

# Home price index prediction for greater Toronto

Queen's University MMAI 823 – Section 2  
Professor: Matthew Thompson

Team Alfred:  
Kelsey Pericak, Shilpa Vishwanath, Mohammad Raza, Shayaan Mehdi, Sky Tang, Nezil Gopinadh, Theebak Sothilingam & Yvan Large



# Single family home price index



## Home price index (HPI)

- CREA
- Sophisticated formula
- Greater Toronto
- 1 unit = 1% from base price (100)
- January 2005 base



## Two main objectives

1. Predict HPI 1 month ahead for greater Toronto area using historical data
2. Learn which features have relationship with/ impact on real estate prices in greater Toronto



## Why does it matter?

- When to buy
- When to sell
- Use trends to plan for future investments/ policy making
- Evaluate potential return on investment
- Real estate agent in-the-know

# Our modeling process

1. Dataset creation
2. Exploratory analysis
3. Data preparation
4. Time sensitive split
  - Train (85%) from Jan. 2006 to Sept. 2018
  - Test (15%) split from Oct. 2018
5. Base model creation
6. Feature selection
7. Hyper parameter tuning
8. Prediction

# 35 features in dataset from 2006 to 2020

|         | Target            | Features, 1 month & 1 year lagged     |                                   |  |                              |
|---------|-------------------|---------------------------------------|-----------------------------------|--|------------------------------|
| Source  | CREA              | CREA                                  | Statistics Canada                 |  | Other                        |
| Toronto | Single family HPI | Single family HPI                     | Retail trade sales                | Population (forecasted)                    | Weather (snow & temperature) |
|         |                   |                                       | Construction starts & completions |  | Month                        |
| Ontario |                   |                                       | Median income                     | Bachelor's degree registrations            |                              |
|         |                   |                                       | Unemployment rate                 |  |                              |
| Canada  |                   | Single family HPI In other big cities | Bank rate                         | Consumer price index in metropolitan areas | S&P 500                      |
|         |                   |                                       | GDP                               |  |                              |



# Examples of features built from hypotheses

|                      |       |   |
|----------------------|-------|---|
| Consumer price index | ----- | ↑ inflation = price ↑                     |
| Bank rate            | ----- | ↓ rate = ↑ demand for mortgage            |
| Weather              | ----- | ↓ temperatures = ↓ interest in moving     |
| Unemployment         | ----- | ↑ rate = ↑ wealth gap OR ↓ buying         |
| Construction         | ----- | ↑ completions = ↑ new homes @ high prices |

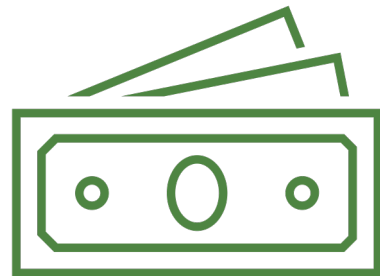
# Building the dataset required some feature engineering



Calculated  
average snowfall  
and temperature  
with daily data



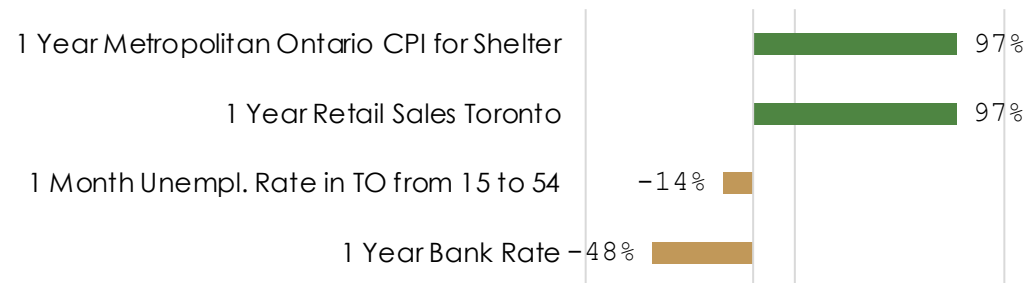
Forecasted  
population  
posted annually



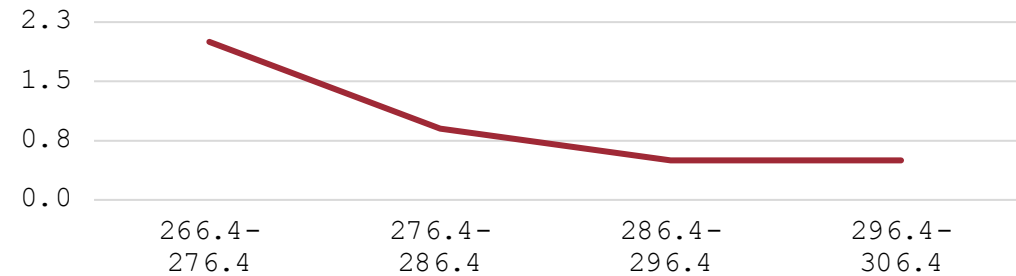
Imputed mean values  
for median family  
income for couples  
and lone parents

