



Lernmodul zum Datenimport und Datenvorbereitung mit pandas

Import und Export

Herstellen einer Datenbankverbindung mit sqlalchemy	<pre>from sqlalchemy import create_engine engine = create_engine('sqlite:///pfad-zur-datenbank') db_connection = engine.connect()</pre>
Import von Daten aus einer Datenbank	<pre># Voraussetzung: Verbindung zur Datenbank hergestellt df = pd.read_sql_query(sql_query, db_connection)</pre>
Import von Daten aus einer csv Datei	<pre>df = pd.read_csv(pfad-zur-csv-datei)</pre>
Mergen von Daten	<pre>merged_df = df1.merge(df2, on='spaltenname-zum-mergen')</pre>
Export von Daten in eine Datenbank	<pre># Voraussetzung: Verbindung zur Datenbank hergestellt df.to_sql('name-für-db-tabelle', db_connection, index = False, if_exists="replace")</pre>
Export von Daten in eine csv Datei	<pre>df.to_csv('pfad-für-csv-datei', index=False)</pre>

Datenanalyse

DataFrame anzeigen	<pre>df.head() # die ersten 5 Zeilen df.tail() # die letzten 5 Zeilen</pre>
Erstellen einer Kopie des DataFrames	<pre>df_copy = df.copy()</pre>
Beschreiben der Daten	<pre>df.shape # (Anzahl der Zeilen, Anzahl der Spalten) df.describe() # Anzeigen der Lageparameter df['spaltenname'].unique() # Anzeigen der einzigartigen Werte df['spaltenname'].value_counts() # Anzeigen der Vorkommen der Werte innerhalb einer Spalte</pre>
Zugriff auf Datenwerte	<pre>df['spaltenname'] # Zugriff auf eine Spalte df['spaltenname'][index] # Zugriff auf einen Wert einer Spalte df.iloc[zeilenindex,spaltenindex] # Index-basierte Auswahl von Daten df.loc[zeilenindex, 'spaltenname'] # label-basierte Auswahl von Daten</pre>
Bedingte Auswahl von Daten	<pre>df.loc[df['spaltenname'] == 'auszuwählender-wert'] # Mehrere Bedingungen mit & oder verknüpfen</pre>
Umbenennen von Spaltennamen	<pre>df.rename(columns = {'alter-name':'neuer-name'}, inplace = True) df = df.rename(columns = {'alter-name':'neuer-name'})</pre>
Datentypen der Spalten	<pre>df['spaltenname'].dtype # Anzeigen des Datentyps einer Spalte df.dtypes # Anzeigen der Datentypen jeder Spalte df.astype({'spaltenname':datatype}) # Konvertieren des Datentyps einer Spalte df['spaltenname'] = pd.to_datetime(df['spaltenname']) # Konvertieren in Datumswerte</pre>
Sortieren der Daten	<pre>df.sort_values('spaltenname') # Sortieren nach einer Spalte df.sort_index() # Sortieren nach den Indexwerten</pre>

Indexing Operatoren

	Python Indexing Operator	pandas iloc Operator	pandas loc Operator
Anwendung	<code>df[Spaltenname][Zeilenindex]</code> <code>df[[Spaltennamen]]</code>	<code>df.iloc[Zeilenindex, Spaltenindex]</code> : für alle Zeilen/Spalten	<code>df.loc[Zeilenindex, Spaltenname]</code> : für alle Zeilen/Spalten
Spalten-Zeilen-Reihenfolge	1. Spalten 2. Zeilen	1. Zeilen 2. Spalten	1. Zeilen 2. Spalten
Umgang mit dem letzten Indexwert	Nur einzelne Indizes	Nicht inklusive	Inklusive
Anwendungsbeispiele	# Alle Zeilen der Spalte 'title' <code>df['title']</code>	# Alle Zeilen der Spalte 'title' (Index 2) <code>df.loc[:, 2]</code> # Ausgabe der Zeilen 2 bis 4 <code>df.loc[2:5]</code> # Ausgabe der ersten 3 Zeilen und Spalten <code>df.loc[:3, :3]</code>	# Alle Zeilen der Spalte 'title' <code>df.loc[:, 'title']</code> # Ausgabe der Zeilen 2 bis 4 <code>df.loc[2:4]</code> # Ausgabe der ersten 3 Zeilen und Spalten <code>df.loc[:3, [Spaltennamen]]</code>

Datenbereinigung

Anzeigen mit fehlenden Werten	<code>df.isnull().sum()</code> # Anzahl der Null Werte je Spalte <code>df[pd.isnull(df.spaltennamen)]</code> # Zeilen mit fehlenden Werten in einer Spalte
Ersetzen von fehlenden Werten	<code>df.spaltenname.fillna('neuer-wert', inplace=True)</code>
Entfernen von Zeilen mit fehlenden Werten in einer bestimmten Spalte	<code>df = df.dropna(subset=['spaltenname'])</code>
Anzeigen von Duplikaten	<code>df.duplicated().sum()</code> # Anzahl der Duplikate <code>df[df.duplicated(subset = 'spaltenname')]</code> # Duplikate innerhalb einer Spalte
Entfernen von Duplikaten	<code>df.drop_duplicates(inplace=True)</code>

Datenmanipulation

Hinzufügen von Zeilen	<code>df = df.append(anzuhängender-dataframe, ignore_index=True)</code>
Löschen von Zeilen	<code>df.drop(zu-löschender-index-wert, axis=0, inplace=True)</code> <code>df.drop(['spaltenname1', 'spaltenname2'], axis=0, inplace=True)</code>
Erzeugen einer Spalte	<code>df['neuer-spaltenname'] = # Neue Werte z. B. mittels map()</code>
Löschen einer Spalte	<code>df.drop('zu-löschender-spaltenname', axis=1, inplace=True)</code> <code>df.drop(['spaltenname1', 'spaltenname2'], axis=1, inplace=True)</code>
Ersetzen von Datenwerten	<code>df.spaltenname.replace('alter-wert', 'neuer-wert')</code>
Mappen von Datenwerten	<code>df['spaltenname'] = df['spaltenname'].map(lambda d: #transformiere den Datenwert hier)</code> <code>df['spaltenname'] = df['spaltenname'].apply(funktionsname)</code>
Gruppieren von Daten	<code>df.groupby('spaltenname')</code> # nach einer Spalte <code>df.groupby(['spaltenname1', 'spaltenname2'])</code> # nach mehreren Spalten
Binning	<code>bins = [# Definition der Grenzen der einzelnen Bins]</code> <code>labels = [# Labels der einzelnen Bins]</code> <code>df['spaltenname'] = pd.cut(df['spaltenname'], bins, labels)</code>
One-hot Encoding	<code>df = pd.concat([df, pd.get_dummies(df['spaltenname'])], axis=1)</code>