

## First Lecture on Pandas

Create Read Update Delete (CRUD)

```
import pandas as pd
import numpy as np

#CREATE
data = {'Name': ['Aishwarya', 'Bhushan',
                'Chetan', 'Dhananjay', 'Eknath', 'Faiz', 'Ganesh', 'NA'],
        'Age': [25, 30, 35, 40, 45, 50, 55, None], #np.nan
        'City': ['Surat', 'Mumbai', 'Pune', 'Nagpur', 'Chennai', 'Delhi', 'Kolka
ta', 'NA'],
        'Status': ['Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'NA']}

df = pd.DataFrame(data)

#READ
df.to_csv(r"C:\Users\YourName\Documents\Parag.csv",
index=False) # Windows
#df.to_csv(r"/Users/YourName/Documents/Parag.csv",
index=False) # Mac/Linux

df.to_csv('Parag.csv', index=False)
df = pd.read_excel("data.xlsx")
k=pd.read_csv('Parag.csv')
print(k)
#print(df)

#print(df.head(2))
#print(df.tail(3))
#print(df.describe(include='all')) # Summary statistics
#print(df['Age'].mean())
#print(df['Status'].value_counts())

print(df.dropna(inplace=False))
#print(df)
#df['Values']=df['Age']+10
#print(df)
#print(k)
```

## 2<sup>nd</sup> Lecture on Pandas

```
import pandas as pd
import numpy as np
```

```
data = {'Name': ['Aishwarya', 'Bhushan',
'Chetan', 'Dhananjay', 'Eknath', 'Faiz', 'Ganesh', 'NA'],
'Age': [25, 30, 35, 40, 45, 50, 55, None], #np.nan
'City': ['Surat', 'Mumbai', 'Pune', 'Surat', 'Chennai', 'Delhi', 'Kolkata', 'NA'],
'Status': ['Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'NA']}
```

```
df = pd.DataFrame(data)
```

```
#print(df.columns)
#print(df.index)
print(df.iloc[1,:])
df.loc[0, "Name"] = 3
df.loc[df["Name"] == "Chetan", "Age"] = 30
df.loc[df["Name"] == "Chetan", "Name"] = 30
print(df)
#print(df.describe(include='all')) # Summary statistics df.info()
print(df['Age'].mean())
print(df['Status'].value_counts())
#print(df.dropna(inplace=False))
df.drop("City", axis=1, inplace=True) #2nd lecture
df.drop(0, axis=0, inplace=True) #2nd lecture
df = df.drop(3) #Similarly, fillna command can be used
print(df)
#df['Values'] = df['Age'] + 10
df['Age'] = df['Age'] + 10 #2nd lecture
#print(df)
```

```
print(df[df["City"] == "Surat"])
print(df[(df["City"] == "Surat") & (df["Status"] == "N")])
print(df[(df["City"] == "Surat") | (df["Status"] == "N")])
print(df[df["City"].str.startswith("S")])
print(df[df["City"].str.endswith("i")])
print(df[df["City"].isin(["Surat", "Mumbai", "Delhi"])])
```

- **Important Commands**

Category	Common Commands	Purpose / Examples
<b>1. Creation / Input-Output (I/O)</b>	<code>DataFrame()</code> , <code>read_csv()</code> , <code>read_excel()</code> , <code>to_csv()</code> , <code>to_excel()</code>	Creating or loading data
<b>2. Inspection / Overview</b>	<code>head()</code> , <code>tail()</code> , <code>info()</code> , <code>describe()</code> , <code>shape</code> , <code>columns</code> , <code>index</code>	Basic summary and structure
<b>3. Selection &amp; Indexing</b>	<code>loc[]</code> , <code>iloc[]</code> , <code>at[]</code> , <code>iat[]</code> , <code>[]</code>	Accessing rows, columns, and cells
<b>4. Filtering &amp; Conditional Selection</b>	Boolean indexing ( <code>df[df['Age'] &gt; 30]</code> ), <code>isin()</code> , <code>query()</code>	Selecting rows by condition
<b>5. Modification / Updating</b>	<code>assign()</code> , <code>replace()</code> , <code>rename()</code> , <code>astype()</code> , <code>map()</code> , <code>apply()</code>	Changing or updating data
<b>6. Missing Data Handling</b>	<code>isna()</code> , <code>notna()</code> , <code>fillna()</code> , <code>dropna()</code> , <code>interpolate()</code>	Handling NaN / missing values
<b>7. Data Cleaning</b>	<code>drop_duplicates()</code> , <code>sort_values()</code> , <code>reset_index()</code> , <code>set_index()</code>	Tidying and sorting data
<b>8. Aggregation / Grouping</b>	<code>groupby()</code> , <code>agg()</code> , <code>mean()</code> , <code>sum()</code> , <code>pivot_table()</code>	Summarizing data
<b>9. Merging / Combining</b>	<code>concat()</code> , <code>merge()</code> , <code>join()</code> , <code>append()</code>	Combining multiple DataFrames
<b>10. Visualization</b>	<code>plot()</code> , <code>hist()</code> , <code>boxplot()</code>	Quick visual analysis (built on Matplotlib)

