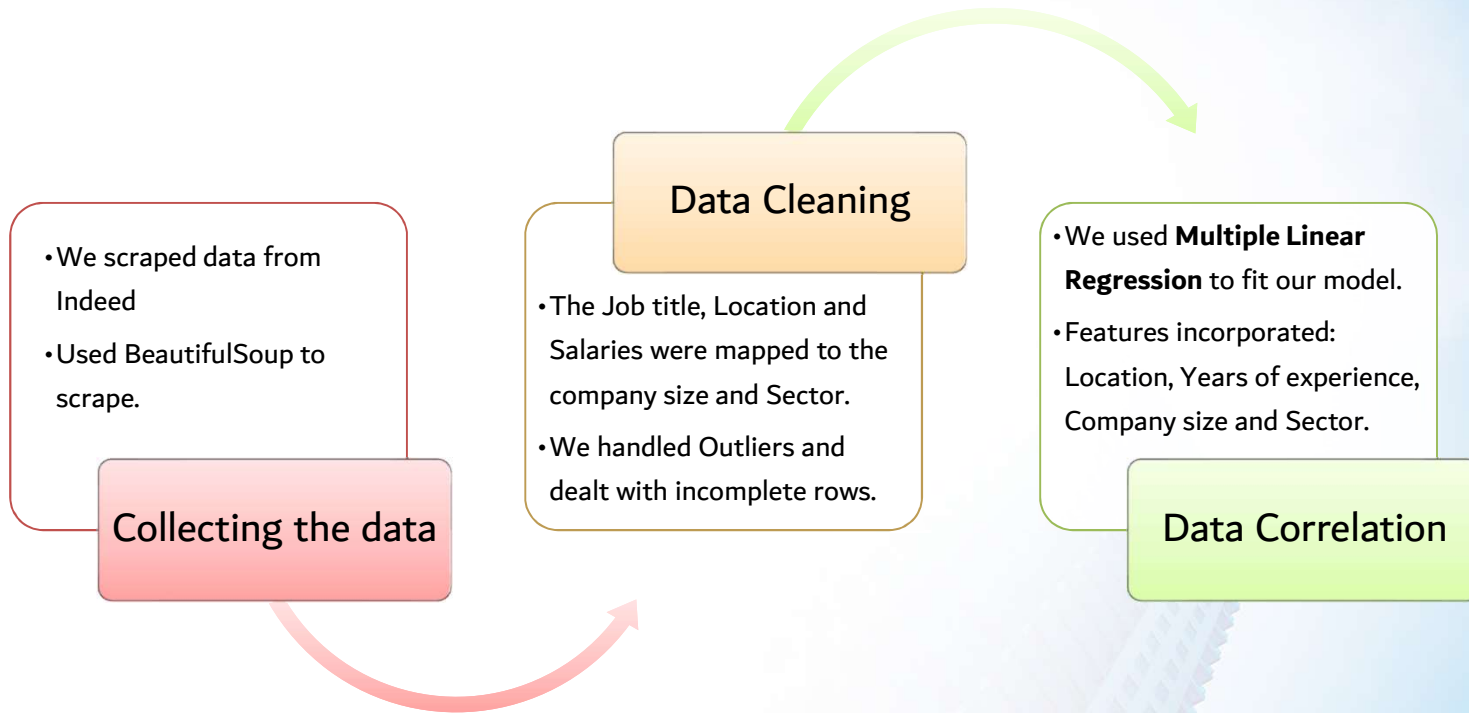
Company-Size: Mid-size

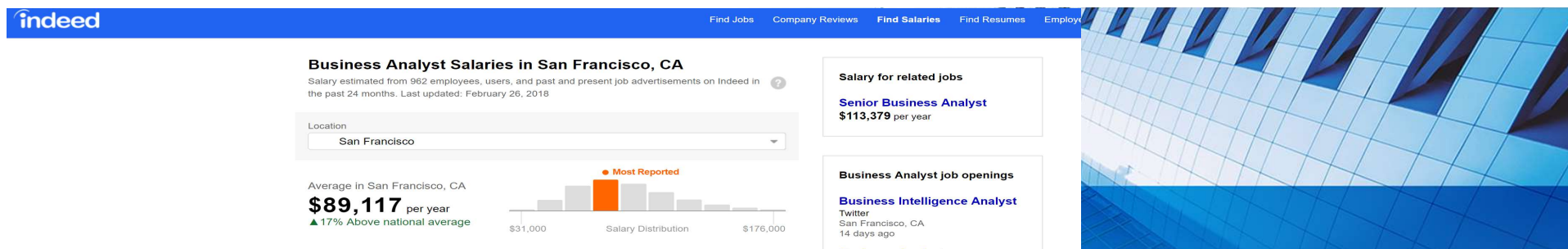
Company-Sector: Software/IT





# Project Pipeline





## Company Name, Avg. Salary, Location from Indeed

	<b>Mason Frank Business Analyst</b> 29 salaries	<b>\$119,634</b> per year
	<b>Salesforce Business Analyst</b> 12 salaries	<b>\$106,648</b> per year
	<b>MUFG Union Bank Business Analyst</b> 5 salaries	<b>\$109,050</b> per year
	<b>FHLBank San Francisco Business Analyst</b> 5 salaries	<b>\$108,796</b> per year
	<b>City and County of San Francisco Business Analyst</b> 28 salaries	<b>\$104,390</b> per year
	<b>San Francisco Department of Public Health Business Analyst</b> 35 salaries	<b>\$102,650</b> per year
	<b>Mason Frank International Business Analyst</b> 12 salaries	<b>\$102,292</b> per year
	<b>Pacific Gas and Electric Company (PG&amp;E) Business Analyst</b> 8 salaries	<b>\$98,386</b> per year
	<b>The San Francisco Department of Public Health Business Analyst</b> 21 salaries	<b>\$101,365</b> per year

## Company Profile from LinkedIn

**in** Search

Twitter, and LinkedIn.

**Company details**

**Website**  
<http://www.symantec.com>

**Headquarters**  
 Mountain View, CA

**Year founded**  
 1982

**Company type**  
 Public Company

**Company size**  
 10,001+ employees

**Specialties**  
 Encryption, Antivirus and Malware protection, Identity Protection and Authentication, Information Protection, Cyber Security Services, Threat Protection, and Cloud Data Protection

# Challenges We Faced

1

Had to manually map the company name with their respective company size and sector by taking information about each company name from LinkedIn

2

