# Yelp Restaurants Recommendation Engine

**---MSCA 31008 Data Mining Principles(Winter 2022)**

Ryan Liao
Jason Lee
Norah Zhang
Milan Toolsidas
Minglun Pan

# Agenda

Executive Summary → Data ETL → EDA/Feature Engineering → Model Mining → Model Evaluation
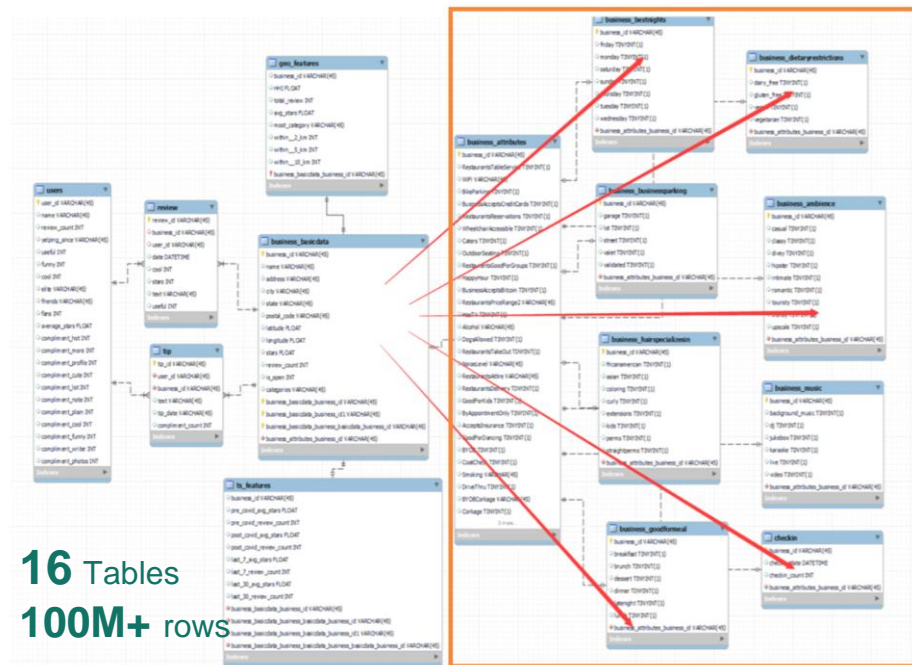
# Executive Summary

- To help users more quickly identify restaurants they would enjoy, our team created a recommendation engine to recommend restaurants based on Yelp's open dataset.

- In this project, we mainly focused on restaurants:

  - Extracted and transformed semi-structured JSON data into structured database

  - Loaded the data onto Research Computing Center, managed and connected to Python through Hive and Spark

  - Leveraged NLP, time series, and geographical analysis to engineer valuable features to add to our dataset

  - Reduced dimensionality through SVD, and generated restaurant recommendations using collaborative filtering, linear models, tree-based models and their ensembles prompted by Facebook.