### 3<sup>rd</sup> Lecture on Pandas

```
import pandas as pd
import numpy as np

data = {'Name': ['Aishwarya', 'Bhushan', 'Chetan', 'Dhananjay',
                'Eknath', 'Faiz', 'Ganesh', None],
        'Age': [25, 30, np.nan, 40, 45, 50, None, 28],
        'City': ['Surat', 'Mumbai', 'Pune', 'Nagpur', 'Chennai',
                None, 'Kolkata', 'Delhi'],
        'Status': ['Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'Y']}

df = pd.DataFrame(data)

print("Original DataFrame:\n", df)

# Detect missing values
print("\nMissing values:\n", df.isnull().sum())

# Fill missing values
df['Age'] = df['Age'].fillna(df['Age'].mean())
df['City'] = df['City'].fillna('Unknown')
df['Name'] = df['Name'].fillna('Anonymous')
#print(df)

# Replace specific values
```

```

df['Status']=df['Status'].replace({'Y': 'Yes', 'N': 'No'})
#print(df)

# Rename columns
df=df.rename(columns={'Name': 'FullName', 'City': 'Location'})
#print(df)

# Change datatype
df['Age'] = df['Age'].astype(int)
print(df)

# Remove duplicates
df=df.drop_duplicates()
print(df)

# Sort values
#df=df.sort_values(by='Age', ascending=False)

#print("\nCleaned DataFrame:\n", df)

```

## 4<sup>th</sup> Lecture on Pandas

```

#7th CPC
import pandas as pd
import matplotlib.pyplot as plt
df=pd.read_csv('/Users/paragthakur/Desktop/Teaching Material/Data
Science for chemical Engineers/Salary.csv', header=2)
print(df.columns); print(df.head()); print(df.info());
print(df.describe())
print(df['1']); print(df.iloc[:,1]);
print(df[df['6']==47600])
print(df[df['6']!=47600])
print(df[(df['1']==35400) | (df['6']==35400)])
print(df[(df['1']<=25000) & (df['6']<=50000)])
print(df[~((df['1']<=25000) & (df['6']<=50000))])
print(df)
df.rename(columns={'1': 'Level 1', '2': 'Level 2', '3': 'Level 3',
                    '4': 'Level 4', '5': 'Level 5', '6': 'Level 6',
                    '7': 'Level 7', '8': 'Level 8', '9': 'Level 9',
                    '10': 'Level 10', '11': 'Level 11', '12': 'Level 12',
                    '13': 'Level 13', '14': 'Level 14', '15': 'Level 15',
                    '16': 'Level 16', '17': 'Level 17', '18': 'Level 18'},
inplace=True)
#print(df.columns)
#print(df)
j=df.T
#print(j)
#print(df.max(axis=1))
(df.sort_values(by='3', ascending=False))
j=df.isna().max()

```

```

df.fillna(j, inplace=True)
#df[['1','3']].plot(kind='line',marker='*')
plt.scatter(df['Level'], df['1'], label='Stage 1', color='blue')
plt.scatter(df['Level'], df['3'], label='Stage 3', color='red')
plt.title("Comparison of Stage 1 and Stage 3 Pay Levels")
plt.xlabel("Stages of yearly Increment")
plt.ylabel("Basic Pay (₹)")
plt.show()
#print(df)
plt.show()

```

## 5<sup>th</sup> Lecture on Pandas

Focus Area	Key Commands
Grouping & Aggregation	groupby(), agg(), mean(), sum()
Combining Data	merge(), join(), concat()
Pivot & Advanced Analysis	pivot_table(), corr(), cov()

```

import pandas as pd

data = {'Department': ['HR', 'IT', 'IT', 'HR', 'Sales', 'Finance', 'Sales'],
        'Employee': ['A', 'B', 'C', 'D', 'E', 'F', 'G'],
        'Salary': [50000, 60000, 70000, 75000, 55000, 52000, 65000],
        'Experience': [3, 5, 7, 8, 2, 4, 6]}

df = pd.DataFrame(data)
#print("Original Data:\n", df)

# Group by department and get average salary
grouped = df.groupby('Department')
#print(grouped)
#print(grouped['Salary'].mean()) # Mean salary by department
#print(grouped['Experience'].sum()) # Total experience by department
#print("\nAverage Salary by Department:\n",
df.groupby('Department')['Salary'].mean())

# Multiple aggregation
agg_data = grouped.agg({'Salary': ['mean', 'max'], 'Experience': 'sum'})
#print("\nAggregated Data:\n", agg_data)

summary = grouped.agg({
    'Salary': ['mean', 'max', 'min'],
    'Experience': ['mean', 'count']})

#print(summary)
#print(summary.sort_values(('Salary', 'mean'), ascending=False))

```

```

#print(grouped['Salary'].mean().sort_values(ascending=False))

# Custom function
def salary_range(x):
    return x.max() - x.min()

#print("\nSalary Range by Department:\n",
df.groupby('Department')['Salary'].apply(salary_range))

# Pivot table
pivot = df.pivot_table(values='Salary', index='Department',
aggfunc=['mean', 'max', 'min'])
print("\nPivot Table:\n", pivot)

# DataFrames for merge
df1 = pd.DataFrame({'EmpID': [1, 2, 3, 4],
                    'Name': ['Aishwarya', 'Bhushan', 'Chetan',
'Dhananjay'],
                    'Department': ['HR', 'Finance', 'IT', 'IT']})

df2 = pd.DataFrame({'EmpID': [3, 4, 5],
                    'Salary': [70000, 80000, 60000]})

print("Left DataFrame:\n", df1)
print("\nRight DataFrame:\n", df2)

# Inner Join
inner_join = pd.merge(df1, df2, on='EmpID', how='inner')
print("\nInner Join:\n", inner_join)

# Left Join
left_join = pd.merge(df1, df2, on='EmpID', how='left')
print("\nLeft Join:\n", left_join)

# Outer Join
outer_join = pd.merge(df1, df2, on='EmpID', how='outer')
print("\nOuter Join:\n", outer_join)

# Concatenation example
df3 = pd.DataFrame({'EmpID': [6, 7],
                    'Name': ['Eknath', 'Faiz'],
                    'Department': ['HR', 'Finance']})

concat_df = pd.concat([df1, df3], ignore_index=True)
print("\nConcatenated DataFrame:\n", concat_df)

```