

```
total_bedrooms      0.047781
population_per_household -0.021991
population           -0.026882
longitude            -0.047466
latitude             -0.142673
bedrooms_per_room   -0.259952
Name: median_house_value, dtype: float64
```

```
housing = strat_train_set.drop("median_house_value", axis=1)
housing_labels = strat_train_set["median_house_value"].copy()
```

```
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy="median")
```

```
housing_num = housing.drop("ocean_proximity", axis=1)
imputer.fit(housing_num)
```

```
SimpleImputer(strategy='median')
```

```
imputer.statistics_
```

```
array([-118.51   ,  34.26   ,  29.   , 2119.   ,  433.   ,
        1164.   ,  408.   ,  3.54155])
```

```
housing_num.median().values
```

```
array([-118.51   ,  34.26   ,  29.   , 2119.   ,  433.   ,
        1164.   ,  408.   ,  3.54155])
```

```
#numpy array
```

```
X = imputer.transform(housing_num)
```

```
#array to dataframe
```

```
housing_tr = pd.DataFrame(X, columns=housing_num.columns,
                          index=housing_num.index)
```

```
housing_cat = housing[["ocean_proximity"]]
print(housing_cat.head(10))
```

```
from sklearn.preprocessing import OrdinalEncoder
ordinal_encoder = OrdinalEncoder()
housing_cat_encoded = ordinal_encoder.fit_transform(housing_cat)
print(housing_cat_encoded[:10])
print(ordinal_encoder.categories_)
```

```
ocean_proximity
12655      INLAND
15502    NEAR OCEAN
2908       INLAND
14053    NEAR OCEAN
20496    <1H OCEAN
1481      NEAR BAY
18125    <1H OCEAN
```

```

5830      <1H OCEAN
17989      <1H OCEAN
4861      <1H OCEAN
[[1.]
 [4.]
 [1.]
 [4.]
 [0.]
 [3.]
 [0.]
 [0.]
 [0.]
 [0.]]
[array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
      dtype=object)]

```

```

from sklearn.preprocessing import OneHotEncoder
cat_encoder = OneHotEncoder()
housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
housing_cat_1hot

```

```

<16512x5 sparse matrix of type '<class 'numpy.float64'>'
  with 16512 stored elements in Compressed Sparse Row format>

```

```

housing_cat_1hot.toarray()
cat_encoder.categories_

```

```

[array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
      dtype=object)]

```

```

from sklearn.base import BaseEstimator, TransformerMixin
rooms_ix, bedrooms_ix, population_ix, households_ix = 3, 4, 5, 6
class CombinedAttributesAdder(BaseEstimator, TransformerMixin):
    def __init__(self, add_bedrooms_per_room = True): # no *args or
**kwargs
        self.add_bedrooms_per_room = add_bedrooms_per_room
    def fit(self, X, y=None):
        return self # nothing else to do
    def transform(self, X):
        rooms_per_household = X[:, rooms_ix] / X[:, households_ix]
        population_per_household = X[:, population_ix] / X[:,
households_ix]
        if self.add_bedrooms_per_room:
            bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
            return np.c_[X, rooms_per_household,
population_per_household, bedrooms_per_room]
        else:
            return np.c_[X, rooms_per_household,
population_per_household]

```

```
attr_adder = CombinedAttributesAdder(add_bedrooms_per_room=False)
housing_extra_attribs = attr_adder.transform(housing.values)

from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
num_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy="median")),
    ('attribs_adder', CombinedAttributesAdder()),
    ('std_scaler', StandardScaler()),
])
housing_num_tr = num_pipeline.fit_transform(housing_num)

from sklearn.compose import ColumnTransformer
num_attribs = list(housing_num)
cat_attribs = ["ocean_proximity"]
full_pipeline = ColumnTransformer([
    ("num", num_pipeline, num_attribs),
    ("cat", OneHotEncoder(), cat_attribs),
])
housing_prepared = full_pipeline.fit_transform(housing)
```