- The ensemble model promoted by Facebook outperforms others, but the traditional Collaborative Filtering has its own advantage in interpretability.
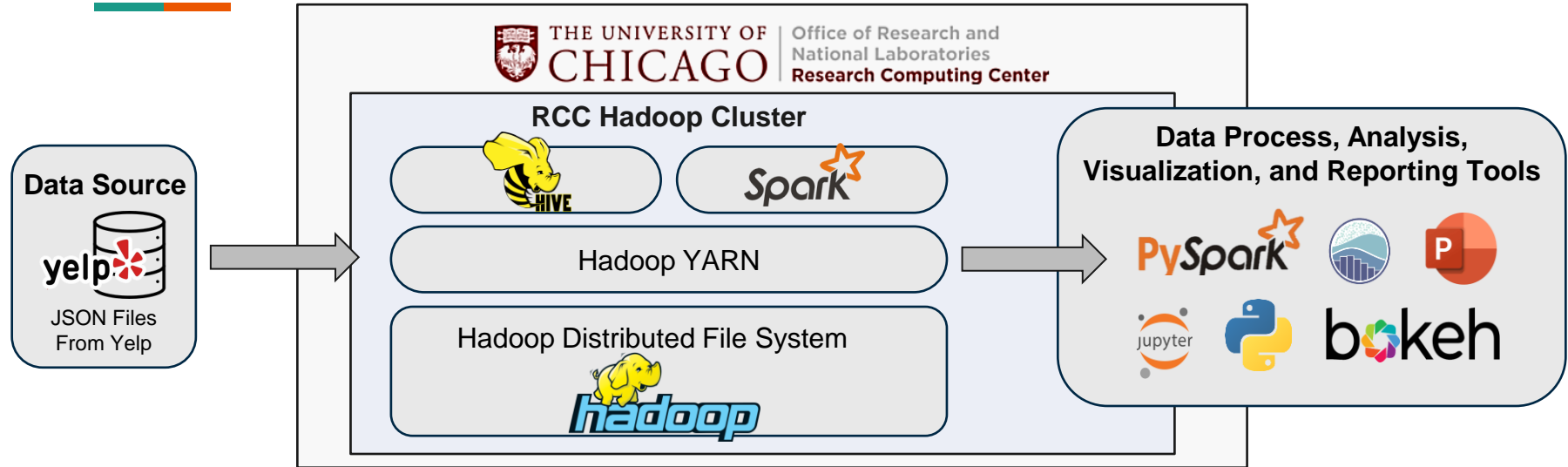
# Database Overview

**Data Extraction:** Source JSON files from Yelp

| | Features | Rows | Complexity |
|---|---|---|---|
| **Business** | 14 | 64K | more than 30 nested variables |
| **Users** | 22 | 2.2M | **3.7 GB** |
| **Review** | 9 | 8.9M | **7 GB** |
| **Tip** | 5 | 938K | |
| **Check in** | 2 | 62K | |

Business Attributes

**16** Tables
**100M+** rows



| Executive Summary | Data ETL | Feature Engineering | Model | Model Evaluation |

# Database Architecture



- **Database Selection:** Considering ease of access and cost, we chose RCC over GCP
- **Data Transformation and Loading:**
  - Steps: Clean, normalize, feature select and upload
  - We created twelve tables focused on restaurants and uploaded in a database in Hadoop

# EDA - NLP

To gain insight into how users could use different words to describe restaurants, here is a simple word cloud Venn diagram showcasing words used in the reviews.

- Blue words are only used in five-star reviews.
- Red words are only used in one-star reviews.
- Gray words in the middle are used in all reviews.



Good Review Words
Bad Review Words
Common Words

# Feature Engineering - NLP

**Goal:** Extract insights from users' reviews to better understand user preferences
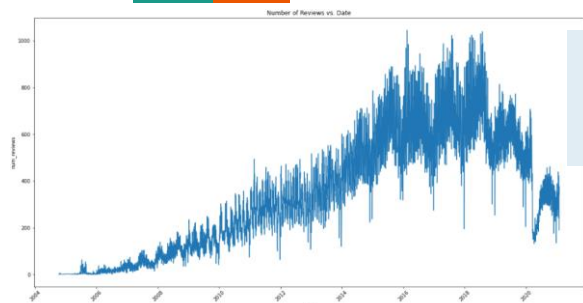
**Methods:**

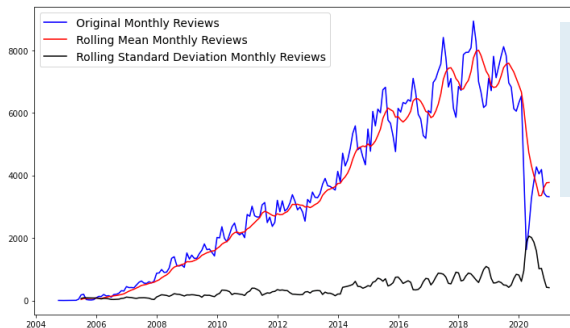| TF- IDF | BERT |
|---|---|
| • TF-IDF is a product of **Term Frequency** and **Inverse Document Frequency**<br>• TD-IDF evaluates how relevant a word is to a document in a collection of documents | • BERT is **pre-trained bidirectionally trained language model** that extract high quality language features from text<br>• BERT could be further **fine-tuned** for specific applications |

