

任务介绍

使用纽约州房地产数据来预测房屋的销售价格。主要包括以下几部分：

1. 数据下载、存储、读取
2. 数据分析
3. 拆分数据集
4. 数据清洗
5. 创建管道与模型训练
6. 模型评估
7. 核心代码

数据下载、存储、读取

- 设置下载地址与保存地址
- 下载tar文件并解压
- 数据读取为DataFrame

```
1 import os
2 import tarfile
3 import urllib.request
4 # download data setting
5 DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml/master/"
6 HOUSING_PATH = os.path.join("datasets", "housing")
7 HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tgz"
```

```
1 # 下载tar文件并解压
2 def fetch_housing_data(housing_url=HOUSING_URL, housing_path=HOUSING_PATH):
3     os.makedirs(housing_path, exist_ok=True)
4     tgz_path = os.path.join(housing_path, "housing.tgz")
5     # stores data on local disk
6     urllib.request.urlretrieve(housing_url, tgz_path)
7     housing_tgz = tarfile.open(tgz_path)
8     housing_tgz.extractall(path=housing_path)
9     housing_tgz.close()
10 #fetch_housing_data()
```

```
1 # 数据读取
2 import pandas as pd
3 import numpy as np
4 # load data
5 def load_housing_data(path = HOUSING_PATH):
6     csv_path = os.path.join(path, 'housing.csv')
7     return pd.read_csv(csv_path)
```

数据分析

- 数据整体情况
 - 通过housing.info()查看数据字段情况
 - 查看数据样例
 - 查看非数字字段“ocean_proximity”的情况
 - housing.describe()查看数据的分布情况

```
1 housing = load_housing_data()
2 housing.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 20640 entries, 0 to 20639
3 Data columns (total 10 columns):
4  #   Column                Non-Null Count  Dtype
5  ---  ---
6  0   longitude              20640 non-null  float64
7  1   latitude               20640 non-null  float64
8  2   housing_median_age     20640 non-null  float64
9  3   total_rooms            20640 non-null  float64
10  4   total_bedrooms        20433 non-null  float64
11  5   population             20640 non-null  float64
12  6   households             20640 non-null  float64
13  7   median_income          20640 non-null  float64
14  8   median_house_value     20640 non-null  float64
15  9   ocean_proximity        20640 non-null  object
16 dtypes: float64(9), object(1)
17 memory usage: 1.6+ MB
```

```
1 housing.head()
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	m
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	45
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	35
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	35
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	34
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	34

```
1 housing["ocean_proximity"].value_counts()
```

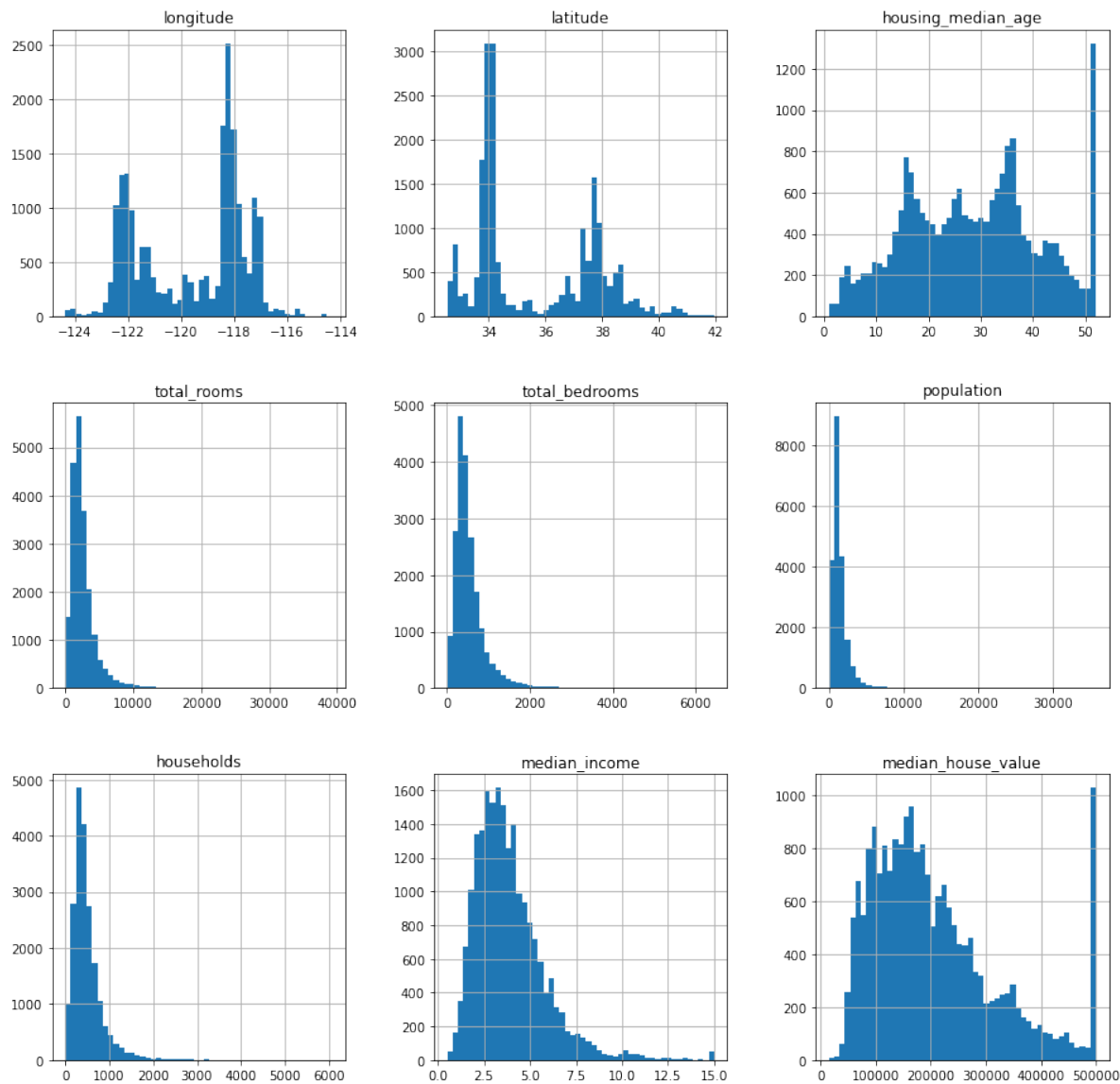
```
1 <1H OCEAN      9136
2 INLAND        6551
3 NEAR OCEAN    2658
4 NEAR BAY      2290
5 ISLAND         5
6 Name: ocean_proximity, dtype: int64
```

```
1 housing.describe()
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000
mean	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.870671
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.534800
75%	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100

- 字段数据分布可视化

```
1 import matplotlib.pyplot as plt
2 housing.hist(bins=50,figsize=(15,15))
3 plt.show()
```



