

Pre-Processing in Structured Data Lab Session- I(b)

TIET, PATIALA

Pandas Introduction

Import pandas as pd

S.No	Feature	Syntax & Examples
1.	Creating Data Frame	<code>df =pd.DataFrame()</code>
2.	Adding Columns	<code>df['Name']=['abc','xyz']</code> <code>df['age']=[38,25]</code>
3.	Loading a Data Frame	<code>df=pd.read_csv(url/path)</code> <code>df=pd.read_csv('C:/Users/jasme/Desktop/titanic.csv')</code>
4.	Navigating Data Frame	<code>df.iloc[row number/slice]</code> <code>df.iloc[4], df.iloc[1:4], df.iloc[:, df.iloc[1:4, 5:8]</code>
5.	Conditional Row Selection	<code>df[condition]</code> <code>df[df['Sex']=='female']</code> or <code>df[(df['Sex']=='female') &(df['Age']>='65')]</code>

Pandas Introduction Contd....

Import pandas as pd

S.No	Feature	Syntax & Examples
6.	Replacing Values	<code>df.replace(old_value,new_value)</code> <code>df.replace("female","Woman")</code> <code>df['Sex'].replace(["female","male"],["woman","man"])</code>
7.	Renaming Columns	<code>df.rename(columns={'Pclass':'Pessanger_Class'})</code>
8.	Mathematical Functions	<code>print(df['Age'].max(), print(df['Age'].min())</code> <code>print(df['Age'].sum(), print(df['Age'].mean())</code> <code>print(df['Age'].count())</code>
9.	Unique Values	<code>print(df['Sex'].unique())</code> <code>print(df['Sex'].nunique())</code> <code>print(df['Sex'].value_counts())</code>
10.	Deleting Columns	<code>df.drop(['Age'],axis=1)</code> <code>df.drop(df.columns[1],axis=1)</code>

Pandas Introduction Contd....

Import pandas as pd

S.No	Feature	Syntax & Examples
11.	Deleting rows/duplicate rows	<code>df[df['Sex']!='male']</code> or <code>df.drop_duplicates()</code>
12.	Grouping rows	<code>df.groupby('Sex').sum()</code>
13.	Looping over Column	<code>[name.upper() for name in df['Name']]</code> Or <code>for name in df['Name']:</code> <code> print(name.upper())</code>
14.	Applying Functions Over all Elements of Column	<code>df['Age'].apply(np.sqrt)</code> <code>df.groupby('Sex').apply(lambda x: x.count())</code>

Data Cleaning- Missing Values

Feature	Syntax/ Example
1. Check Missing Values	<pre>print(pd.isna(df)) print(pd.isna(df['Pclass']))</pre>
2. Deleting Missing Values Rows	<pre>df.dropna()</pre>
3. Replacing Missing value with scalar value	<pre>df['Age'].fillna(0) or df.fillna(0)</pre>
4 Fill NA forward and backward	<pre>df.fillna(method='pad') df.fillna(method='bfill')</pre>
5. Central Tendency Imputation (mean, median, most_frequent, constant)	<pre>from sklearn.impute import SimpleImputer mean_imputer=SimpleImputer(missing_values=np.nan,strategy="mean") mean_imputer.fit_transform((df['Age'].values).reshape(1,-1))</pre>
6. Nearest Neighbor Imputer	<pre>from fancyimpute import KNN X_filled_knn = KNN(k=1).fit_transform((df['Age'].values).reshape(1,-1))</pre>

Data Cleaning- Noisy Data

Binning:

```
import pandas as pd
df=pd.read_csv('C:/Users/jasme/Desktop/titanic.csv')
import numpy as np
bins=np.linspace(min(df['Age']),max(df['Age']),4)
group_names=['child','young','old']
df['Age']=pd.cut(df['Age'],bins,group_names,include_lowest=True)
```

Data Cleaning- Handling Outliers

S.No	Feature	Syntax / Example
1.	Detecting Outliers (Box Plots, Scatter Plots)	<pre>import seaborn as sns sns.boxplot(x=df['Age'])</pre>
2.	Removing Outlier (using z score)	<pre>from scipy import stats import numpy as np df=df.fillna(method='pad') z = np.abs(stats.zscore(df['Age'])) print(np.where(z > 2.2))</pre>
3.	Removing Outlier (using z score)	<pre>Q1 = df.quantile(0.25) Q3 = df.quantile(0.75) IQR = Q3 - Q1 print(df < (Q1 - 1.5 * IQR)) (df > (Q3 + 1.5 * IQR))</pre>

Data Transformation- Scaling

S.No	Feature	Synatx Example
1.	MinMax Scaler	<pre>from sklearn.preprocessing import MinMaxScaler scaler = MinMaxScaler() df['Age']=scaler.fit_transform(df['Age'].values.reshape(-1, 1)).flatten()</pre>
2	Standard Scaler (z-score)	<pre>from sklearn.preprocessing import StandardScaler scaler = StandardScaler() df['Age']=scaler.fit_transform(df['Age'].values.reshape(-1, 1)).flatten()</pre>
3	Robust Scaler (IQR)	<pre>from sklearn.preprocessing import RobustScaler scaler = RobustScaler() df['Age']=scaler.fit_transform(df['Age'].values.reshape(-1, 1)).flatten()</pre>
4	MaxAbs Scaler	<pre>from sklearn.preprocessing import MaxAbsScaler scaler =MaxAbsScaler() df['Age']=scaler.fit_transform(df['Age'].values.reshape(-1, 1)).flatten()</pre>

Data Transformation- Encoding

S.No	Feature	Synatx Example
1.	Label Encoder	<pre>from sklearn import preprocessing le = preprocessing.LabelEncoder() le.fit([1, 2, 2, 6])</pre>
2	One Hot Encoder	<pre>from sklearn.preprocessing import OneHotEncoder onehotencoder = OneHotEncoder() onehotencoder.fit_transform(data)</pre>