- Generated over 20 GB of features from TF-IDF and BERT
- Quasi Map-Reduce process
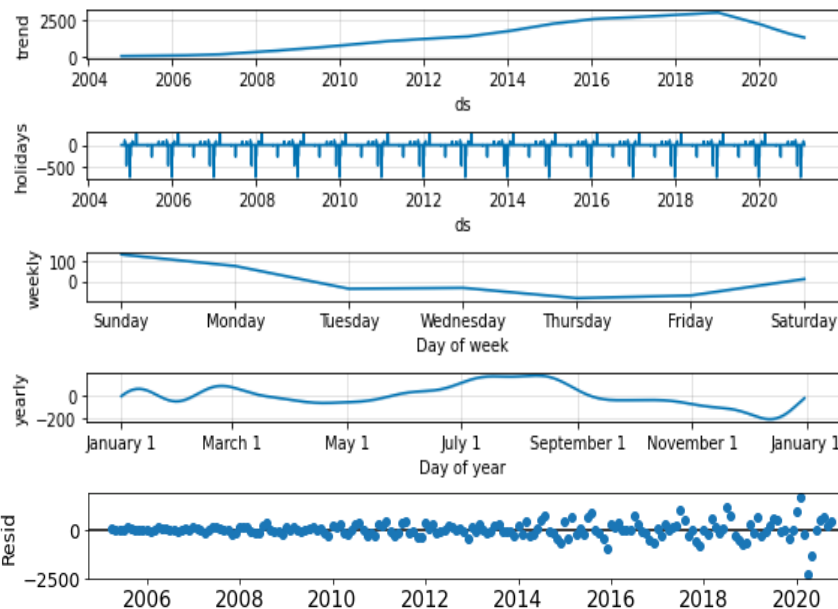- Combined the sparse matrix output into final model

# EDA - Time Series



Preliminary view of review counts show cyclic pattern
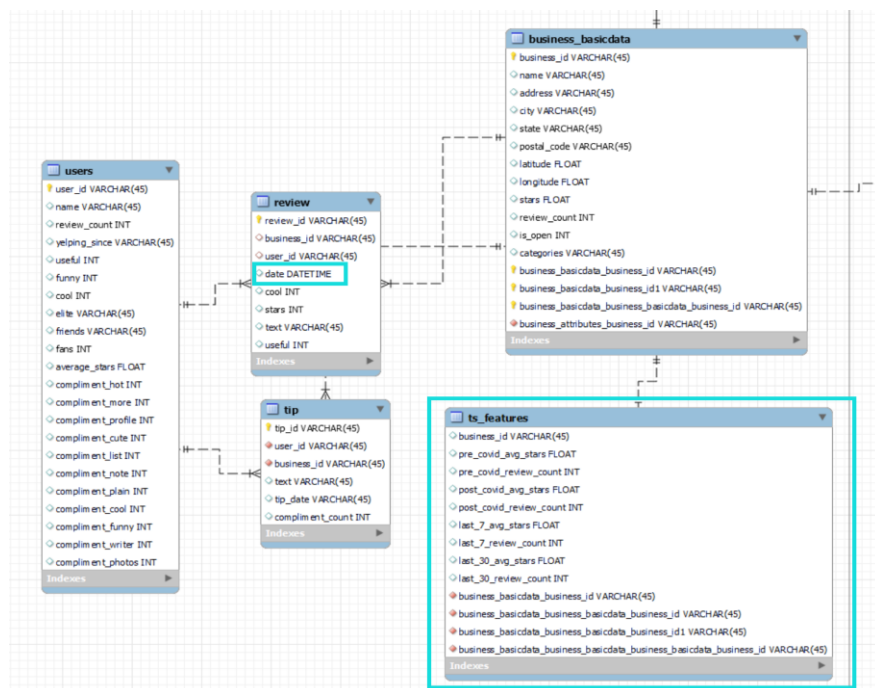


Sample of the largest category shows cyclic pattern in rolling means



The feature shows seasonality after decomposition:
- Notice disturbance due to COVID
- Number of reviews peak in summers and on Sundays

# Feature Engineering - Time Series



**Goal:**

Extract seasonality and other features from time and turn them into single dimention variables
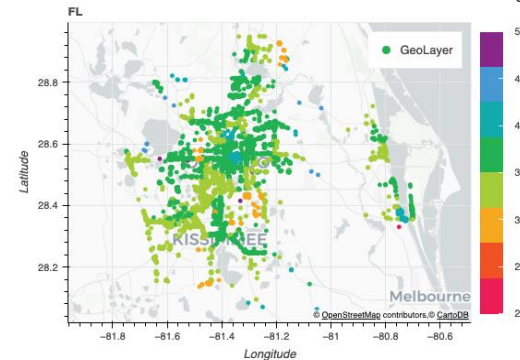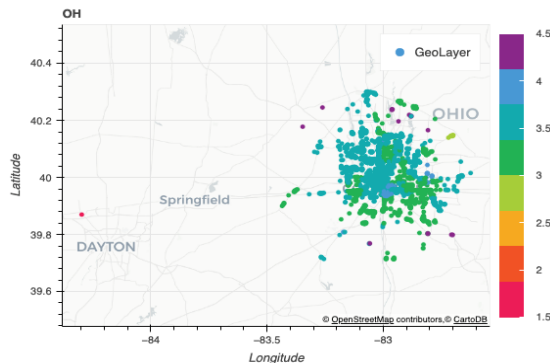
**Method:**

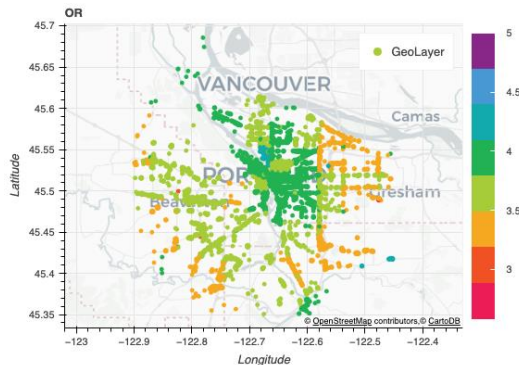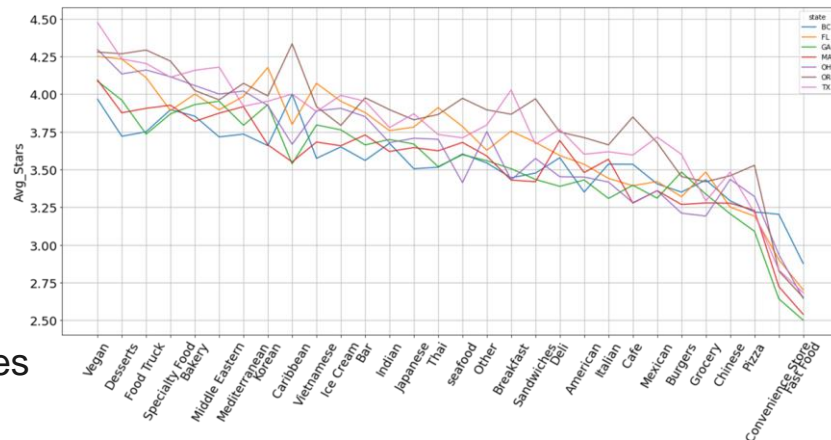- Calculated features for average stars and number of reviews before and after COVID
- Created features for average stars and number of reviews in last 7 days and last 30 days

# EDA – Geographic info

At state level:

- Ratings seem to have patterns among the states
- Higher ratings tend to be in in more centralized areas within the state
- Some states display preference of certain categories



Rating of Categories By State

# Feature Engineering - Geographic Info



**Goal:**

Convert (Latitude,Longitude) -> Multiple 1D data about the business

**Method:**

Calculate distance between each Business, then calculate aggregative informations

# Feature Engineering - Geographic Info

Geo data clustered by states



Geo data clustered by city



**Challenge:**
- Program time complexity is $O(N^2)$
- Memory issues (kernel dies)
- Running forever...

**Solution:**
- Calculate distance separatly within each states to reduced data-size by 85.6%
- Save and load intermediate results on hard drive to avoid memory issues
- Introduce Numpy vectorization to speed up calculation by 93.5%

