

# Feature Engineering 101



Topic - 13

## Binarization

### Binarization

```
In [1]: import pandas as pd  
import numpy as np
```

```
In [2]: import matplotlib.pyplot as plt

        from sklearn.model_selection import train_test_split

        from sklearn.tree import DecisionTreeClassifier

        from sklearn.metrics import accuracy_score
        from sklearn.model_selection import cross_val_score

        from sklearn.preprocessing import KBinsDiscretizer
        from sklearn.compose import ColumnTransformer
```

```
In [3]: df = pd.read_csv('train.csv', usecols=['Age', 'Fare', 'Survived'])
```

```
In [4]: df.dropna(inplace=True)
```

```
In [5]: df.shape
```

```
Out[5]: (714, 3)
```

```
In [6]: df.head()
```

```
Out[6]:
```

	Survived	Age	Fare
0	0	22.0	7.2500
1	1	38.0	71.2833
2	1	26.0	7.9250
3	1	35.0	53.1000
4	0	35.0	8.0500

```
In [7]: X = df.iloc[:,1:]
        y = df.iloc[:,0]
```

```
In [8]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [9]: X_train.head(2)
```

```
Out[9]:
```

	Age	Fare
328	31.0	20.5250
73	26.0	14.4542

```
In [10]: clf = DecisionTreeClassifier()
```

```
In [11]: clf.fit(X_train, y_train)
        y_pred = clf.predict(X_test)
```

In [12]: accuracy\_score(y\_test,y\_pred)

Out[12]: 0.6433566433566433

In [13]: np.mean(cross\_val\_score(DecisionTreeClassifier(),X,y,cv=10,scoring='accuracy'))

Out[13]: 0.6289319248826291

In [14]: kbin\_age = KBinsDiscretizer(n\_bins=15,encode='ordinal',strategy='quantile')  
kbin\_fare = KBinsDiscretizer(n\_bins=15,encode='ordinal',strategy='quantile')

In [15]: trf = ColumnTransformer([  
 ('first',kbin\_age,[0]),  
 ('second',kbin\_fare,[1])  
])

In [16]: X\_train\_trf = trf.fit\_transform(X\_train)  
X\_test\_trf = trf.transform(X\_test)

In [17]: trf.named\_transformers\_['first'].bin\_edges\_

Out[17]: array([array([ 0.42, 6. , 16. , 19. , 21. , 23. , 25. , 28. , 30. ,  
 32. , 35. , 38. , 42. , 47. , 54. , 80. ])  
dtype=object)

In [18]: trf.named\_transformers\_['second'].bin\_edges\_

Out[18]: array([array([ 0.42, 6. , 16. , 19. , 21. , 23. , 25. , 28. , 30. ,  
 32. , 35. , 38. , 42. , 47. , 54. , 80. ])  
dtype=object)

In [19]: output = pd.DataFrame({  
 'age':X\_train['Age'],  
 'age\_trf':X\_train\_trf[:,0],  
 'fare':X\_train['Fare'],  
 'fare\_trf':X\_train\_trf[:,1]  
})

In [20]: output['age\_labels'] = pd.cut(x=X\_train['Age'],  
 bins=trf.named\_transformers\_['first'].bin\_edges\_[0].to  
output['fare\_labels'] = pd.cut(x=X\_train['Fare'],  
 bins=trf.named\_transformers\_['second'].bin\_edges\_[0].to

In [21]: output.sample(5)

Out[21]:

	age	age_trf	fare	fare_trf	age_labels	fare_labels
821	27.0	6.0	8.6625	4.0	(25.0, 28.0]	(8.158, 10.5]
230	35.0	10.0	83.4750	13.0	(32.0, 35.0]	(76.292, 108.9]
784	25.0	6.0	7.0500	0.0	(23.0, 25.0]	(0.0, 7.25]
660	50.0	13.0	133.6500	14.0	(47.0, 54.0]	(108.9, 512.329]

	age	age_trf	fare	fare_trf	age_labels	fare_labels
621	42.0	12.0	52.5542	12.0	(38.0, 42.0]	(51.479, 76.292]

```
In [22]: clf = DecisionTreeClassifier()
         clf.fit(X_train_trf,y_train)
         y_pred2 = clf.predict(X_test_trf)
```

```
In [23]: accuracy_score(y_test,y_pred2)
```

```
Out[23]: 0.6363636363636364
```

```
In [24]: X_trf = trf.fit_transform(X)
         np.mean(cross_val_score(DecisionTreeClassifier(),X,y,cv=10,scoring='accuracy'))
```

```
Out[24]: 0.6303012519561815
```

```
In [25]: def discretize(bins, strategy):
         kbin_age = KBinsDiscretizer(n_bins=bins, encode='ordinal', strategy=strategy)
         kbin_fare = KBinsDiscretizer(n_bins=bins, encode='ordinal', strategy=strategy)

         trf = ColumnTransformer([
             ('first', kbin_age, [0]),
             ('second', kbin_fare, [1])
         ])

         X_trf = trf.fit_transform(X)
         print(np.mean(cross_val_score(DecisionTreeClassifier(), X, y, cv=10, scoring='accuracy')) )

         plt.figure(figsize=(14, 4))
         plt.subplot(121)
         plt.hist(X['Age'])
         plt.title("Before")

         plt.subplot(122)
         plt.hist(X_trf[:,0], color='red')
         plt.title("After")

         plt.show()

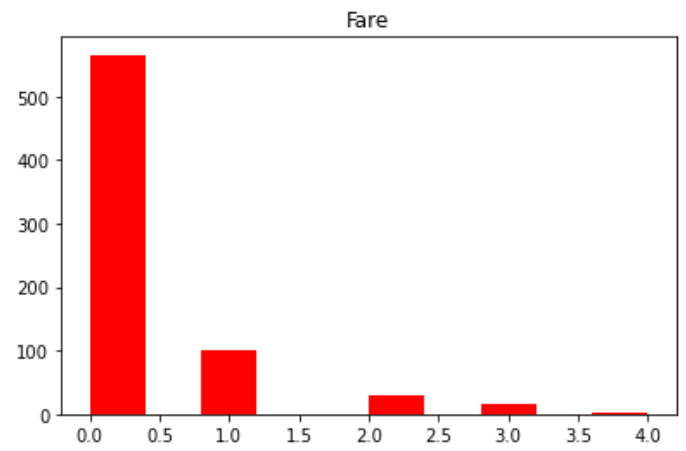
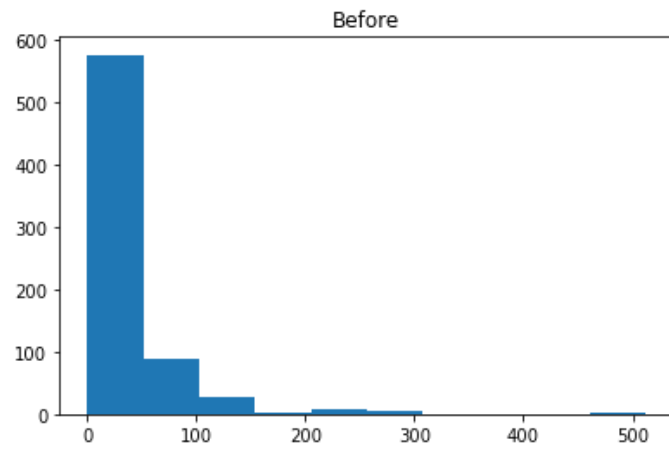
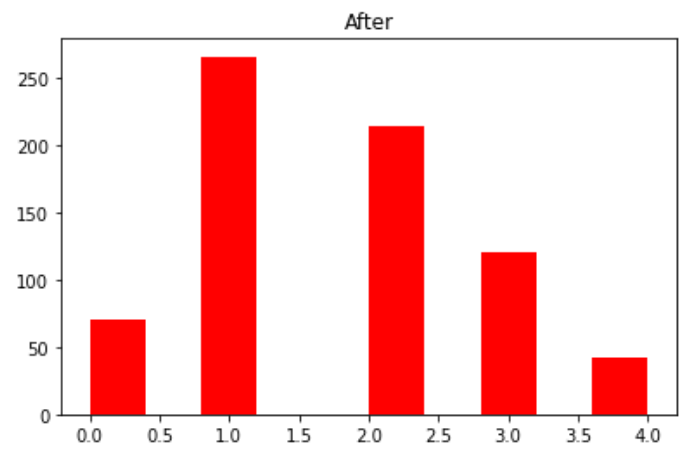
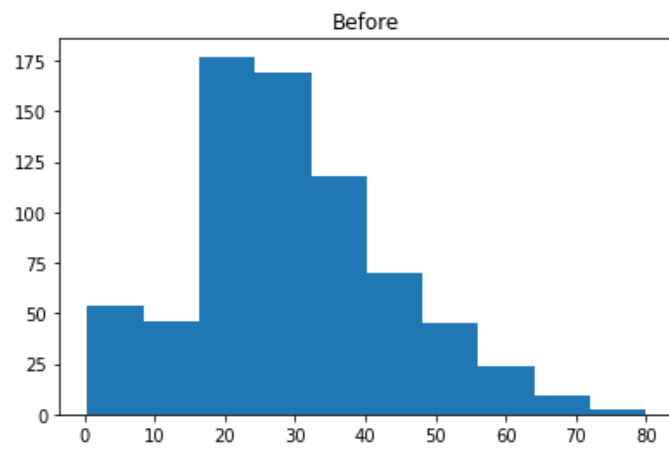
         plt.figure(figsize=(14, 4))
         plt.subplot(121)
         plt.hist(X['Fare'])
         plt.title("Before")

         plt.subplot(122)
         plt.hist(X_trf[:,1], color='red')
         plt.title("Fare")

         plt.show()
```

```
In [26]: discretize(5, 'kmeans')
```

```
0.6288928012519561
```



In [ ]: