

Topic Selection

At first, we find and come up with some interesting topics, including:

- Black Friday prom
- Housing price prediction
- Movie rating/prediction
- TED topic prediction
- Car license plate recognition
- Agriculture
- Google travel rating

Black Friday. Housing Price Movie. TED.

Car Plate. Agriculture. Google travel rating

NY: {B, LC, M}

1.02^v
↓
x 1.03



After discussion, we selected the Housing Price Prediction as the final topic. As this is a popular and familiar topic for us, and there are abundant resources we can refer to on the web. Then we come up with some interesting features that are related to the housing price.

Black Friday.

Housing Price

Movie:

TZD.

- Iowa - HP.
- Sale Price, Lot Area, Sale Condition.
- Year Remod. Add., Bathrooms, Bedrooms.
- Kitchen, Neighborhood.
- housestyle, lotconfig

~~Car Plate.~~

Agriculture.

~~Google travel rating~~

$$\left\{ \text{NY} : \begin{matrix} 1.02 \\ \downarrow \\ \{B, Lc, M\} \\ \times 1.03 \end{matrix} \right.$$

General Procedure

- Data visualization. Select some features manually from the data and visualize the data to get an intuitive recognition of the data set.
- Data preprocessing. Do preprocessing like normalization and filling in the blank on the data

to improve the validation accuracy.

- Model selection.Try to use some models we studied or from external resources to train the data, and use K-fold validation methods to evaluate the error on the test data to select the best model.
- Evaluate the model.Select the best model and retrain the model with some extreme value removed to check whether it improves the performance.
- Generate the summary.

Data Visualization

We selected 9 features that are most relevant to house price from on data set.

- bedrooms
- bathrooms
- sqft_lot
- floors
- waterfront
- condition
- grade
- sqft_above
- sqft_basement

1. Distribution. For every feature, we first compute each distinct value and the number of samples equal to the value in the training data.If there are no more than 20 distinct values, we plot a line graph, otherwise we plot a histogram with 20 blocks.This gives us an intuitive distribution of the feature values.
2. Correlation with the price.For every feature, we plot a scatter graph showing the feature value and the price.Usually all features are in the positive correlation with the price attribute.We can briefly draw a conclusion from the graph which feature has more weight in deciding the price.
3. Correlation of each feature.Calculate and plot the heatmap of each pair of attributes.And may exclude the features with the least correlation with the price in future training.



