



PAYCHEX®

Payroll | Benefits | HR | Insurance

Final Presentation December 6, 2021

```
paused = false;
S_Resourcesound D;

if (skill == sk_nightmare)
    skill = sk_nightmare;

// This was quite noisy with SPECIAL and common
// Supposedly hacks to make the latest edition work
// It might not work properly
if (episode < 1)
    episode = 1;

if (gameMode == retail)
{
    if (episode > 4)
        episode = 4;
    else if (gameMode == shareware)
        episode = 1;
}

if (episode > 3)
    episode = 3;
}

if (map < 1)
    map = 1;
if (Cmap > 89
    && CgameMode != commercial)
    map = 91;
M_CloseRandom();

if (skill == sk_nightmare || respawnmons)
    respawnmons = true;
else
    S8805();
    respawnmons = false;

if (lastparm || (skill == sk_nightmare & gameSkill == sk_nightmare))
    for (i=S_SARG_PAIN1; i<=S_SARG_PAIN2; i++)
        statedm[i].dmg = 1;
mobility(MT_BRUISINGSHOT).speed = 20*FRACUNIT;
mobility(MT_HEADSHOT).speed = 20*FRACUNIT;
mobility(MT_HOODSHOT).speed = 20*FRACUNIT;
ACCESS POINT
if (skill == sk_nightmare || gameSkill == sk_nightmare)
    for (i=S_SAWG_PAIN1; i<=S_SAWG_PAIN2; i++)
        statedm[i].dmg = 1;
mobility(MT_BRUISINGSHOT).speed = 15*FRACUNIT;
mobility(MT_HEADSHOT).speed = 10*FRACUNIT;
mobility(MT_HOODSHOT).speed = 10*FRACUNIT;

// focus players to be initialized upon first level load
for (i=0; i<MAXPLAYERS; i++)
    player[i].initialized = PST_REBORN;
```

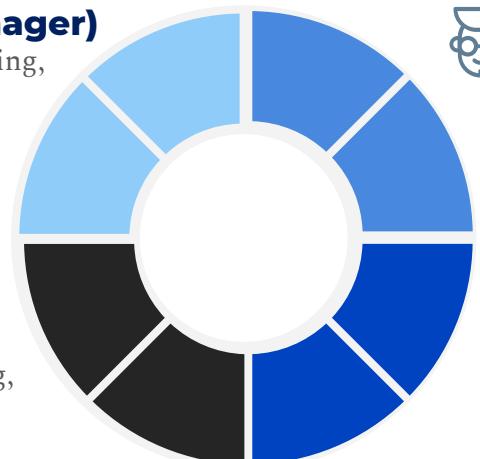
Team member

**Shijing Li (Project Manager)**

Data Collection/Preprocessing,
Visualization

**Yangxin Fan**

Data Collection, EDA, Modeling

**Yuan Wang**

EDA, Feature Engineering,
Modeling

**Lingyu Ye**

Data Collection/Preprocessing,
Visualization

Agenda

01 Introduction

02 Data Overview

03 Modeling

04 Future Work

1. INTRODUCTION

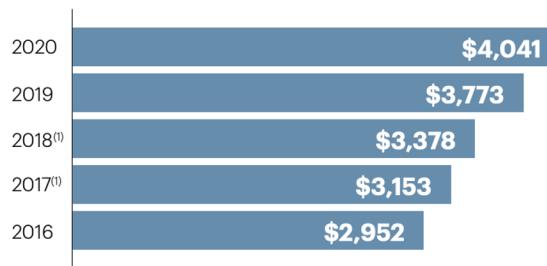


Who is Paychex?

#1 Ranked

401(k) Plan Recordkeeper

Total Revenue (\$Millions)
For the fiscal year ended May 31



- ✓ Leading provider of integrated human capital management solutions for payroll, benefits, human resources, and insurance services
- ✓ **680,000** business clients
- ✓ **41%** Return on Equity

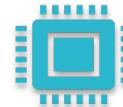
Goal & Vision

Goal:



Identify

Identify correlations between Paychex revenue and external features by time



Build

Build various time-series, hybrid, and deep learning-based models to predict Paychex revenue by time



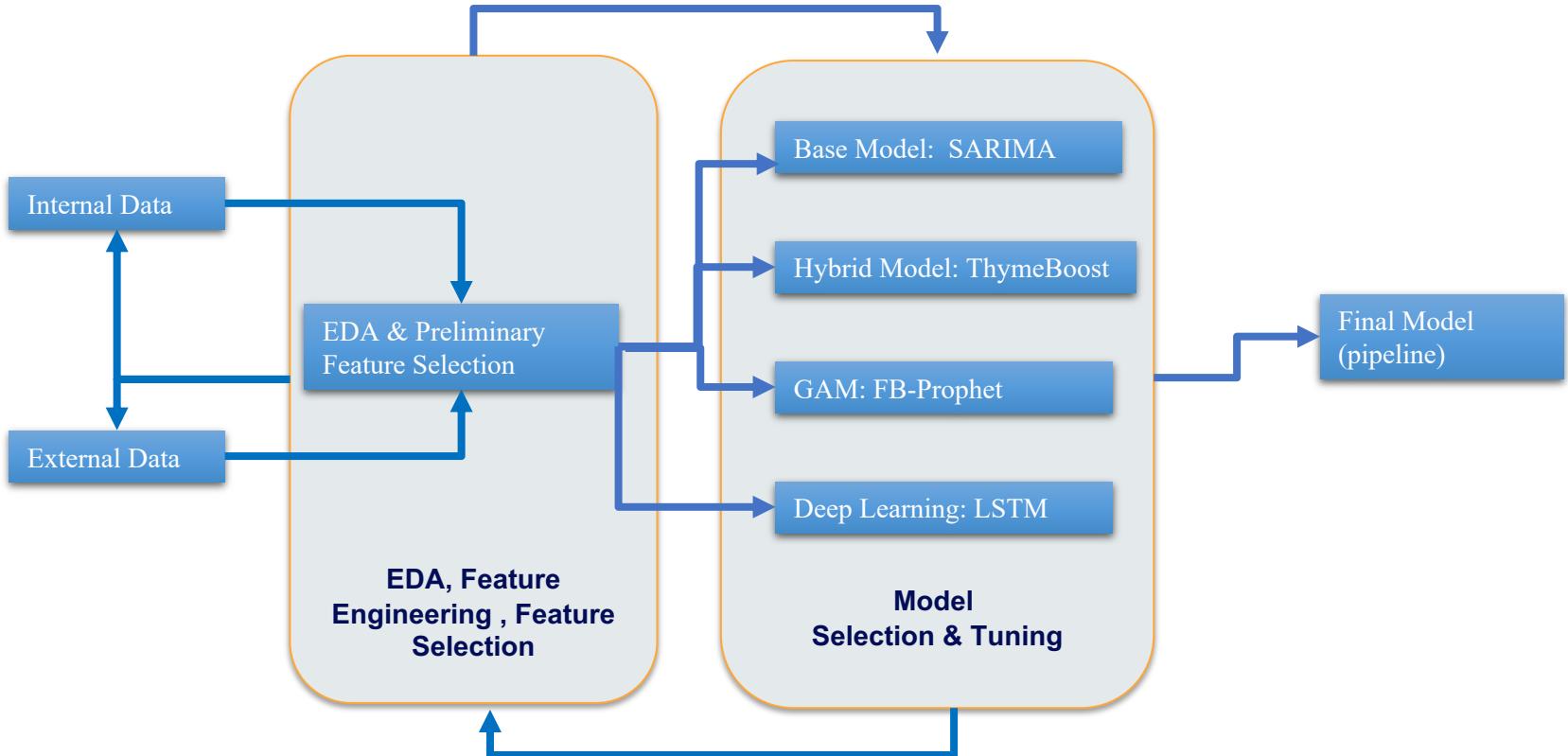
Deliver

Deliver constructive suggestions for Paychex to increase revenue by identifying new "growth"

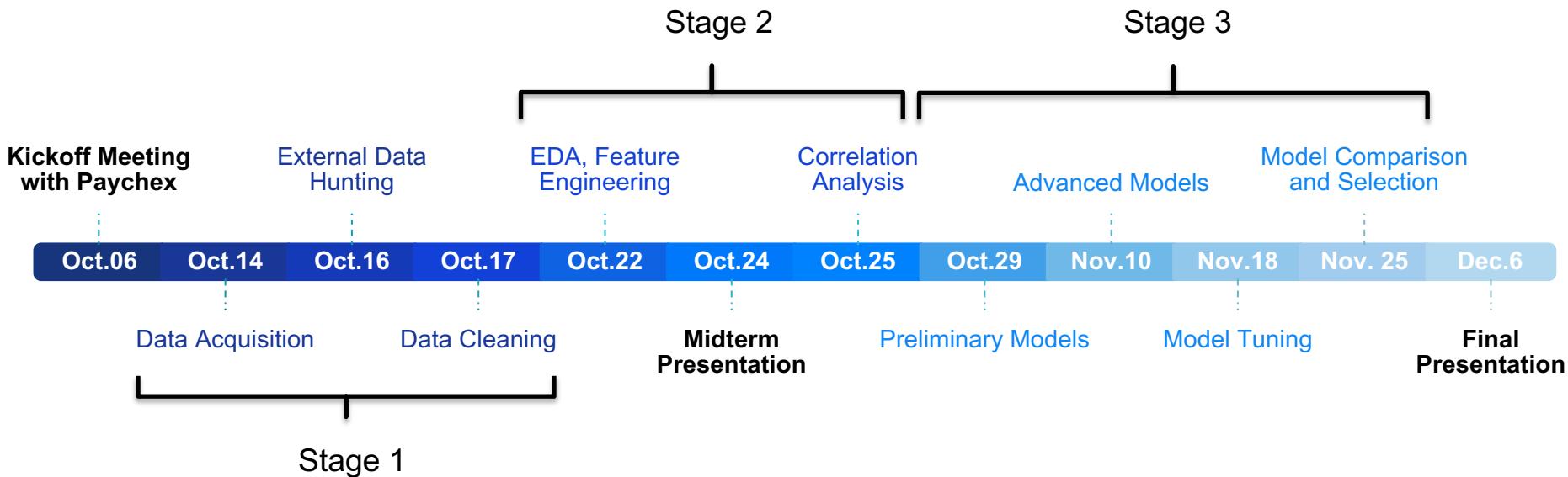
Vision:

Identify the external activities that have a significant impact on Paychex's internal activity to build prediction models that will forecast Pyachex's future monthly revenue gain or loss. Produce the best fit model by using performance measurements.

Method Overview



Milestones



2. Data Overview



Internal Data



Monthly Revenue Dataset:

Column Name	Description
UNIQUESYSTEMCLIENTID	Client ID
EECOUNT	Employee count
BUSTYPELEVEL1	Industry type (NAICS2)
STATECODE	State Code
TOTALREV	Total monthly revenue
PRODUCTCOUNT	Total monthly product count
PAYROLLREVENUE	Total monthly revenue in Payroll
PEOREVENUE	Total monthly revenue in PEO
ASOREVENUE	Total monthly revenue in ASO
TIMEREVENUE	Total monthly revenue in Time
RETIREMENTREVENUE	Total monthly revenue in Retirement
HNBREVENUE	Total monthly revenue in Health & Benefit
WRKSCOMPREVENUE	Total monthly revenue in Work Compensation

Dataset Description

FY17-FY21

823,963 rows, 247 columns
(No missing values & duplicates)

Client Information

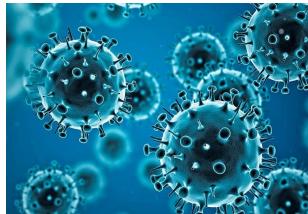
Monthly total revenue and product count of all Paychex's service products

Monthly revenue of 7 different service products

External Data



PAYCHEX



BUREAU OF
ECONOMIC
ANALYSIS

U.S. BUREAU OF
LABOR
STATISTICS

Centers for
Disease Control
and Prevention

- Monthly
- GDP per capita
 - Person income per capita
 - Inflation rate (quarterly)

- Monthly
- Job opening Rate
 - Hiring rate
 - Separation rate
 - Unemployment rate

- Monthly
- Case number
 - Death number
 - Vaccination number

GDP Per Capita

State	2018-Q1	2018-Q2	2018-Q3	2018-Q4
Alabama	40963.28	41091.58	41063.68	41
Alaska	71920.29	72419.56	72300.30	72
Arizona	43801.17	43874.85	43792.51	43
Arkansas	38205.09	38494.65	38597.23	38
California	66231.62	66856.95	67369.28	67
Colorado	60117.41	60159.96	60089.51	60
Connecticut	69444.61	69778.74	70191.01	70

Job Turnover and Unemployment Rate

Month	Job openings rate	Hiring rate	Total separations rate
June-17	4.00%	3.70%	3.60%
July-17	4.00%	3.80%	3.70%
August-17	4.00%	3.70%	3.60%
September-17	4.00%	3.60%	3.60%
October-17	4%	3.80%	3.60%

COVID-19 Case Number

State	20-Jan	20-Feb	20-Mar
Alabama	0	0	6
Alaska	0	0	1
American Samoa	0	0	0
Arizona	6	23	6
Arkansas	0	0	3

Data Preprocessing

Data cleaning

1. Check missing and duplicate values
2. Standardize data granularity

Data transformation

1. Transform internal & external data to usable formats
2. Merge internal & external datasets

Data reduction

1. Aggregate monthly revenue of all industries
2. Feature selection

```
[9] def revenue_table(df):
    strings = df.columns.tolist()
    #Total Revenue
    total_cols = [string for string in strings if "TOTALREV" in string]
    #PAYROLLREV
    payroll_cols = [string for string in strings if "PAYROLLREV" in string]
    #PEOREV
    peo_cols = [string for string in strings if "PEOREV" in string]
    #ASOREV
    aso_cols = [string for string in strings if "ASOREV" in string]
    #TMRREV
    timer_cols = [string for string in strings if "TMRREV" in string]
    #RETIREMENTREV
    retire_cols = [string for string in strings if "RETIREMENTREV" in string]
    #HNBREV
    hnb_cols = [string for string in strings if "HNBREV" in string]
    #WRKRSCOMPREV
    wrkscomp_cols = [string for string in strings if "WRKRSCOMPREV" in string]

    other = df[total_cols].sum().sum() - (df[payroll_cols].sum().sum() + df[peo_cols].sum().sum()
    + df[timer_cols].sum().sum() + df[retire_cols].sum().sum() + df[hnb_cols].sum().sum()
    + df[wrkscomp_cols].sum().sum())
    rev_by_category = [df[total_cols].sum().sum(), df[payroll_cols].sum().sum(),
    df[peo_cols].sum().sum(), df[aso_cols].sum().sum(),
    df[timer_cols].sum().sum(), df[retire_cols].sum().sum(),
    df[hnb_cols].sum().sum(), df[wrkscomp_cols].sum().sum(),
    other]
    label = ['Total', 'Payroll', 'PEO', 'ASO', 'Time', 'Retire', 'H&B', 'WorkComp', 'Others']
    return rev_by_category
```

```
table_rev = pd.DataFrame(
    columns=['Total', 'Payroll', 'PEO', 'ASO', 'Time', 'Retire', 'H&B', 'WorkComp', 'Others'],
    index=[2018, 2019, 2020, 2021, 2022])
table_rev.loc[2018] = revenue_table(px_18)
table_rev.loc[2019] = revenue_table(px_19)
table_rev.loc[2020] = revenue_table(px_20)
table_rev.loc[2021] = revenue_table(px_21)
table_rev.loc[2022] = revenue_table(px_22)
# Add percentage for convenience
table_rev['Payroll %'] = table_rev['Payroll']/table_rev['Total']
table_rev['PEO %'] = table_rev['PEO']/table_rev['Total']
table_rev['ASO %'] = table_rev['ASO']/table_rev['Total']
table_rev['Time %'] = table_rev['Time']/table_rev['Total']
table_rev['Retire %'] = table_rev['Retire']/table_rev['Total']
table_rev['H&B %'] = table_rev['H&B']/table_rev['Total']
table_rev['WorkComp %'] = table_rev['WorkComp']/table_rev['Total']
table_rev['Others %'] = table_rev['Others']/table_rev['Total']
pd.options.display.float_format = '{:.2f}'.format
table_rev
```

