



# Deep Learning

## Housing Price Prediction

**U Kang**  
**Seoul National University**



# In This Lecture

- House pricing data
- Time series prediction with LSTM network



# Outline

- ➡ ☐ **Problem Definition**
- ☐ Preprocessing Codes



# Dataset (1)

- Zillow Economic Data
- Housing and economic data from a variety of public and proprietary sources by Zillow



# Dataset (2)

- We will focus on these files
  - 'City\_time\_series.csv'
  - 'cities\_crosswalk.csv'
- In EDA part, we will explore each file
- In prediction part, we will use  
'ZHVIPerSqft\_AllHomes' column in  
'City\_time\_series.csv'
  - Mean value of the all homes per square feet



# Problem Definition

- We will extract time series house pricing data and predict future price
- Input will be series of past house prices
- Output will be series of future house prices



# Method

- Input data are time series
- We will use LSTM as main predictor
- We will apply techniques for treating time series data
  - Stationarity
  - Making data for supervised learning



# Outline

☒ Problem Definition

 ☐ **Preprocessing Codes**





# Import libraries

- We will use following libraries
  - ❑ Numpy
  - ❑ Pandas
  - ❑ Sklearn
  - ❑ Matplotlib
  - ❑ Seaborn # for prettifying graph
  - ❑ Plotly # for drawing graph



# EDA (1)

- First file: 'City\_time\_series.csv'
- In EDA part, we will focus on next columns
  - 'ZHVIPerSqft\_AllHomes'
  - 'MedianListingPricePerSqft\_AllHomes'
  - 'MedianRentalPricePerSqft\_AllHomes'
  - 'ZHVI\_2bedroom'
  - 'ZHVI\_3bedroom'
  - 'ZHVI\_4bedroom'
  - 'RegionName'



# EDA (2)

## ■ ZHVI means 'Zillow Home Value Index'

```
df_city_time_seris = pd.read_csv('./data/City_time_series.csv')  
df_city_time_seris.head()
```

	Date	RegionName	InventorySeasonallyAdjusted_AllHomes	InventoryRaw_AllHomes	MedianListingPricePerSqft_1Bedro
0	1996-04-30	abbottstownadamspa	NaN	NaN	I
1	1996-04-30	aberdeenbinghamid	NaN	NaN	I
2	1996-04-30	aberdeenhafordmd	NaN	NaN	I
3	1996-04-30	aberdeenmonroems	NaN	NaN	I
4	1996-04-30	aberdeenmoorenc	NaN	NaN	I

5 rows × 81 columns



# EDA (3)

- Second file: 'cities\_crosswalk.csv'
- 'Unique\_City\_ID' column will be used later to draw nation wide map

```
df_cities_crosswalk = pd.read_csv('./data/cities_crosswalk.csv')  
df_cities_crosswalk.head()
```

	Unique_City_ID	City	County	State
0	oak_grovechristianky	Oak Grove	Christian	KY
1	jarvisburgcurritucknc	Jarvisburg	Currituck	NC
2	mcminnvilleyamhillor	McMinnville	Yamhill	OR
3	union_townshiperiepa	Union Township	Erie	PA
4	oshkoshwinnebago wi	Oshkosh	Winnebago	WI

