



# Big Data Salary Prediction

With Linear Regression

# Introduction

## Objection

- Predicting annual salary for big data career by data collected from employment website.

## Goal- Create the optimal model to predict salary

- Help people in job industry to find expected salary before negotiating.
- For people who already in data related industry, they can use this model to check if they get reasonable salary



# Data Processed

**Target Variable:**

**Annual Salary**

Web Scraping tool:  
Beautifulsoup & Selenium  
Website:



## Features Variables:

About Company:

- Headquarter\_State
- Company Annual Revenue
- Company Size
- Founded year
- Industry

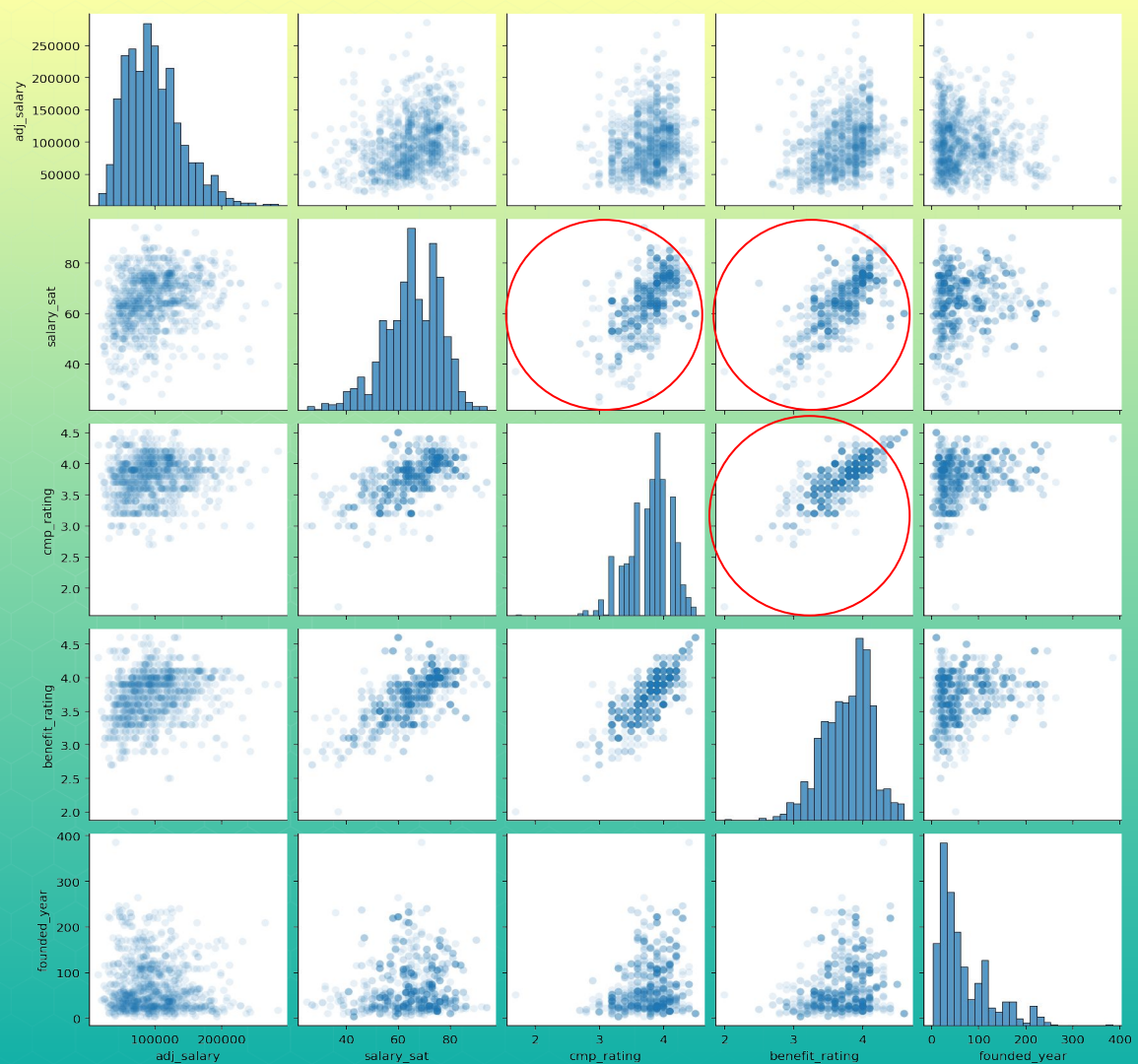
Company Review:

- Salary Satisfaction %
- Company Overall Rating
- Company Benefit Rating

Job:

- Job title
- Interview Process Length
- Interview Difficulty
- Payment Type

# Pairs Plot for Numerical Variables



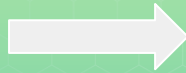
# Preliminary Linear Regression

**Linear Regression** with all features variables

$$R^2 = 0.5184$$

$$R^2 = 0.4831$$

$$R^2 = 0.5603$$



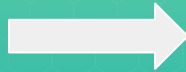
$$\text{Avg } R^2 = \mathbf{0.5206}$$

**Polynomial Regression** with all features variables

$$R^2 = 0.5252$$

$$R^2 = 0.4886$$

$$R^2 = 0.5699$$

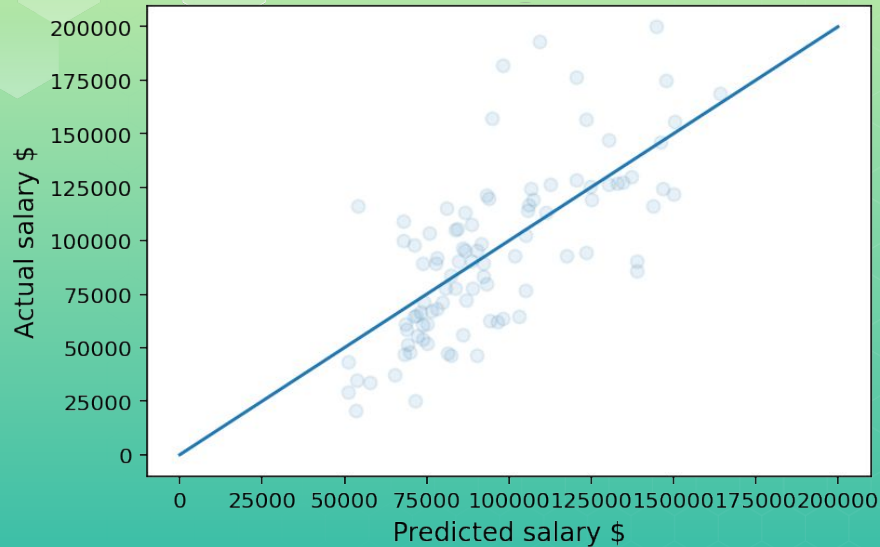


$$\text{Avg } R^2 = \mathbf{0.5279}$$

# LASSO Regression VS. Ridge Regression

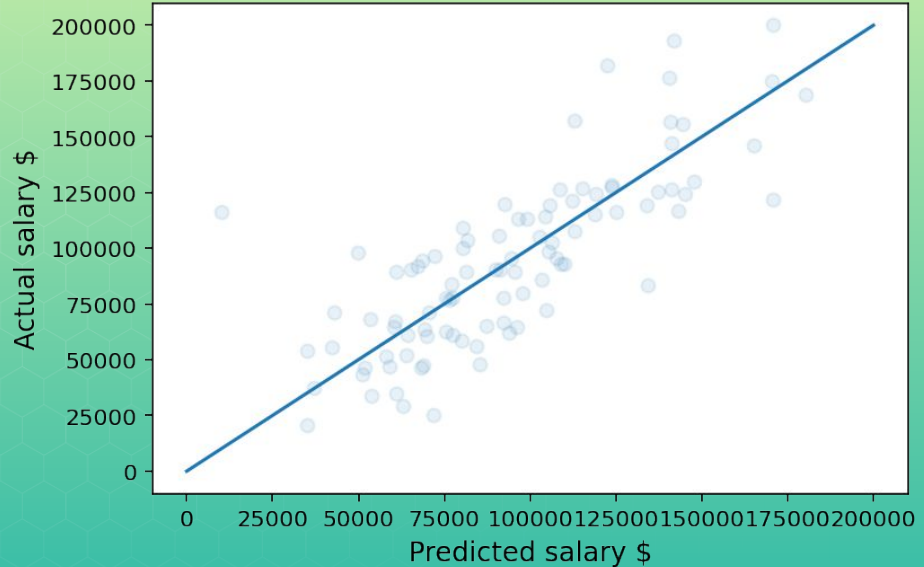
## Regularization

### Lasso Regression



**$R^2 = 0.51$**   
**MAE= 22,105.59**

### Ridge Regression



**$R^2 = 0.61$**   
**MAE= 18,331.80**

# LASSO Regression

## Variable Selection

Variables	Coefficient
title.analyst salary_sat	-411.457
title.analyst founded_year	-16.27
title.senior salary_sat	157.12
title.senior founded_year	3.40
title.scientist salary_sat	354.76
Headquarters_State_DC founded_year	-1.33
Headquarters_State_Other salary_sat	-10.44
revenue_\$1B to \$5B (USD) founded_year	59.50
revenue_more than \$10B (USD) salary_sat	7.82

Variables	Coefficient
inter_len_About a month founded_year	9.84
inter_len_About two weeks salary_sat	63.99
inter_len_More than one month founded_year	-25.67
salary_unit_per year salary_sat	400.26
Energy, Mining & Utilities founded_year	13.75
Government & Public Administration salary_sat	-8.95
salary_sat^2	5.243
salary_sat founded_year	-0.49
founded_year^2	-0.02

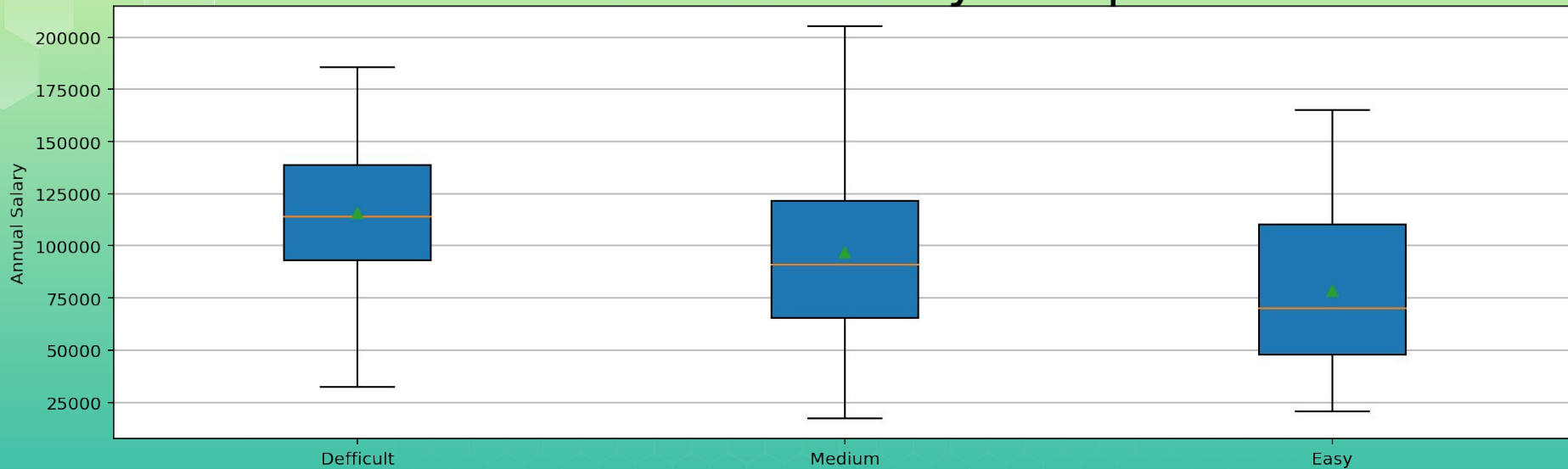


# Remove Non-Statistically Significant Variables to optimal Ridge Model

Removed Variable	Current R <sup>2</sup>	Updated R <sup>2</sup>	R <sup>2</sup> Improvement
Company Size	0.61	0.64	0.03
Company Overall rating	0.64	0.65	0.01
Benefit Rating	0.65	0.66	0.01
<b>Interview Difficulty</b>	0.66	0.55	<b>-0.11</b>



# Interview Difficulty box plot



# Final Model & Conclusions



## The optimal model:

- Ridge Regression Model



## Including features:

- Headquarter\_State
- Company Annual Revenue
- Founded year
- Salary Satisfaction %
- Job title (scientist, analyst, senior)
- Interview Process Length
- Interview Difficulty
- Payment Type



