



# MSBANA DATA SCIENCE JOB RECOMMENDER

*By –*

*Vaishali Pawar*

*Rohan Sharma*

*Manoj Tomar*

*Aditya Nawal*

# AGENDA


---

- Problem Definition – What are we trying to solve
- Data description and Feature Engineering
- Solution Approach
  - Personalized job Finder
  - Skills Recommender
- Walkthrough of the Power BI Dashboard
- Overall Insights and Recommendations



Problem Definition

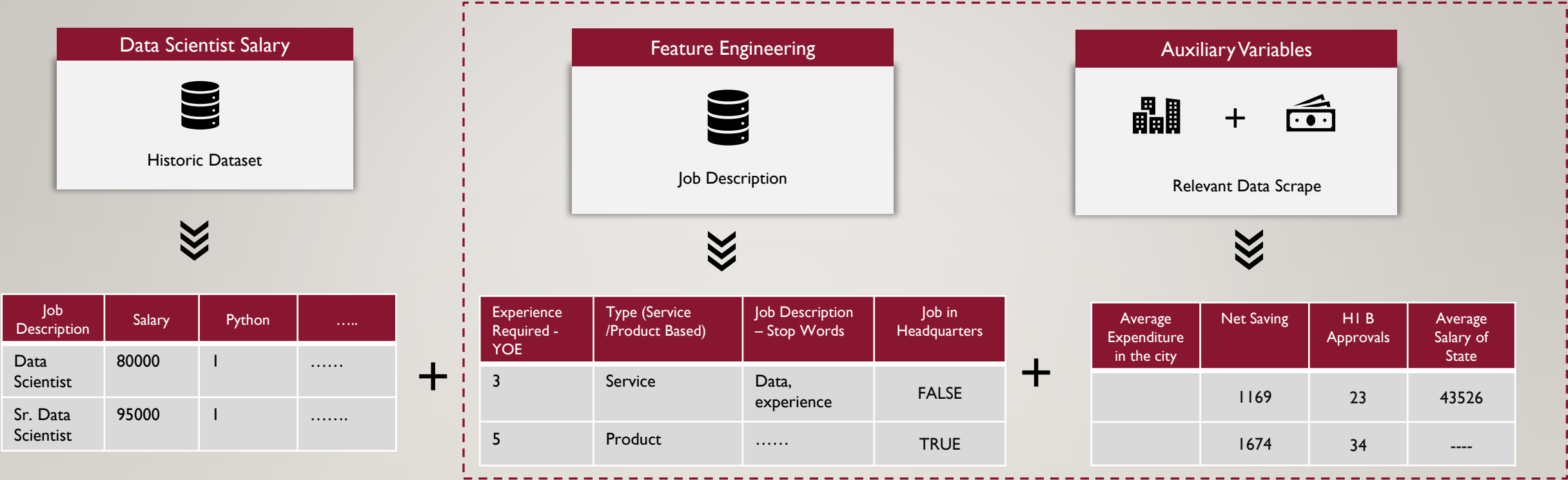
MSBANA Data Science Job Recommender Would help students secure Data Scientist Jobs

| Current State  | Gap  | Desired Future State  |
|--|--|---|
| <ul style="list-style-type: none"><li>MS Business Analytics students find it difficult to secure Data Scientist jobs upon completion of the program</li><li>The OBAIS Department is of the opinion that - It would be helpful for the students to know beforehand about the tools needed, availability of opportunities and job titles, and the salaries associated with the roles</li><li>Understanding what would be best for students to pursue and what are the best suited job profiles, is important in their job search</li></ul> | <p>Identify the “best” Data Science jobs and have an understanding about the skills required for these jobs</p>  | <ul style="list-style-type: none"><li>The Power BI dashboard helps the MSBANA students to find Data Science jobs based on their current skillset, and recommends additional courses that can be helpful to land desired jobs</li><li>Creating a similarity matrix on existing and auxiliary features gives best suited jobs for a candidate and suggests similar jobs he/she can apply to</li><li>Feature engineering and adding auxiliary variables to the existing dataset helps us generate insights on the popular tools and technologies required to acquire Data Science jobs</li></ul> |
|  | Key Question   |   |
|  | <p>What are the jobs in Data Science domain and what is the average skillset required for these jobs</p>   |   |