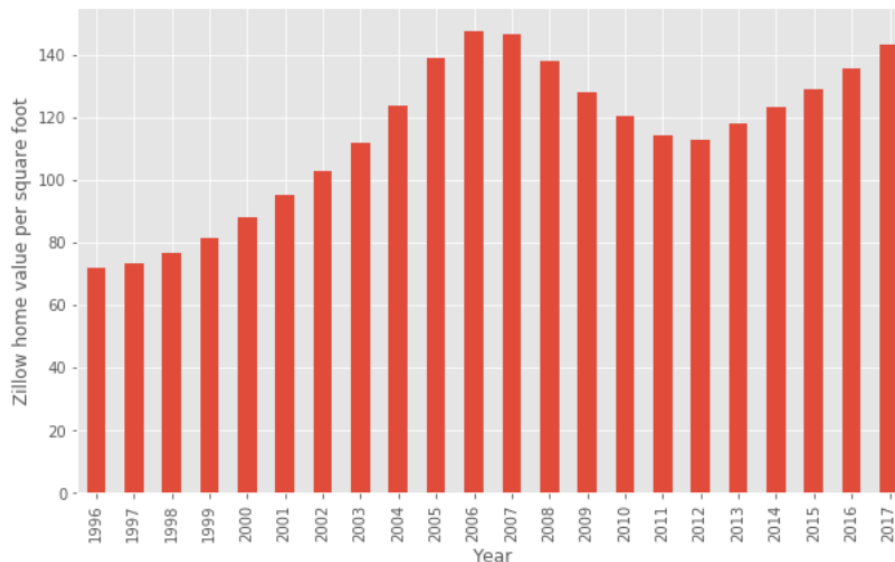


# EDA (4)

- Value of all homes per square in different years
  - Mean of the value of all homes per square foot

```
1 df_city_time_seris.Date = pd.to_datetime(df_city_time_seris.Date)
2 df_city_time_seris.groupby(df_city_time_seris.Date.dt.year)['ZHVIPerSqft_AllHomes'].mean().plot(
3     kind='bar', figsize=(10, 6))
4 plt.suptitle('Mean of the value of all homes per square foot in different year', fontsize=12)
5 plt.ylabel('Zillow home value per square foot')
6 plt.xlabel('Year')
7 plt.show()
```

Mean of the value of all homes per square foot in different year



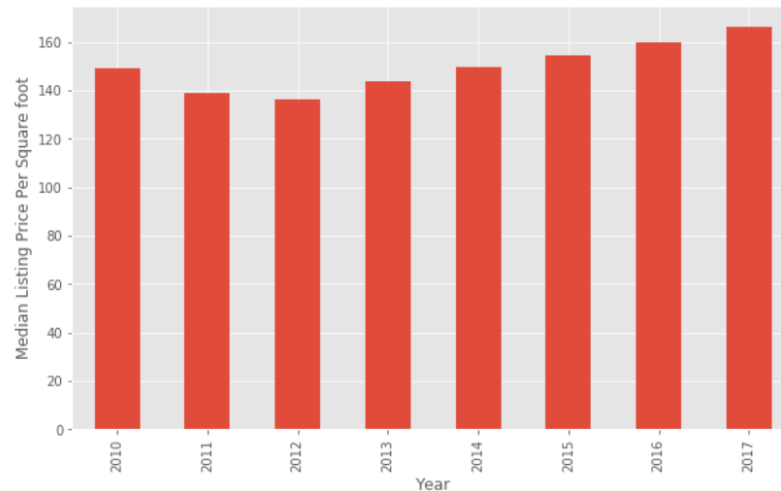


# EDA (5)

- Median of list prices per square foot in different years

```
1 df_city_time_seris_without_null = df_city_time_seris.dropna(  
2     subset=['MedianListingPricePerSqft_AllHomes'], how='any')  
3 df_city_time_seris_without_null \  
4     .groupby(df_city_time_seris_without_null.Date.dt.year)['MedianListingPricePerSqft_AllHomes'] \  
5     .mean().plot(kind='bar', figsize=(10, 6))  
6 plt.suptitle('Median of list prices per square foot in different year', fontsize=24)  
7 plt.ylabel('Median Listing Price Per Square foot')  
8 plt.xlabel('Year')  
9 plt.show()
```

Median of list prices per square foot in different year



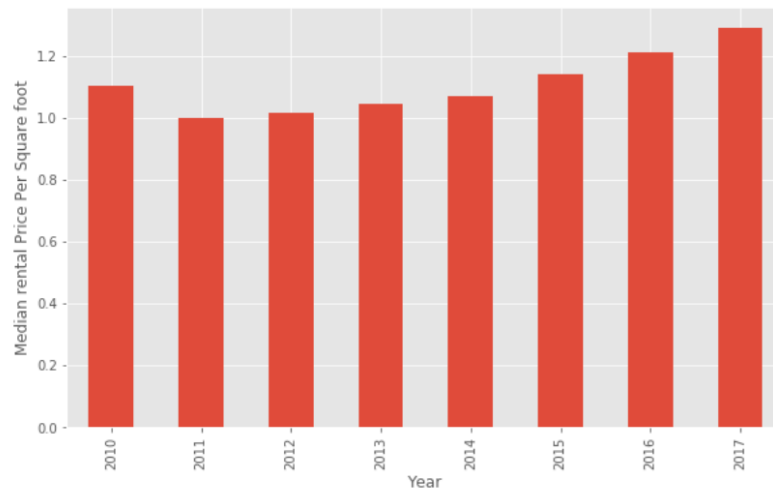


# EDA (6)

## ■ Median of rental prices per square foot in different years

```
1 df_city_time_seris_without_null_rent = df_city_time_seris.dropna(  
2     subset=['MedianRentalPricePerSqft_AllHomes'], how='any')  
3 df_city_time_seris_without_null_rent.groupby(df_city_time_seris_without_null_rent.Date.dt.year) \  
4     ['MedianRentalPricePerSqft_AllHomes'].mean().plot(kind='bar', figsize=(10, 6))  
5 plt.suptitle('Median of rental prices per square foot in different year', fontsize=24)  
6 plt.ylabel('Median rental Price Per Square foot')  
7 plt.xlabel('Year')  
8 plt.show()
```

Median of rental prices per square foot in different year



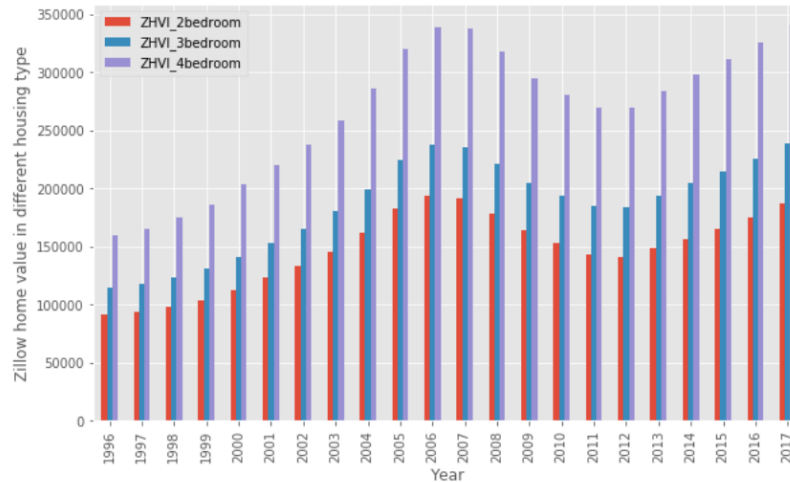


# EDA (7)

- Zillow's different home values in different years
  - Values vary based on housing types, e.g. 2,3,4 bed rooms

```
1 df_city_time_seris.groupby(df_city_time_seris.Date.dt.year) \  
2   [['ZHVI_2bedroom', 'ZHVI_3bedroom', 'ZHVI_4bedroom']].mean().plot(kind='bar', figsize=(10, 6))  
3 plt.suptitle("Zillow's different home value in different year", fontsize=24)  
4 plt.ylabel('Zillow home value in different housing type')  
5 plt.xlabel('Year')  
6 plt.show()
```