**R<sup>2</sup>:** 0.666



**Mean Absolute Error:** 16831.94

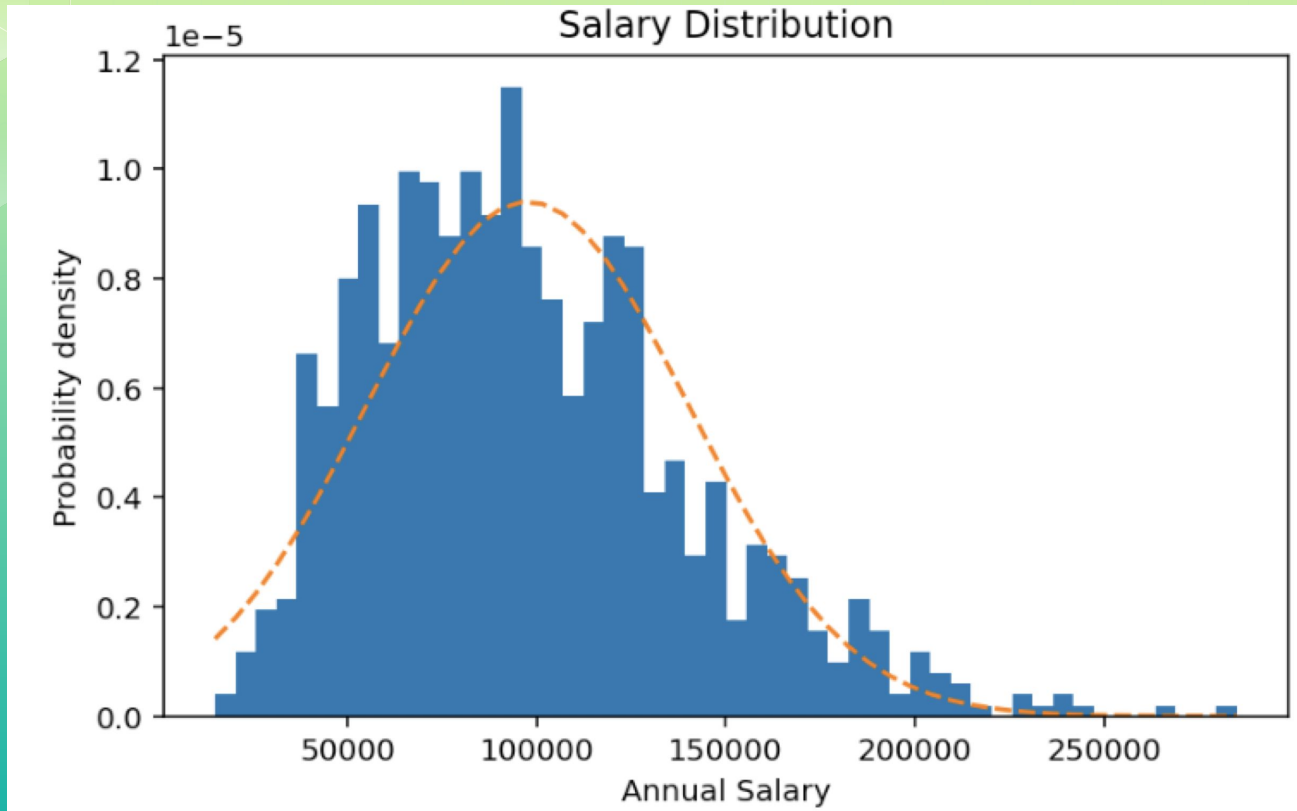


**THANK YOU**

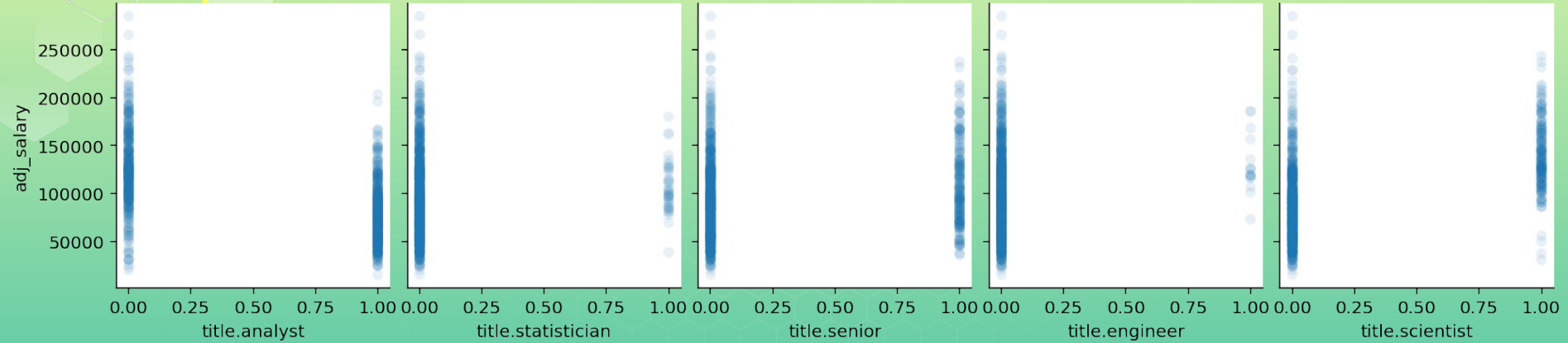


# APPENDIX

# Distribution of Annual Salary

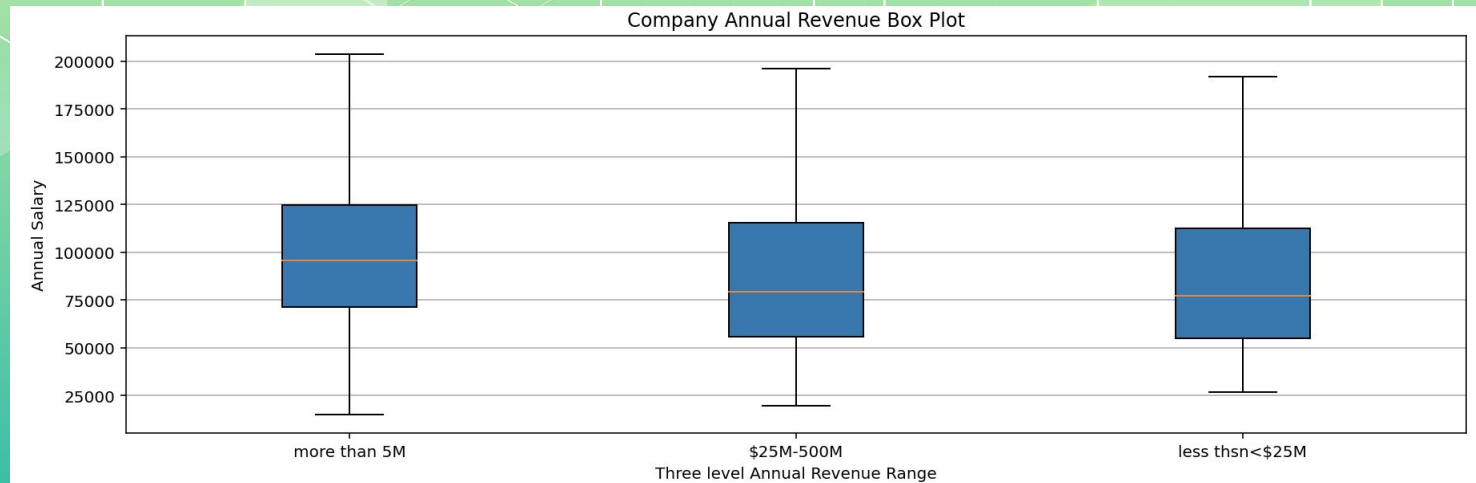


# Pairs Plot for Title Variables



**Remove title.statistician and title.engineer from ridge model:  
R<sup>2</sup> score increase 0.004 to 0.666.**

# BoxPlot of Company Annual Revenue

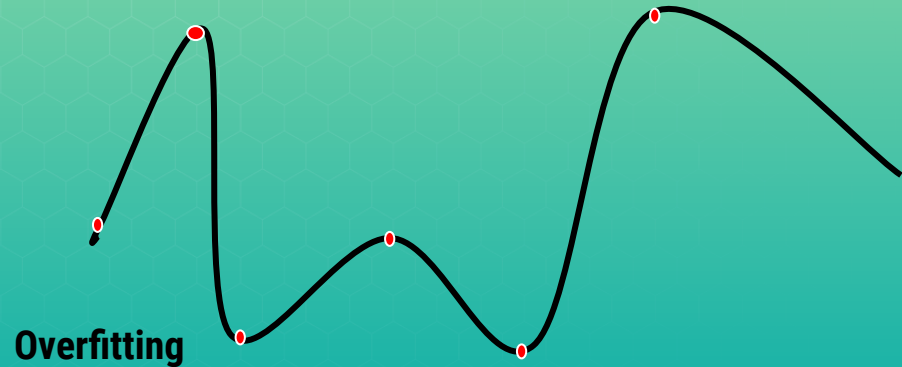




# Polynomial Regression

$R^2$  for training data is **0.916**

$R^2$  for testing data is **- 3.3**



('title.senior', -5762.426346233193),  
(**'title.scientist', 985.8822957477951**),  
(**'Headquarters\_State\_CT', 6188.745788723318**),  
( 'Headquarters\_State\_DC', 202.08306782503792),  
( 'Headquarters\_State\_FL', 1117.9844778638997),  
( 'Headquarters\_State\_GA', -220.64875407813315),  
( 'Headquarters\_State\_IL', -921.3283836144151),  
( 'Headquarters\_State\_MA', -754.7972169554096),  
( 'Headquarters\_State\_MD', -526.2022470509401),  
( 'Headquarters\_State\_MN', 641.3332858633435),  
( 'Headquarters\_State\_NC', -470.1012289344508),  