**Limited Internal
Datasets (4yrs)**

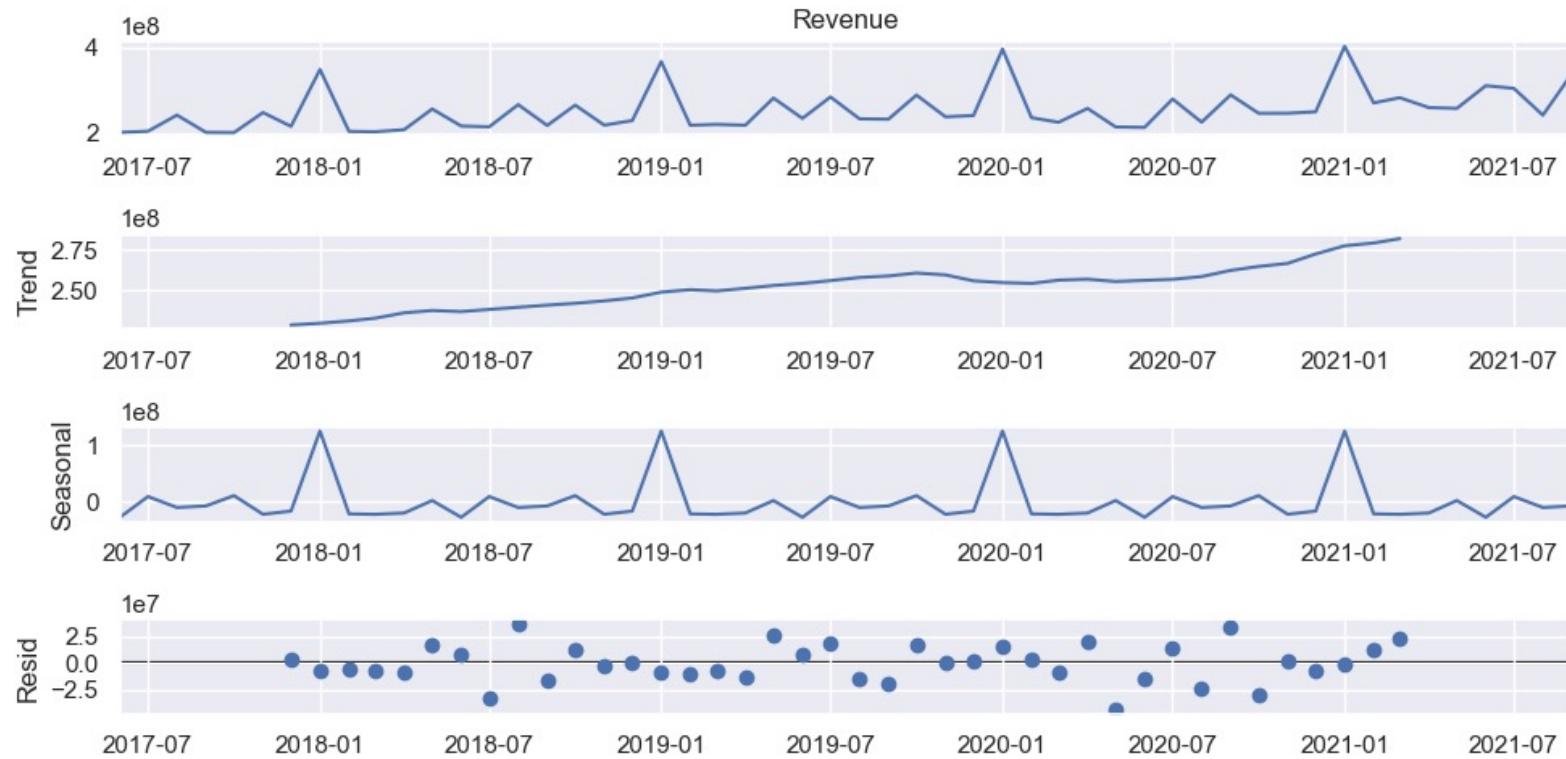
**External Datasets
With Different
Granularities**

**No Boundary for
External
Factors Hunting**

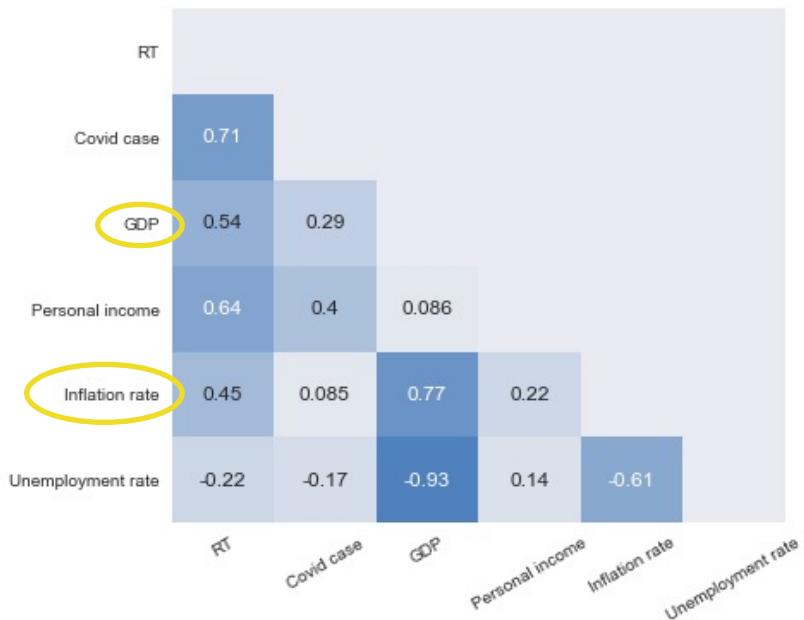
**Internal Dataset
Structured
Differently**