# Insights about single family home price index

## Interesting Correlations with Target



## Average 1 Month Lag Bank Rate (%) vs HPI Ranges in 2020



**1.5X HPI**

Target Increase  
From 2015 to 2020

**>17%**

Growth in avg annual HPI  
in 2016 & 2017

**Seasonal**

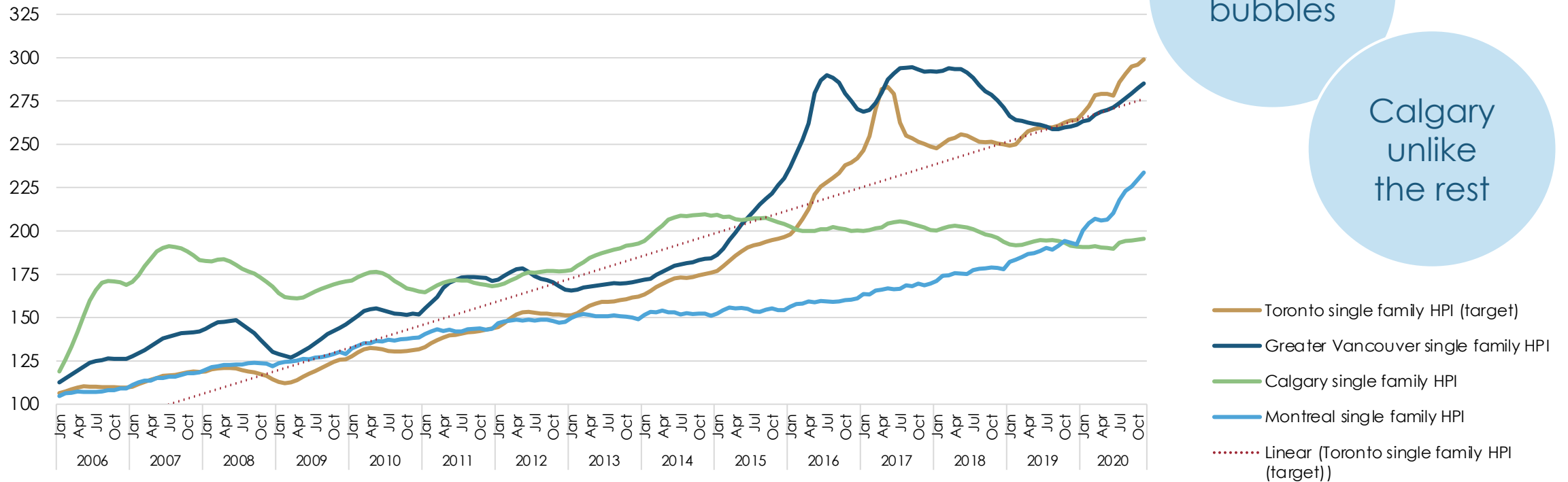
Troughs or stability in HPI  
during colder months

**19**

Base features with >.9  
correlation with target

# Comparing HPI indices for single family

HPI with no lag



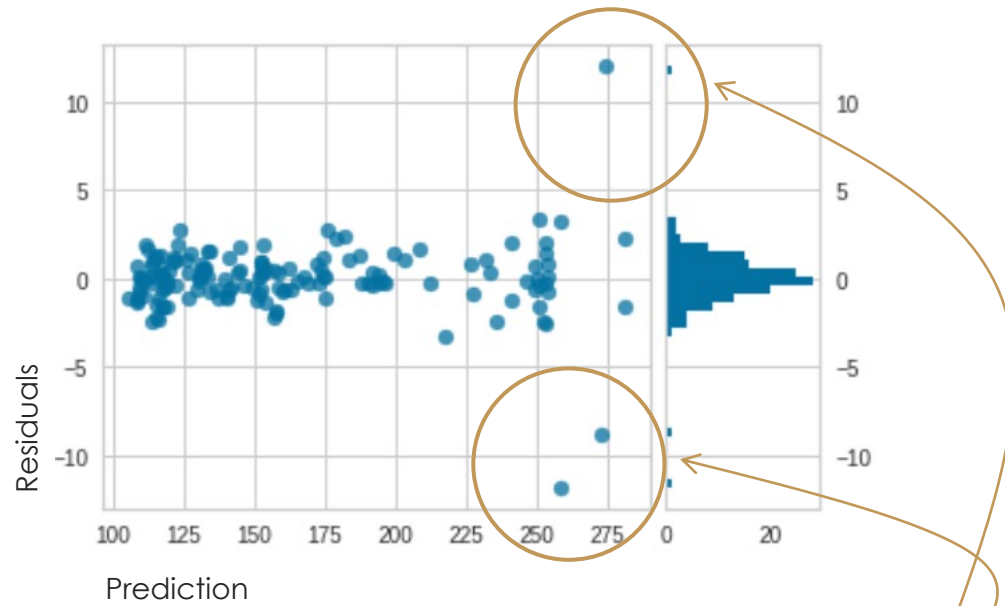


# Three main data cleaning steps beyond dataset creation and before feature selection

1. Dropped Calgary
2. Eliminated bubble time range
3. Capped construction outliers

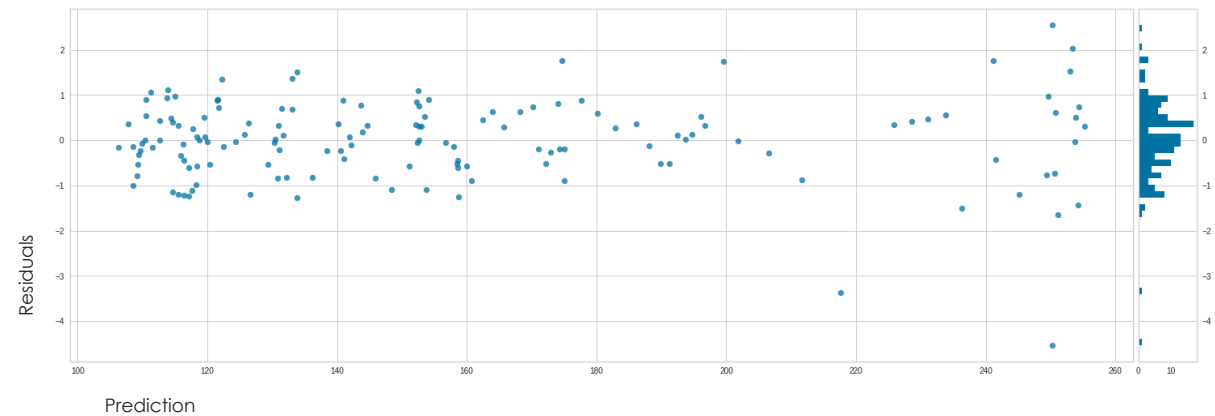
# Normality testing and outlier detection

All data

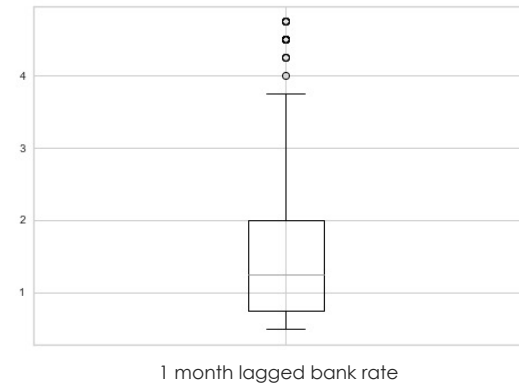
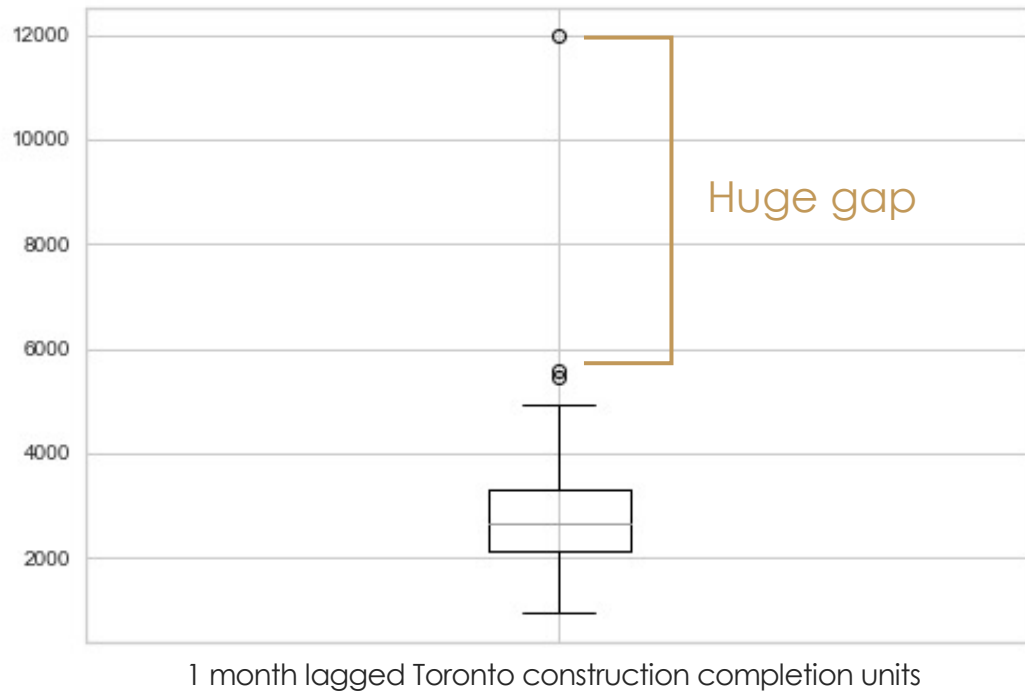


Housing bubble

Cleaned data



# Features with outliers



# Bayesian model selected

No feature selection  
or tuning

|                                 | ARD  | Linear Regression | Bayesian |
|---------------------------------|------|-------------------|----------|
| R <sup>2</sup> on training data | 0.99 | 0.99              | 0.99     |
| R <sup>2</sup> on test data     | 0.97 | 0.76              | 0.74     |
| RMSE on test data               | 2.76 | 8.0               | 8.43     |