**Results:**
- Reduced total run time by 99.1%
- Generated following aggregated features:

  - HHI (To measure category diversity)
  - Total review, average ratings within 2 km$^2$
  - Most frequent category within 2 km$^2$

# Recommendation System Model - SVD

# Recommendation System – Collaborative Filtering

**ITEM-BASED** COLLABORATIVE FILTERING

**Decomposed Business (Item) Matrix**

| feature1 | feature2 | feature... |
|----------|----------|------------|
| ... | ... | ... |
| ... | ... | ... |
| ... | ... | ... |

- Focuses on **relationship between pairs of items**
- Matches each users' rated items to similar items

$$Similarity(\vec{A}, \vec{B}) = \frac{\vec{A} \cdot \vec{B}}{||\vec{A}|| * ||\vec{B}||}$$

| Executive Summary | Data ETL | Feature Engineering | Model | Model Evaluation |

# Recommendation System – Collaborative Filtering

Use items already rated by user that are **most similar** to missing item we want to generate rating for

$$rating(U, I_i) = \frac{\sum_j rating(U, I_j) * s_{ij}}{\sum_j s_{ij}}$$

Generate a recommendation based on a **weighted sum of the ratings of other similar products**

Time & Memory Comparison



| Executive Summary | Data ETL | Feature Engineering | Model | Model Evaluation |

# Recommendation System – XGB+LR

**Structure**



**Step 1**  **Embedding**



$$\#Features = \#Tree * 2^{\max Depth}$$

**Step 2**  **OneHot + LR**

| Executive Summary | Data ETL | Feature Engineering | Model | Model Evaluation |

# Model Evaluation

## RMSE

| Model | Value |
|-------|-------|
| SVD+CF | 56.47% |
| XGBOOST+LR ✓ | 44.03% |
| XGBOOST | 44.71% |
| RIDGE+LR | 46.57% |

## NDCG

| Model | Value |
|-------|-------|
| SVD+CF | 92.23% |
| XGBOOST+LR ✓ | 97.01% |
| XGBOOST | 97.14% |
| RIDGE+LR | 95.91% |

### NDCG

$$DCG_k = \Sigma_{i=1}^{k} \frac{rel(i)}{\log_2(i+1)}$$

| i | rel(i) | log(i+1) | rel(i)/log(i+1) |
|---|--------|----------|-----------------|
| 1 = A | 0.5 | 1 | 0.5 |
| 2 = B | 0.9 | 1.59 | 0.57 |
| 3 = C | 0.3 | 2 | 0.15 |
| 4 = D | 0.6 | 2.32 | 0.26 |
| 5 = E | 0.1 | 2.59 | 0.04 |

Executive Summary | Data ETL | Feature Engineering | Model | Model Evaluation

# Key Takeaways



**Agile Development**

**Three Principles**

**Tech Stacks**

# Future Work

- System
  - Improve the **stability** of daily recommendation system considering the memory and computation
  - Deal with the frequent down-time of Spark Service
  - Deep dive into the relationships between features, such as Network Analysis
- Model
  - Fine-tune the NLP Embedding models and XGBoost models
  - Implement **Gradient Descending Method** to find the proper matrix decomposition comparing to SVD

# Thank you!
# Questions?

# Appendix

# Resources

**RCC Hadoop Hive Database Name:** dmp_yelp_rs

**Link to GitHub Repository:** https://github.com/MinglunPan/MSCA31008-Data-Mining-Principles

**Link to Original Yelp Data:** https://www.kaggle.com/yelp-dataset/yelp-dataset

# Team Introduction

| Team Member | Role | Description |
| --- | --- | --- |
| Jason Lee | Script | Note-taking for meetings |
| Milan Toolsidas | Quality Control | Checking if we meet the criteria of the assignments, formatting the PowerPoints and reports |
| Minglun Pan | Leader | Creating agendas for each meeting, making sure other roles' expectations are meet |
| Norah Zhang | Facilitator | Planning meetings, creating Zoom meetings when needed, making sure everyone's involved, timekeeping |
| Ryan Liao | Devil's Advocate | Providing constructive critical opinions, observing team dynamics |