Challenges

Exploratory Data Analysis – Data Decomposition



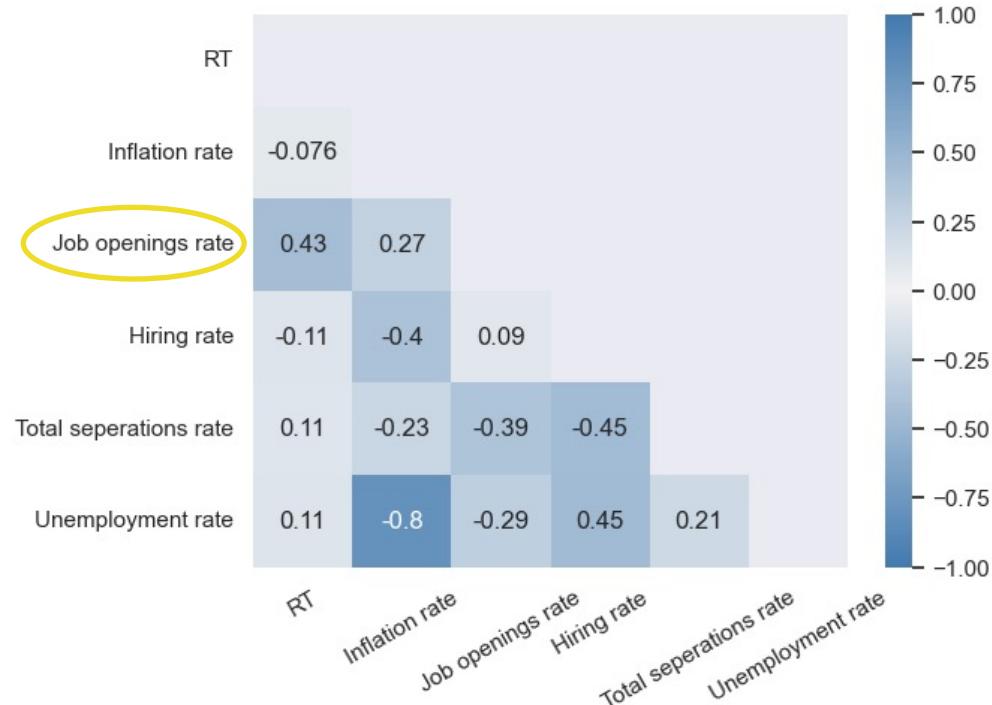
Exploratory Data Analysis – Covid Correlation Analysis



- Correlation between RT (residual + trend) Seasonally Adjusted Revenue, COVID-19 and Economic external factors
- Quarterly data from 2020 Q1 to 2021 Q3

Conclusion: GDP & Inflation Rate are **highly** correlated with Seasonality Adjusted Revenue.
 $(r = 0.81)$

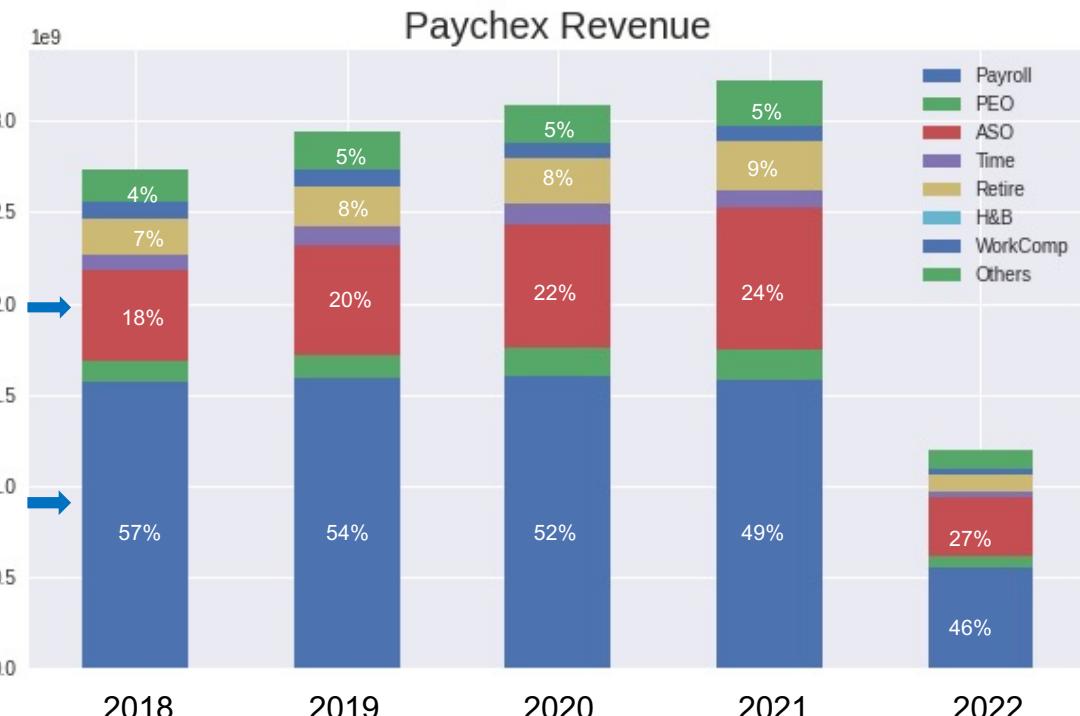
Exploratory Data Analysis – Correlation Analysis



- Correlation between RT (residual + trend) **Seasonally Adjusted Revenue** and **BLS** external factors
- Monthly data from 06/2017 - 08/2021

Conclusion: Job Opening Rate is moderately correlated with Seasonality Adjusted Revenue ($r = 0.43$)

EDA – Data Visualization

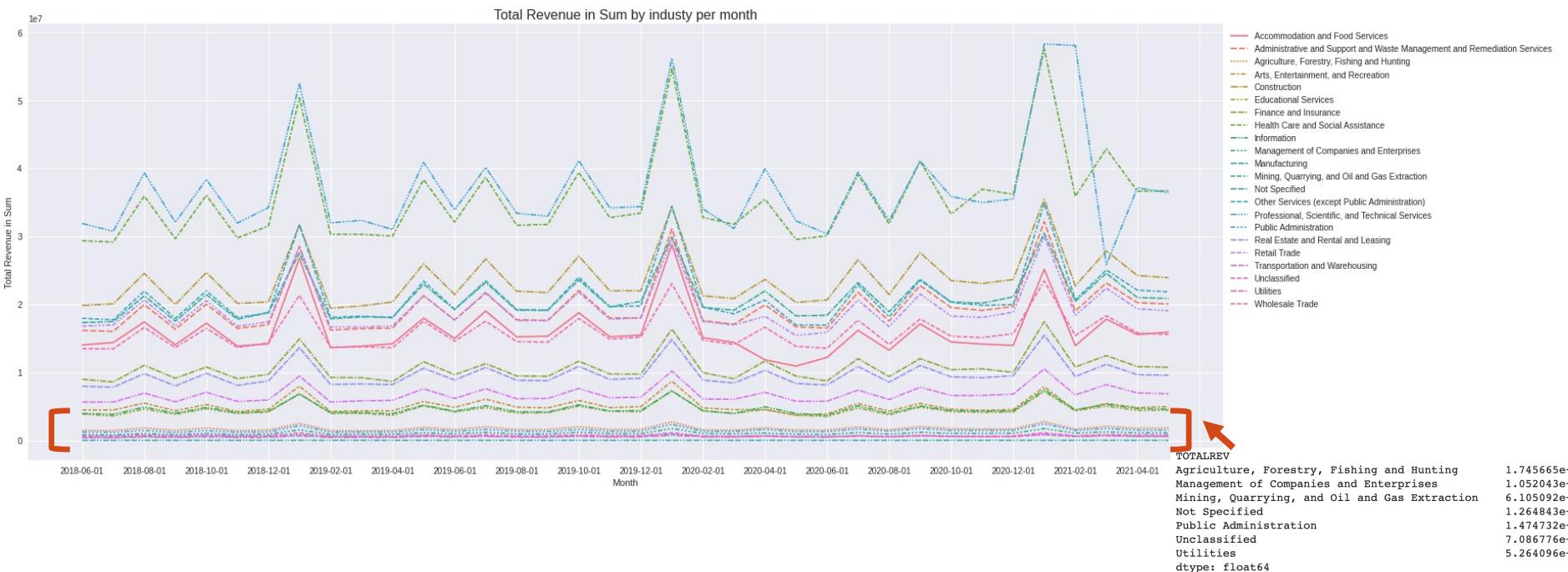


PEO: Professional Employer Organization
 ASO: Administrative Services Organization
 H&B: Health & Benefits

- The **total revenue** is increasing over the years.
- The revenue percentage for the **Payroll** service product is decreasing.
- The revenue percentage for the **ASO** service product is increasing.

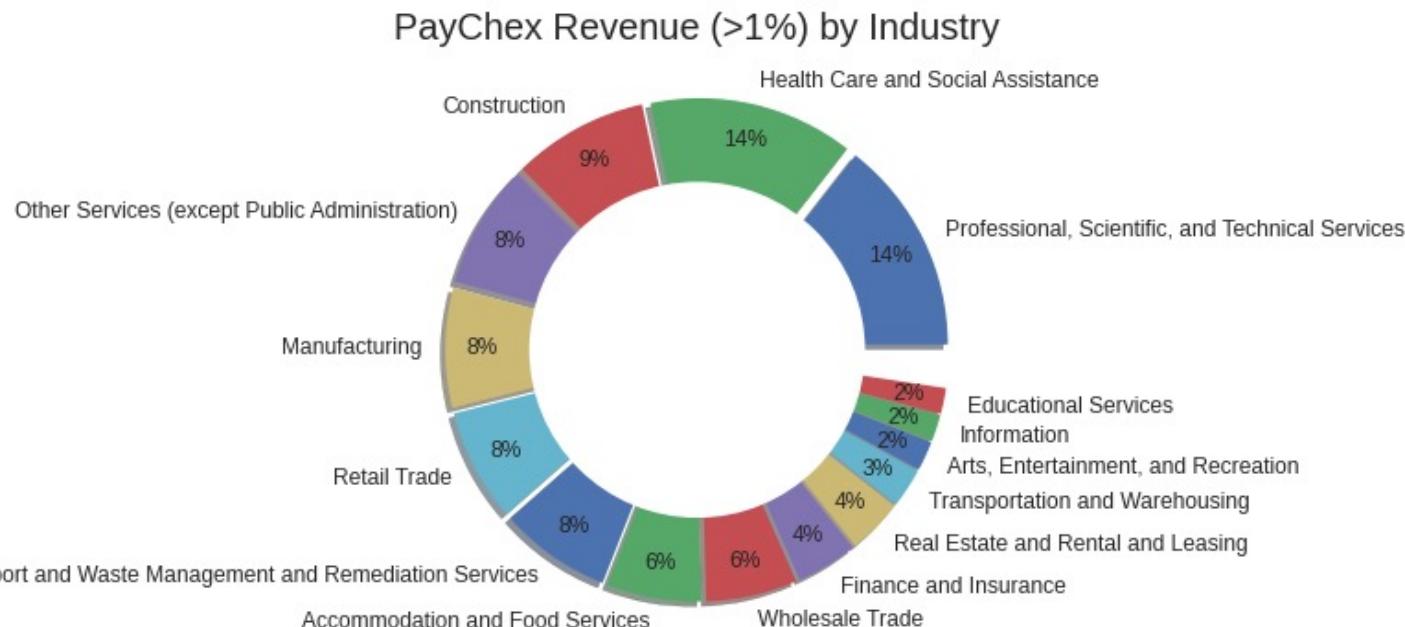
* 2022 fiscal year is incomplete (3 months data), but the trend of each portion stays consistent.

EDA – Data Visualization



- Over the years, there is no significant change in revenue contribution among different industries.
- We can use Total Revenue instead of each industry for the future revenue prediction.

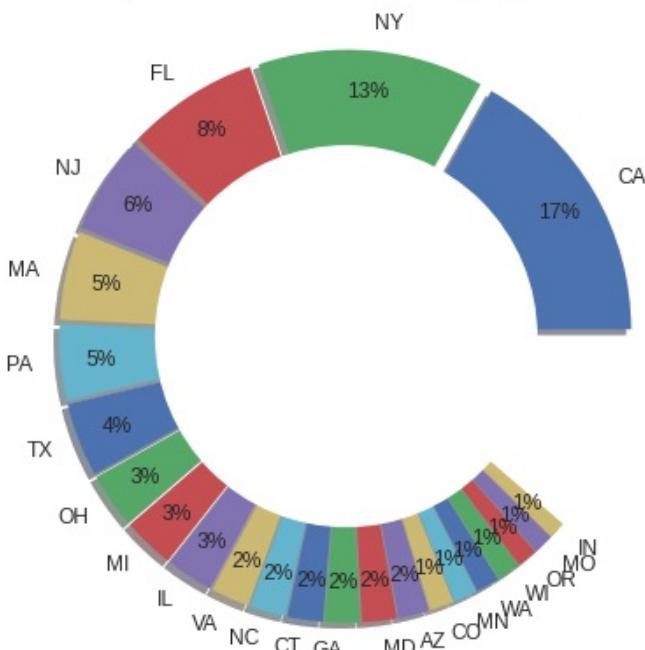
EDA – Data Visualization



- The top-3 industries are Professional services, Health care, and Construction.
- One trillion infrastructure bill passed will lead a high potential opportunity in construction industry.

EDA – Data Visualization

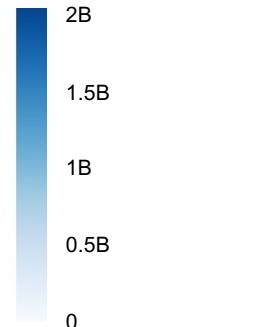
PayChex Revenue (>1%) by State



2018-2022 Paychex Revenue by State

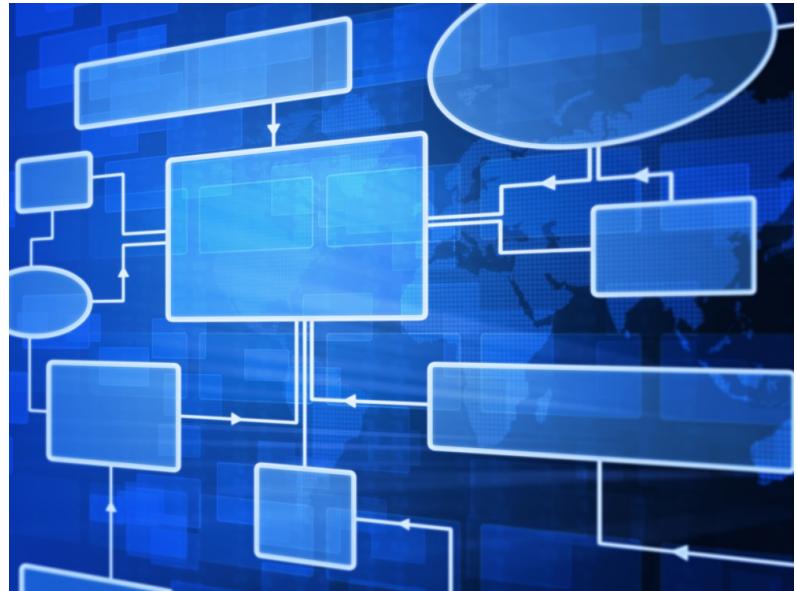


2018-2022 Paychex Revenue by State



- Top-5 states with the highest revenue: **CA, NY, FL, NJ, and MA**
- **TX, IL** have large potential to obtain more revenue since they are 2nd and 5th largest economy state in the U.S.

3. Modeling



Performance Measurement

- **MAE:** Mean absolute error (MAE) is the average of the absolute errors.