Zillow's different home value in different year







# EDA (8)

- Median of the value of all homes per square foot in different states

```
1 # let's replace the regionName column value with State name from cities_crosswalk.csv
2 df_city_time_seris['RegionName'] = df_city_time_seris['RegionName'] \
3     .map(df_cities_crosswalk.set_index('Unique_City_ID')['State'])
4 # group regionName with ZHVIpersqft mean value
5 df_regi_zhvi_sq_mean = df_city_time_seris.groupby(df_city_time_seris.RegionName)['ZHVIPerSqft_AllHomes'] \
6     .mean().reset_index(name = "ZHVIpersqft_mean")
7 # drop null values
8 df_regi_zhvi_sq_mean = df_regi_zhvi_sq_mean.dropna(subset=['ZHVIpersqft_mean'], how='any')
```



# EDA(9)

## ■ Draw graph across the states

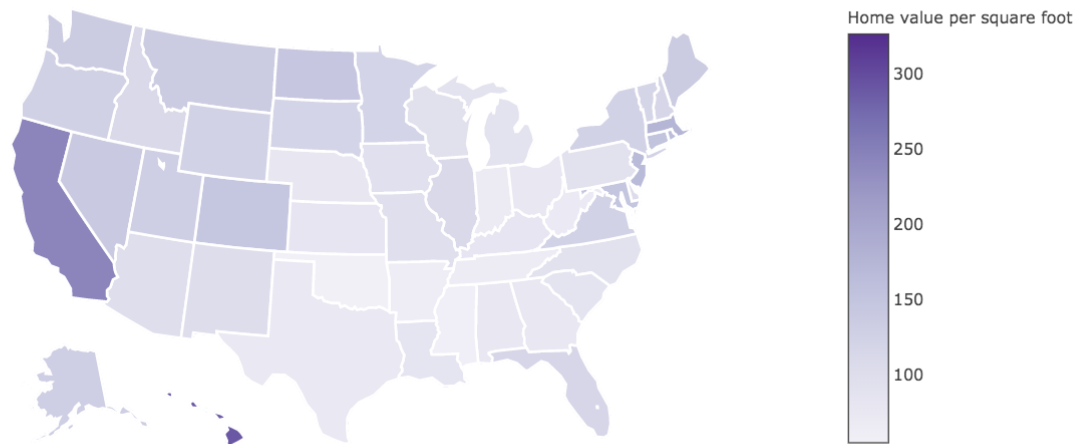
```
1 scl = [  
2     [0.0, 'rgb(242,240,247)'],  
3     [0.2, 'rgb(218,218,235)'],  
4     [0.4, 'rgb(188,189,220)'],  
5     [0.6, 'rgb(158,154,200)'],  
6     [0.8, 'rgb(117,107,177)'],  
7     [1.0, 'rgb(84,39,143)']]  
8  
9  
10 # define our data for plotting  
11 data = [ dict(  
12     type='choropleth',  
13     colorscale = scl,  
14     autocolorscale = False,  
15     locations = df_regi_zhvi_sq_mean['RegionName'], # location (states)  
16     z = df_regi_zhvi_sq_mean['ZHVIpersqft_mean'].astype(float), # Zillow Home value per square foot  
17     locationmode = 'USA-states', # let's define the location mode to USA_states  
18     text = 'Median home value per square foot',  
19     marker = dict(  
20         line = dict (   
21             color = 'rgb(255,255,255)',  
22             width = 2  
23         ) ),  
24     colorbar = dict(  
25         title = "Home value per square foot")  
26     ) ]  
27  
28 layout = dict(  
29     title = 'Median of the value of all homes per square foot in different states<br>(Hover for breakdown)  
30     geo = dict(  
31         scope='usa',  
32         projection=dict( type='albers usa' ),  
33         showlakes = True,  
34         lakecolor = 'rgb(255, 255, 255)',  
35     )  
36  
37  
38 fig = dict( data=data, layout=layout )  
39 # let's plot  
40 py.iplot( fig, filename='d3-cloropleth-map' )
```



# EDA (10)

## ■ Result

Median of the value of all homes per square foot in different states  
(Hover for breakdown)





