

# 직장인을 위한 파이썬 데이터 분석

# Pandas Cheat Sheet

#### **Pandas**

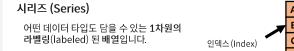


Pandas는 panel data analysis(패널 자료)의 약자로, 파이 썬을 활용한 정형 데이터 분석에서 가장 많이 활용되는 라이 브러리입니다.

#### 라이브러리 로딩

import pandas as pd

#### Pandas 데이터 구조



s = pd.Series([1, 2, 3, 4], index=['A', 'B', 'C', 'D'])

#### 데이터 프레임 (DataFrame)

인덱스 (Index) 와 열 (Columns) 로 이루어진 2차원의 라벨링(labeled) 된 데이터 구조입니다.

열 (Columns) 인구(만) 한국 서울 인덱스 중국 (Index) 도쿄 14.0

data = {'국가': ['한국', '국가': ['한국', '중국', '일본'], '수도': ['서울', '베이징', '도쿄'], '인구': [9828094, 19612368, 14049146]}

df = pd.DataFrame(data, columns=['국가', '수도', '인구'])

# 슬라이싱 (Slicing) / 인덱싱 (Indexing)

#### 값 선택

```
s['A']
                                  시리즈 (Series) 값 선택
df[1:]
                                  데이터 프레임 (DataFrame) 값 선택
          수도
    중국
         베이징 19612368
> 2 일본
          도쿄 14049146
```

#### 값 변경

s['A'] = 6'A' 인덱스에 위치한 값을 6으로 변경

#### loc / iloc

df.loc(0, ['국가']) > '한국'	라벨(인덱스 값)을 통한 선택
> 연구 df.iloc(0, 0) > '한국'	위치를 통한 선택

#### Boolean 인덱싱

s[(s < 1)   (s > 2)] > C 3 > D 4	값이 1 미만이거나 2 초과인 값 선택
df[df['인구']>12000000] > 국가 수도 인구 > 1 중국 베이징 19612368 > 2 일본 도쿄 14049146	'인구' 열 값이 12,000,000 이상인 데이터 선택

#### 데이터 삭제

s.drop(['A', 'B'])	행 방향 데이터 삭제 (axis = 0)
df.drop('국가', axis=1)	열 방향 데이터 삭제 (axis = 1)
ut.urob( 5/1 , axis=1)	들 당성 데이터 국제 (dxl2 - 1)

#### 정렬

df.sort_index(by='국가')	행 / 열 인덱스 기준 정렬
<pre>df.reset_index(drop=True)</pre>	인덱스 초기화 (기존 인덱스 열 삭제)

# 파일 읽기 / 쓰기

#### CSV 파일

```
df = pd.read csv('file.csv')
                           # 파일 읽기
df.to_csv('file.csv')
                            # 파일 저장
```

#### EXCEL 파일

```
# xlsx 파일 읽기
df = pd.read_excel('file.xlsx')
df.to_excel('file.xlsx', sheet_name='Sheet1') # 파일 저장
                                             # xls 파일 읽기
xlsx = pd.ExcelFile('file.xls')
df = pd.read excel(xlsx, 'Sheet1')
```

#### 데이터 베이스 (SQL) 연결

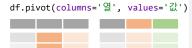
```
from sqlalchemy import create_engine
engine = create_engine('sqlite:///:memory:')
df=pd.read sql('SELECT * FROM my table;', engine) # sql 문으로 로딩
df=pd.read sql table('my table', engine)
                                                # 테이블명으로 로딩
df=pd.read_sql_query('SELECT * FROM my_table;', engine)
                                                # sal 파일 저장
df.to sql('my dataframe', engine)
```

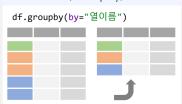
※ read\_sql() 은 read\_sql\_table() 과 read\_sql\_query() 를 묶어서 사용하는 것과 같습니다.

#### 데이터 프레임 결합

# 피벗(Pivot) 테이블

그룹 연산 (Groupby)





#### 데이터 프레임 결합

# Merge (겹치는 값을 기준으로 데이터 결합)

```
pd.merge(df1, df2, how='left', on='Key 열')
                                                Left Merge
pd.merge(df1, df2, how='right', on='Key 열')
                                                Right Merge
pd.merge(df1, df2, how='inner', on='Key 열')
                                                Inner Merge
pd.merge(df1, df2, how='outer', on='Key 열')
                                                Outer Merge
```

#### Concat (행 / 열 방향 데이터 결합)

pd.concat([df1,df2],	axis=0)	행 방향 결합
pd.concat([df1,df2],	axis=1)	열 방향 결합

## 결측 / 중복 처리

#### 결측 값 처리

df.dropna()	결측값 행 삭제
df.fillna('대체값')	결측값 대체

# 중복 값 처리

중복 행 삭제 df.drop\_duplicates()

#### 데이터 정보 확인

# 기본 정보 확인

df.info()	전체 데이터 정보
df.columns	열 (Columns) 이름
df.index	인덱스 (Index) 이름

# 요약 통계 확인

<pre>df.describe()</pre>	요약 통계 정보
df['열이름'].value_counts()	값 개수 집계
<pre>df.sum() / df.mean() / df.min() / df.max()</pre>	합계 / 평균 / 최소 / 최대값

# 시각화