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n} = \frac{\sum_{i=1}^n |e_i|}{n}$$

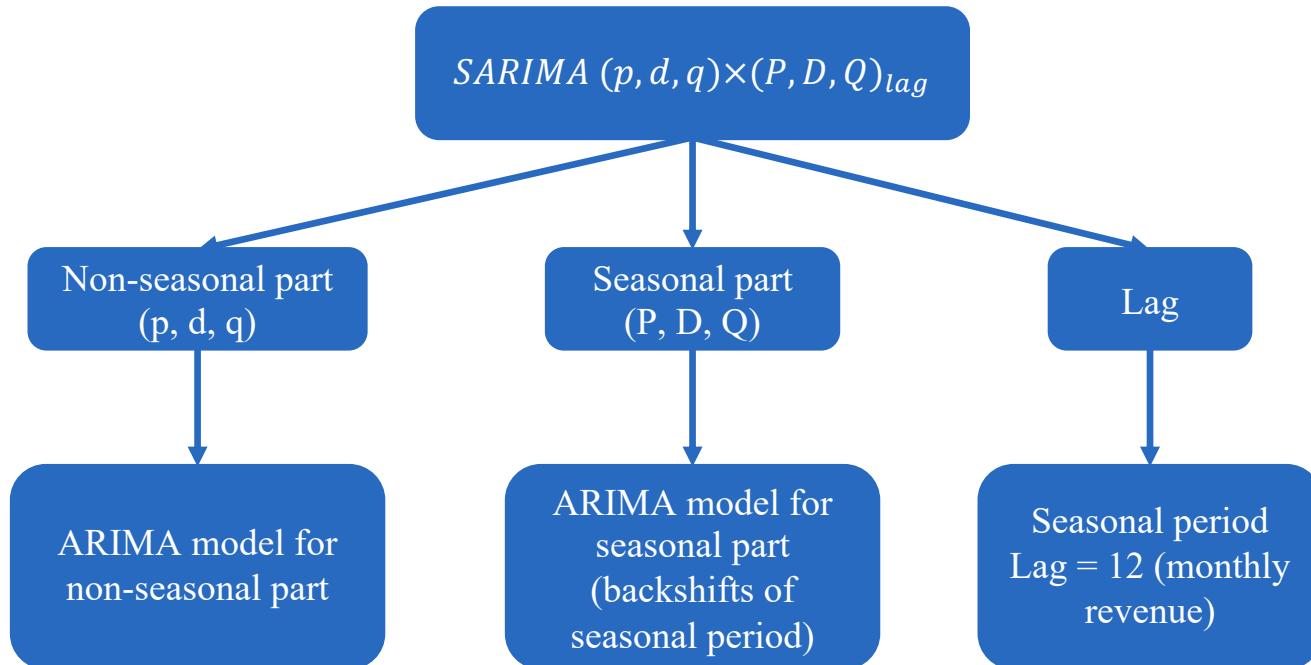
- **MAPE:** Mean absolute percentage error (MAPE) is a measure of prediction accuracy of a forecasting method in statistics.

$$\text{MAPE} = \frac{1}{n} \sum_{t=1}^n \left| \frac{V(t) - P(t)}{V(t)} \right| * 100$$

- **RMSE:** Root Mean Square Error (RMSE) is the standard deviation of the errors/residuals.

$$\text{RMSE} = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

Model - SARIMA



- (p, d, q) (P, D, Q) are integers ≥ 0 and refer to the order of the autoregressive, integrated, and moving average of the model respectively

Model - SARIMA

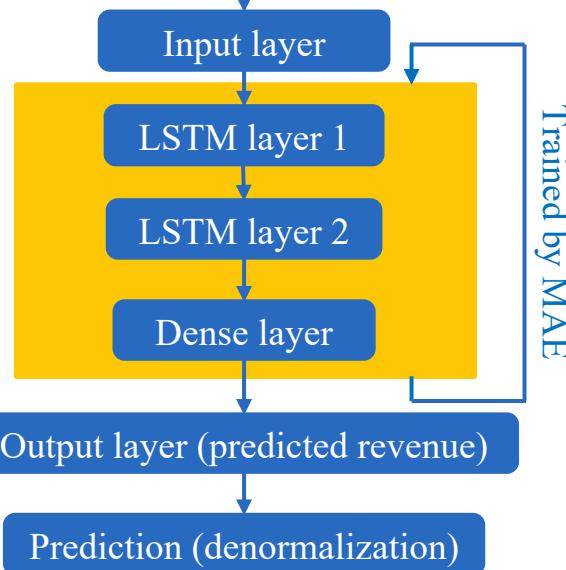


- External feature used: **Job openings rate** (correlated with seasonally adjusted revenue)
- Final model chosen by **AIC cost**: $SARIMA(0, 1, 1) \times (1, 1, 1)_{12} + Job\ Openings\ Rate$

Model - LSTM

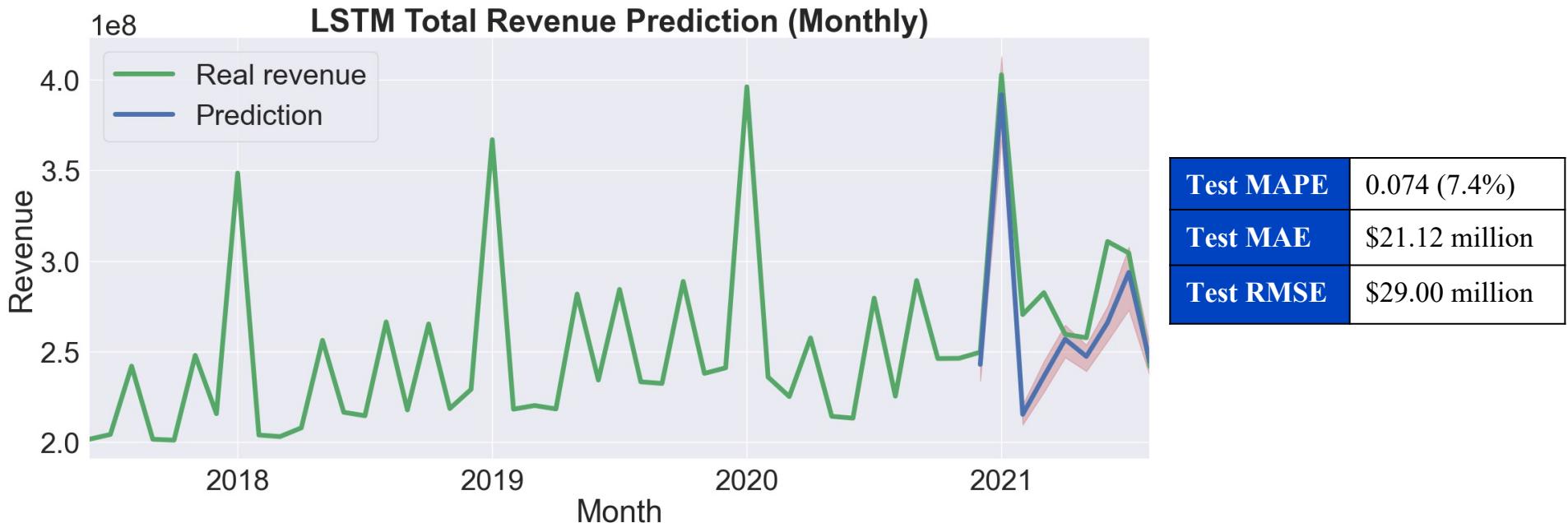
Min-Max normalization of revenue and external factors

Training samples (look back period = 12 months)