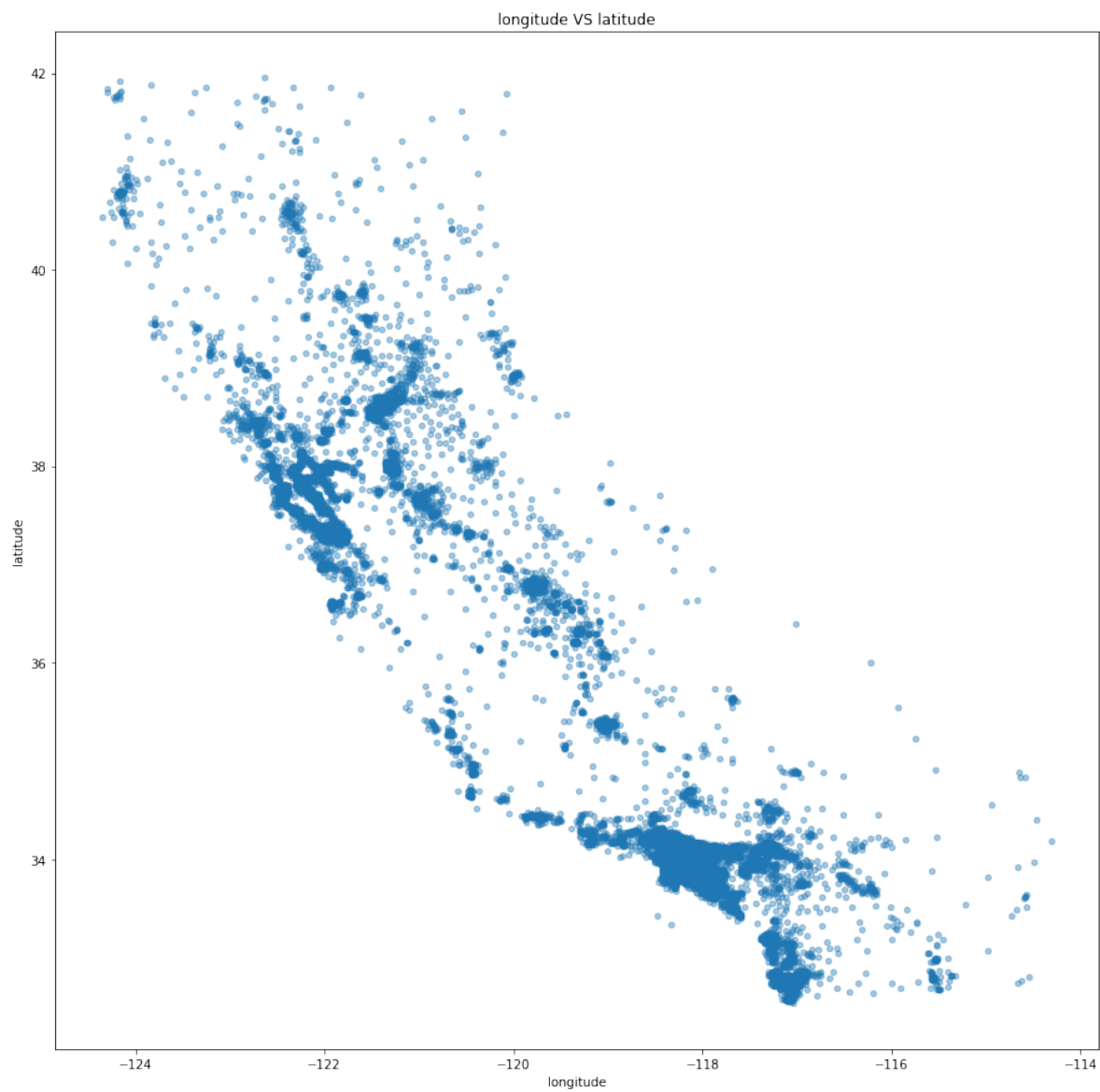
• 可视化分析

- 查看按照经纬度数据的分布情况
- 经纬度与房价分布情况
- 经纬度、收入与房价分布情况
- 主要特征之间的相关性可视化
- 各主要特征与房价的相关性可视化

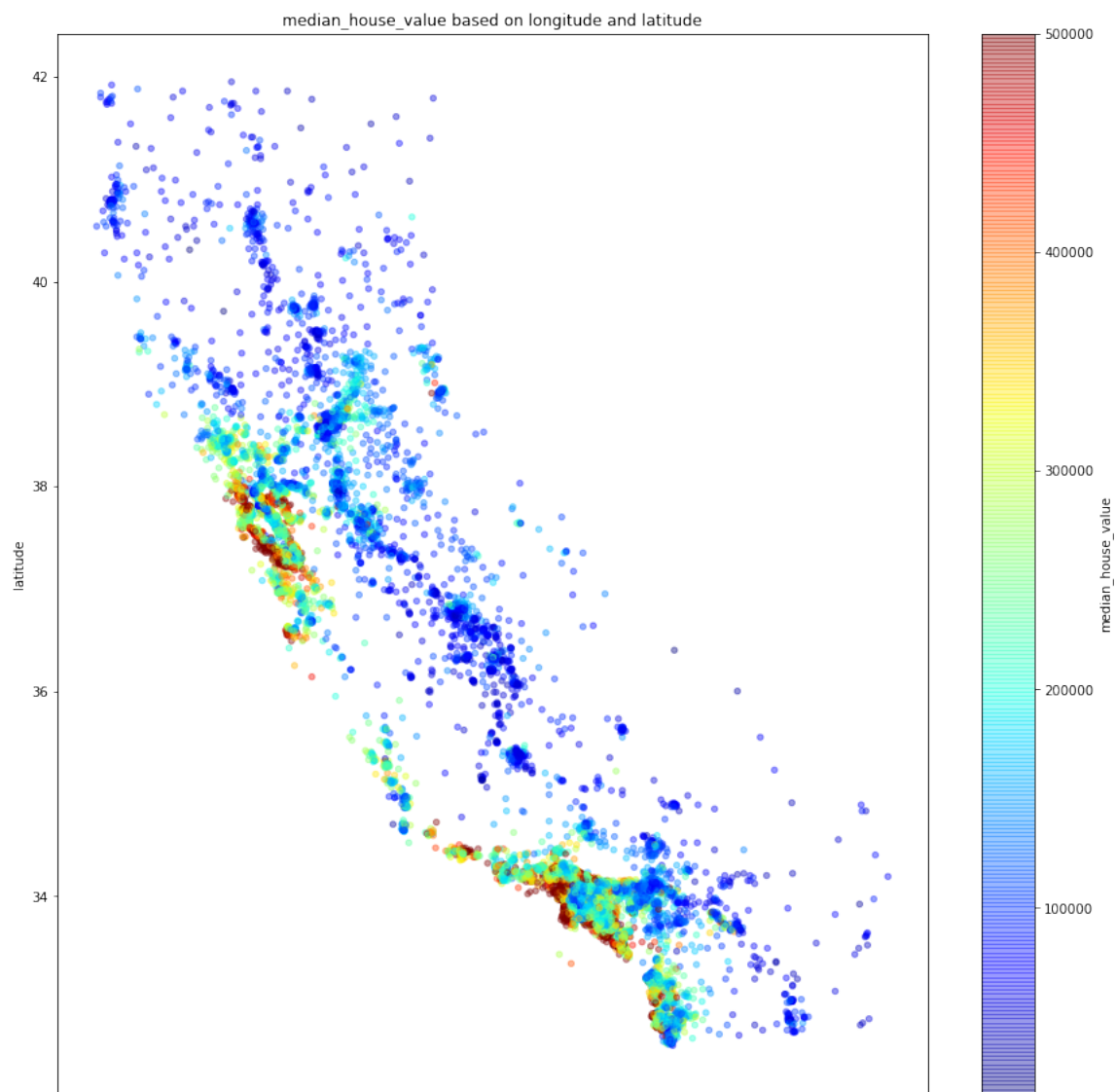
```

1 # 数据可视化
2 housing = strat_train_set.copy()
3 # 经纬度
4 housing.plot(figsize=(15,15),kind='scatter',x='longitude',y='latitude',alpha=0.4,title="longitude VS latitude")
5 plt.show()

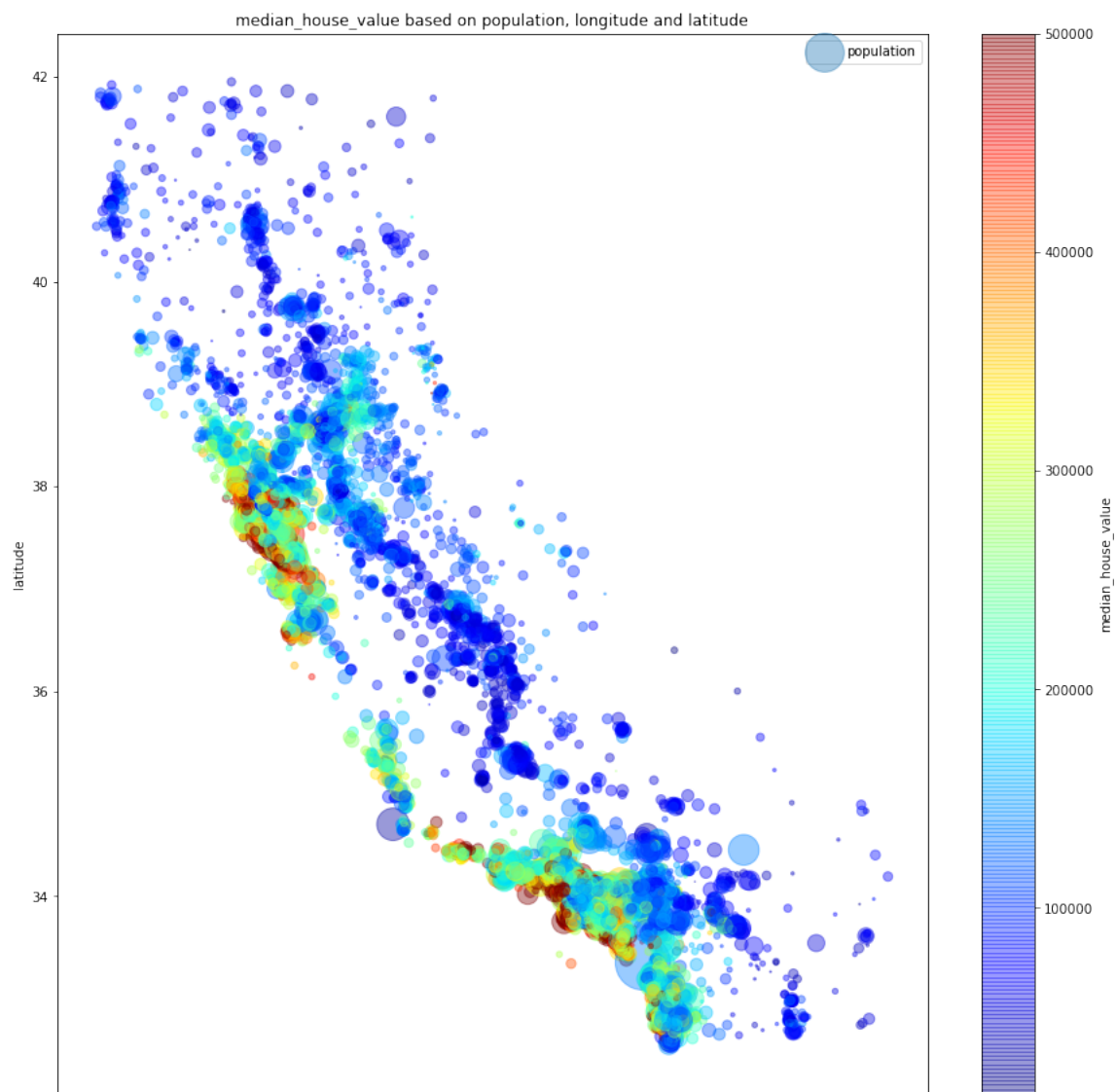
```



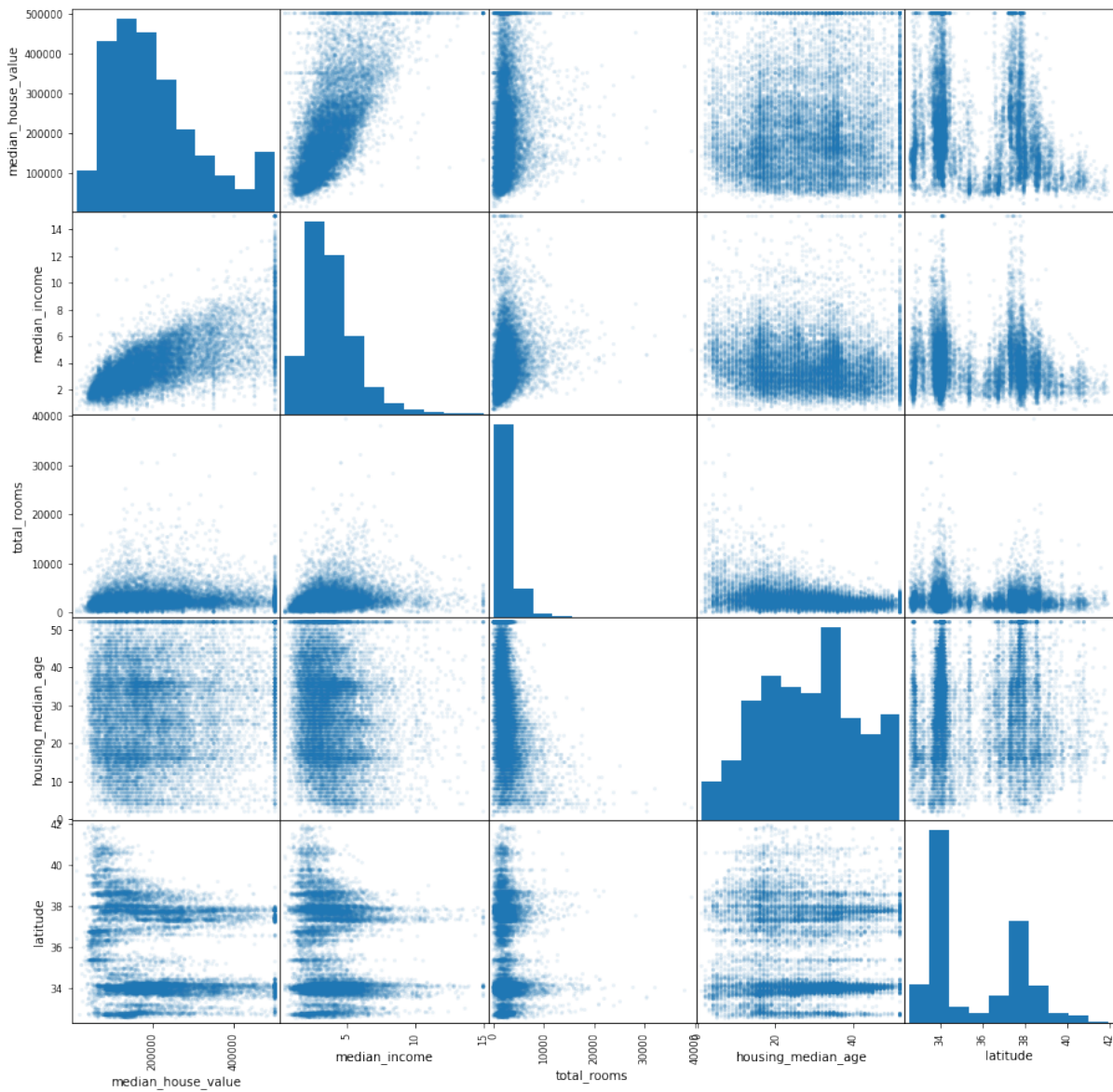
```
1 # 经纬度、房价
2 housing.plot(figsize=(15,15),kind='scatter',x='longitude',y='latitude',alpha=0.4,c = 'median_house_value',
3             cmap=plt.get_cmap("jet"),colorbar = True,title="median_house_value based on longitude and latitude")
4 plt.show()
5
```



```
1 # 经纬度、收入、房价
2 housing.plot(figsize=(15,15),kind='scatter',x='longitude',y='latitude',alpha=0.4,
3             s=housing['population']/20,label='population', c = 'median_house_value',cmap=plt.get_cmap("jet"),colorbar = True,
4             title="median_house_value based on population, longitude and latitude")
5 plt.show()
```



```
1 # 主要特征之间的相关性
2 from pandas.plotting import scatter_matrix
3 attributes = ["median_house_value", "median_income", "total_rooms", "housing_median_age", "latitude"]
4 scatter_matrix(housing[attributes], figsize=(15,15), alpha=0.1)
5 plt.show()
6
```



```

1 # 计算相关系数
2 corr_matrix = housing.corr()
3 corr_matrix['median_house_value'].sort_values(ascending=False)

```

```

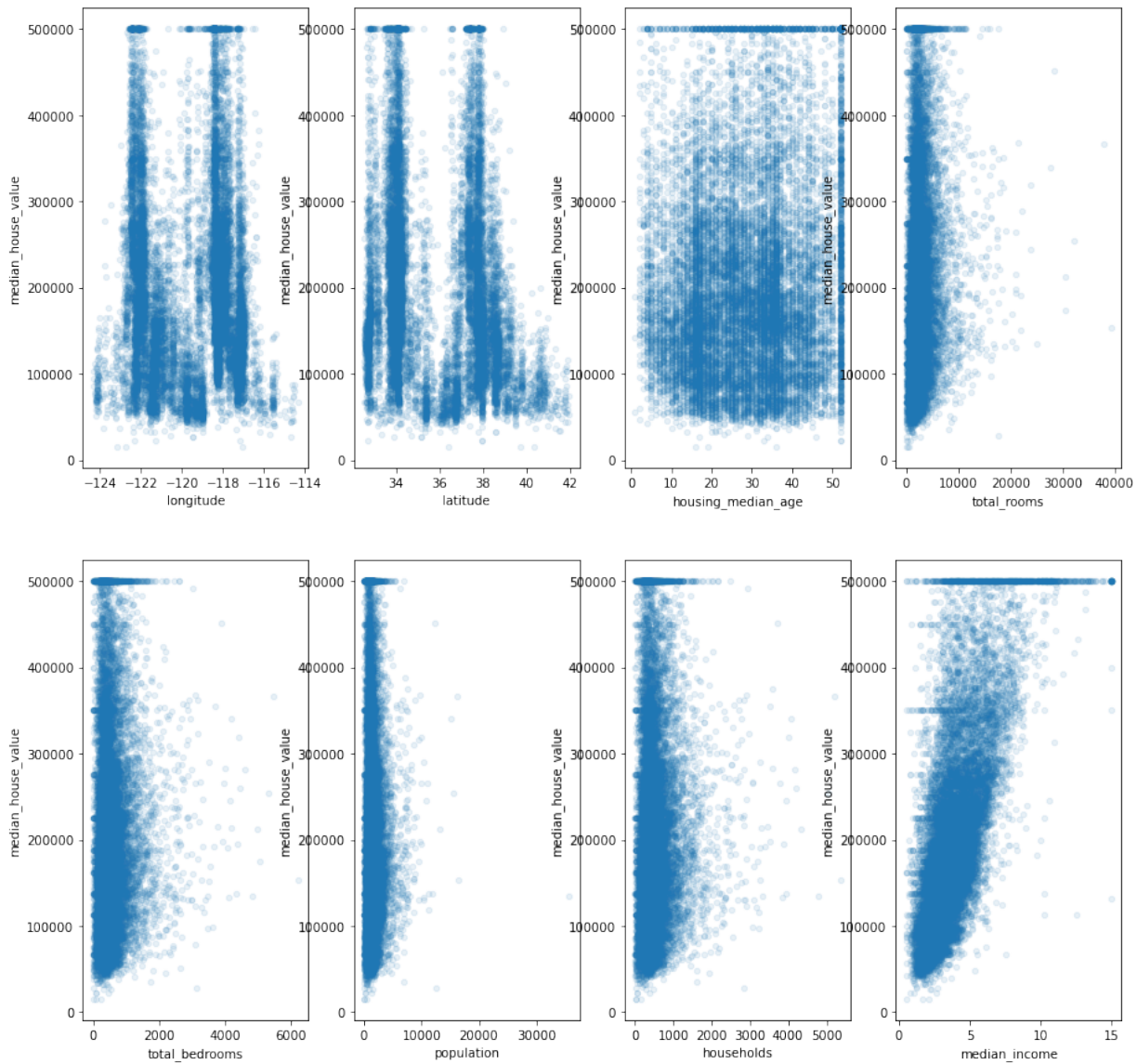
1 median_house_value    1.000000
2 median_income         0.687160
3 total_rooms           0.135097
4 housing_median_age    0.114110
5 households            0.064506
6 total_bedrooms        0.047689
7 population            -0.026920
8 longitude             -0.047432
9 latitude              -0.142724
10 Name: median_house_value, dtype: float64

```

```

1 # 主要特征与房价的相关性
2 fig, axes = plt.subplots(2, 4, figsize=(15, 15))
3 axes = axes.ravel()
4 columns = list(housing.columns)
5 columns.remove("median_house_value")
6 columns.remove("ocean_proximity")
7 for index in range(len(columns)):
8     housing.plot(figsize=(15, 15), kind='scatter', x=columns[index], y='median_house_value', alpha=0.1, ax=axes[index])
9 plt.show()
10

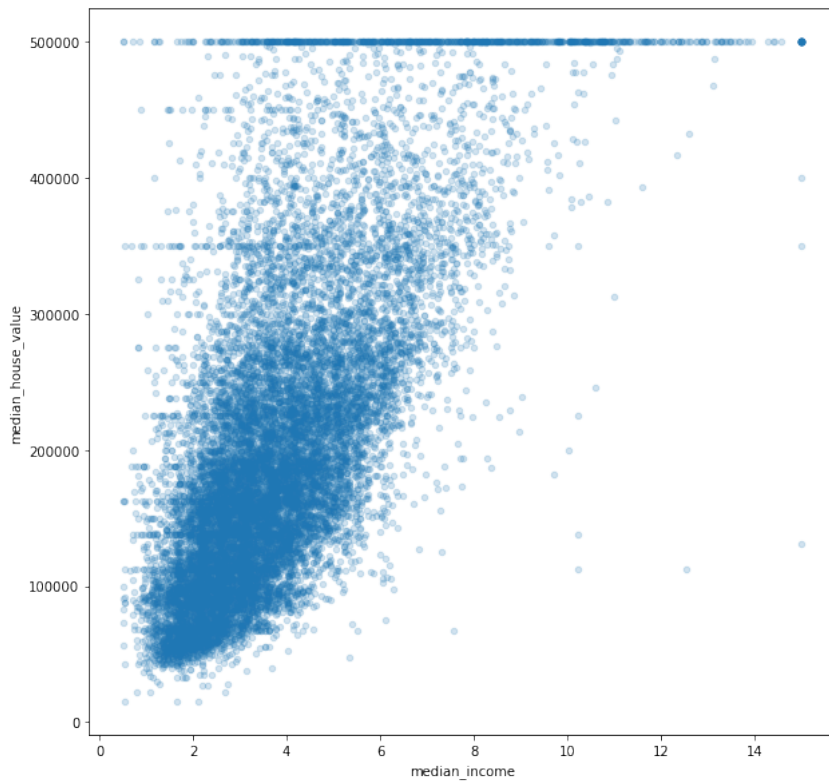
```



```
1 | columns
```

```
1 | ['longitude',
2 |  'latitude',
3 |  'housing_median_age',
4 |  'total_rooms',
5 |  'total_bedrooms',
6 |  'population',
7 |  'households',
8 |  'median_income']
```

```
1 | # 收入与房价的相关性
2 | housing.plot(figsize=(10,10),kind='scatter',x='median_income',y='median_house_value',alpha=0.2)
3 | plt.show()
```

数据拆分

- 生成收入等级
 - 将收入进行缩放，并取整分级
 - 增加一列income_cat临时存储收入的分级
- 数据拆分
 - 使用StratifiedShuffleSplit函数，构建随机拆分器
 - 根据收入等级来将数据分组，在每组随机抽样
 - 将数据分成测试集和验证集
- 删除临时保存的收入等级数据

```
1 # 分层抽样
2 from sklearn.model_selection import StratifiedShuffleSplit
3 housing['income_cat'] = np.ceil(housing['median_income']/1.5)
4 housing['income_cat'].where(housing['income_cat']<5,5,inplace=True)
5 split = StratifiedShuffleSplit(n_splits=1,test_size=0.2,random_state=42)
6 for train_index,test_index in split.split(housing,housing['income_cat']):
7     strat_train_set = housing.iloc[train_index]
8     strat_test_set = housing.iloc[test_index]
9 strat_train_set.drop(["income_cat"],axis =1, inplace = True)
10 strat_test_set.drop(["income_cat"],axis =1, inplace = True)
11
```

数据清洗

- 数值型数据
 - 缺失数据填充
 - 使用SimpleImputer，通过中位数策略进行填充
 - 补充特征，添加户型平均室数和人均房屋数
 - 使用标准差对数据进行缩放
- 文本型数据
 - 独热编码进行数字化

```

1 # 中位数补充缺失值
2 housing = strat_train_set.drop("median_house_value", axis=1) # drop labels for training set
3 housing_labels = strat_train_set["median_house_value"].copy()
4 from sklearn.impute import SimpleImputer
5 imputer = SimpleImputer(strategy="median")
6 housing_num = housing.drop("ocean_proximity", axis=1)
7 imputer.fit(housing_num)
8 housing_num.median()
9 X = imputer.transform(housing_num)
10 X.shape
11

```

```

1 (16512, 8)

```

```

1 # 查看补充后数据
2 housing_tr = pd.DataFrame(X, columns=housing_num.columns, index=housing_num.index)
3 # columns=housing_num.columns
4 housing_tr.head()

```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
17606	-121.89	37.29	38.0	1568.0	351.0	710.0	339.0	2.7042
18632	-121.93	37.05	14.0	679.0	108.0	306.0	113.0	6.4214
14650	-117.20	32.77	31.0	1952.0	471.0	936.0	462.0	2.8621
3230	-119.61	36.31	25.0	1847.0	371.0	1460.0	353.0	1.8839
3555	-118.59	34.23	17.0	6592.0	1525.0	4459.0	1463.0	3.0347

```

1 # 处理文本数据: oneHot_encoder
2 housing_cat = housing[['ocean_proximity']]
3 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
4 # encoder = LabelEncoder()
5 # housing_cat_encoded = encoder.fit_transform(housing_cat)
6 # # housing_cat_encoded
7 oneHot_encoder = OneHotEncoder()
8 housing_cat_lhot = oneHot_encoder.fit_transform(housing_cat)
9 housing_cat_lhot.shape
10 # housing_cat_lhot[1].toarray()
11 # from sklearn.preprocessing import La
12

```

```

1 array([[1., 0., 0., 0., 0.]])

```

```

1 # 处理文本数据: LabelBinarizer
2 from sklearn.preprocessing import LabelBinarizer
3 encoder = LabelBinarizer(sparse_output=True)
4 housing_cat_lhot = encoder.fit_transform(housing_cat)
5 housing_cat_lhot

```

```

1 <16512x5 sparse matrix of type '<class 'numpy.int64'>'
2 with 16512 stored elements in Compressed Sparse Row format>

```

```

1 # 自定义转换器
2 ## 将Dataframe 转换为Numpy数组
3 from sklearn.base import BaseEstimator, TransformerMixin
4 class DataFrameSelector(BaseEstimator,TransformerMixin):
5     def __init__(self,attributes):
6         self.attributes = attributes
7     def fit(self,X,y=None):
8         return self
9     def transform(self,X):
10        return X[self.attributes].values

```

```

1 # 添加新的特征
2 from sklearn.preprocessing import FunctionTransformer
3 rooms_ix, bedrooms_ix, population_ix, household_ix = [
4     list(housing.columns).index(col)
5     for col in ("total_rooms", "total_bedrooms", "population", "households")]
6
7 def add_extra_features(X, add_bedrooms_per_room=True):
8     rooms_per_household = X[:, rooms_ix] / X[:, household_ix]
9     population_per_household = X[:, population_ix] / X[:, household_ix]
10    if add_bedrooms_per_room:
11        bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
12        return np.c_[X, rooms_per_household, population_per_household,
13                    bedrooms_per_room]
14    else:
15        return np.c_[X, rooms_per_household, population_per_household]
16
17 attr_adder = FunctionTransformer(add_extra_features, validate=False,
18                                 kw_args={"add_bedrooms_per_room": False})
19 housing_extra_attribs = attr_adder.fit_transform(housing.values)
20 housing_extra_attribs = pd.DataFrame(
21     housing_extra_attribs,
22     columns=list(housing.columns)+["rooms_per_household", "population_per_household"])
23 housing_extra_attribs.head()

```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	ocean_proximity
0	-121.89	37.29	38	1568	351	710	339	2.7042	<1H OCEAN
1	-121.93	37.05	14	679	108	306	113	6.4214	<1H OCEAN
2	-117.2	32.77	31	1952	471	936	462	2.8621	NEAR OCEAN
3	-119.61	36.31	25	1847	371	1460	353	1.8839	INLAND
4	-118.59	34.23	17	6592	1525	4459	1463	3.0347	<1H OCEAN

创建管道与模型训练

通过管道方式对数据清洗和训练过程进行封装

- 数值类特征管道
- 文本类特征管道
- 合并管道
- 训练管道
- 模型选择
 - 线性回归
 - 决策树

```

1 # 数据清洗管道-数字类
2 from sklearn.pipeline import FeatureUnion,Pipeline
3 from sklearn.preprocessing import StandardScaler
4 # 增加两列, 9-》11
5 num_pipeline = Pipeline([
6     ('imputer',SimpleImputer(strategy="median")),
7     ("attribs_adder",FunctionTransformer(add_extra_features,validate=False)),
8     ('stander',StandardScaler())
9 ])
10 #数据清洗管道-文本类
11 #ColumnTransformer vs FeatureUnion
12 # https://blog.csdn.net/fendouaini/article/details/109211421
13 from sklearn.compose import ColumnTransformer
14 num_attribs = list(housing_num.columns)
15 # 从一列, 变为5列

```

```

16 cat_attribs = ["ocean_proximity"]
17 full_pipeline = ColumnTransformer([
18     ('num', num_pipeline, num_attribs),
19     ('cat', OneHotEncoder(), cat_attribs)
20 ])

```

```

1 housing_prepared = full_pipeline.fit_transform(housing)
2 housing_prepared.shape

```

```

1 (16512, 16)

```

```

1 # 搭建训练管道
2 from sklearn.linear_model import LinearRegression,
3 from sklearn.metrics import mean_squared_error
4 def train_model(model):
5     train_pipeline = Pipeline([
6         ('datas', full_pipeline),
7         ('model', model)
8     ])
9     return train_pipeline.fit(housing, housing_labels)
10

```

```

1 # 模型选择-线性回归
2 lr_model = train_model(LinearRegression())

```

```

1 # 模型选择-决策树
2 from sklearn.tree import DecisionTreeRegressor
3 dt_model = train_model(DecisionTreeRegressor())
4

```

模型评估

- 预测数据分布
- 标准差
- 模型保存与加载

```

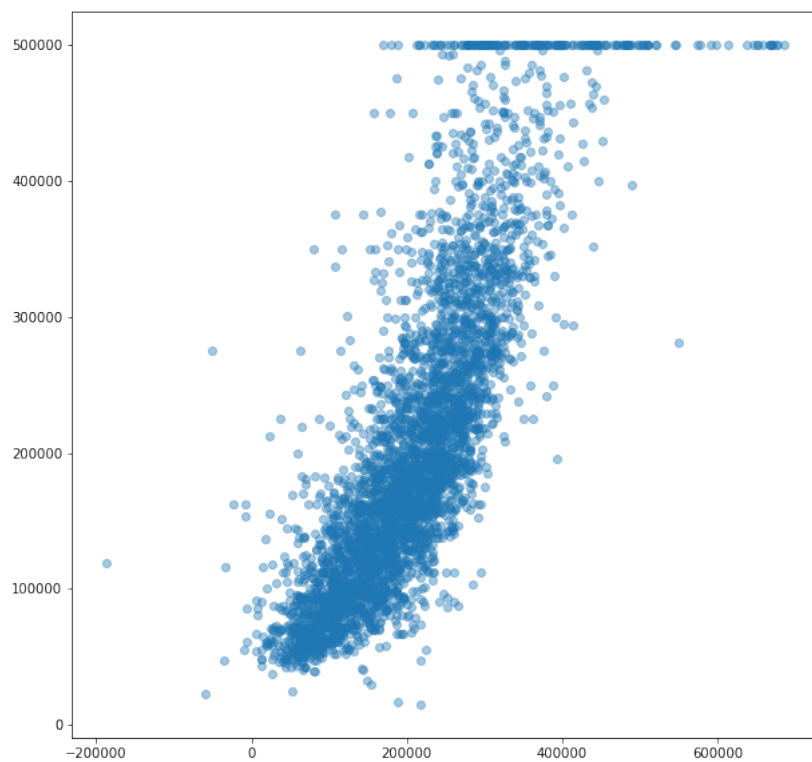
1 # 画出预测值与真实值的分布情况，返回方差
2 test_data = strat_test_set.drop("median_house_value", axis=1) # drop labels for training set
3 test_labels = strat_test_set["median_house_value"].copy()
4 def test(model):
5     test_data = strat_test_set.drop("median_house_value", axis=1) # drop labels for training set
6     test_labels = strat_test_set["median_house_value"].copy()
7     predictions = model.predict(test_data)
8     plt.figure(1, figsize=(10, 10))
9     plt.scatter(x=predictions, y=test_labels, alpha=0.4)
10    plt.xlabel = 'predictions'
11    plt.ylabel = 'ground truth'
12    plt.show()
13    mse = mean_squared_error(test_labels, predictions)
14    return np.sqrt(mse)

```

```

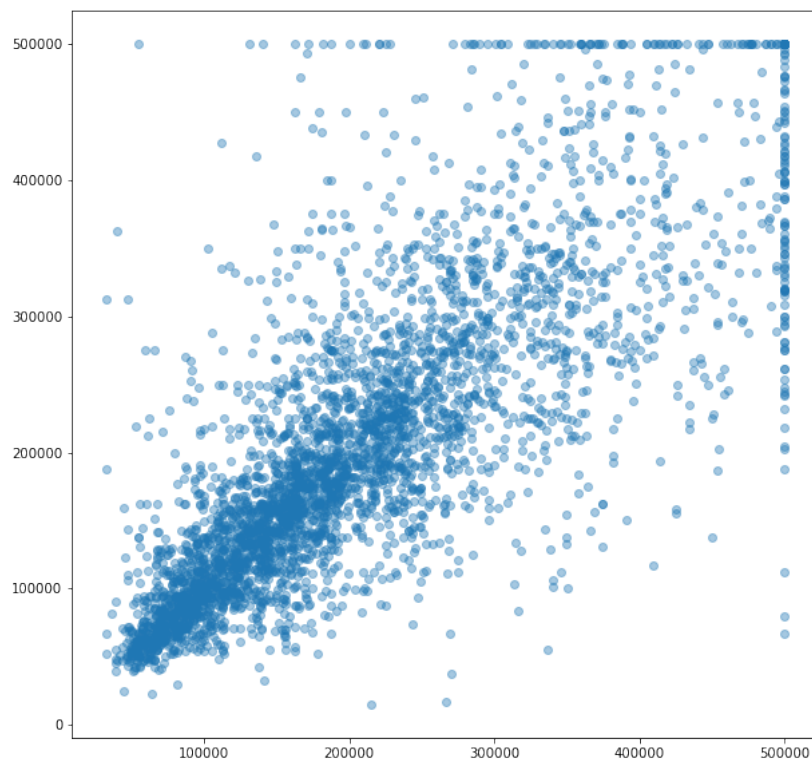
1 test(lr_model)
2

```



```
1 | 66911.98070857547
```

```
1 | test(dt_model)
```

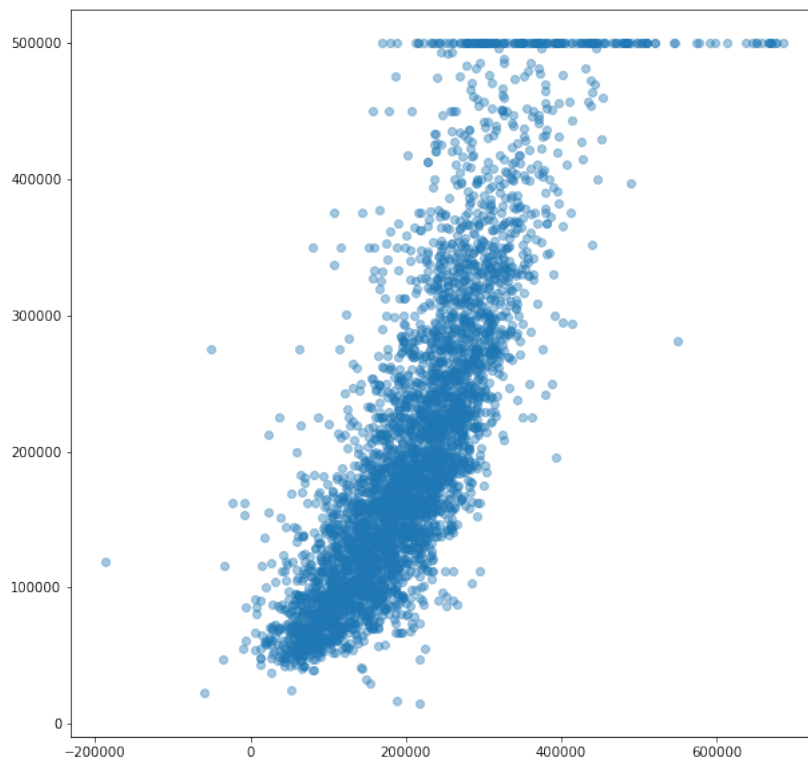


```
1 | 70537.35639423419
```

```
1 | # 模型保存
2 | import joblib
3 | joblib.dump(lr_model, 'lr_model.pkl')
```

```
1 | ['lr_model.pkl']
```

```
1 | # 模型加载
2 | lr = joblib.load('lr_model.pkl')
3 | test(lr)
```



```
1 | 66911.98070857547
```

核心代码

```
1 | import os
2 | import tarfile
3 | import urllib.request
4 | import pandas as pd
5 | import numpy as np
6 | import matplotlib.pyplot as plt
7 | from sklearn.preprocessing import OneHotEncoder, FunctionTransformer, StandardScaler
8 | from sklearn.impute import SimpleImputer
9 | from sklearn.pipeline import Pipeline
10 | from sklearn.compose import ColumnTransformer
11 | from sklearn.linear_model import LinearRegression
12 | from sklearn.tree import DecisionTreeRegressor
13 | from sklearn.metrics import mean_squared_error
14 |
15 | # download data setting
```

```

16 DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml/master/"
17 HOUSING_PATH = os.path.join("datasets", "housing")
18 HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tgz"
19
20 # 下载tar文件并解压
21 def fetch_housing_data(housing_url=HOUSING_URL, housing_path=HOUSING_PATH):
22     os.makedirs(housing_path, exist_ok=True)
23     tgz_path = os.path.join(housing_path, "housing.tgz")
24     # stores data on local disk
25     urllib.request.urlretrieve(housing_url, tgz_path)
26     housing_tgz = tarfile.open(tgz_path)
27     housing_tgz.extractall(path=housing_path)
28     housing_tgz.close()
29
30 # load data
31 def load_housing_data(path = HOUSING_PATH):
32     csv_path = os.path.join(path, 'housing.csv')
33     return pd.read_csv(csv_path)
34
35
36 #fetch_housing_data()
37 housing = load_housing_data()
38
39 # 数据拆分
40 from sklearn.model_selection import StratifiedShuffleSplit
41 housing['income_cat'] = np.ceil(housing['median_income']/1.5)
42 housing['income_cat'].where(housing['income_cat']<5,5.0,inplace=True)
43 split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
44 for train_index, test_index in split.split(housing, housing['income_cat']):
45     strat_train_set = housing.iloc[train_index]
46     strat_test_set = housing.iloc[test_index]
47 strat_train_set.drop(["income_cat"], axis = 1, inplace = True)
48 strat_test_set.drop(["income_cat"], axis = 1, inplace = True)
49
50 # 准备管道数据
51 housing_train = strat_train_set.drop("median_house_value", axis=1) # drop labels for training set
52 housing_labels = strat_train_set["median_house_value"].copy()
53 housing_num = housing_train.drop("ocean_proximity", axis=1)
54
55 # 添加新的特征
56 rooms_ix, bedrooms_ix, population_ix, household_ix = [
57     list(housing.columns).index(col)
58     for col in ("total_rooms", "total_bedrooms", "population", "households")]
59
60 def add_extra_features(X, add_bedrooms_per_room=True):
61     rooms_per_household = X[:, rooms_ix] / X[:, household_ix]
62     population_per_household = X[:, population_ix] / X[:, household_ix]
63     if add_bedrooms_per_room:
64         bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
65         return np.c_[X, rooms_per_household, population_per_household,
66                     bedrooms_per_room]
67     else:
68         return np.c_[X, rooms_per_household, population_per_household]
69
70 # 构造数值特征管道
71 num_pipeline = Pipeline([
72     ('imputer', SimpleImputer(strategy="median")),
73     ('attribs_adder', FunctionTransformer(add_extra_features, validate=False)),
74     ('stander', StandardScaler())
75 ])
76
77 # 数值类特征
78 num_attribs = list(housing_num.columns)
79 # 从一列, 变为5列
80 cat_attribs = ["ocean_proximity"]
81
82 # 所有特征管道
83 full_pipeline = ColumnTransformer([
84     ('num', num_pipeline, num_attribs),
85     ('cat', OneHotEncoder(), cat_attribs)
86 ])
87
88 # 搭建训练管道
89 def train_model(model):
90     train_pipeline = Pipeline([
91         ('datas', full_pipeline),
92         ('model', model)
93     ])
94     return train_pipeline.fit(housing_train, housing_labels)

