

ch2 지도학습 - knn - 실습

- Machine Learning with sklearn @ DJ,Lim
- date : 21/10

실습해 보기

- titanic 데이터 셋을 활용하여 knn 모델을 구현한다.
- 가장 높은 일반화 성능을 갖는 k의 값은 무엇인지 찾아보자.

01 데이터 준비

02 머신러닝 모델 만들고 예측하기

01 데이터 준비

In [15]:



```
import pandas as pd
from sklearn.model_selection import train_test_split
import numpy as np
```

In [16]:



```
dat = pd.read_csv("train.csv")
dat
```

Out[16]:

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	2	1	1Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833
2	3	1	3Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
3	4	1	1Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000
4	5	0	3Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500
...
886	887	0	2Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000
887	888	1	1Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000
888	889	0	3Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500
889	890	1	1Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000
890	891	0	3Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500

891 rows × 12 columns

In [17]:

```
dat.columns
```

Out[17]:

```
Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',  
      'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],  
      dtype='object')
```

In [18]:

```
dat.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 891 entries, 0 to 890  
Data columns (total 12 columns):  
#   Column          Non-Null Count  Dtype    
---  ---            -  
0   PassengerId      891 non-null    int64    
1   Survived         891 non-null    int64    
2   Pclass           891 non-null    int64    
3   Name             891 non-null    object    
4   Sex              891 non-null    object    
5   Age              714 non-null    float64   
6   SibSp            891 non-null    int64    
7   Parch            891 non-null    int64    
8   Ticket           891 non-null    object    
9   Fare             891 non-null    float64   
10  Cabin            204 non-null    object    
11  Embarked         889 non-null    object    
dtypes: float64(2), int64(5), object(5)  
memory usage: 83.7+ KB
```

데이터 선택 및 나누기

In [19]:

```
X = dat[ ['Pclass' , 'SibSp' ] ]  
y = dat['Survived']  
  
# 90% : 학습용, 10% : 테스트용  
X_train, X_test, y_train, y_test = train_test_split(X, y,  
                                                    test_size=0.1,  
                                                    random_state=0)  
  
X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[19]:

```
((801, 2), (90, 2), (801,), (90,))
```

In [20]:

```
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier(n_neighbors=2)
model.fit(X_train, y_train)

### 예측시키기
pred = model.predict(X_test)
```

In [21]:

```
(pred == y_test).sum() / len(pred)
```

Out[21]:

0.6

In [22]:

```
print("테스트 세트의 정확도 : {:.2f}".format(np.mean(pred == y_test)))
```

테스트 세트의 정확도 : 0.60

결측치 처리 및 레이블 인코딩

In [23]:

```
dat.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column        Non-Null Count  Dtype  
---  -
 0   PassengerId   891 non-null    int64  
 1   Survived      891 non-null    int64  
 2   Pclass        891 non-null    int64  
 3   Name          891 non-null    object  
 4   Sex           891 non-null    object  
 5   Age           714 non-null    float64 
 6   SibSp         891 non-null    int64  
 7   Parch         891 non-null    int64  
 8   Ticket        891 non-null    object  
 9   Fare          891 non-null    float64 
10   Cabin         204 non-null    object  
11   Embarked      889 non-null    object  
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

In [24]:



```
dat.head()
```

Out[24]:

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	
1	2	1	1Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	
2	3	1	3Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	
3	4	1	1Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	
4	5	0	3Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	

In [25]:



```
mapping = { "male":1, 'female':2 }  
dat['Sex_num'] = dat['Sex'].map(mapping)  
dat.head()
```

Out[25]:

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	
1	2	1	1Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	
2	3	1	3Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	
3	4	1	1Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	
4	5	0	3Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	

In [26]:

```
val_mean = dat['Age'].mean()
dat['Age'] = dat['Age'].fillna( val_mean )
dat.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     891 non-null    int64
1   Survived        891 non-null    int64
2   Pclass          891 non-null    int64
3   Name            891 non-null    object
4   Sex             891 non-null    object
5   Age             891 non-null    float64
6   SibSp           891 non-null    int64
7   Parch           891 non-null    int64
8   Ticket          891 non-null    object
9   Fare            891 non-null    float64
10  Cabin           204 non-null    object
11  Embarked        889 non-null    object
12  Sex_num         891 non-null    int64
dtypes: float64(2), int64(6), object(5)
memory usage: 90.6+ KB
```