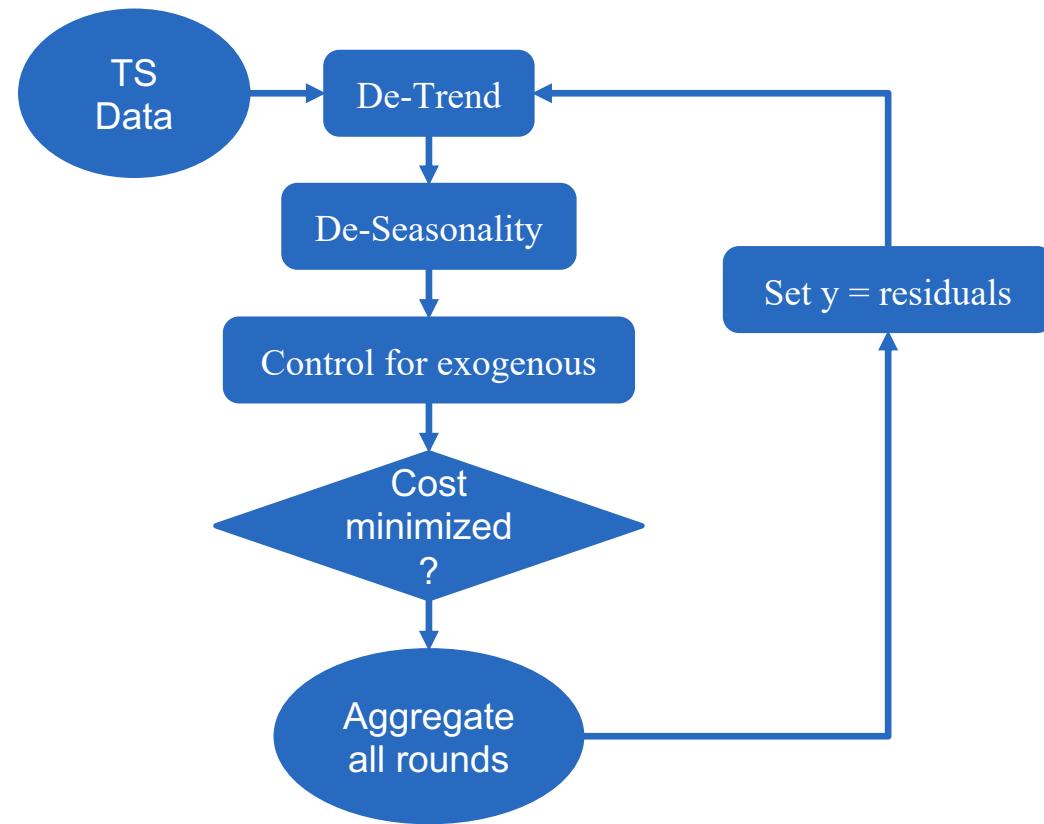
- **Training period:** 06/2018 - 11/2020 (30 months)
- **Testing period:** 12/2020 - 08/2021 (9 months)

Model - LSTM



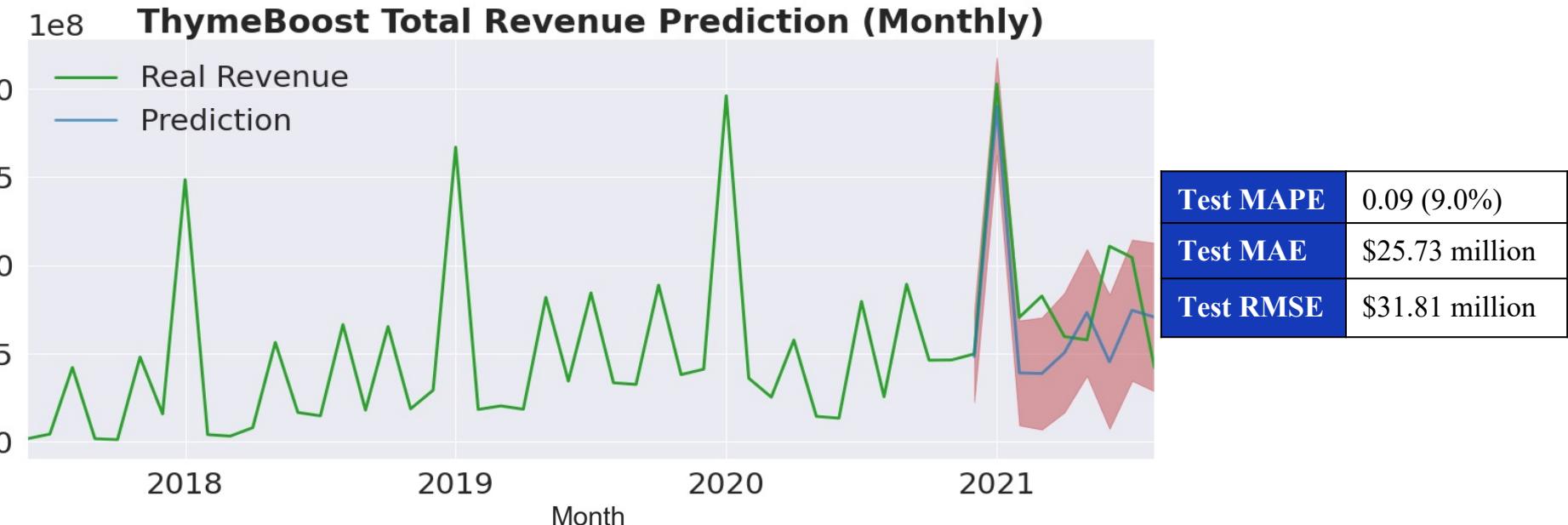
- No external features since including external features worsens model performance
- Ensembled LSTM model (100 LSTM models)

Model - ThymeBoost



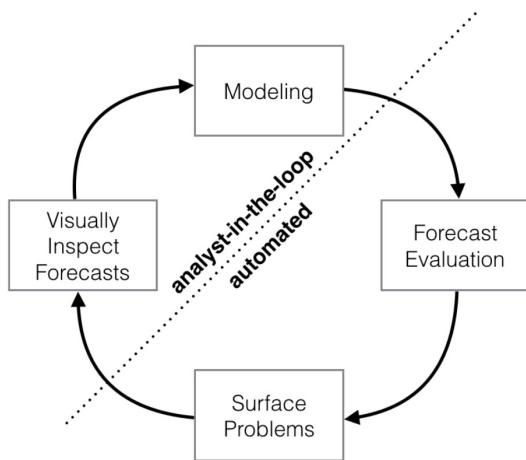
- Traditional time series decomposition for trend, seasonality
- **Gradient boost model for residual**
- Cost minimize

Model - ThymeBoost



- No external features
- Trend estimator: Linear
- Seasonal estimator: Fourier
- Cost: AIC

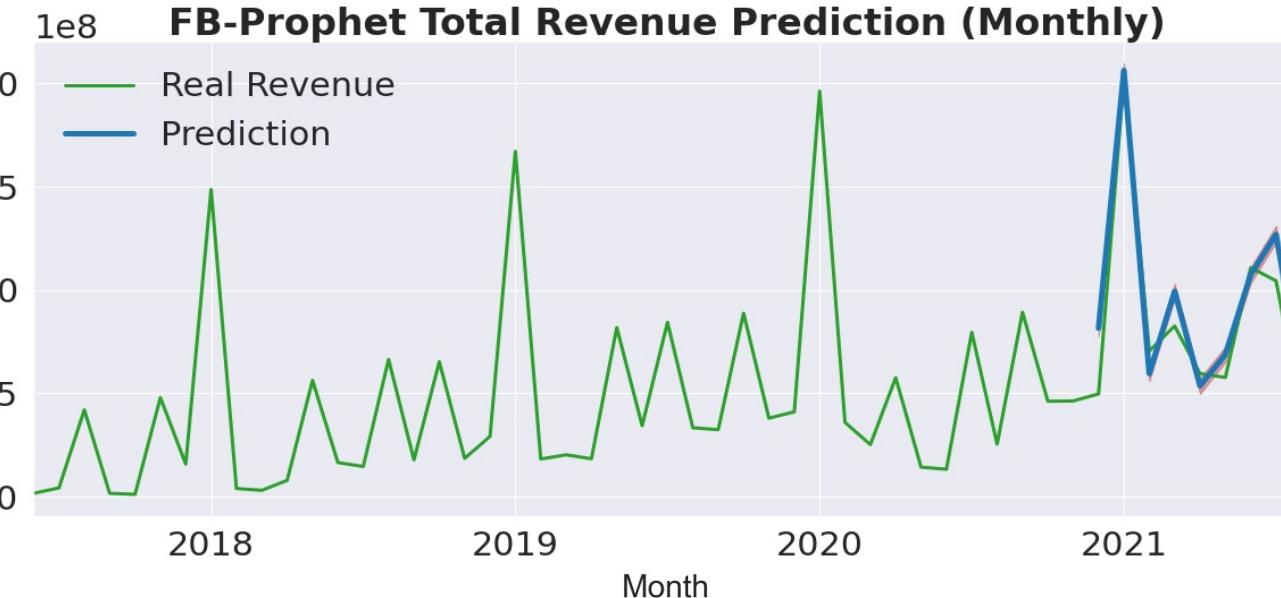
Model – FB-Prophet



$$y(t) = g(t) + s(t) + h(t) + e(t)$$

- $g(t)$ is a trend function which models the non-periodic changes.
 - $s(t)$ represents a periodic changes
 - $h(t)$ is a function that represents the effect of holidays which occur on irregular schedules. ($n \geq 1$ days)
 - $e(t)$ represents error changes that are not accommodated by the model.