## Solution Approach | Why not Clustering or Regression?

- The Data Scientist Salary data had ~765 records, out of which we found that there were duplicates in the data. The Final row count after data cleaning was 469
- For any statistical or machine learning algorithm to give a good accuracy, we need a good amount of training data. With a smaller dataset, the recommendations that the model would suggest would be biased
- Moreover, the problem statement background is to help MSBANA students secure Data Scientist jobs. Our solution approach is more focused on what the students should or should not do, to land Data Scientist jobs
- While developing this tool and analyzing the problem statement, we found very interesting insights that would be helpful for us in our job hunt
- This tool is scalable and with more historical data, we would be able to give personalized recommendations for all the MSBANA students to help them streamline their job search process

# Data Description | Engineering Features and adding auxiliary variables helped in finding meaningful insights



Final Dataset

| Job description    | Salary | Python | Years of Experience | Type (Service/Product based) | Job Description – Stop Words | Median Rental Income | Average Salary of state | Standardized Salary | FICO Score | Job in Headquarter |
|--------------------|--------|--------|---------------------|------------------------------|------------------------------|----------------------|-------------------------|---------------------|------------|--------------------|
| Data Scientist     | 80000  | 1      | 3                   | Service                      | Data, experience             | 1169                 | 93156                   |                     | 717        | FALSE              |
| Sr. Data Scientist | 95000  | 1      | 5                   | Product                      | .....                        | 1674                 | 89139                   |                     | 734        | TRUE               |

Metadata →

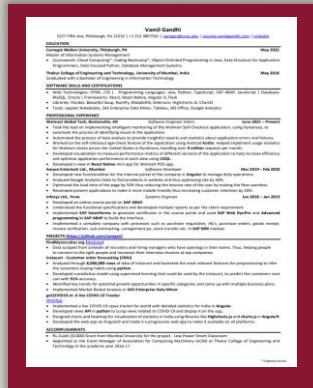


**Data Description** | **Definitions of Engineered and Auxiliary Variables**

| Auxiliary/ Engineered Features    Definition |   |
|--|---|
| Experience Required -YOE                     | Years of Experience required for the job  |
| Type (Service /Product Based)                | Is the Company Service based or Product based   |
| Job Description – Stop Words                 | Relevant information after removing Connecting words  |
| Job in Headquarters                          | If Job Location is the Headquarter of the Company   |
| Net Saving                                   | Average salary of job – Average cost of Living for the city   |
| HI Approvals                                 | Number of HI approvals for the company – Important field to consider for the International Students |
| Average Salary of State                      | Average salary of the state   |

# Solution Approach | Our Approach helps students land best jobs and also improve their skillset

## Sample Resume



What are the best 10 jobs you should apply, based on current skillset?



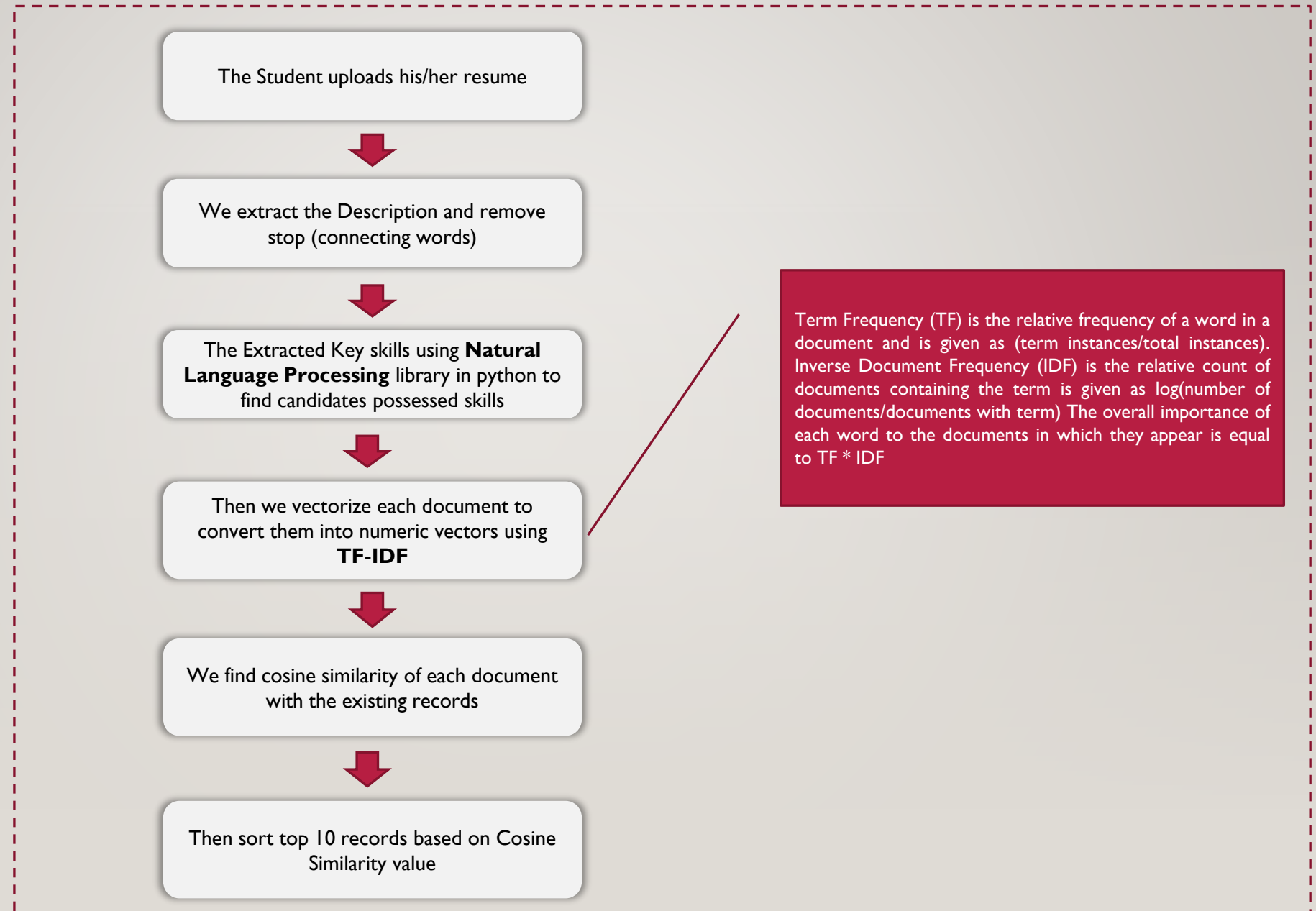
- The student uploads his/her resume to the app
- The algorithm scans the resume and extracts important tools and technologies mentioned in the resume
- Then the extracted tools and technologies are matched with the records in the database
- Based on the matches found, the algorithm suggests 10 'best' jobs that the student can apply for

What courses you should take, to boost your skillset?



- The 2<sup>nd</sup> part of the algorithm, compares the extracted tools and technologies from the resume with the popular tools and technologies required for Data Science jobs
- The tools and technologies that are missing in the resume are highlighted
- The algorithm also suggests the BANA and Data Camp courses the student should take to boost their skills which increases their chances to land Data Science jobs

## Solution Approach | Step 1 - How do we suggest 10 best jobs you should apply ?





## Solution Approach | Step 2 - What courses you should take, to boost your skillset?



The Student uploads his/her resume



We extract the Description and remove stop (connecting words)



The Extracted Key skills using **Natural Language Processing** library in python to find candidates possessed skills



Same as Step 1

Compare skills required in the best jobs recommended and skills from the Student's resume

Create a flag –

1. **Required** – if the skill is missing in Resume
2. **Not Required** – If the skill is in Resume but not required for the job
3. **Satisfied** – If the skill is required for the job and is present in the student's resume

Recommend courses from BANA curriculum

# Insights | Streamlining Insights for groups would help students find best suited jobs

## Product Companies that need 0 to 3 YOE

Words to use: Data, Learning, Team, Support, Scientist  
Best Possible Net Saving PA: 156K  
Best Job : Lead DE @Credit Sesame  
Top3 Skills: Python, SQL, Excel

## Product Companies that need 0 to 3 YOE & Sponsor H1B Visas

Words to use: Data, Learning, Complex, Skills, Team  
Best Possible Net Saving PA: 124K  
Best Jobs: Data Scientist @ Pfizer  
Top3 Skills: Python, SQL, Excel



## Product Companies that need >5 YOE

Words to use: Data, Learning, Team, Support, Scientist  
Best Possible Net Saving PA: 220K  
Best Job : Director Data Science @ Liberty Mutual  
Top3 Skills: Python, SQL, Excel

## Product Companies that need >5 YOE & Sponsor H1B Visas

Words to use: Data, Development, Solutions, Science  
Best Possible Net Saving PA: 220K  
Best Jobs: Director Data Science @ Liberty Mutual  
Top3 Skills: Python, SQL, Excel

## Insights | Streamlining Insights for groups would help students find best suited jobs

### All Companies that need 0 to 3 YOE

Words to use: Data, Learning, Team, Support, Models  
Best Possible Net Saving PA: 156K  
Best Job : Lead DE @Credit Sesame  
Top3 Skills: Python, SQL, Excel

### All Companies that need 0 to 3 YOE & Sponsor H1B Visas

Words to use: Data, Learning, Statistical, Support  
Best Possible Net Saving PA: 147K  
Best Jobs: Lead Data Scientist @ Visa  
Top3 Skills: Python, SQL, Excel



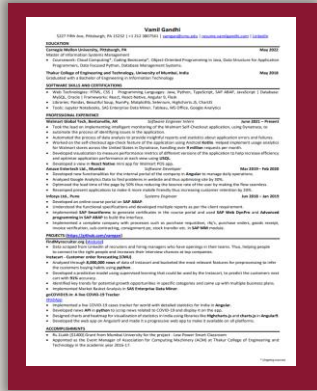
### All Companies that need >5 YOE

Words to use: Data, Learning, Team, Support, Scientist  
Best Possible Net Saving PA: \$ 220K  
Best Job : Director Data Science @ Liberty Mutual  
Top3 Skills: Python, SQL, Excel

### All Companies that need >5 YOE & Sponsor H1B Visas

Words to use: Data, Learning, Complex, Skills, Team  
Best Possible Net Saving PA: 220K  
Director Data Science @ Liberty Mutual  
Top3 Skills: Python, SQL, Excel

# Insights | Manoj needs to work on Scikit, AWS and SAAS to support his candidature for top 5 jobs



### Resume Analysis

**Filters**

Job Title: All

Type: ☐ Product Based ☐ Service Based

Years: 0 to 10

Avg Salary(K): 15.50 to 254.00

H1 B Approvals: 0 to 889

**Top 10 Recommended Jobs to Apply - For Manoj**

| Job Title                           | Job Location | Type          | Estimated Net Savings | Similarity % | Tableau      | Tensor       | SQL          | Spark        | scikit       |
|-------------------------------------|--------------|---------------|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Wave Partners Senior Data Scientist | IL           | Service Based | 95,190.16             | 0.12         | Satisfied    | Not required | Satisfied    | Not required | Not required |
| Wave Partners Data Scientist        | IL           | Service Based | 47,690.16             | 0.12         | Satisfied    | Not required | Satisfied    | Not required | Not required |
| Residential Data Engineer           | IL           | Service Based | 59,190.16             | 0.11         | Not required | Not required | Satisfied    | Required     | Not required |
| Point Lead Data Scientist           | OH           | Product Based | 69,252.46             | 0.11         | Not required | Not required | Satisfied    | Not required | Required     |
| de International Data Scientist     | CA           | Product Based | 60,713.15             | 0.10         | Satisfied    | Not required | Not required | Not required | Not required |
| Residential Senior Data Scientist   | IL           | Service Based | 96,690.16             | 0.10         | Satisfied    | Not required | Satisfied    | Required     | Required     |
| e Data Scientist - Research         | MA           | Service Based | 12,375.28             | 0.10         | Satisfied    | Not required | Satisfied    | Not required | Not required |
| ady Data Scientist                  | NY           | Service Based | 72,406.00             | 0.10         | Satisfied    | Not required | Satisfied    | Required     | Not required |
| eger Group Senior Data Scientist    | TX           | Service Based | 77,222.59             | 0.10         | Satisfied    | Not required | Not required | Not required | Not required |
| qma Data Scientist - Alpha Insights | NY           | Product Based | 129,406.00            | 0.10         | Not required | Not required | Not required | Not required | Not required |

**Skills Missing From Resume - For Top Jobs**

| Keras        | Python    | pytorch      | scikit       | Spark        | SQL          | Tableau      | Tensor       | Hadoop       | AWS          | Mongo        | SAAS         |
|--------------|-----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Not required | Satisfied | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Required     |
| Not required | Satisfied | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Not required | Not required |
| Not required | Satisfied | Not required | Required     | Not required | Satisfied    | Not required | Not required | Not required | Not required | Not required | Not required |
| Not required | Satisfied | Not required | Not required | Not required | Not required | Satisfied    | Not required | Not required | Required     | Not required | Not required |
| Not required | Satisfied | Not required | Not required | Not required | Satisfied    | Not required | Not required | Not required | Not required | Not required | Not required |

Python SQL Excel Google Ana. Spark Tensor Scikit SAAS

AWS Hadoop pytorch Keras Mongo Tableau BI Flink

Source :Kaggle (Scraped Glassdoor) & Manoj\_resume.pdf

Update Frequency : Data Dump

Data is available till 2019

❖ Hovering over the buttons shows the target BANA course / Data Camp course to gain particular skill, users can click the buttons in the dashboard to be redirected to the course links.

## References for Auxiliary Variables

- H1B data : US CIS Gov. website: <https://www.uscis.gov/>
- <https://wise.com/us/blog/cost-of-living-in-the-usa>
- <https://advisorsmith.com/data/coli/>
- <https://meric.mo.gov/data/cost-living-data-series>