```

```

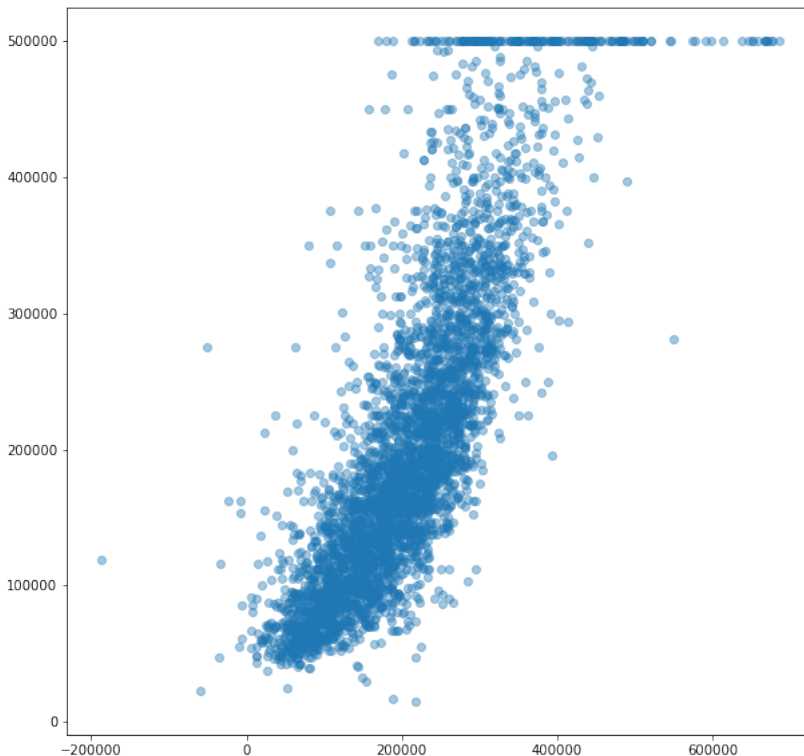
95 # 模型选择-线性回归
96 lr_model = train_model(LinearRegression())
97 # 模型选择-决策树
98 dt_model = train_model(DecisionTreeRegressor())
99
100 # 画出预测值与真实值的分布情况, 返回方差
101 test_data = strat_test_set.drop("median_house_value", axis=1) # drop labels for training set
102 test_labels = strat_test_set["median_house_value"].copy()
103 def test(model):
104     test_data = strat_test_set.drop("median_house_value", axis=1) # drop labels for training set
105     test_labels = strat_test_set["median_house_value"].copy()
106     predictions = model.predict(test_data)
107     plt.figure(1,figsize=(10,10))
108     plt.scatter(x=predictions,y=test_labels,alpha=0.4)
109     plt.xlabel = 'predictions'
110     plt.ylabel = 'ground truth'
111     plt.show()
112     mse = mean_squared_error(test_labels,predictions)
113     return np.sqrt(mse)
114 print('线性回归标准差: '+str(test(lr_model)))
115 print('决策树标准差: '+str(test(dt_model)))
116
117 # 模型保存
118 import joblib
119 joblib.dump(lr_model, 'lr_model.pkl')
120 # 模型加载
121 # lr = joblib.load('lr_model.pkl')
122 # test(lr)

```

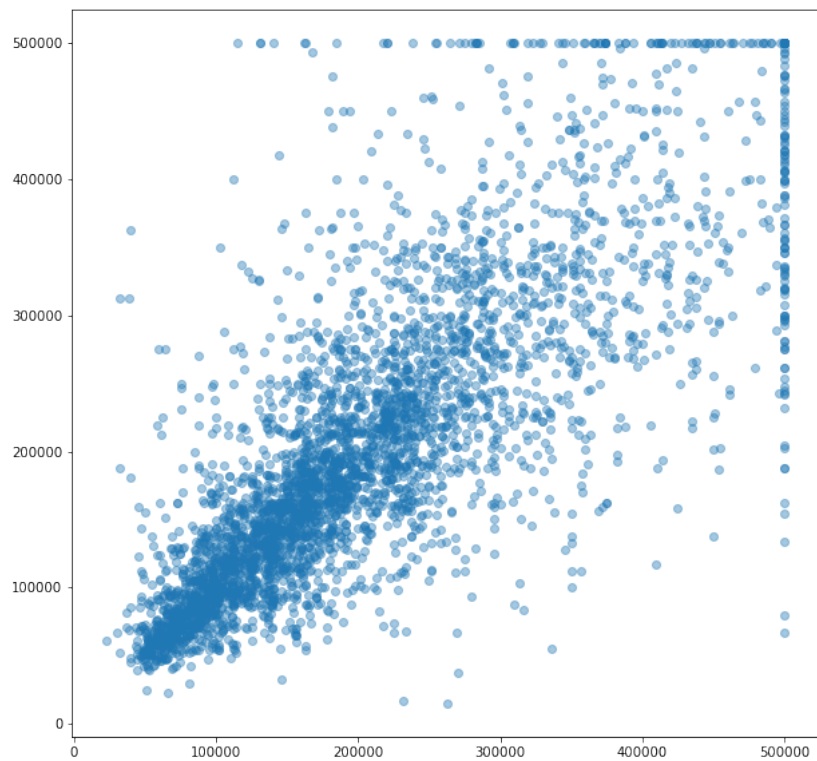
```

1 /root/anaconda3/envs/ai/lib/python3.8/site-packages/pandas/core/frame.py:4163: SettingWithCopyWarning:
2 A value is trying to be set on a copy of a slice from a DataFrame
3
4 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
5 return super().drop(

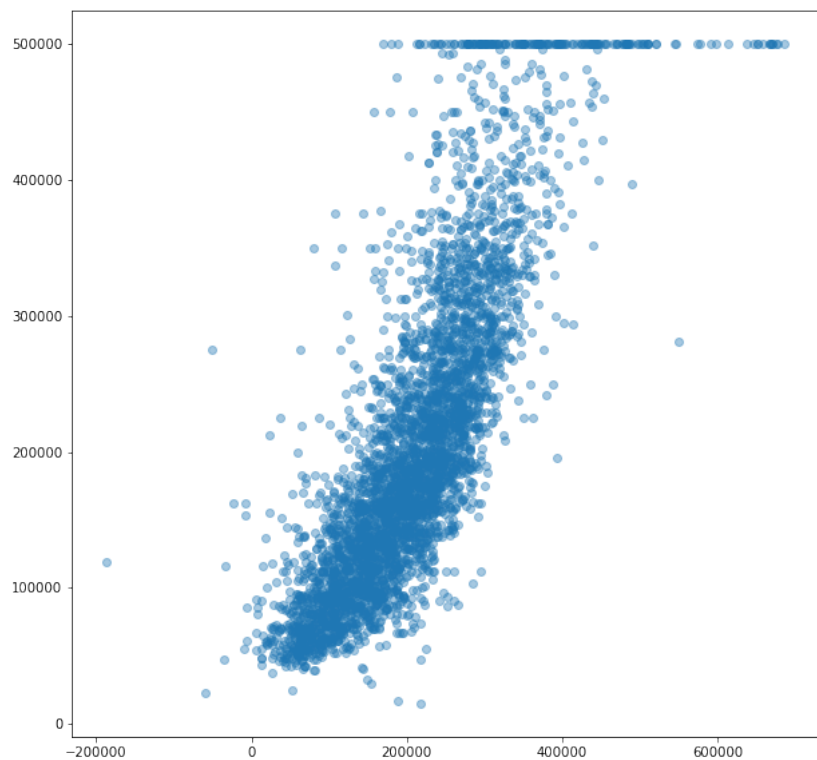
```



1 | 线性回归标准差: 66911.98070857547



1 | 决策树标注差: 70375.68286769741



1 | 66911.98070857547