Model – FB-Prophet + peak indication



Test MAPE	0.049 (4.9%)
Test MAE	\$13.33 million
Test RMSE	\$15.98 million

- External features used : Job Openings Rate, Hiring Rate, Total Separations Rate, Unemployment Rate, Inflation Rate
- Holidays/Peak events** indication applied both on **train** and **test periods**

Model Comparison

MODEL	Inputs	MAPE	Threshold	Interpretation
FB-Prophet	TS data, external features	4.9%	<5%	Highly Accurate
LSTM	TS data	7.4%	<10%	Good
SARIMA	TS data, external features	8.0%	<10%	Good
ThymeBoost	TS data	9.0%	<10%	Good

Conclusion

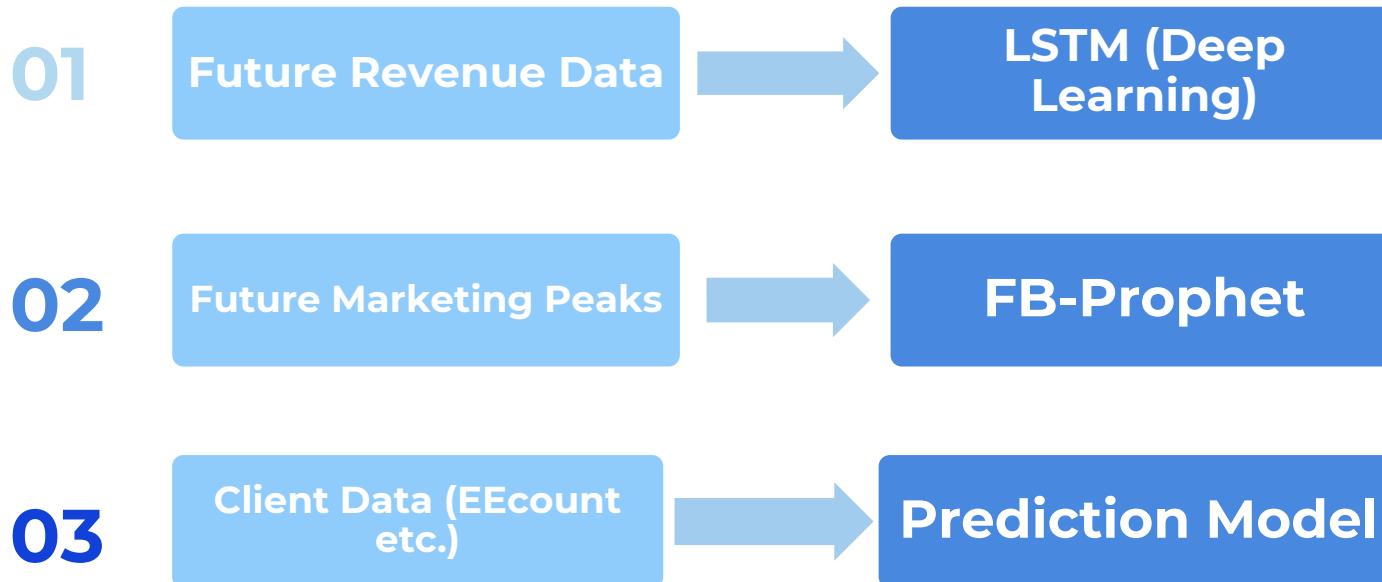
- The increase in GDP, Inflation rate and, Job opening rate may indicate an increase in Paychex revenue
- FB-Prophet has the best performance with 4.9% of MAPE
- Potential revenue growth opportunities in Texas and Illinois. (enhance marketing activities)
- Potential opportunities in construction industry



4. Future Work



Future works



Acknowledgement

- Paychex Client

Ryan Hayden
Joel Page
Oseoba Airewele
Mike Childress

- Advising Professor

Prof. Cantay Caliskan
Prof. Ajay Anand



Thanks!



Appendix – Project Directory

General			Code			Paycheck project documents		
+ New	Upload	Sync	Correlation Analysis.ipynb	2 minutes ago	Fan, Yangxin	Final Presentation.pptx	A few seconds ago	Wang, Yuan
			EDA_Yangxin.ipynb	October 26	Fan, Yangxin	Kickoff_email.docx	October 6	Li, Shijing
			LSTM model.ipynb	19 minutes ago	Fan, Yangxin	Meeting_1_Kickoff Questions_UofR Capston...	October 6	Li, Shijing
			PX_0_yuan.ipynb	18 minutes ago	Wang, Yuan	Meeting round 3.pptx	November 22	Ye, Lingyu
			PX_v1_Yuan.ipynb	18 minutes ago	Wang, Yuan	Meeting round 4.pptx	November 22	Li, Shijing
			SARIMA model.ipynb	19 minutes ago	Fan, Yangxin	PaycheX project meeting round 5.pptx	November 18	Fan, Yangxin
			Time series analysis_Yangxin.ipynb	October 26	Fan, Yangxin	Project_Charter_Paychex2_Starship.docx	October 15	Fan, Yangxin
			Documents > General	External Data			Project_Charter_Template_CAPSTONE_V7.d...	October 15
			Code	October 26	Fan, Yangxin	Script for midterm presentation.docx	October 26	Wang, Yuan
			External Data	October 18	Wang, Yuan	Timeline.xlsx	October 19	Wang, Yuan
			Paycheck Data	October 15	Wang, Yuan	Timeline-Yuan's MacBook Pro.xlsx	October 19	Wang, Yuan
			Paycheck project documents	October 15	Wang, Yuan	Weekly Meeting Notes.docx	October 6	Fan, Yangxin
			Final_Report_Paychex2.docx	About an hour ago	Fan, Yangxin	Documents > General > Paycheck Data		
			SpectralAnalysis_TimeSeries_MachineLearni...	October 28	Wang, Yuan	FY18 PX Clients - Monthly Rev Buckets.csv	November 10	Li, Shijing
						FY19 PX Client Data.csv	October 15	Fan, Yangxin
						FY20 PX Clients - Monthly Rev Only (SUM)...	October 18	Li, Shijing
						FY20 PX Client Data.csv	October 15	Fan, Yangxin
						FY20 PX Clients - Monthly Rev Only (SUM)...	October 15	Fan, Yangxin
						FY21 PX Client Data.csv	October 15	Fan, Yangxin
						FY21 PX Clients - Monthly Rev Only (SUM)...	October 15	Fan, Yangxin
						FY22 PX Clients - Monthly Rev Buckets.csv	November 10	Li, Shijing
						px_total_rev.csv	October 20	Wang, Yuan
						U of R-DrB Data.csv	November 5	Wang, Yuan