학습, 테스트 데이터 셋 나누기

In [28]:

```
X = dat[ ['Pclass' , 'SibSp' , 'Sex_num' , 'Age' ] ]
y = dat['Survived']

# 90% : 학습용 , 10% : 테스트용
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.1,
                                                    random_state=0)

X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[28]:

```
((801, 4), (90, 4), (801,), (90,))
```

In [29]:

```
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier(n_neighbors=2)
model.fit(X_train, y_train)
### 예측시키기
pred = model.predict(X_test)
print("테스트 세트의 정확도 : {:.2f}".format(np.mean(pred == y_test)))
```

테스트 세트의 정확도 : 0.76

k값에 따른 정확도 확인

In [30]:



```
tr_acc = []
test_acc = []
k_nums = range(1, 22, 2)# 1,3,5~21

for n in k_nums:
    # 모델 선택 및 학습
    model = KNeighborsClassifier(n_neighbors=n)
    model.fit(X_train, y_train)

    # 정확도 구하기
    acc_tr = model.score(X_train, y_train)
    acc_test = model.score(X_test, y_test)

    # 정확도 값 저장.
    tr_acc.append(acc_tr)
    test_acc.append(acc_test)

    print("k : ", n)
    print("학습용셋 정확도 {:.3f}".format(acc_tr) )
    print("테스트용셋 정확도 {:.3f}".format(acc_test) )
```

```
k : 1
학습용셋 정확도 0.885
테스트용셋 정확도 0.756
k : 3
학습용셋 정확도 0.863
테스트용셋 정확도 0.800
k : 5
학습용셋 정확도 0.850
테스트용셋 정확도 0.800
k : 7
학습용셋 정확도 0.839
테스트용셋 정확도 0.733
k : 9
학습용셋 정확도 0.830
테스트용셋 정확도 0.722
k : 11
학습용셋 정확도 0.826
테스트용셋 정확도 0.744
k : 13
학습용셋 정확도 0.824
테스트용셋 정확도 0.756
k : 15
학습용셋 정확도 0.815
테스트용셋 정확도 0.744
k : 17
학습용셋 정확도 0.805
테스트용셋 정확도 0.756
k : 19
학습용셋 정확도 0.792
테스트용셋 정확도 0.711
k : 21
학습용셋 정확도 0.790
테스트용셋 정확도 0.722
```


In [31]:



```
import seaborn as sns
print(sns.__version__)
```

0.11.1

In [32]:



```
# tr_acc = []
# test_acc = []
dat = { "tr_acc":tr_acc, "test_acc":test_acc }
data_df = pd.DataFrame(dat, index=range(1, 22, 2))
data_df
```

Out[32]:

	tr_acc	test_acc
1	0.885144	0.755556
3	0.862672	0.800000
5	0.850187	0.800000
7	0.838951	0.733333
9	0.830212	0.722222
11	0.826467	0.744444
13	0.823970	0.755556
15	0.815231	0.744444
17	0.805243	0.755556
19	0.791511	0.711111
21	0.790262	0.722222

In [37]:

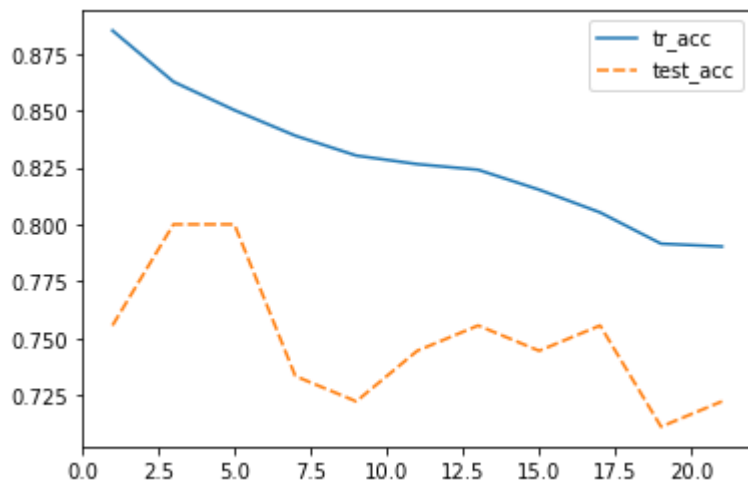


```
import matplotlib.pyplot as plt
```

In [38]:



```
sns.lineplot(data=data_df, palette="tab10")  
plt.show()
```



In []:

