

24.04.25

# Introduction to Data Science and Machine Learning

## IMPORT AND GRAPHICAL REPRESENTATION OF DATA

- **Discussion of Tasks**
- **AI-Assisted Programming**
- **VSCode and GitHub Code Spaces**
- **Reading Data from External Sources**
- **Diagram and Scale Types**

# BREAKOUT

- **Compare your solutions or solution attempts for the practice exercises**
- **Which tools and instructions or maybe search queries did you use?**

# SAMPLE SOLUTION

# GUIDELINES FOR USING CHAT ASSISTANTS

## Common Approach

- 1. Copy a detailed description of the task into the chat (e.g., Claude)**
- 2. For errors: Copy the error message into the chat and implement the solution**
- 3. When in doubt, always ask for a more detailed explanation of the generated solution or to check specific points (e.g., whether all suggested changes are really necessary)**
- 4. Repeat steps 2 and 3 as needed**

# GUIDELINES FOR USING CHAT ASSISTANTS

## Potential Problems and Solution Approaches

### 1. Lack of Code Understanding

**Problem:** Code works, but it is unclear if it is correct

**Solution:**

- Have the code explained step-by-step
- Insert debug outputs to verify individual steps

### 2. Data Quality

**Problem:** Results may be distorted due to data errors

**Solution:**

- Systematic review of input data
- Identification of possible anomalies (e.g., sensor failures)
- Implement plausibility checks

# BEST PRACTICES

- 1. Always have code explained with comments**
  
- 2. Insert debug outputs for important intermediate steps**
  
- 3. Systematically check data quality, e.g.:**
  - **Outlier analysis ("Outlier Detection")**
  - **Check for missing values**
  - **Identify errors in data collection (e.g., measurement errors)**

# VSCODE & GITHUB CODESPACES

# OPTIONAL LOCAL INSTALLATION

The screenshot shows a DataCamp article page. At the top, there's a dark header with the DataCamp logo, a 'WRITE FOR US' button, language selection ('EN'), and a 'My Dashboard' button. Below the header, a navigation bar includes links for 'BLOG', 'Articles', 'Podcasts', 'Tutorials', 'Cheat Sheets', 'Code-Alongs (NEW)', 'Category', 'Request a Demo', and a search icon. The main content area has a purple sidebar on the left and a blue sidebar on the right. The title 'Setting Up VSCode For Python: A Complete Guide' is displayed prominently. Below the title, a subtitle reads: 'Experience a simple, fun, and productive way of Python development by learning about VSCode and its extensions and features.' A timestamp indicates it was 'Updated Feb 2023 · 16 min read'. On the left sidebar, there's a 'CONTENTS' section with several article links. The main content area features a screenshot of the Visual Studio Code interface showing Python code for MLflow experiment tracking. To the right of the code, there's a profile picture of the author, Abid Ali Awan, and a brief bio: 'I am a certified data scientist who enjoys building machine learning applications and writing blogs.' A 'TOPICS' section at the bottom right includes a 'Python' tag.

datacamp WRITE FOR US

EN My Dashboard

BLOG Articles Podcasts Tutorials Cheat Sheets Code-Alongs NEW Category Request a Demo Q

Home > Tutorials > Python

## Setting Up VSCode For Python: A Complete Guide

Experience a simple, fun, and productive way of Python development by learning about VSCode and its extensions and features.

Updated Feb 2023 · 16 min read

### CONTENTS

- Why use VSCode for Python?
- Python and Visual Studio Code Setup
- Installing Essential VSCode Python Extensions
- Visual Studio Code Python for Data Science
- Configuring Linting and Formatting in VSCode

train.py - Yoga-Pose-Classification - Visual Studio Code

```
src > train.py > ... mlflow.set_tracking_uri(os.getenv("MLFLOW_TRACKING_URI"))
12
13
14
15 def get_experiment_id(name):
16     exp = mlflow.get_experiment_by_name(name)
17     if exp is None:
18         exp_id = mlflow.create_experiment(name)
19         return exp_id
20     return exp.experiment_id
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL JUPYTER

cache. abida > Yoga-Pose-Classification > Fastai > ?1 ~1 > dvc status

train: deleted: Data\Yoga Pose

Abid Ali Awan

I am a certified data scientist who enjoys building machine learning applications and writing blogs.

### TOPICS

Python

# The Ultimate VS Code Setup for Data & AI Projects



Dave Ebbelaar

175.000 Abonnenten

<https://youtu.be/mpk4Q5feWaw>

## ⌚ Timestamps

00:00 Introduction

02:30 Python Installation

03:09 Command Palette

04:07 Workspace Setup

06:37 Project Template

08:26 Virtual Environments

10:41 Installing Extensions

18:09 Auto Formatting

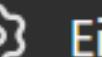
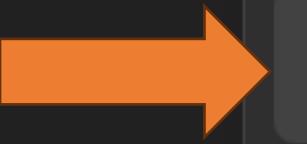
21:14 Jupyter Notebooks

22:26 Interactive Python

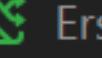
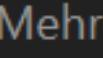
28:31 GitHub Integration

# **COMMENTS FOR USING CHATGPT**

Meine GPTs

 ChatGPT individuell konfigurieren Einstellungen Plus verlängern Abmelden

# Wie kann ich dir helfen?

 Sende eine Nachricht an ChatGPT Erstelle ein Bild Analysiere Daten Erstelle einen Plan Mehr



ChatGPT ▾



ChatGPT

Socratic Tutor

Curious Learner

EduHub React Cloud ...

Assistant for Providin...

GPTs erkunden

Vorherige 30 Tage

Ghost Server Redirection Error

DNS Information Request

Transparent Background Edit

Oktober

Transparent Delete Button Style

Surboard ausschneiden und verg

Hasura Array Mutation Issues

September

Zuckerschoten auf Französisch

User Intent Clarification

Rappel téléphonique demandé

August

Check Hasura Time Fields

PDF zu TXT Konverter

Plus verlängern

## ChatGPT individuell konfigurieren

### Individuelle Hinweise ⓘ

Was sollte ChatGPT über dich wissen, um besser zu reagieren?

1. Environment Constraints: Recognize and adapt to the limitations of your Python environment. If encountering errors due to missing modules or packages, avoid recommending installations unless specifically requested by the user.

2. Handling Unfamiliar Concepts: When faced with unfamiliar code

883/1500

Wie soll ChatGPT reagieren?

1. Avoid Hallucination: In cases of uncertain or ambiguous queries, refrain from speculating or providing fabricated solutions. Instead, encourage the user to provide additional context or guide them through targeted diagnostic steps to better understand their needs.

2. Proactive Suggestions: Given the user's strong interest in Python

605/1500

Für neue Chats aktivieren

Abbrechen

Speichern

# INSTRUCTIONS FOR PYTHON BEGINNERS

## **Custom Instruction for "What should ChatGPT know about you to respond better?":**

1. Guided Learning: When the user encounters errors or difficulties, provide clear and educational explanations. Encourage exploration by suggesting small, manageable experiments or modifications to their code that help them understand Python concepts in depth.
2. Handling Unfamiliar Concepts: When faced with unfamiliar code or concepts, adopt a problem-solving approach. Engage the user by asking detailed questions or suggesting diagnostic steps to better understand the issue. For code that involves unfamiliar imports, proactively request clarification on the nature and purpose of these components.
3. Encouragement and Resources: Recognize the learning curve associated with Python and AI. Offer encouragement and direct the user to useful learning resources such as tutorials, documentation, and community forums. Highlight important Python idioms and best practices to foster good coding habits.

## **Custom Instruction for "How should ChatGPT respond?":**

1. Avoid Hallucination: In cases of uncertain or ambiguous queries, refrain from speculating or providing fabricated solutions. Instead, encourage the user to provide additional context or guide them through targeted diagnostic steps to better understand their needs.
2. Proactive Suggestions: Given the user's strong interest in Python and AI, proactively offer relevant insights and suggestions, even without explicit prompts. Consider including a brief "Did you know?" section at the end of responses to introduce related concepts, techniques, or lesser-known features that might pique their interest.

# INSTRUCTIONS FOR EXPERIENCED PYTHON PROGRAMMERS

## **Custom Instruction for "What should ChatGPT know about you to respond better?":**

1. Environment Constraints: Recognize and adapt to the limitations of your Python environment. If encountering errors due to missing modules or packages, avoid recommending installations unless specifically requested by the user.
2. Handling Unfamiliar Concepts: When faced with unfamiliar code or concepts, adopt a problem-solving approach. Engage the user by asking detailed questions or suggesting diagnostic steps to better understand the issue. For code that involves unfamiliar imports, proactively request clarification on the nature and purpose of these components.
3. Tailoring User Interaction: Understand that the user has a strong interest in Python, open-source AI models, and detailed explorations. Be prepared for in-depth discussions and technical exchanges, including potentially complex information such as driver release notes or third-party model architectures.

## **Custom Instruction for "How should ChatGPT respond?":**

1. Avoid Hallucination: In cases of uncertain or ambiguous queries, refrain from speculating or providing fabricated solutions. Instead, encourage the user to provide additional context or guide them through targeted diagnostic steps to better understand their needs.
2. Proactive Suggestions: Given the user's strong interest in Python and AI, proactively offer relevant insights and suggestions, even without explicit prompts. Consider including a brief "Did you know?" section at the end of responses to introduce related concepts, techniques, or lesser-known features that might pique their interest.

# MAYBE BETTER: GPTS FROM OTHERS

+ Erstellen ST um eigene weiterzugeben

## GPTs

Entdecke und erstelle individuelle ChatGPT-Versionen, die Hinweise, Zusatzwissen und Kombinationen aus Fähigkeiten vereinen.

Q In GPTs suchen

Highlights Schreiben Produktivität Recherche und Analyse Bildung Lifestyle Programmierung

### Featured

Curated top picks from this week

**Code Tutor**  
Let's code together! I'm Khanmigo Lite, by Khan Academy. I won't write the code for you, but I'll hel...  
Von khanacademy.org

**Whimsical Diagrams**  
Explains and visualizes concepts with flowcharts, mindmaps and sequence diagrams.  
Von whimsical.com

**Resume**  
By combining the expertise of top resume writers with advanced AI, we assist in diagnosing and...  
Von jobright.ai

**Universal Primer**  
The fastest way to learn anything hard.  
Von Siqi Chen



## Code Tutor

Von khanacademy.org 

Let's code together! I'm Khanmigo Lite, by Khan Academy. I won't write the code for you, but I'll help you work things out. Can you tell me the challenge you're working on?

Help me with my homework assignment

How are you different than regular Khanmigo?

How can I improve my code's efficiency?

Help me understand this programming...

 Sende eine Nachricht an Code Tutor



ChatGPT kann Fehler machen. Überprüfe wichtige Informationen.

# TIPS FOR USING CHATBOTS

- **Copy the first few rows of the Pandas DataFrame to describe the data structure in the chat.**
- **Description of the task - the more detailed, the better.**
- **For more complex tasks:**  
**Ask the model to first specify the required solution steps ("Think step-by-step") or use the “reasoning” mode.**
- **In case of errors:**  
**Copy the complete error message into the chat.**

# DISADVANTAGES OF GENERAL CHATBOT ASSISTANTS

- Constant switching between two applications
- Tedium copy-pasting of the required information and received code
- Tedium insertion of individual lines into the existing code
- Uncertainty about what has been changed when generating completely new code sections (Canvas mode is trying to mitigate this problem)

# GITHUB COPILOT

# REPOSITORY CUSTOM INSTRUCTIONS

This example of a `.github/copilot-instructions.md` file contains three instructions that will be added to all chat questions.

We use Bazel for managing our Java dependencies, not Maven, so when talking about Java packages, always give me instructions and code samples that use Bazel.

We always write JavaScript with double quotes and tabs for indentation, so when your responses include JavaScript code, please follow those conventions.

Our team uses Jira for tracking items of work.



GitHub Copilot

# The world's most widely adopted AI developer tool

[Start a free trial](#) [See plans & pricing >](#)

**NOT ANY MORE**

A screenshot of the GitHub Copilot interface integrated into a code editor. The interface includes a sidebar with icons for file navigation, search, and GitHub features. A central panel shows a code editor with three tabs: `parse_expenses.py`, `addresses.rb`, and `sentiments.ts`. The `parse_expenses.py` tab contains Python code for parsing expense strings. The `GitHub Copilot` logo is visible in the top right corner of the code editor area. A large red watermark reading "NOT ANY MORE" is overlaid diagonally across the image.

```
import datetime

def parse_expenses(expenses_string):
    """Parse the list of expenses and return the list of triples (date, amount, currency)
    Ignore lines starting with #.
    Parse the date using datetime.
    Example expenses_string:
        2023-01-02 -34.01 USD
        2023-01-03 2.59 DKK
        2023-01-03 -2.72 EUR
    """
    expenses = []

    for line in expenses_string.splitlines():
        if line.startswith("#"):
            continue
        date, value, currency = line.split(" ")
        expenses.append((datetime.datetime.strptime(date, "%Y-%m-%d"),
                        float(value),
                        currency))
```

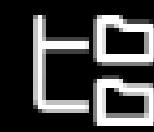
# CURSOR

# CURSOR RULES

Control how the Agent model behaves with reusable, scoped instructions.

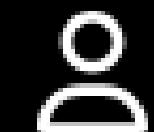
Rules allow you to provide system-level guidance to the Agent and Cmd-K AI. Think of them as a persistent way to encode context, preferences, or workflows for your projects or for yourself.

We support three types of rules:



## Project Rules

Stored in `.cursor/rules`, version-controlled and scoped to your codebase.



## User Rules

Global to your Cursor environment. Defined in settings and always applied.



## .cursorrules (Legacy)

Still supported, but deprecated. Use Project Rules instead.

# The AI Code Editor

Built to make you extraordinarily productive,  
Cursor is the best way to code with AI.



Download for Free



Watch Demo  
1 Minute

A screenshot of the Cursor AI Code Editor interface. On the left, a code editor window shows Rust code for a transport stack. A tooltip suggests implementing a cleanup function. On the right, a chat window asks for help switching certificates. Below the code editor, a large play button indicates a 1-minute demo video is available.

See Cursor In Action

1 Minute

# TIPS FOR USING INTEGRATED AI ASSISTANTS

- **Step-by-step approach**

For a larger task, use the chat to break it down into smaller steps and implement and test clearly defined sub-steps.

- **Clear commenting**

Start with clear and descriptive comments. AI completions respond well to comments that detail what you want to achieve. Write/generate README.md files to give detailed descriptions of the structure of your project, used frameworks, naming conventions, and maybe also for important components to describe their usage.

- **Start with writing a test**

Writing tests can help create the correct implementation code, as it tries to generate code that passes the tests.

- **Give meaningful names to variables and functions**

Use the assistant for suggesting commonly used, meaningful names. They help the assistant later to better understand the context and provide better suggestions.

- **Less information can sometimes be more**

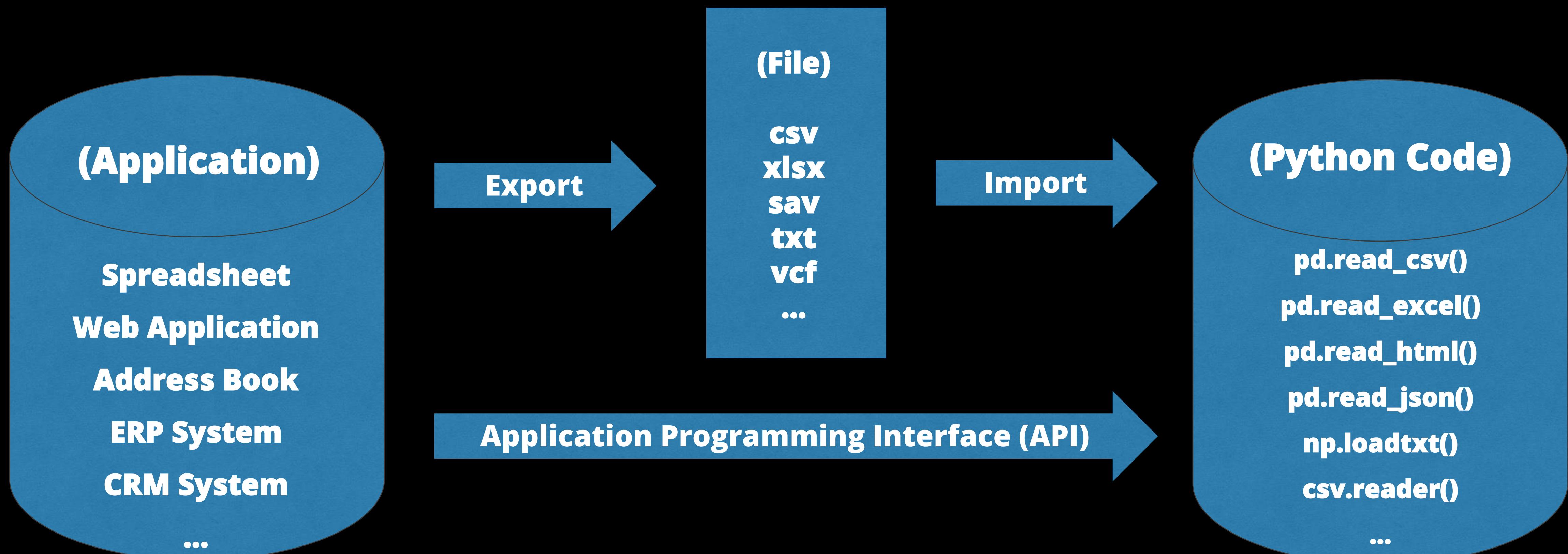
For example, use @codebase only in exceptional cases.

# TIPPS ZUR SUCHE MIT GOOGLE

- **English**
- **Mention the name of the programming language or package you are using ("python", "matplotlib", "seaborn", ...)**
- **Complete questions with a sensible word order yield better results**

# DATA IMPORT

# DATA IMPORT



# ENCODING

**Method for representing texts/characters in computer systems: mapping bit sequences to characters**

## Relevant for:

- **Text/CSV files (.txt, .csv)**
- **XML files**
- **HTML files**
- **JSON files (when stored as text)**

## Typical Issues:

- **Umlauts (ä, ö, ü) and special characters (€, §, °)**
- **Mojibake – incorrectly rendered characters (e.g., "Ã¼" instead of "ü")**
- **Different standards across operating systems**

## Common encoding formats:

- **UTF-8 (Unicode) – Standard for web applications, supports all languages**
- **ASCII – Basic character set (128 characters, English only)**
- **Latin-1 (ISO-8859-1) – Extended ASCII for Western European languages**
- **Windows-1252 – Microsoft's version of Latin-1**

## Best Practices:

- **Use UTF-8 as the default**
- **Explicitly specify encoding when importing files (`encoding='utf-8'`)**
- **Check the source file's encoding before importing (pay attention to BOM - Byte Order Mark)**

# PANDAS DATAFRAMES

- **Simple Data Structure**

A data structure called DataFrame that represents data in tabular form, similar to a database table or an Excel spreadsheet.

- **Powerful Functions**

Offers powerful functions for data manipulation, filtering, aggregation, and visualization.

- **Supports different Data Formats**

Allows importing data from CSV, Excel, JSON, SQL databases, and more.

- **Integration mit anderen Bibliotheken**

Integrates with libraries like NumPy, Matplotlib, Seaborn, and Scikit-learn, enabling advanced analysis and visualization on imported data.

- **Not Ideal for Very Large Datasets**

Suitable only for datasets smaller than 1 GB in file size. For larger datasets, alternatives such as the Datasets library from Hugging Face are recommended.

# PROCEDURE FOR IMPORTING WITH AI ASSISTANCE

- **Instruction that includes the full filename as well as the directory path or the link where the file can be found.**
- **For text files (such as CSV), include a snippet from the beginning of the file in the instruction to indicate the format.**

# **EXAMPLE: IMPORT FROM GITHUB**

# BREAKOUT

- Download the files "kiwo.csv", "umsatzdaten\_gekuerzt.csv", and "wetter.csv", and save them in your workspace.
- The files can be found at:  
<https://github.com/opencampus-sh/einfuehrung-in-data-science-und-ml>
- Import the file "wetter.csv" as a Pandas DataFrame.

# **GRAPHICAL REPRESENTATIONS**

# DIAGRAM TYPES

← Python Graph Gallery    CHART TYPES    TOOLS ▾    ALL    BEST    RELATED ▾    ABOUT    [Subscribe](#)

## Distribution



Violin    Density    Histogram    Boxplot    Ridgeline    Beeswarm

## Correlation



Scatterplot    Heatmap    Correlogram    Bubble    Connected Scatter    2D Density

## Ranking



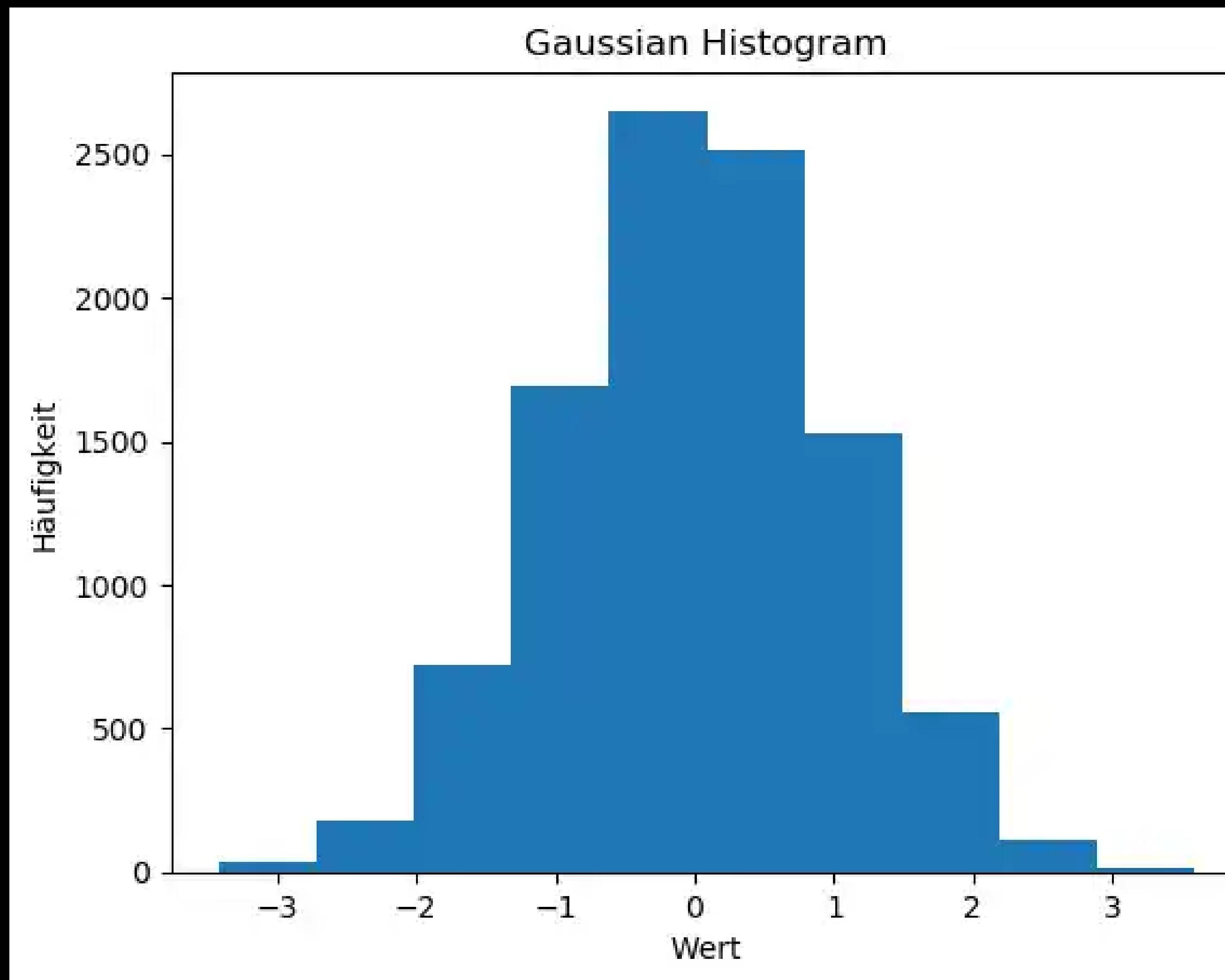
<https://python-graph-gallery.com/>

# TYPES OF SCALES

- **Nominal scale (categorical)**  
→ [Traffic light colors, federal state]
- **Ordinal scale**  
→ [English grade, test response on a scale from good - medium - poor]
- **Interval scale**  
→ [Temperature in Celsius, IQ score]
- **Ratio scale**  
→ [Speed, income]

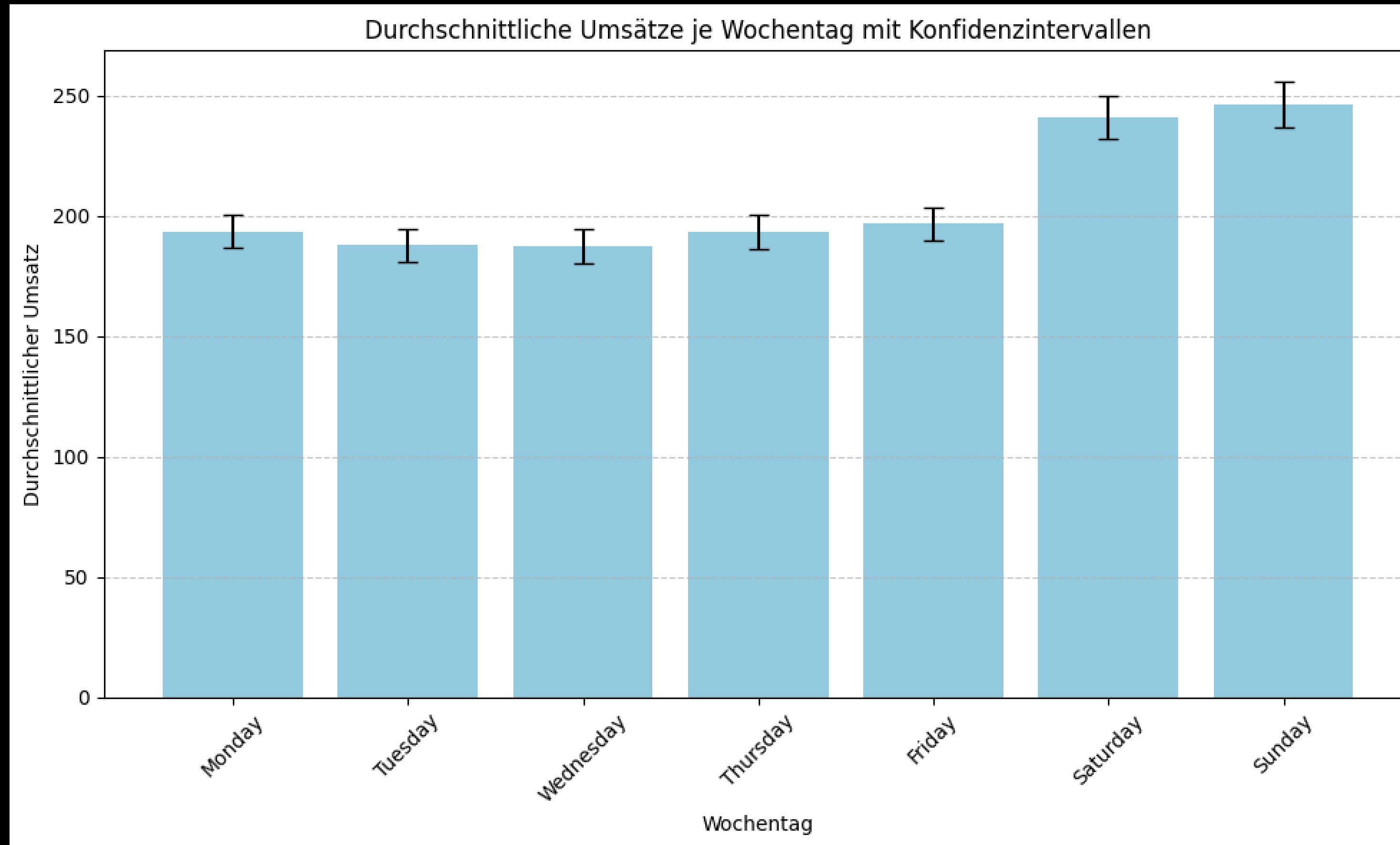
# COMMON CHART TYPES

# HISTOGRAM



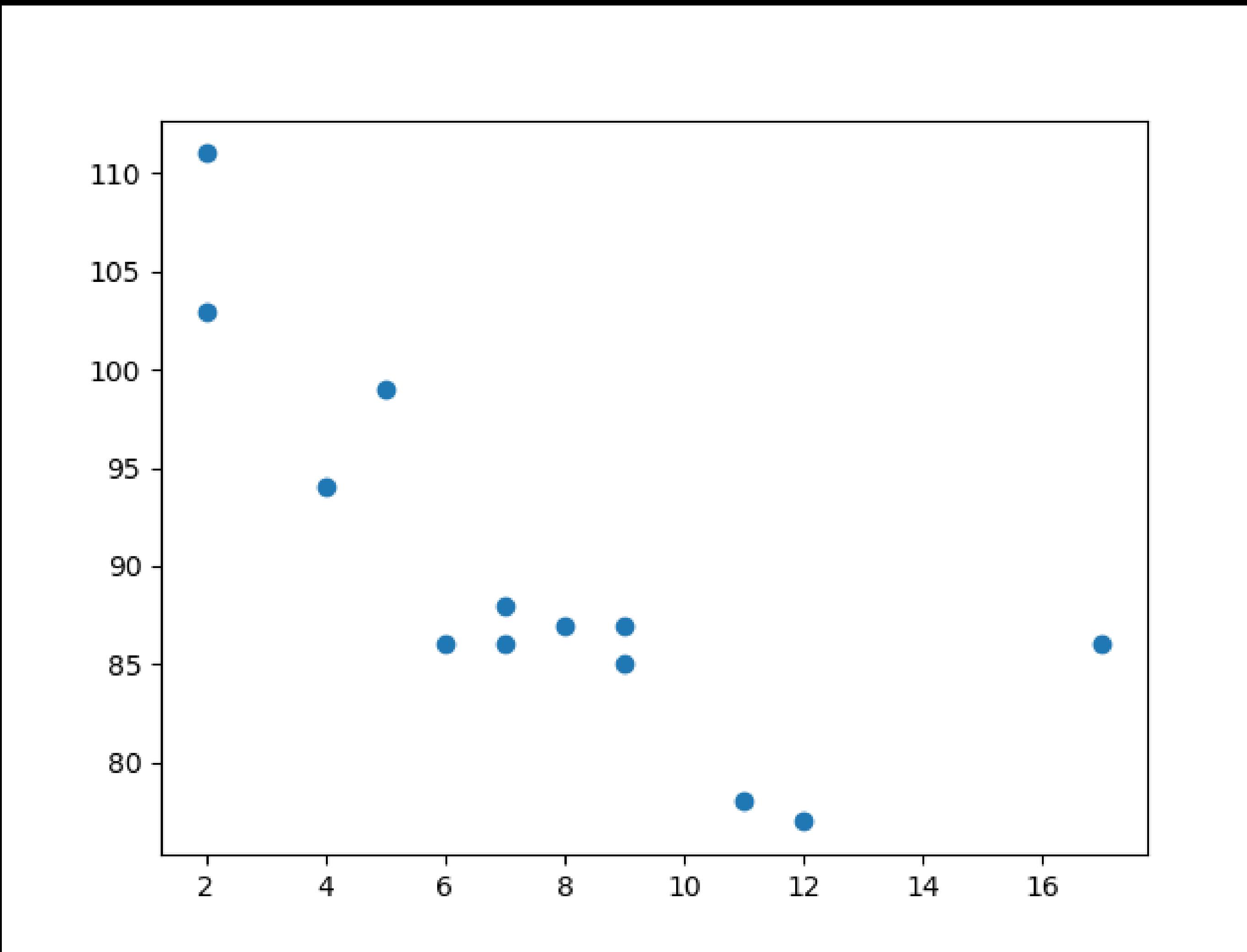
**Representation of the distribution of an at least ordinal-scaled variable**

# BAR CHART



**Representation of all variable types possible; commonly used to illustrate the relationship with an at least interval-scaled variable.**

# SCATTERPLOTT



**Representation of the relationship between two at least ordinal-scaled variables; more meaningful for interval-scaled variables."**

# PROJECT DATASET

- **Sales data of various product groups from a bakery branch for the period from 01.07.2013 to 30.07.2018**
- **Weather data for the period from 01.07.2013 to 30.07.2019**
- **Available at:**  
[https://raw.githubusercontent.com/opencampus-sh/einfuehrung-in-data-science-und-ml/main/umsatzdaten\\_gekuerzt.csv](https://raw.githubusercontent.com/opencampus-sh/einfuehrung-in-data-science-und-ml/main/umsatzdaten_gekuerzt.csv)

 Search or jump to... / Pull requests Issues Marketplace Explore

 + 

 [opencampus-sh/einfuehrung-in-data-science-und-ml](#)

Public

Edit Pins Watch 8 Fork 9 Star 12

<> [Code](#) [Issues](#) [Pull requests](#) [Actions](#) [Projects](#) [Wiki](#) [Security](#) [Insights](#) [Settings](#)