( 'Headquarters\_State\_NJ', 682.1501291217144),  
( 'Headquarters\_State\_NY', -818.0260283338703),  
( 'Headquarters\_State\_Other', -251.65588938198744),  
( 'Headquarters\_State\_TN', 167.62429886666996),  
( 'Headquarters\_State\_TX', -30.94786531332454),  
( 'Headquarters\_State\_VA', 428.4371344457495),  
( 'Headquarters\_State\_WA', 1849.0129016708506),  
( 'Headquarters\_State\_WI', 770.5408505907744),  
( 'revenue\_\$1B to \$5B (USD)', 314.7564727284416),  
( 'revenue\_\$1M to \$5M (USD)', -3360.7248734854857),  
( 'revenue\_\$25M to \$100M (USD)', -27.73111185644302),  
( 'revenue\_\$500M to \$1B (USD)', 619.3579974454169),  
( 'revenue\_\$5B to \$10B (USD)', 4664.067912251381),  
( 'revenue\_\$5M to \$25M (USD)', -1884.482517891447),  
( 'revenue\_less than \$1M (USD)', 144.22199985563157),  
( 'revenue\_more than \$10B (USD)', (-460.11154569072164)

( 'inter\_diff\_Easy', 322.4279199996231),  
( 'inter\_diff\_Medium', -1078.3239247232077),  
( 'inter\_diff\_none', 1735.9507493741849),  
( 'inter\_len\_About a month', -712.5190726615895),  
( 'inter\_len\_About a week', 1977.6141850854697),  
( 'inter\_len\_About two weeks', 656.94997482089),  
( 'inter\_len\_More than one month', 1168.941275478176),  
( 'inter\_len\_none', -1703.1611623214235),  
( 'salary\_unit\_per month', 1644.850613584671),  
( 'salary\_unit\_per year', -230.91775638753856),  
(**'Aerospace & Defense', 2268.9577127856664**),  
( 'Education', -183.2852442346466),  
( 'Energy, Mining & Utilities', 19.364111537963936),  
( 'Financial Services', -42.16054278952453),  
( 'Government & Public Administration', 332.03296850658626),  
( 'Healthcare', 260.42133738771645),  
( 'Human Resources & Staffing', -1488.6333514954922),  
( 'Information Technology', 591.7546380215108),  
( 'Insurance', -843.4599369498455),  
(**'Manufacturing', 2164.3388443918657**),  
( 'Pharmaceutical & Biotechnology', 23.213216077430843),  
( 'Retail & Wholesale', -1467.122205810767),  
( 'Telecommunications', -697.4622294364681),  
( 'salary\_sat', -445.44321887897524),  
( 'founded\_year', -136.9720554654632)

# Further Improvements

- **Add more feature variables: work experience, skill requirement**
- **Collect more observations**