# Project Workflow

| | | W6 | | | | | | | W7 | | | | | | | W8 | | | | | | | W9 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | M | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| 1 | Data Clearning | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1.1 | Prepration & Cleaning | | | 1 | | | 2 | | | | | | | | | | | | | | | | | | | | | | |
| 1.2 | Loading & Sharing | | | | | | | 1 | | | | | | | | | | | | | | | | | | | | | |
| 2 | Feature Engineering | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2.1 | NLP | | | | | | | | | 1 | | 2 | | | | | | | | | 3 | | | | | | | | |
| 2.2 | Time Series | | | | | | | | | | 1 | | | | | | | | | | 2 | | | | | | | | |
| 2.3 | Geo | | | | | | | | | | | | 1 | | | | | | | | 2 | | | | | | | | |
| 3 | Programs Design | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 3.1 | Framework Design(Data Structures, etc.) | | 1 | | 2 | | | | | | | | | | | | | | | | | | | | | | | | |
| 3.2 | Framework Implementation | | | | | | 1 | | | | | | | | | | | | | | | | | | | | | | |
| 3.3 | Integration & Unit Test | | | | | | | 1 | | | | | | | | | | 2 | | 3 | | | | | | | | | |
| 4 | Experimental Design & Model | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 4.1 | User Match Algorithm | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 4.1.1 | Propensity Score | | | | | | | | | | | 1 | | 2 | | | | | | | | | | | | | | | |
| 4.1.2 | Clustering | | | | | | | | | | | 1 | | 2 | | | | | | | | | | | | | | | |
| 4.2 | Recommendation System | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 4.2.1 | SVD & CF | | | | | | | | | | | | 1 | | | 2 | | 3 | | | 4 | | | | | | | | |
| 4.2.2 | GBDT & LR | | | | | | | | | | | | 1 | | | 2 | | 3 | | 4 | | | | | | | | | |
| 4.2.3 | Attention & DL* | | | | | | | | | | | | | | | | | | | 1 | | 2 | | | | | | | |
| 5 | Metrics Design & Implementation | | 1 | | | | | | | | | | | 2 | 3 | | | | | | 4 | | | | | | | | |
| 6 | Documents & Slides & Rehearsal | | | | | | | 1 | | | | 2 | | | | 3 | | | | 4 | | 5 | | | | | | | |

25

# Project Workflow

# EDA- Geograph



| user_state | | | |
|---|---|---|---|
| BC | 0.0 | 200.0 | 1192.0 |
| CO | 231.0 | 0.0 | 468.0 |
| FL | 1205.0 | 720.0 | 0.0 |
| GA | 1135.0 | 646.0 | 10451.0 |
| MA | 2897.0 | 1666.0 | 10003.0 |
| OH | 276.0 | 358.0 | 2367.0 |
| OR | 7193.0 | 1406.0 | 3378.0 |

| user_state | BC | CO | FL | GA | MA | OH | OR | TX |
|---|---|---|---|---|---|---|---|---|
| BC | 3.6 | 4.1 | 3.8 | 3.8 | 3.8 | 4.1 | 4.1 | 4 |
| CO | 4 | 3.8 | 3.8 | 3.9 | 3.9 | 3.9 | 4.1 | 4 |
| FL | 4.1 | 4.1 | 3.8 | 3.9 | 3.9 | 4 | 4.1 | 4.1 |
| GA | 4 | 4.1 | 3.8 | 3.7 | 3.8 | 3.8 | 4.1 | 4 |
| MA | 4 | 4 | 3.8 | 3.8 | 3.7 | 3.9 | 4.1 | 4 |
| OH | 4.2 | 4 | 3.9 | 3.8 | 3.9 | 3.7 | 4.1 | 4 |
| OR | 4 | 4 | 3.8 | 3.8 | 3.9 | 3.9 | 3.9 | 4 |
| TX | 4 | 4 | 3.8 | 3.9 | 3.9 | 4 | 4.1 | 3.8 |
| WA | 4.2 | 3.8 | 3.7 | 3.7 | 4 | 3.8 | 3.9 | 3.7 |

business_state

# EDA - Time Series