Had to manually match the salary with the Years of experience taking references from Glassdoor.com

# Our Data Looks like

	Company Name	Salary	Years of Experience	Location	Company Size	Sector	Job Title
1	3DI INC	167000	7 to 10	Los Angeles	Small-size	Software/IT	Business Analyst
2	3DI INC	123683	4 to 6	Los Angeles	Small-size	Software/IT	Business Analyst
3	3DI INC	80366	1 to 3	Los Angeles	Small-size	Software/IT	Business Analyst
4	3EDGEUSAGROUP LLC	121100	7 to 10	New York	Small-size	LAW	Business Analyst
5	3EDGEUSAGROUP LLC	101727.5	4 to 6	New York	Small-size	LAW	Business Analyst
6	3EDGEUSAGROUP LLC	82355	1 to 3	New York	Small-size	LAW	Business Analyst
7	4C Connect Inc.	149950	7 to 10	Atlanta	Small-size	Accounting	Business Analyst
8	4C Connect Inc.	112475	4 to 6	Atlanta	Small-size	Accounting	Business Analyst
9	4C Connect Inc.	75000	1 to 3	Atlanta	Small-size	Accounting	Business Analyst
10	A.T.KearneyInc.	150000	7 to 10	Atlanta	Mid-size	Consulting	Business Analyst
11	A.T.KearneyInc.	110731.5	4 to 6	Atlanta	Mid-size	Consulting	Business Analyst
12	A.T.KearneyInc.	71463	1 to 3	Atlanta	Mid-size	Consulting	Business Analyst
13	Absolute Opportunities	151920	7 to 10	Washington	Small-size	Education	Business Analyst
14	Absolute Opportunities	127960	4 to 6	Washington	Small-size	Education	Business Analyst
15	Absolute Opportunities	104000	1 to 3	Washington	Small-size	Education	Business Analyst
16	Acadia Technologies	140000	7 to 10	Atlanta	Small-size	Software/IT	Business Analyst
17	Acadia Technologies	97500	4 to 6	Atlanta	Small-size	Software/IT	Business Analyst
18	Acadia Technologies	55000	1 to 3	Atlanta	Small-size	Software/IT	Business Analyst
19	Acadia Technologies						



# Why Multiple linear regression?



Location



Years of Experience

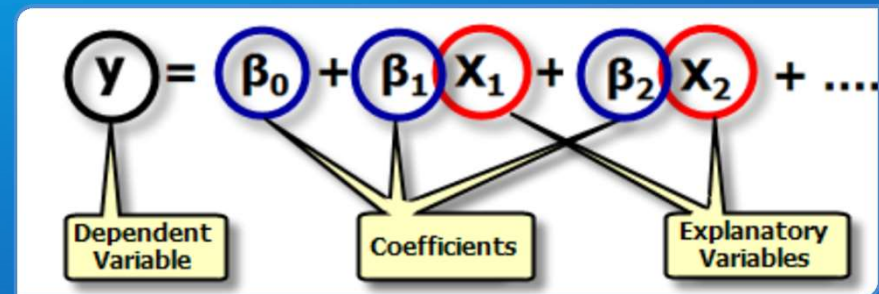


Company Size



Company Sector

- Independent Variables( $X_1, X_2, \dots, X_n$ ): Location, Yrs. of Experience, Company Size, Sector.
- Dependent Variable( $y$ ): Salary



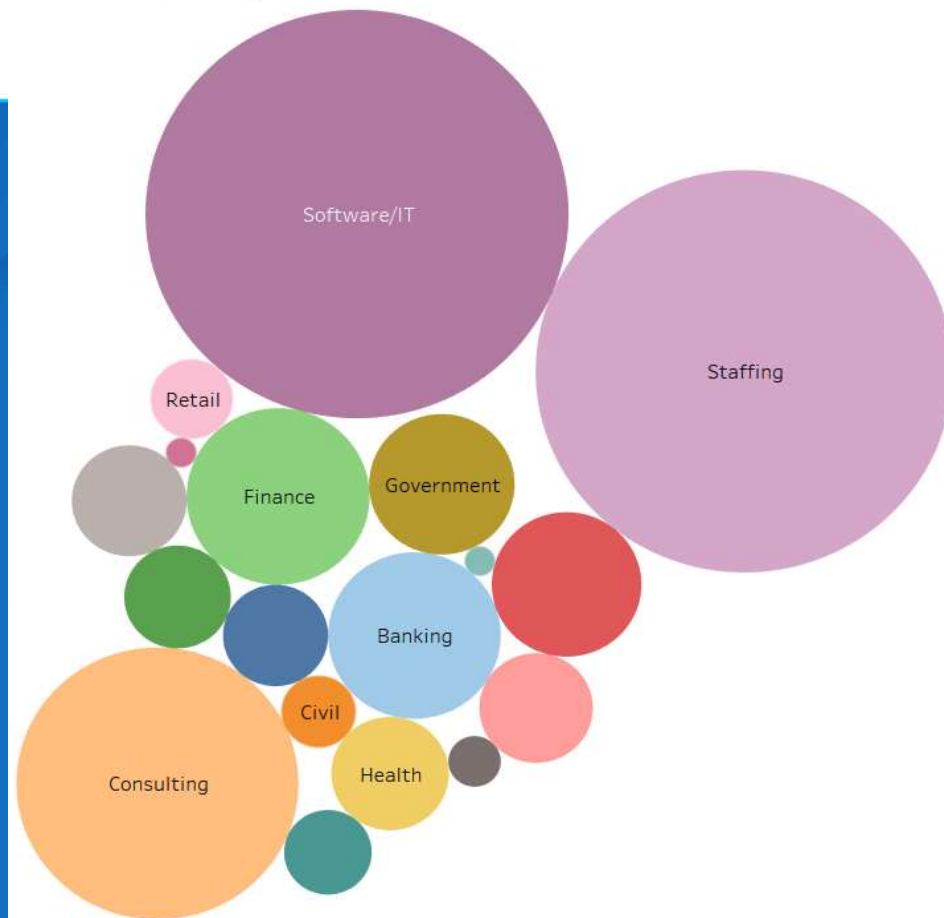


## Location and Yrs of Experience wise Salary

7 to 10 Berkley 189,000	7 to 10 Boston 176,971	7 to 10 Austin 162,043	7 to 10 Houston 159,685	7 to 10 Chicago 158,211	1 to 3 Berkley 120,000	1 to 3 San Francisco 104,341	1 to 3	
7 to 10 San Francisco 185,971	7 to 10 Los Angeles 174,155	7 to 10 San Jose 161,000	7 to 10 Seattle 156,137		1 to 3 New York 97,667	1 to 3 San Jose 96,170	1 to 3 Boston 89,807	
7 to 10 New York 183,115	7 to 10 Mountain View 167,143	7 to 10 Orlando 160,400			7 to 10 Cupertino	1 to 3 Austin 87,194	1 to 3 Los Angeles 86,231	1 to 3 Seattle 83,749
4 to 6 Berkley 154,500	4 to 6 New York 140,199	4 to 6 Los Angeles 129,943	4 to 6 Seattle 119,943	4 to 6 Houston 119,326	1 to 3 Orlando 80,153		1 to 3	1 to 3
4 to 6 San Francisco 145,049	4 to 6 Mountain View 134,780	4 to 6 San Jose 129,200	4 to 6 Chicago 119,117				4 to 6 Cupertino 91,725	1 to 3 Chicago 79,511
4 to 6 Austin 142,838	4 to 6 Boston 133,102	4 to 6 Orlando 120,276						



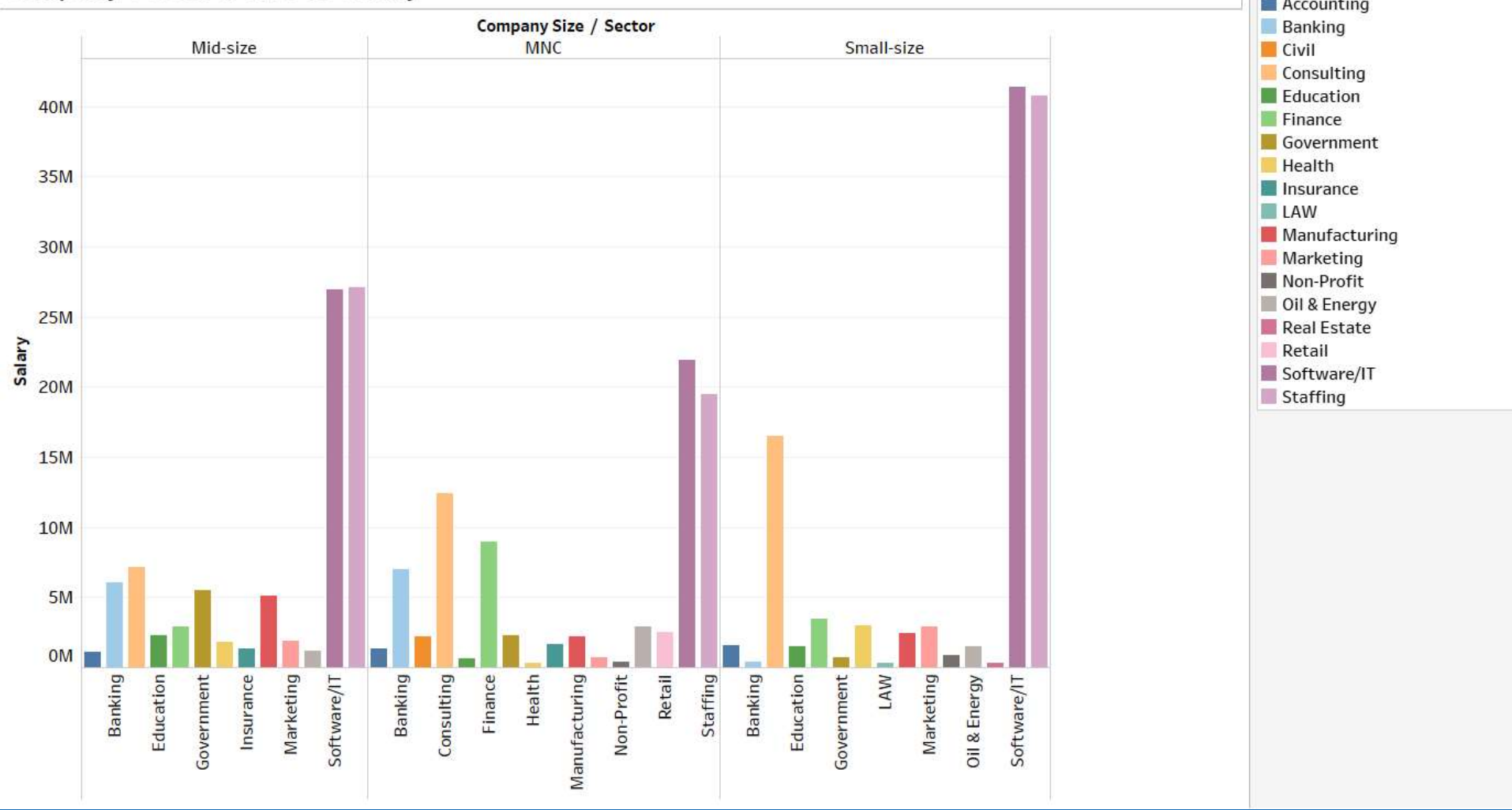
## Sector Vs Salary



### Sector

- Accounting
- Banking
- Civil
- Consulting
- Education
- Finance
- Government
- Health
- Insurance
- LAW
- Manufacturing
- Marketing
- Non-Profit
- Oil & Energy
- Real Estate
- Retail
- Software/IT
- Staffing

Company Sector & Size Vs Salary





# Company Location

In [10]: *# Mapping the Location*

```
city = {"Austin": 0, "Berkley":1, "Boston": 2, "Chicago": 3, "Cupertino": 4, "Houston": 5, "Los Angeles": 6, "Mountain View":7, "New York": 8, "Orlando": 9, "San Francisco":10, "San Jose": 11, "Seattle": 12, "Cleveland": 14, "MenloPark": 15, "OAKLAND": 16, "PaloAlto": 17, "REDWOODCITY": 18, "Sacramento": 19, "SanMateo": 21, "Sunnyvale": 22, "Washington": 23, "SantaClara": 24, "Minneapolis": 25}
data["Location"] = data["Location"].map(city)
data.head()
```

Out[10]:

	Company Name	Salary	Years of Experience	Location	Company Size	Sector	Job Title
0	3DI INC	167000.0	2	6	Small-size	Software/IT	Business Analyst
1	3DI INC	123683.0	1	6	Small-size	Software/IT	Business Analyst
2	3DI INC	80366.0	0	6	Small-size	Software/IT	Business Analyst
3	3EDGEUSAGROUP LLC	121100.0	2	8	Small-size	LAW	Business Analyst
4	3EDGEUSAGROUP LLC	101727.5	1	8	Small-size	LAW	Business Analyst

# Company Size

Out[234]:

	Company Name	Salary	Years of Experience	Location	Company Size	Sector	Job Title
0	3DI INC	167000.0	2	6	Small-size	Software/IT	Business Analyst
1	3DI INC	123683.0	1	6	Small-size	Software/IT	Business Analyst
2	3DI INC	80366.0	0	6	Small-size	Software/IT	Business Analyst
3	3EDGEUSAGROUP LLC	121100.0	2	8	Small-size	LAW	Business Analyst
4	3EDGEUSAGROUP LLC	101727.5	1	8	Small-size	LAW	Business Analyst

In [235]: *# Mapping the Company Size*

```
size = {'Small-size': 0, 'Mid-size': 1, 'MNC': 2}
data['Company Size'] = data['Company Size'].map(size)
data.head()
```

Out[235]:

	Company Name	Salary	Years of Experience	Location	Company Size	Sector	Job Title
0	3DI INC	167000.0	2	6	0	Software/IT	Business Analyst
1	3DI INC	123683.0	1	6	0	Software/IT	Business Analyst
2	3DI INC	80366.0	0	6	0	Software/IT	Business Analyst
3	3EDGEUSAGROUP LLC	121100.0	2	8	0	LAW	Business Analyst
4	3EDGEUSAGROUP LLC	101727.5	1	8	0	LAW	Business Analyst

# Sector

In [12]: *# Mapping the Sector*

```
sector = {"Oil & Energy": 0, "Accounting": 1, "Banking": 2, "Civil": 3, "Consulting":4, "Education":5,
          "Government":7, "Health": 8, "Insurance": 9, "LAW": 10, "Manufacturing": 11, "Marketing": 12,
          "Software/IT": 14, "Staffing": 15, "Real Estate": 16, "Retail": 17}
data["Sector"] = data["Sector"].map(sector)
data.head()
```

Out[12]:

	Company Name	Salary	Years of Experience	Location	Company Size	Sector	Job Title
0	3DI INC	167000.0	2	6	0	14	Business Analyst
1	3DI INC	123683.0	1	6	0	14	Business Analyst
2	3DI INC	80366.0	0	6	0	14	Business Analyst
3	3EDGEUSAGROUP LLC	121100.0	2	8	0	10	Business Analyst
4	3EDGEUSAGROUP LLC	101727.5	1	8	0	10	Business Analyst





## Split the data for training and testing

```
In [30]: feature_col = ["Location", "Years of Experience", "Company Size", "Sector" ]
X = data[feature_col]
y = data["Salary"]

xtrain, xtest, ytrain, ytest = train_test_split(X, y, test_size=0.10, random_state = 52)

train_x = xtrain
train_y = ytrain
test_x = xtest
test_y = ytest

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
train_x = sc.fit_transform(train_x)
test_x = sc.transform(test_x)
```



## Train the model

```
In [31]: regr = linear_model.LinearRegression()

# Train the model using the training sets
clf = regr.fit(train_x, train_y)

# Make predictions using the testing set
Sal_y_pred = regr.predict(test_x)

#Accuracy:
clf.score(test_x, test_y)
```

```
Out[31]: 0.70769752746380643
```

Model Accuracy : 70.76%

# Predicting Salary

## Prediction

```
In [15]: # Location: Mountain View, Yrs of exp: 4 to 6, size: MNC, Sector: Consulting  
         regr.predict([[7,1,2,4]])
```

```
Out[15]: array([ 142767.42008853])
```



# Save the model in Pickle

```
In [ ]: # Pickling the model
```

```
In [32]: import pickle

with open("python_lin_reg_model.pkl", "wb") as file_handler:
    pickle.dump(regr, file_handler)
with open("python_lin_reg_model.pkl", "rb") as file_handler:
    loaded_pickle = pickle.load(file_handler)
```

# Future Enhancements



- Currently we have data for **Business Analyst** only.
- Addition of Job titles like **Project Manager, Software Engineer, Operations Manager, Data Scientist.**
- Currently we have collected the salaries in **26 cities.**
- Addition of more cities.



Thank You!