

# Income estimation

Based on skills

Can we calculate salaries for jobs in any industry based only on skills and location?

# Datasets

Salary data from Bureau of Labor Statistics:

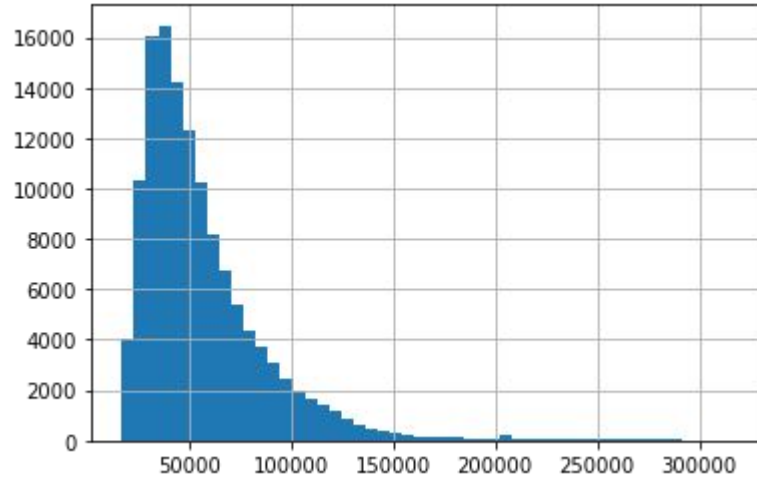
- Occupation
- State
- Annual Salary

Skills data from Onet

There are two scales, importance and level

For this exercise we choose importance with a scale of 1-5

# Target feature



	Annual	Hourly
Mean	57,448	27.34
Min	17,300	8.32
Max	315830	151.84

# Features

Features

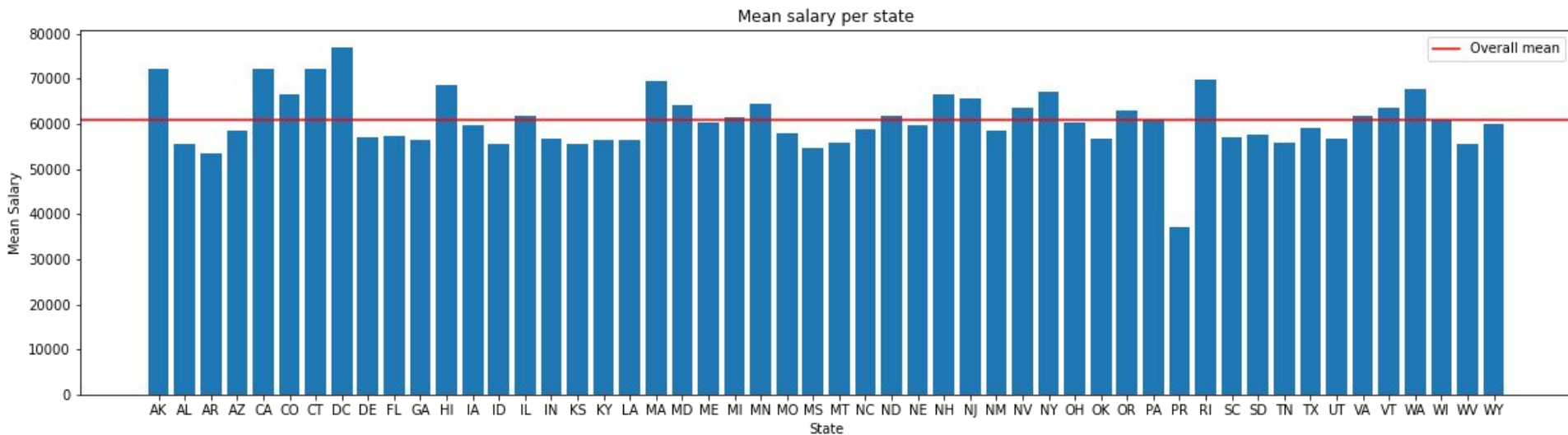
states

788 occupations

35 skills

379 locations

# States



# Skills scale

1 to 5

Just looking at these histograms, we can see that Installation, Equipment Maintenance, Equipment Selection, and Repairing are highly skewed to the right. Most of their values are concentrated around 1.

Most of the distributions above are bimodal. Time Management and Service Orientation have the most normal distributions. Technology design and Programming have most of their values concentrated below 2. Science has a high value at 1 which means that science is not ranked as important for many occupations! We also see this features in Troubleshooting, Programming and Operations and Control.

**Relationship between target feature (salary) and variables (skills)**

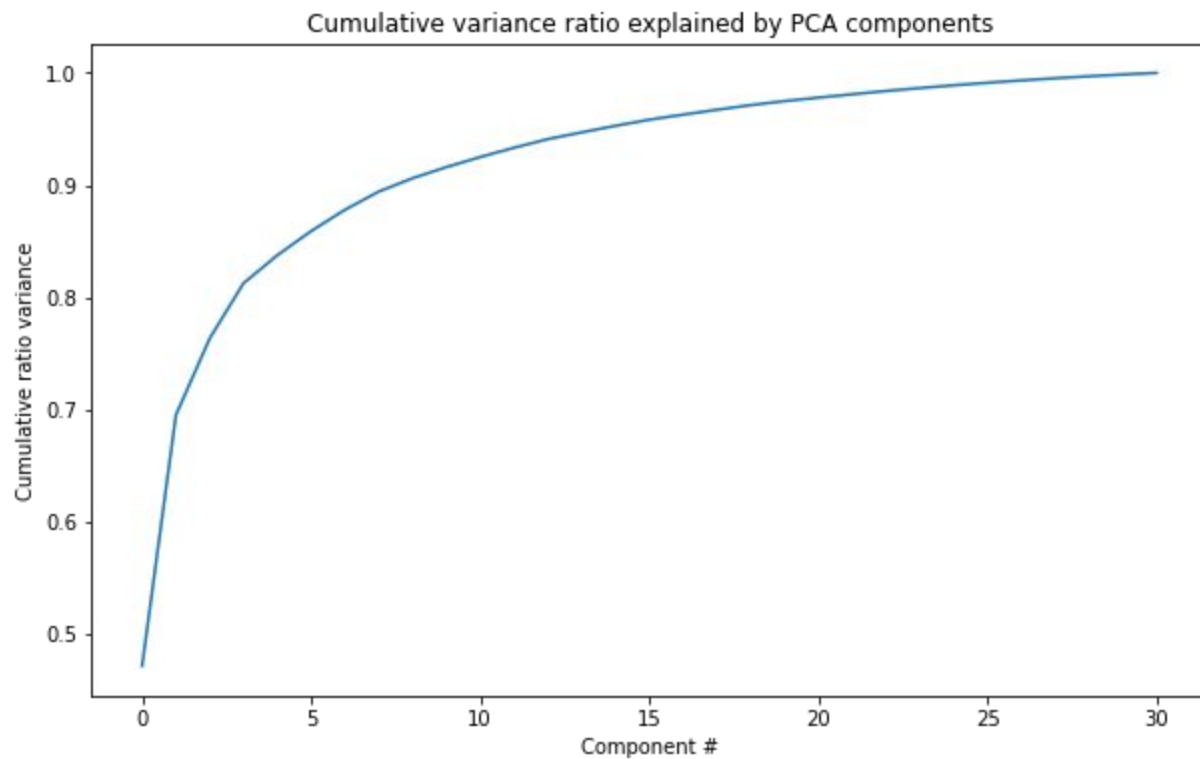


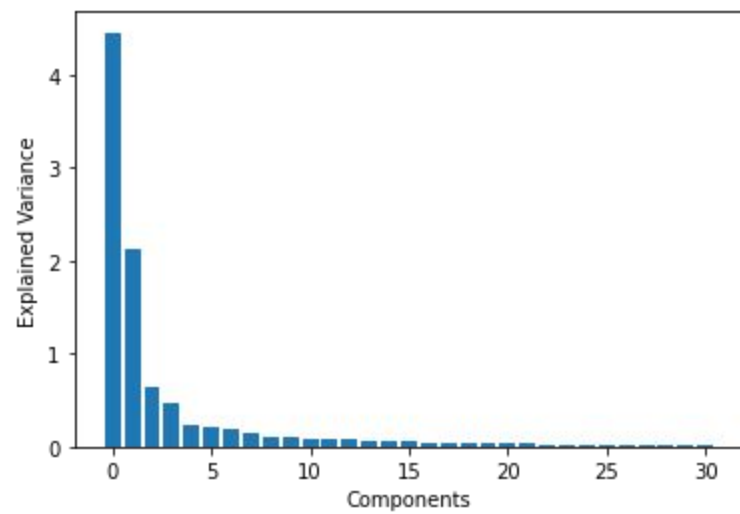
**Relationship between target feature (salary) and variables (skills)**

## **Correlation between each feature**

Heatmaps, pairplots

# PCA





# Baseline model

Uses mean as the best estimate

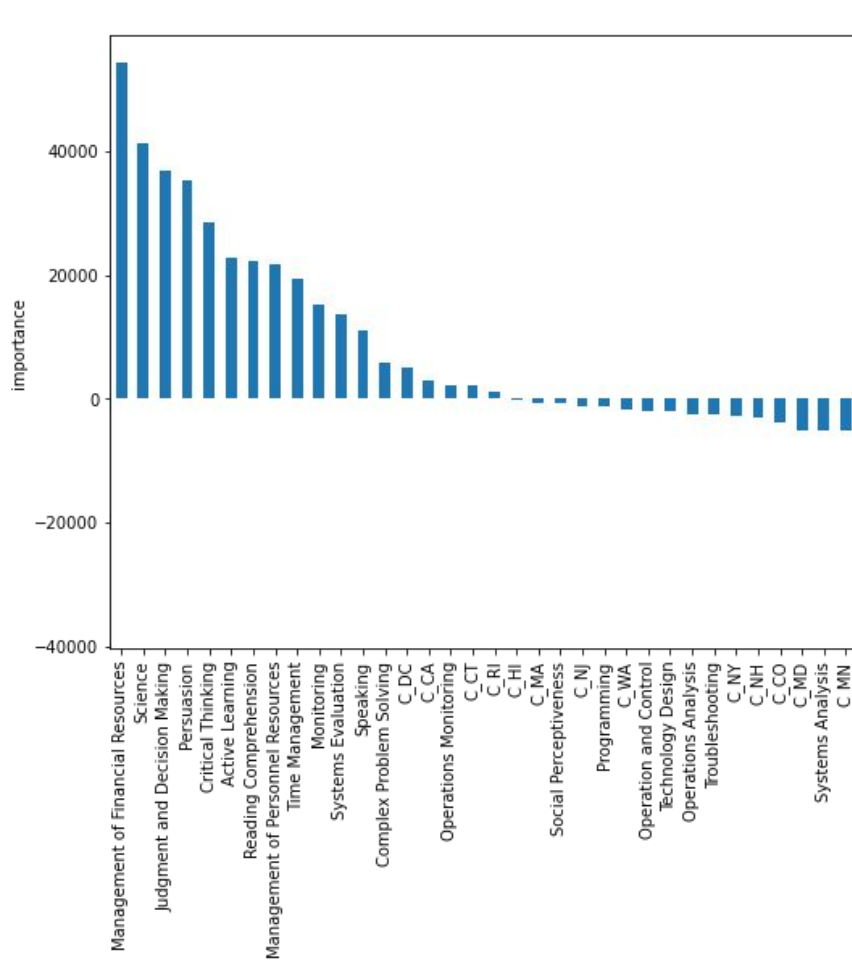
Get basically zero for  $R^2$  error: the model does a terrible job

Fit time = 0.5

Score time = 0.08

# Regression Models

	<b>OLS</b>	<b>Ridge</b>	<b>SVR</b>
<b>R-squared</b>	0.606	0.606	0.586
<b>Fit time</b>	12.88	0.93	91.27
<b>Score time</b>	0.58	0.066	0.0198
<b>MAE test</b>	\$14,251		



