

Pandas

1) What is the max temperature in new york in the month of January.

in windspeed there are some missing data.

Cleaning messy data

`df.fillna(0, inplace=True)`

2) On which days did it rain

3) What was average speed of wind during the month.

Sol.

1) `df['Temp'].max()` ^{How to find the date.}

2) `df['EST'][df['Event']=='Rain']`

3) `df['windspeed'].mean()`

`a = df['windspeed'].sum()`

`b = a/30`

Dataframe

Some it is used to represent data with rows and columns.

- 1) Creating dataframe ✓
- 2) Dealing with rows and columns
- 3) ops: min, max, std, describe.
- 4) conditional selection
- 5) set index.

```
weather-data = {  
    'day': ['1/1', '1/2', '1/3']  
    'temp': [32, 35, 28]  
    'winds': [6, 7, 2]  
    'event': ['Rain', 'sunny', 'snow']  
}  
  
df = pd.DataFrame(weather-data)
```

① import pandas as

```
df = pd.read_csv('weather-data.csv')
```

creating dataframe by using
Dictionary.

wa

② rows, column = df.shape.

indexing-slicing.

0 1 2 3 4 5
1 5 1

df[2:5] → Column (2 to 4)

df.columns → shows all column

df.day or df['event']

type of dataframe

df[['day', 'event', 'temp']]

③

df.describe()

↑ satisfaction.

SQL

Select data in dataframe.

df[df.temp >= 32]

↑ == df.temp.max()

Select

df[['day']][df.temp == df['temp'].max(),
'temp']

df.index → represent range index

df.index inplace = True

df.set_index('day') → set the

index in day format / column.

df.loc['1/3/17']

↳ ~~see~~ shows the

df.reset_index(inplace = True)

ways of creating Data frame.

dt = pd.read_csv('weather.csv')
excel .xlsx

pd.DataFrame(weather) ← Dictionary.

Read / write

read_csv
read_excel
- hdf
- SQL
- json
- html
- stata
- sar

df = pd.read_csv('stock.csv', header = 1)
header = None, names = ['ticker', 'eps'])

df = pd.read_csv('stock.csv', nrows = 3) → ↑ create the header.
↑ represent only 3 rows.
, na_values = ['not available', 'Na']

Fungitac
Sertaconazole Nitrate


```
df = pd.read_csv('stock.csv', na_values = ["not available", 'n.a', -1])
```

Supply a dictionary

```
df = pd.read_csv("stock.csv", na_values = {  
    'eps' : ["not available", "n.a."],  
    'revenue' : ["not available", "n.a.", -1],  
})
```

→ write csv file.

```
df.to_csv('new.csv') ← creat new file in directory.
```

→ index = false) ← to remove index

→ columns = ['ticker', 'eps']) ← only two
column
in new.csv

- 1) fill
- 2) inter
- 3) drop

```
df = pd.  
type (
```

```
1/1  
df = pd.  
2017
```

```
df.se
```

Index

How to handle missing data. (panda-5)

1) fillna to fill missing value

2) interpolate to guess on missing values

3) dropna to drop missing value rows

```
df = pd.read_csv("weather.csv")
```

convert type into str → date.

type(df.day[0]) → str

↓
1/1/2017
df = pd.read_csv("weather.csv", parse_dates=["day"])
↓
2017-01-01
→ pandas.tseries.Timestamp

Index
df.set_index('day', inplace=True)

↑ modify the original dataframe
otherwise new dataframe.

Fungitac
Sertaconazole Nitrate

df.
new-df = df.fillna(0)

NaN \rightarrow 0

different value for column

new-df = df.fillna({

'temp': 0,

'wind': 0,

'event': 'no event'

})

new-df = df.fillna(method='bfill',
axis='column')

copying value horizontally

df.fillna(method='ffill', limit=1)

copy value only once.

[ffill]

new-df = df.fillna(method="ffill")

carry forward

4/7	32	sum
1/8	32	sum
4/9	32	snow

forward

backward

`new-df = df.interpolate()`

linear interpolation

1/1	32
4/1	30
5/1	28

$\frac{32+28}{2}$

←

missing 2/1 and 3/1

`df.interpolate(method="time")`

1/1	32
4/1	29
5/1	28

the closest time `df = df.neindex(idx)`

`new-df = df.dropna()`

`new-df` ↑ drop the rows where missing any data

`df.dropna(thresh=1)`

↑ keep the rows if there one parameter

`dt = pd.date_range("01-01-2017", "01-11-17")`

`idx = pd.DatetimeIndex(dt)`

`df = df.reindex(idx)`

`dt` ↑ How to get the dataframe?

Fungitac

Sertaconazole Nitrate

handling missing data (pandas)

replace function

```
df = pd.read_csv('weather.csv')
```

```
new_df = df.replace(-999, np.NaN)
```

↑ ↑
value replace value

```
df.replace([-999, -888], np.NaN)
```

```
new_df = df.replace({
```

'temp' : -999, *this dictionary*

'wind' : -999, *column wise*

'event' : '0'

```
}, np.NaN)
```

```
new_df
```

mapping

```
df.replace({
```

-999 : np.NaN

~~"noevent"~~ : "sunny"

} *in data from*

*this mapping using for replace
for this point (define) value.*

if there are given with unite.

Now chop off the unite

Temp	
32 f	61 mph
0 f	2 mph
28 f	5 mph

```
new_df = df.replace(['[A-Za-z]', ''], regex = True)
```

so provide a dictionary: → it remove event column

```
df.replace({
```

```
    'temp' : '[A-Za-z]',
```

```
    'wind' : '[A-Za-z]',
```

```
})
```

creating data frame :

```
df = pd.DataFrame({
```

```
    'score' : ['excep', 'avg', 'poor', 'avg', 'excep'],
```

```
    'student' : ['nod', 'Maya', 'prnt', 'tom', 'Jalilah']
```

```
})
```

Fungitac
Sertaconazole Nitrate

df.replace(['poor', 'avg', 'good', 'excl'], [1, 2, 3, 4])

Group By. panda-7

1) find maximum temperature in each of the cities.

2) find average wind speed per city.

g = df.groupby('city')

for city, city_df in g:

print(city) ^{How to work}

print(city_df)

print city ~~and~~ city info.

g.get_group('Mumbai')

A.1)

g.max() ← represent the

A.2)

g.mean()

g.describe()

Concat data frame

us - weather

india-weather = pd.DataFrame({

'city': ['mumbai', 'delhi', 'bangalore']

'temp': [32, 45, 50]

'humidity': [80, 60, 78]

})

df = pd.concat([india-weather, us-weather])

~~ignore index = True~~

keys = ['India', 'US']

0	mumbai	0
1	delhi	1
2	bangalore	2
0	NY	3
1	chicago	4
2	orlando	5

retrieve the data

df.loc["india"]

Fungitac

Sertaconazole Nitrate