

데이터 전처리

In [107...]

```
import pandas as pd
import glob
import numpy as np
import seaborn as sns

sns.set(style="whitegrid", color_codes=True)

# 데이터 불러오기
DATA_DIR = 'titanic'
data_files = sorted(glob.glob(DATA_DIR + "/*.csv"), reverse=True)
print("불러온 파일:", data_files)

# train, test 각각 읽기
train_df = pd.read_csv(data_files[0])
test_df = pd.read_csv(data_files[1])

n_train = len(train_df)
n_test = len(test_df)

# 두 데이터 합치기
df = pd.concat([train_df, test_df], sort=False).reset_index(drop=True)
print("통합 후 데이터 크기:", df.shape)
```

불러온 파일: ['titanic/train.csv', 'titanic/test.csv']
통합 후 데이터 크기: (1309, 13)

In [108...]

```
pd.options.display.float_format = '{:.2f}'.format
print(df.isnull().sum() / len(df) * 100)
```

```
Unnamed: 0      0.00
PassengerId    0.00
Survived       0.00
Pclass          0.00
Name            0.00
Sex             0.00
Age            20.09
SibSp          0.00
Parch          0.00
Ticket         0.00
Fare           0.08
Cabin          77.46
Embarked       0.15
dtype: float64
```

In [109...]

```
df[df["Embarked"].isnull()]
```

Out[109...]

	Unnamed: 0	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
61	61	62	1	1	Icard, Miss. Amelie	female	38.00	0	0	113572	80.00	B28	NaN
829	829	830	1	1	Stone, Mrs. George Nelson (Martha Evelyn)	female	62.00	0	0	113572	80.00	B28	NaN

In [110...]

```
df[df["Fare"].isnull()]
```

Out[110...]

	Unnamed: 0	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1043	152	1044	0	3	Storey, Mr. Thomas	male	60.50	0	0	3701	NaN	NaN	S

In [111...]

```
# Age: Pclass 기준 평균으로 채움
df["Age"] = df["Age"].fillna(df.groupby("Pclass")["Age"].transform("mean"))

# Embarked: 결측을 'S'로 채움
df["Embarked"] = df["Embarked"].fillna("S")

# Fare: 결측을 3등급의 평균으로 채움
mean_fare_pclass3 = df.loc[df["Pclass"] == 3, "Fare"].mean()
df["Fare"] = df["Fare"].fillna(mean_fare_pclass3)
```

In [112...]

```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1309 entries, 0 to 1308
Data columns (total 13 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Unnamed: 0    1309 non-null   int64  
 1   PassengerId  1309 non-null   int64  
 2   Survived     1309 non-null   int64  
 3   Pclass       1309 non-null   int64  
 4   Name         1309 non-null   object  
 5   Sex          1309 non-null   object  
 6   Age          1309 non-null   float64 
 7   SibSp        1309 non-null   int64  
 8   Parch        1309 non-null   int64  
 9   Ticket       1309 non-null   object  
 10  Fare         1309 non-null   float64 
 11  Cabin        295 non-null   object  
 12  Embarked     1309 non-null   object  
dtypes: float64(2), int64(6), object(5)
memory usage: 133.1+ KB

```

```

In [113]: object_columns = ["PassengerId", "Pclass", "Name", "Sex",
                  "Ticket", "Cabin", "Embarked"]
numeric_columns = ["Age", "SibSp", "Parch", "Fare"]

for col in object_columns:
    df[col] = df[col].astype(object)
for col in numeric_columns:
    df[col] = df[col].astype(float)

df["SibSp"] = df["SibSp"].astype(int)
df["Parch"] = df["Parch"].astype(int)

```

```

In [114]: def merge_and_get(ldf, rdf, on=None, how="inner", index=None):
    if index is True:
        return pd.merge(ldf, rdf, how=how, left_index=True, right_index=True)
    else:
        return pd.merge(ldf, rdf, how=how, on=on)

one_hot_df = merge_and_get(df, pd.get_dummies(df["Sex"], prefix="Sex"), index=True)
one_hot_df = merge_and_get(one_hot_df, pd.get_dummies(df["Pclass"], prefix="Pclass"), index=True)
one_hot_df = merge_and_get(one_hot_df, pd.get_dummies(df["Embarked"], prefix="Embarked"), index=True)

# 원본 범주형 제거
one_hot_df = one_hot_df.drop(columns=["Sex", "Pclass", "Embarked"])

# 학습에 의미 없는 텍스트 컬럼 제거
drop_cols = ["PassengerId", "Name", "Ticket", "Cabin"]
one_hot_df = one_hot_df.drop(columns=drop_cols, errors="ignore")

```

```

In [115]: # 위에서 기록한 n_train, n_test 사용
x_train = one_hot_df.iloc[:n_train, :].copy()
x_test = one_hot_df.iloc[n_train:, :].copy()

y_train = train_df["Survived"].reset_index(drop=True)
y_test = test_df["Survived"].reset_index(drop=True)

# 입력에서는 Survived 제거
x_train = x_train.drop(columns=['Survived'], errors='ignore')
x_test = x_test.drop(columns=['Survived'], errors='ignore')

print("Train 크기:", x_train.shape)
print("Test 크기:", x_test.shape)

```

Train 크기: (891, 13)
Test 크기: (418, 13)

Linear Regression

```

In [116]: import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# 모델 학습
lr = LinearRegression()
lr.fit(x_train, y_train)

# 예측
y_pred = lr.predict(x_test)

# cutoff 값 여러 개 시도

```

```
thresholds = np.arange(0.1, 0.9, 0.01)

results = []
print("Linear Regression Thresholds")
for t in thresholds:
    y_bin = (y_pred >= t).astype(int)

    acc = accuracy_score(y_test, y_bin)
    prec = precision_score(y_test, y_bin)
    rec = recall_score(y_test, y_bin)
    f1 = f1_score(y_test, y_bin)

    results.append([t, acc, prec, rec, f1])
    print(f't={t:.2f} | acc={acc:.3f} | prec={prec:.3f} | rec={rec:.3f} | f1={f1:.3f}')

# F1이 가장 높은 cutoff 찾기
best_row = result_df.loc[result_df["F1"].idxmax()]
best_t = best_row["Threshold"]

print("\nBest cutoff")
print(best_row)
```

Linear Regression Thresholds					
t=0.10	acc=0.471	prec=0.414	rec=0.962	f1=0.579	
t=0.11	acc=0.510	prec=0.431	rec=0.930	f1=0.589	
t=0.12	acc=0.531	prec=0.442	rec=0.918	f1=0.597	
t=0.13	acc=0.543	prec=0.448	rec=0.905	f1=0.600	
t=0.14	acc=0.553	prec=0.453	rec=0.886	f1=0.600	
t=0.15	acc=0.557	prec=0.456	rec=0.886	f1=0.602	
t=0.16	acc=0.574	prec=0.467	rec=0.886	f1=0.611	
t=0.17	acc=0.586	prec=0.474	rec=0.867	f1=0.613	
t=0.18	acc=0.608	prec=0.489	rec=0.854	f1=0.622	
t=0.19	acc=0.612	prec=0.493	rec=0.842	f1=0.621	
t=0.20	acc=0.617	prec=0.496	rec=0.835	f1=0.623	
t=0.21	acc=0.627	prec=0.504	rec=0.835	f1=0.629	
t=0.22	acc=0.636	prec=0.512	rec=0.835	f1=0.635	
t=0.23	acc=0.641	prec=0.516	rec=0.829	f1=0.636	
t=0.24	acc=0.648	prec=0.522	rec=0.829	f1=0.641	
t=0.25	acc=0.651	prec=0.524	rec=0.816	f1=0.639	
t=0.26	acc=0.656	prec=0.529	rec=0.810	f1=0.640	
t=0.27	acc=0.663	prec=0.536	rec=0.810	f1=0.645	
t=0.28	acc=0.672	prec=0.545	rec=0.810	f1=0.651	
t=0.29	acc=0.677	prec=0.549	rec=0.810	f1=0.655	
t=0.30	acc=0.691	prec=0.564	rec=0.810	f1=0.665	
t=0.31	acc=0.701	prec=0.574	rec=0.810	f1=0.672	
t=0.32	acc=0.708	prec=0.583	rec=0.797	f1=0.674	
t=0.33	acc=0.708	prec=0.585	rec=0.785	f1=0.670	
t=0.34	acc=0.725	prec=0.605	rec=0.785	f1=0.683	
t=0.35	acc=0.730	prec=0.613	rec=0.772	f1=0.683	
t=0.36	acc=0.734	prec=0.619	rec=0.772	f1=0.687	
t=0.37	acc=0.739	prec=0.627	rec=0.766	f1=0.689	
t=0.38	acc=0.749	prec=0.640	rec=0.766	f1=0.697	
t=0.39	acc=0.746	prec=0.644	rec=0.734	f1=0.686	
t=0.40	acc=0.749	prec=0.650	rec=0.728	f1=0.687	
t=0.41	acc=0.746	prec=0.649	rec=0.715	f1=0.681	
t=0.42	acc=0.754	prec=0.661	rec=0.715	f1=0.687	
t=0.43	acc=0.758	prec=0.671	rec=0.709	f1=0.689	
t=0.44	acc=0.758	prec=0.671	rec=0.709	f1=0.689	
t=0.45	acc=0.763	prec=0.683	rec=0.696	f1=0.690	
t=0.46	acc=0.763	prec=0.686	rec=0.690	f1=0.688	
t=0.47	acc=0.763	prec=0.688	rec=0.684	f1=0.686	
t=0.48	acc=0.763	prec=0.688	rec=0.684	f1=0.686	
t=0.49	acc=0.763	prec=0.690	rec=0.677	f1=0.684	
t=0.50	acc=0.761	prec=0.688	rec=0.671	f1=0.679	
t=0.51	acc=0.763	prec=0.695	rec=0.665	f1=0.680	
t=0.52	acc=0.768	prec=0.705	rec=0.665	f1=0.684	
t=0.53	acc=0.770	prec=0.709	rec=0.665	f1=0.686	
t=0.54	acc=0.775	prec=0.719	rec=0.665	f1=0.691	
t=0.55	acc=0.780	prec=0.736	rec=0.652	f1=0.691	
t=0.56	acc=0.782	prec=0.745	rec=0.646	f1=0.692	
t=0.57	acc=0.778	prec=0.741	rec=0.633	f1=0.683	
t=0.58	acc=0.780	prec=0.746	rec=0.633	f1=0.685	
t=0.59	acc=0.778	prec=0.748	rec=0.620	f1=0.678	
t=0.60	acc=0.780	prec=0.758	rec=0.614	f1=0.678	
t=0.61	acc=0.782	prec=0.768	rec=0.608	f1=0.678	
t=0.62	acc=0.785	prec=0.779	rec=0.601	f1=0.679	
t=0.63	acc=0.782	prec=0.786	rec=0.582	f1=0.669	
t=0.64	acc=0.789	prec=0.818	rec=0.570	f1=0.672	
t=0.65	acc=0.785	prec=0.821	rec=0.551	f1=0.659	
t=0.66	acc=0.773	prec=0.812	rec=0.519	f1=0.633	
t=0.67	acc=0.768	prec=0.851	rec=0.468	f1=0.604	
t=0.68	acc=0.766	prec=0.849	rec=0.462	f1=0.598	
t=0.69	acc=0.761	prec=0.845	rec=0.449	f1=0.587	
t=0.70	acc=0.758	prec=0.861	rec=0.430	f1=0.574	
t=0.71	acc=0.763	prec=0.893	rec=0.424	f1=0.575	
t=0.72	acc=0.761	prec=0.903	rec=0.411	f1=0.565	
t=0.73	acc=0.758	prec=0.901	rec=0.405	f1=0.559	
t=0.74	acc=0.756	prec=0.900	rec=0.399	f1=0.553	
t=0.75	acc=0.756	prec=0.912	rec=0.392	f1=0.549	
t=0.76	acc=0.751	prec=0.922	rec=0.373	f1=0.532	
t=0.77	acc=0.751	prec=0.922	rec=0.373	f1=0.532	
t=0.78	acc=0.749	prec=0.934	rec=0.361	f1=0.521	
t=0.79	acc=0.749	prec=0.949	rec=0.354	f1=0.516	
t=0.80	acc=0.749	prec=0.965	rec=0.348	f1=0.512	
t=0.81	acc=0.737	prec=0.962	rec=0.316	f1=0.476	
t=0.82	acc=0.732	prec=0.960	rec=0.304	f1=0.462	
t=0.83	acc=0.730	prec=0.959	rec=0.297	f1=0.454	
t=0.84	acc=0.725	prec=0.978	rec=0.278	f1=0.433	
t=0.85	acc=0.718	prec=0.976	rec=0.259	f1=0.410	
t=0.86	acc=0.711	prec=0.974	rec=0.241	f1=0.386	
t=0.87	acc=0.706	prec=0.973	rec=0.228	f1=0.369	
t=0.88	acc=0.701	prec=0.971	rec=0.215	f1=0.352	
t=0.89	acc=0.689	prec=0.967	rec=0.184	f1=0.309	

Best cutoff

Threshold 0.38

```
Accuracy    0.75
Precision   0.64
Recall      0.77
F1          0.70
Name: 28, dtype: float64
```

Logistic Regression

```
In [117]:  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, classification_report  
import numpy as np  
  
# 로지스틱 회귀 모델 학습  
log_reg = LogisticRegression(max_iter=1000)  
log_reg.fit(train_processed, y_train)  
  
# 예측  
y_pred_proba = log_reg.predict_proba(test_processed)[:, 1] # 생존(1)일 확률  
y_pred = (y_pred_proba >= 0.5).astype(int)  
  
# 성능 평가  
acc = accuracy_score(y_test, y_pred)  
prec = precision_score(y_test, y_pred)  
rec = recall_score(y_test, y_pred)  
f1 = f1_score(y_test, y_pred)  
  
print("Logistic Regression 성능:")  
print(f"Accuracy : {acc:.3f}")  
print(f"Precision: {prec:.3f}")  
print(f"Recall   : {rec:.3f}")  
print(f"F1 Score : {f1:.3f}")
```

Logistic Regression 성능:
Accuracy : 0.770
Precision: 0.704
Recall : 0.677
F1 Score : 0.690

두 모델을 비교해 봅시다. 어느 모델이 성능이 나은가요?

두 모델은 전반적으로 비슷한 수준의 분류 성능을 보이지만 세부적인 성능 지표와 안정성 측면에서 차이를 보인다. 선형 회귀는 F1score가 0.70으로 로지스틱 회귀보다 약간 높게 나타났다. 특히 Recall이 0.77로 매우 높아 실제 생존자를 놓치지 않고 잘 예측하는 경향을 보인다. 이는 민감도가 중요한 상황에서 유리하게 작용할 수 있다. 반면 Accuracy와 Precision은 로지스틱 회귀가 더 높은 수치를 기록했다. 로지스틱 회귀는 모델이 생존이라고 예측한 경우 실제로 생존일 확률이 더 높다는 것을 의미한다. 또한 F1score은 선형 회귀와 거의 비슷한 수준을 유지하면서도 전체적인 지표 간의 균형이 더 안정적이다. 종합적으로 봤을 때 선형 회귀는 재현율이 높아 민감도가 중요한 상황에서 적합하지만 로지스틱 회귀는 정확도와 정밀도가 높고 모델의 안정성이 더 우수하여 일반적인 분류 문제에 더 적합한 모델이라고 볼 수 있다. 따라서 로지스틱 회귀가 보다 일관적이고 신뢰할 수 있는 성능을 보이는 모델로 평가된다.

왜 성능 차이가 날까요?

두 모델의 성능 차이는 학습 원리와 손실 함수의 구조적 차이에서 비롯된다. 선형 회귀는 평균제곱오차를 최소화하여 연속적인 값을 예측하도록 설계된 모델이다. 따라서 분류 문제에 적용할 경우 예측 결과가 확률이 아닌 단순 실수 값으로 출력되며 이를 0과 1로 구분하기 위해 임의적인 cutoff설정이 필요하다. 이러한 특성 때문에 선형 회귀는 예측값의 범위가 제한되지 않고 분류 기준에 따라 결과가 크게 달라질 수 있다. 반면, 로지스틱 회귀는 시그모이드 함수를 사용해 예측값을 자연스럽게 0과 1 사이의 확률로 변환한다. 또한 회귀 문제가 아닌 분류 문제에 특화된 로그 손실을 최소화하여 클래스 간의 경계를 보다 안정적이고 일관되게 학습할 수 있다. 결국 두 모델의 성능 차이는 선형 회귀가 분류 문제에 부적합한 손실 함수를 사용한다는 점과 로지스틱 회귀가 확률 기반의 최적화 구조를 통해 분류에 최적화되어 있다는 점에서 발생한다. 즉, 로지스틱 회귀는 분류 문제에서 해석하기 쉽고 안정적인 결과를 제공하는 반면 선형 회귀는 연속 예측에 적합하지만 분류 문제에서는 근본적인 한계를 갖는다.

왜 선형회귀를 분류기로 쓰면 안 될까요?

선형회귀는 연속적인 값을 예측하기 위한 회귀 분석 기법으로 수치 예측 문제에는 적합하지만 생존 여부 예측과 같은 이진 분류 문제에 적용할 경우 근본적인 한계가 존재한다. 첫째, 선형 회귀의 예측값은 $-\infty$ 에서 $+\infty$ 까지의 연속적인 실수 범위를 가지므로 확률로 해석할 수 없다. 예측값이 -0.3이나 1.5와 같이 비현실적인 값으로 출력될 수 있으며 이는 확률적인 의미를 갖지는 못한다. 둘째, 선형 회귀는 출력이 정규화되어 있지 않기 때문에 분류를 위해 임의의 기준값을 설정해야 한다. 그러나 이 기준값의 변화에 따라 모델의 성능이 크게 달라질 수 있어 결정 경계가 불안정하고 일관성이 부족하다. 셋째, 선형 회귀의 손실 함수인 평균제곱오차는 연속형 예측 오차를 최소화하는데 초점을 맞추고 있어 클래스 확률을 최적화해야하는 분류 문제에는 부적합하다. 마지막으로, 선형 회귀는 MSE손실을 사용하기 때문에 이상치에 매우 민감하다. 극단적인 값이 존재할 경우 전체 회귀선이 크게 왜곡되어 예측 결과가 불안정해질 수 있다.