# Extract feature (1)

- Among many features in 'City\_time\_series.csv', we will use average of 'ZHVIPerSqft\_AllHomes' per month as feature
- Then, our feature is series of floating point numbers containing average house prices over months

```
1 df_city_time_series = pd.read_csv('./data/City_time_series.csv', parse_dates=['Date'])
2 # drop null values in ZHVIPerSqft_AllHomes because we are interested in this column
3 df_city_time_series = df_city_time_series.dropna(subset=['ZHVIPerSqft_AllHomes'])
4 df_city_time_series.head()
```



# Extract feature (2)

- Corresponding codes are below
- This will plot average price graph over years

```
1 # the ZHVIPerSqft_AllHomes column has many value in same date but for different location.
2 # For this notebook we are not interested in location. We mean all the value in same date
3 df_zhvi_sqft_all = df_city_time_series.set_index('Date') \
4     .groupby(pd.Grouper(freq='d')).mean().dropna(how='all') \
5     .ZHVIPerSqft_AllHomes
6
7 fig, ax = plt.subplots(figsize=(15, 10))
8 ax.scatter(df_zhvi_sqft_all.index, df_zhvi_sqft_all)
9 # change x axis year location interval to 1 year. So that it displays data in interval of 1 year
10 ax.xaxis.set_major_locator(mdates.YearLocator(1))
11 # Add the title to the graph
12 plt.title('Zillow Home Value Index in Per Square foot in different year', fontsize=18)
13 # add xlabel
14 plt.xlabel('Year', fontsize=18)
15 # add ylabel
16 plt.ylabel('Zillow Home Value Index in Per Square foot', fontsize=18)
17 # beautify the x axis date presentation
18 fig.autofmt_xdate()
19 # And finally show the plot in a new window.
20 plt.show()
21
```



# Stationarity (1)

- So far, we have a series of prices
- Before feeding in these numbers to LSTM, we need to make this time series stationary
- Stationarity in time series means statistical properties of a process generating a time series do not change over time
- We can give stationarity by using difference between current and next elements in series



# Stationarity (2)

- The following functions can be helpful with respect to stationarity

```
# create a differenced series  
# this is to make time series stationary  
# why stationarity? => https://towardsdatascience.com/stationarity-in-time-series-analysis-90c94f27322  
def difference(dataset, interval=1):  
    diff = list()  
    for i in range(interval, len(dataset)):  
        value = dataset[i] - dataset[i - interval]  
        diff.append(value)  
    return pd.Series(diff)
```



# Make supervised data (1)

- RNN is a network for supervised dataset
- This means that we need both feature and 'label'
- We only have feature (e.g. differenced series data) and not 'label'
- We can convert this series to feature and 'label'





# Make supervised data (2)

- We can frame current value as feature and next value as label
- By doing this, we will predict the very next value from past values

```
# frame a sequence as a supervised learning problem  
# this methods will create a column and column value will be 1 shift from the data.  
# it will make our data to supervised so that we can feed into network  
def timeseries_to_supervised(data, lag=1):  
    df = pd.DataFrame(data)  
    columns = [df.shift(i) for i in range(1, lag+1)]  
    columns.append(df)  
    df = pd.concat(columns, axis=1)  
    df.fillna(0, inplace=True)  
    return df
```



# Make supervised data (3)

## ■ After timeseries\_to\_supervised

- 1<sup>st</sup> column: shifted value
- 2<sup>nd</sup> column: orig. value

	0	0
0	0.000000	-0.034709
1	-0.034709	-0.041769
2	-0.041769	-0.003739
3	-0.003739	-0.011523
4	-0.011523	0.022334
5	0.022334	0.051720
6	0.051720	0.136330
7	0.136330	0.175334
8	0.175334	0.238790
9	0.238790	-0.037459
10	-0.037459	0.308243
11	0.308243	0.216467



# Scaler (1)

- Scaler can convert arbitrary range of values to given range
- This is a kind of standardization method and also can improve training quality



# Scaler (2)

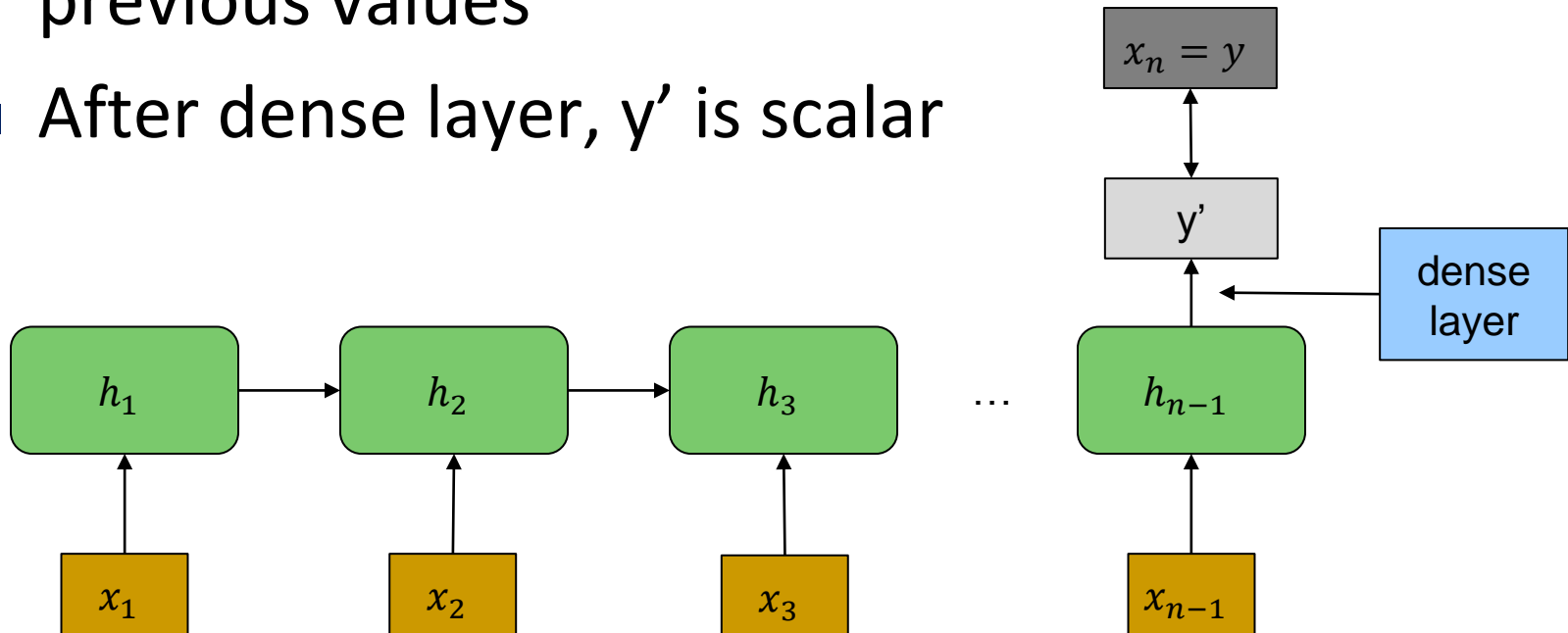
- We will use minmax scaler from sklearn

```
# scale train and test data to [-1, 1]  
def scale(train, test):  
    # fit scaler  
    scaler = MinMaxScaler(feature_range=(-1, 1))  
    scaler = scaler.fit(train)  
    # transform train  
    train_scaled = scaler.transform(train)  
    # transform test  
    test_scaled = scaler.transform(test)  
    return scaler, train_scaled, test_scaled
```



# Model Suggestion

- Our model gets monthly data as input
- As you can see, we predict last value given previous values
- After dense layer,  $y'$  is scalar





# Questions?