

결측값(missing values)

결측(missing value)은 데이터에서 특정 위치 값이 없는 경우를 의미합니다. 결측값은 여러 가지 이유로 발생할 수 있습니다. 예를 들어, 설문조사 응답자나 질문에 답하지 않은 고객, 센서나 일시적으로 작동하지 않아 데이터가 수집되지 못한 경우, 데이터 입력 오류 등 여러 상황에서 결측값이 생길 수 있습니다.

head(5)

	PassengerId	Survived	Pclass		Name	Sex	Age	SibSp	Parch		Ticket	Fare	Cabin	Embarked
0	1	0	3		Braund, Mr. Owen Harris	male	22.0	1	0		A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley	Florence Briggs Th.	female	38.0	1	0		PC 17599	71.2833	C85	C
2	3	1	3		Heikkinen, Miss. Laina	female	26.0	0	0		STON/O2 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0			113803	53.1000	C123	S
4	5	0	3		Allen, Mr. William Henry	male	35.0	0	0		37450	8.0500	NaN	S

df.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 284 non-null object
11 Embarked 889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

결측값 확인하기

결측값이 어디에 있는지 확인하는 것이 첫 단계입니다. Pandas는 isnull() 또는 isna() 메서드를 사용하여 결측값을 확인할 수 있습니다.

#isnull()
#isna()은 결측값
df.isnull()
df.isna()

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	False	False	False	False	False	False	False	False	False	False	True	False
1	False	False	False	False	False	False	False	False	False	False	True	False
2	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	True	False
...
886	False	False	False	False	False	False	False	False	False	False	True	False
887	False	False	False	False	False	False	False	False	False	False	False	False
888	False	False	False	False	False	True	False	False	False	False	True	False
889	False	False	False	False	False	False	False	False	False	False	False	False
890	False	False	False	False	False	False	False	False	False	False	True	False

891 rows × 12 columns

#isna()
df.isna()

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	False	False	False	False	False	False	False	False	False	False	True	False
1	False	False	False	False	False	False	False	False	False	False	True	False
2	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	True	False
...
886	False	False	False	False	False	False	False	False	False	False	True	False
887	False	False	False	False	False	False	False	False	False	False	False	False
888	False	False	False	False	False	True	False	False	False	False	True	False
889	False	False	False	False	False	False	False	False	False	False	False	False
890	False	False	False	False	False	False	False	False	False	False	True	False

891 rows × 12 columns

결측값 개수 확인하기

각 열에 결측값이 몇 개인지 확인하려면 sum() 메서드를 사용합니다.

각 열의 결측값 개수 확인
df.isnull().sum()

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	0	0	0	0	0	177	0	0	0	0	687	2
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0
...
886	0	0	0	0	0	0	0	0	0	0	0	0
887	0	0	0	0	0	0	0	0	0	0	0	0
888	0	0	0	0	0	0	0	0	0	0	0	0
889	0	0	0	0	0	0	0	0	0	0	0	0
890	0	0	0	0	0	0	0	0	0	0	0	0

891 rows × 12 columns

결측값 제거하기

결측값을 제거하는 방법은 dropna() 메서드를 사용합니다.

결측값이 있는 행 제거
df.dropna()은 결측값
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True)
df.dropna(inplace=True