## Basic Linear Regressions

Management of Financial Resources  
Science

Judgment and Decision Making

Persuasion

Critical Thinking

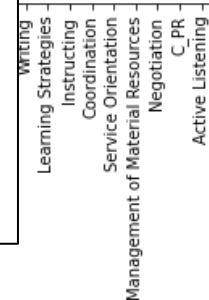
Active Learning

Reading Comprehension

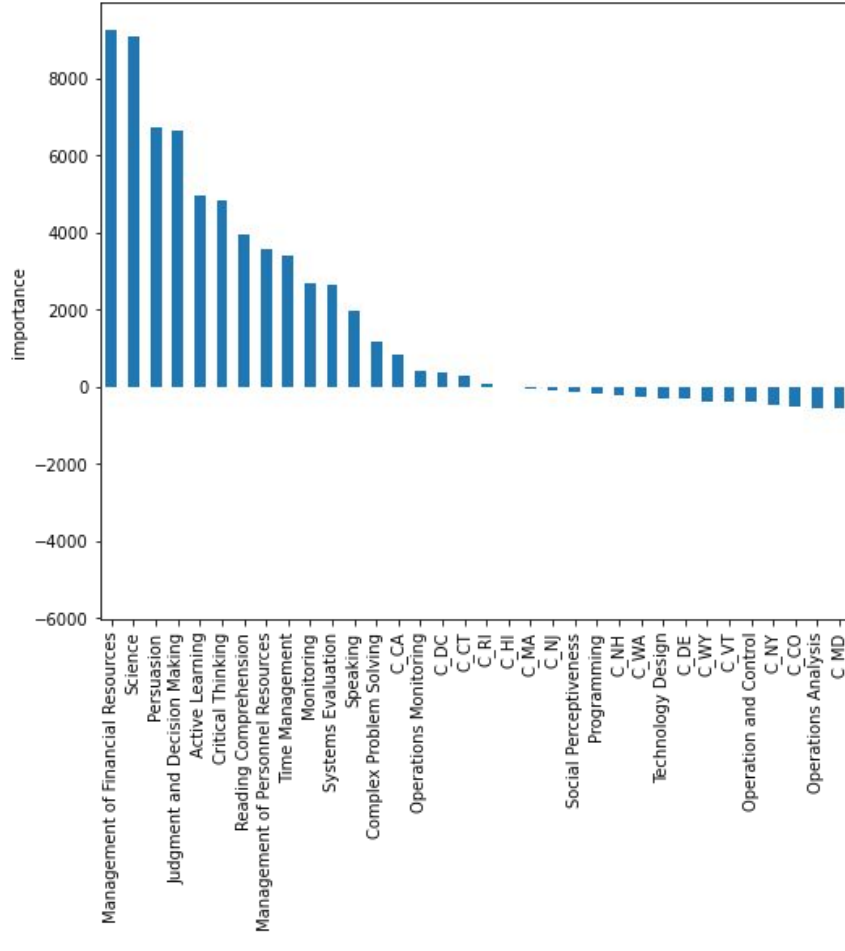
Management of Personnel Resources

Time Management

Monitoring



Feature importances



## Ridge Regression

Management of Financial Resources

Science

Persuasion

Judgment and Decision Making

Active Learning

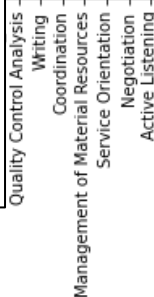
Critical Thinking

Reading Comprehension

Management of Personnel Resources

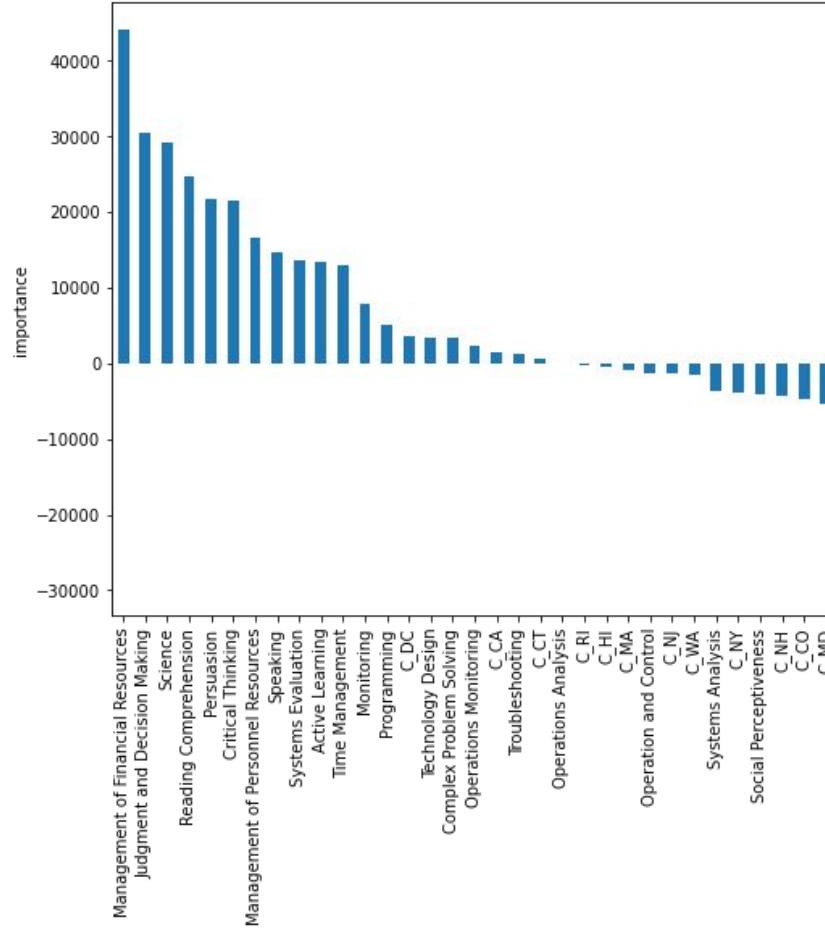
Time Management

Monitoring





Feature importances



## Support Vector Regression(LinearSVR)

Management of Financial Resources

Judgment and Decision Making

Science

Reading Comprehension

Persuasion

Critical Thinking

Management of Personnel Resources

Speaking

Systems Evaluation

Active Learning

Management

Service Orientation  
Active Listening  
C\_PR

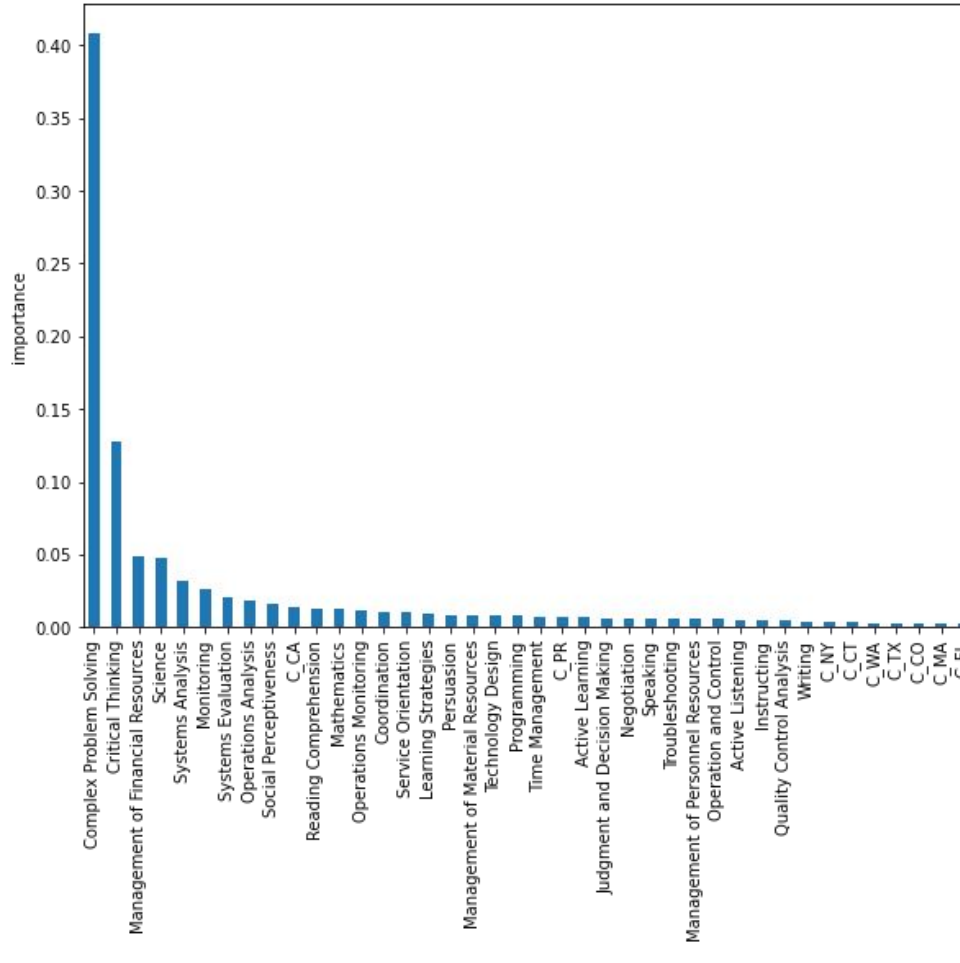
# Random Forest Model

(0.8875267860497988, 0.0031824517006334503)

Fit time: 92.8323

Score time: 0.5231

Best random forest regressor feature importances



## Random Forest Model

Complex Problem Solving

Critical Thinking

Management of Financial Resources

Science

Systems Analysis

Monitoring

Systems Evaluation

Operations Analysis

Social Perceptiveness

C\_CA

C\_IL  
C\_MT  
C\_VT  
C\_WY  
C\_DE

features