Feature selection  
and tuning

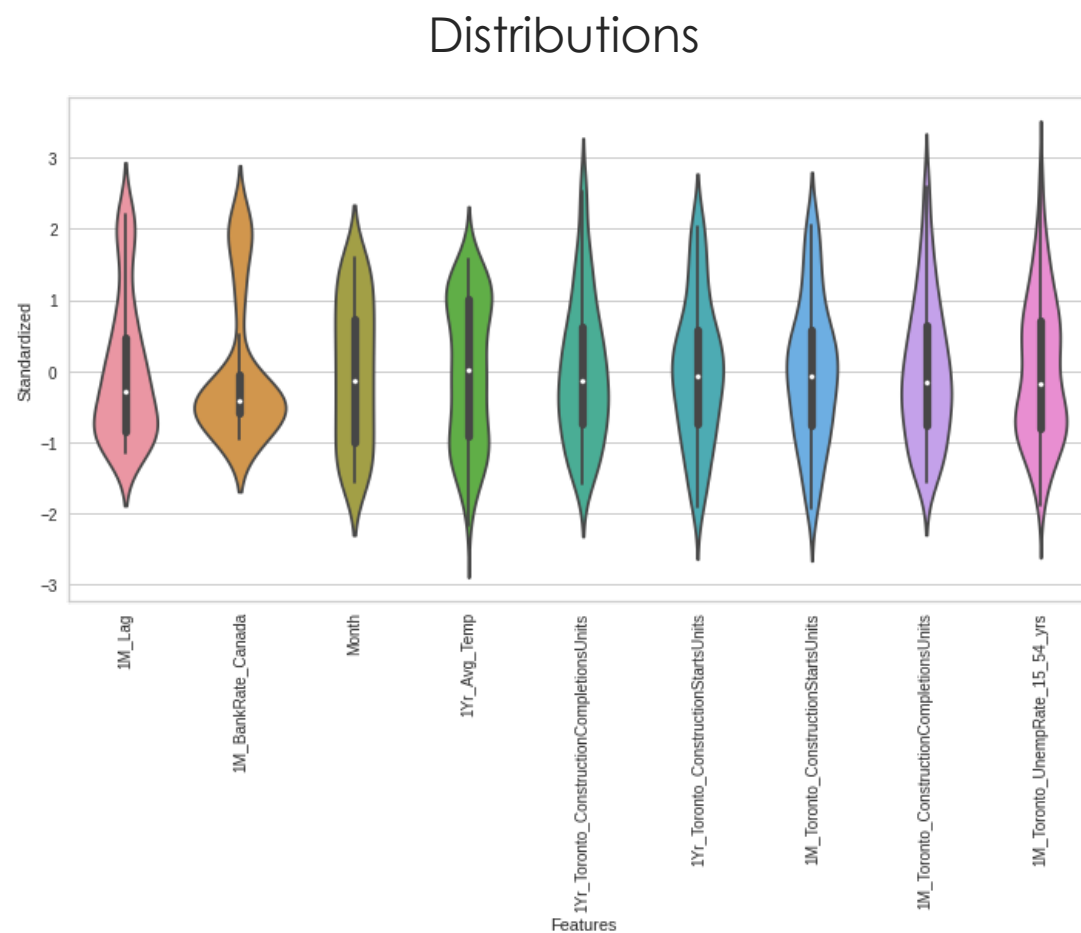
| Eliminated   | Eliminated   | Selected  |
|--|--|---|
| <ul style="list-style-type: none"><li>• Only 2 features<ul style="list-style-type: none"><li>• Last month HPI</li><li>• Weather</li></ul></li><li>• 0.980 R<sup>2</sup></li><li>• Skeptic during economic shifts</li></ul> | <ul style="list-style-type: none"><li>• Gave coefficients for all features selected</li><li>• High performance</li><li>• 0.984 R<sup>2</sup> on test</li><li>• Few hyperparameters to tune over time</li></ul> | <ul style="list-style-type: none"><li>• Gave coefficients for all features selected</li><li>• High performance</li><li>• 0.984 R<sup>2</sup> on test</li><li>• Many hyperparameters</li></ul> |

# Multicollinearity removed

Correlations

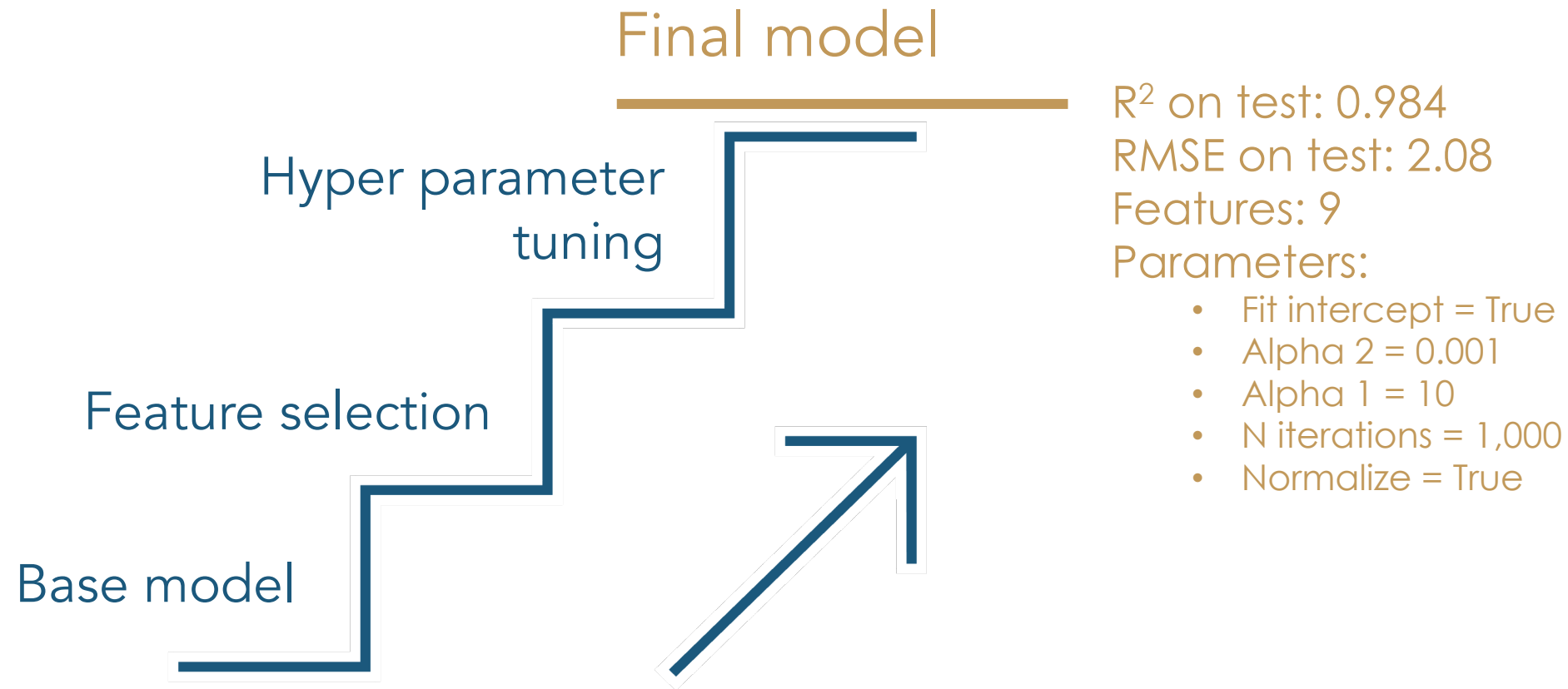


# Selected features are nicely distributed





# Bayesian regression evolution



# Parsimonious equation (rounded & standardized)

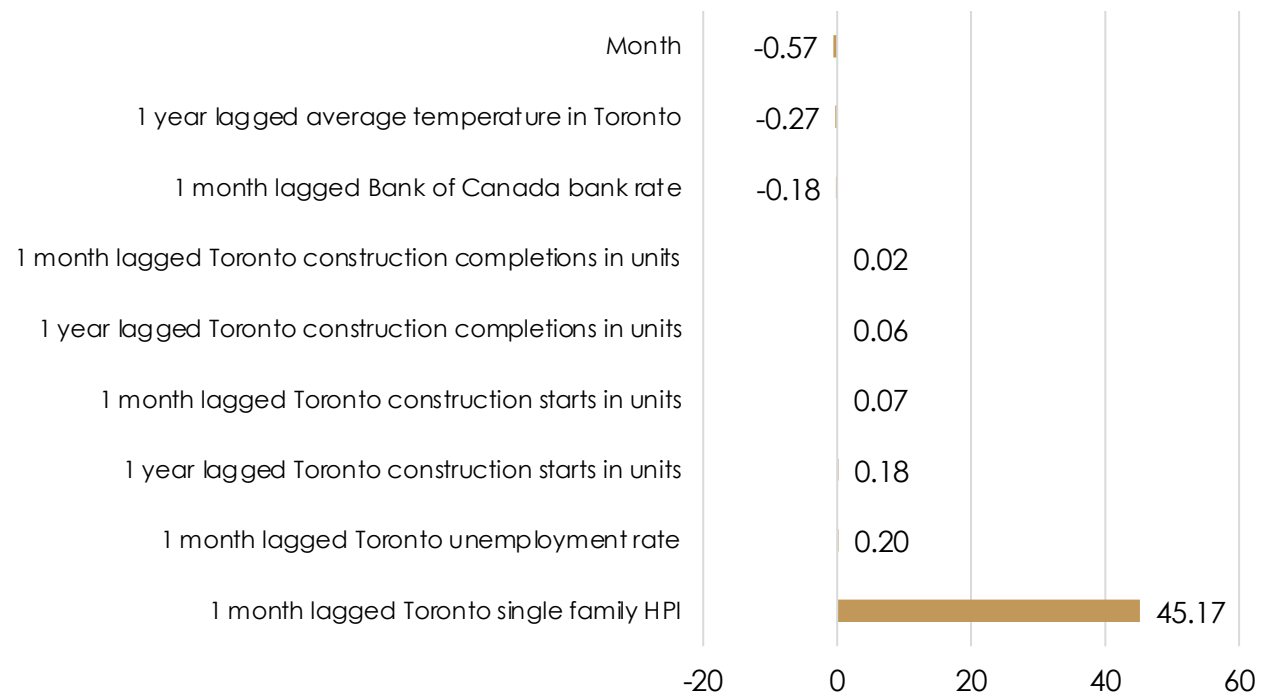


## Monthly HPI Prediction

No change in  $R^2$  or RMSE  
whether standardized or not  
(normalized in both cases)

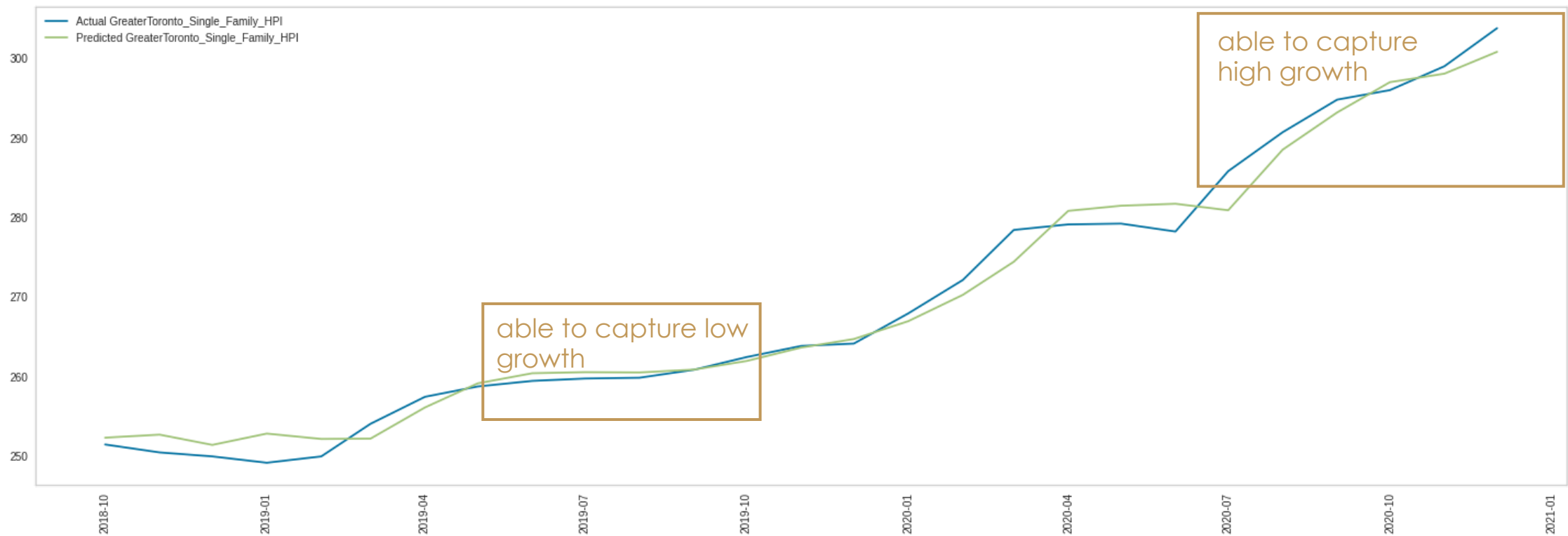
=

Coefficient with intercept of 158.36



# Model output on test data

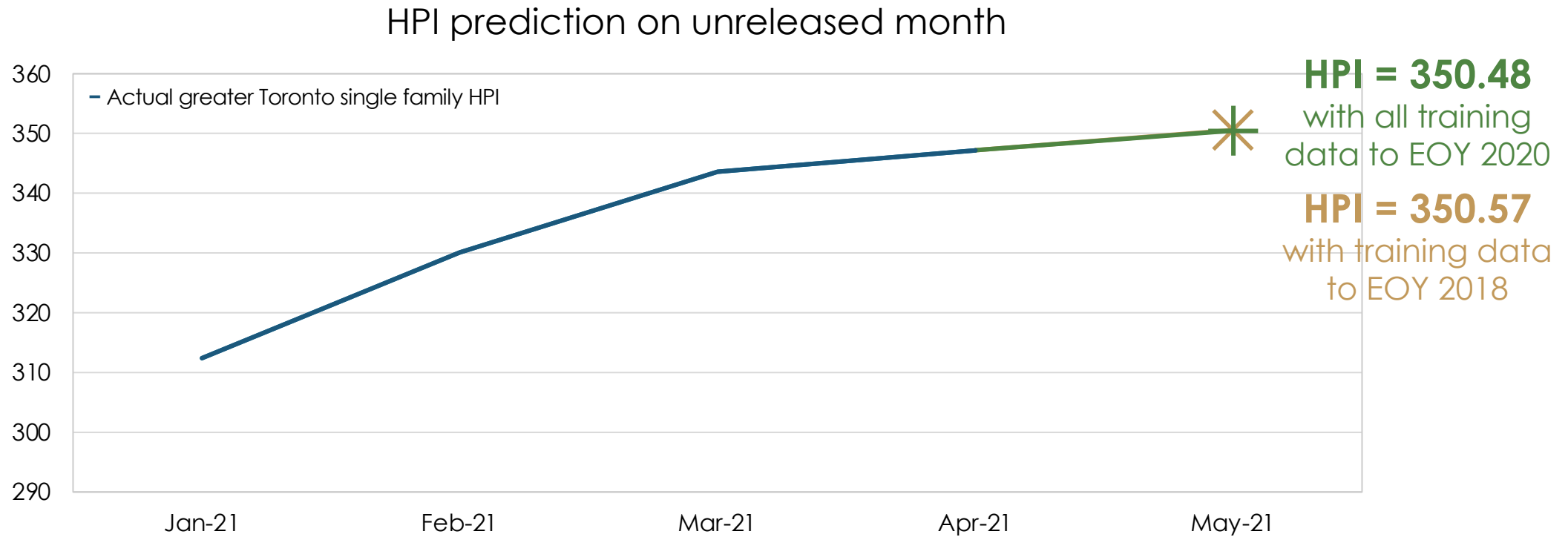
Prediction vs Actual Results



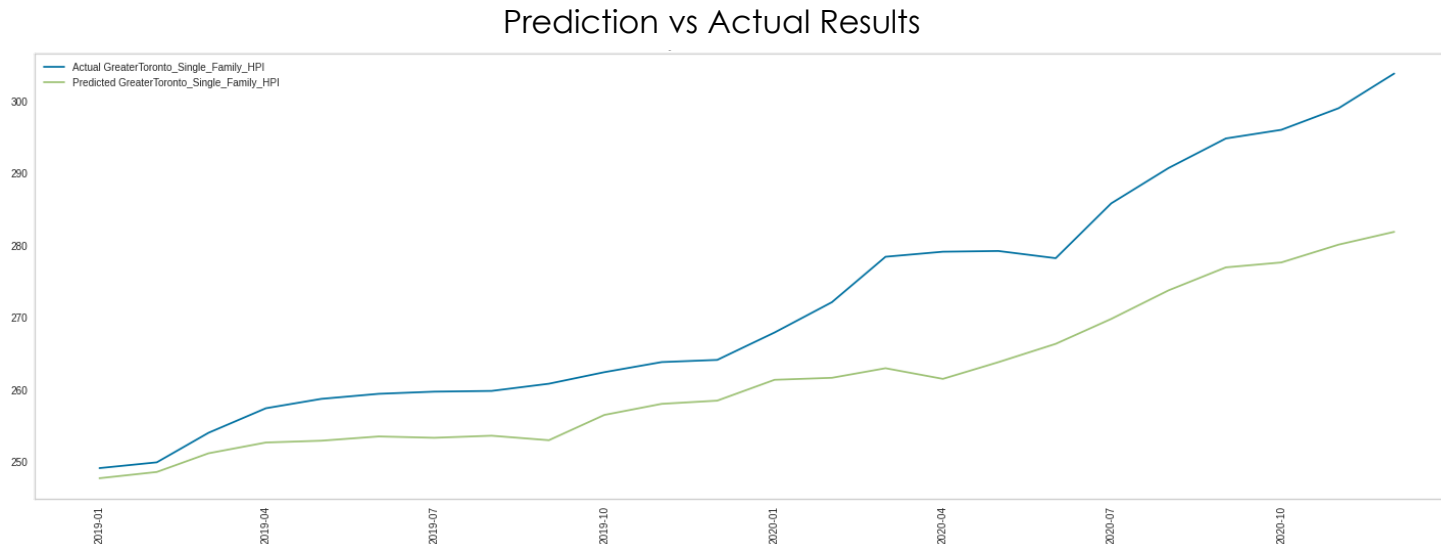
# Prediction for May (unreleased) using model on normalized but not standardized data

| Feature   | Coefficient | Value  | Product       |
|---|-------------|--------|---------------|
| Intercept   | -1.12       | 1      | -1.12         |
| 1 month lagged Toronto single family HPI                | 1.01        | 347.20 | 350.32        |
| Month   | -0.16       | 5      | -0.82         |
| 1 month lagged bank of Canada bank rate                 | -0.13       | 0.50   | -0.07         |
| 1 month lagged average temperature in Toronto           | -0.03       | 7.89   | -0.22         |
| 1 year lagged construction completion units in Toronto  | 6.71e-05    | 2,942  | 0.20          |
| 1 month lagged construction starts units in Toronto     | 7.31e-05    | 2,802  | 0.21          |
| 1 year lagged construction starts units in Toronto      | 0.00019     | 2,558  | 0.49          |
| 1 month lagged construction completion units in Toronto | 2.55e-05    | 2,822  | 0.08          |
| 1 month lagged unemployment in Toronto                  | 0.18        | 8.50   | 1.51          |
|   |             |        | <b>350.57</b> |

# 2021 single family HPI with 2 predictions showing little to no data drift



# Model with no self lag and standardized data



$R^2$  on test: 0.432  
RMSE on test: 12.02  
Features: all  
Dropped:

- 2016 & 2017



# Recommendations

## **Consumers**

Strategize on best-time to refinance, purchase or sell a home

## **Banks**

Anticipate higher demand for mortgages when lending rate lowers

## **Home Repair/Renovation**

Adjust prices proportional to HPI and plan for material demand

## **Policy Makers**

Education of consumer impacts

## **Mortgage/Real Estate Agents**

Optimize marketing spend and personnel for high seasons

## **Developers**

Adjust prices proportional to HPI index

## **Property Tax**

Evaluate if re-assessment is required in the future

## **Policy Makers**

Proactive action to various changes in factors (such as employment rate)

# Model implications

- Retraining required if **base year for index changes** for either or both HPI and CPI
- Monitor **data drift during abnormal times** and retrain if model predicting poorly
- Can **only predict one month** in advance since using prior month's data
- Relies on **CREA publishing prior month's data** in a timely manner

# Model improvements

- **Additional features:** for example, acquire immigration data from CIC and learn how new immigrants moving to Toronto drives real estate demand
- Using **web scraping** for **sentiment analysis** to determine what public feels about housing trends or to capture policy changes which could influence housing bubbles
- Use similar model to predict HPIs for other cities and home types

Thank you!

# Appendix A: All columns

```
Index(['Date', 'Target_GreaterToronto_Single_Family_HPI', '1Yr_MetroOntario_CPI_All',  
'1M_MetroOntario_CPI_all', '1Yr_MetroOntario_CPI_Shelter', '1M_MetroOntario_CPI_Shelter', 'Month',  
'1Yr_Lag', '1M_Lag', '1Yr_GreaterVancouver_Single_Family_HPI',  
'1M_GreaterVancouver_Single_Family_HPI', '1Yr_Calgary_Single_Family_HPI',  
'1M_Calgary_Single_Family_HPI', '1Yr_Montreal_Single_Family_HPI', '1M_Montreal_Single_Family_HPI',  
'1M_RealEstate_GDP_Trading_AdjustedConstant_BasePrice',  
'1Yr_RealEstate_GDP_Trading_AdjustedConstant_BasePrice', '1Yr_RetailSales_Toronto',  
'1M_RetailSales_Toronto', '1M_Avg_Temperature', '1Yr_Avg_Temp', '1M_Snow_onGround',  
'1Yr_Snow_onGround', '1Yr_Toronto_Median_Annual_Income_Couple_Families',  
'1Yr_Toronto_Median_Annual_Income_Lone_Parent_Families',  
'1M_Toronto_ConstructionCompletionsUnits', '1Yr_Toronto_ConstructionCompletionsUnits',  
'1M_Toronto_ConstructionStartsUnits', '1Yr_Toronto_ConstructionStartsUnits',  
'1M_Toronto_UnempRate_15_54_yrs', '1Yr_BankRate_Canada', '1M_BankRate_Canada', '1M_Population',  
'1Y_S&P500_Close', '1M_S&P500_Close', '1Y_Bachelors_Degree_Ontario_NewRegEducation',  
'1Y_Bachelors_Degree_Ontario_NewRegAll'], dtype='object')
```

# Appendix B: Coefficients trained to 2020

0 intercept

{'fit\_intercept': False, 'alpha\_2': 1e-05, 'alpha\_1': 0.0001}

[('1M\_Lag', 1.008367052199725), ('1M\_BankRate\_Canada', -0.13658420440000074), ('Month', -0.17928942913176657), ('1Yr\_Toronto\_ConstructionCompletionsUnits', 6.61131129536024e-05), ('1M\_Toronto\_ConstructionCompletionsUnits', -0.00014693483912440138), ('1Yr\_Avg\_Temp', -0.018867444940951345), ('1Yr\_Toronto\_ConstructionStartsUnits', 5.027050709300945e-05), ('1M\_Toronto\_ConstructionStartsUnits', 2.974158240514532e-05), ('1M\_Toronto\_UnempRate\_15\_54\_yrs', 0.1852511975687902)]

# Appendix C: Coefficients no self lag on all standardized data

Intercept 150.36969696969837

[('1Yr\_MetroOntario\_CPI\_All', 4.3359339530660534), ('1M\_MetroOntario\_CPI\_all', 1.058312664546516), ('1Yr\_MetroOntario\_CPI\_Shelter', 0.7337013540851673), ('1M\_MetroOntario\_CPI\_Shelter', 2.525820407753989), ('Month', 5.042981919409141), ('1Yr\_GreaterVancouver\_Single\_Family\_HPI', 2.0421408473418263), ('1M\_GreaterVancouver\_Single\_Family\_HPI', 15.499695971193294), ('1Yr\_Montreal\_Single\_Family\_HPI', 3.030558814575963), ('1M\_Montreal\_Single\_Family\_HPI', -0.3703416158921996), ('1M\_RealEstate\_GDP\_Trading\_AdjustedConstant\_BasePrice', 13.097832233141274), ('1Yr\_RealEstate\_GDP\_Trading\_AdjustedConstant\_BasePrice', 0.19156564556744485), ('1Yr\_RetailSales\_Toronto', 2.245223467007025), ('1M\_RetailSales\_Toronto', -1.3489615333798972), ('1M\_Avg\_Temperature', 0.5475206285935481), ('1Yr\_Avg\_Temp', -1.8959158054863605), ('1M\_Snow\_onGround', 0.015421834192089808), ('1Yr\_Snow\_onGround', 0.22586881853668353), ('1M\_Toronto\_ConstructionCompletionsUnits', 0.15368612655643293), ('1Yr\_Toronto\_ConstructionCompletionsUnits', 0.0616036075097103), ('1M\_Toronto\_ConstructionStartsUnits', -0.008778186200790662), ('1Yr\_Toronto\_ConstructionStartsUnits', 0.09451417591770342), ('1M\_Toronto\_UnempRate\_15\_54\_yrs', 0.6475051730984084), ('1Yr\_BankRate\_Canada', -2.290543693106973), ('1M\_BankRate\_Canada', 1.325407759711066), ('1M\_Population', -74.56530972420299), ('1Y\_S&P500\_Close', -2.021108263069027), ('1M\_S&P500\_Close', 3.151273147348706), ('1Y\_Bachelors\_Degree\_Ontario\_NewRegAll', 1.339805027762643), ('Year', 68.19038720103602)]



# Appendix D: Normality no self lag