 main ▾ [einfuehrung-in-data-science-und-ml / umsatzdaten\\_gekuerzt.csv](#) Go to file ...

 steffen74 Datensatz Sommersemester 2022 Latest commit fd081db 9 minutes ago 

 1 contributor

10910 lines (10910 sloc) | 319 KB Raw Blame   

 Search this file...

1	Datum	Warengruppe	Umsatz
2	2013-07-01	1	148.828353112183
3	2013-07-02	1	159.79375714468
4	2013-07-03	1	111.885593514353
5	2013-07-04	1	168.864940979931
6	2013-07-05	1	171.280754117955
7	2013-07-06	1	174.552359998476
8	2013-07-07	1	92.6377553788373

# PRODUCT GROUPS

- 1    **Brot**
- 2    **Brötchen**
- 3    **Croissant**
- 4    **Konditorei**
- 5    **Kuchen**
- 6    **Saisonbrot**

# WEATHER DATA

- **For the period from 01.07.2013 to 30.07.2019**
- **Available at:**  
<https://raw.githubusercontent.com/opencampus-sh/einfuehrung-in-data-science-und-ml/main/wetter.csv>
- **Variables:**
  - **Average cloud cover during the day (0: min to 8: max)**
  - **Average temperature in Celsius**
  - **Average wind speed in m/s**
  - **Weather code (a list of descriptions can be found, for example, here: [http://www.seewetter-kiel.de/seewetter/daten\\_symbole.htm](http://www.seewetter-kiel.de/seewetter/daten_symbole.htm))**

# BREAKOUT

Create one of the following types of charts and use the dataset '`wetter.csv`' for this::

- **Scatterplot**
- **Histogram**
- **Bar Chart**

 Search or jump to... / Pull requests Issues Marketplace Explore

 + 

 [opencampus-sh/einfuehrung-in-data-science-und-ml](#)

Public

Edit Pins Watch 8 Fork 9 Star 12

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

 main ▾ [einfuehrung-in-data-science-und-ml / wetter.csv](#) Go to file ...

 steffen74 Project Data Latest commit c61a127 on 20 Apr 2021 History

1 contributor

2602 lines (2602 sloc) 64.2 KB Raw Blame   

Search this file...

1	Datum	Bewoelkung	Temperatur	Windgeschwindigkeit	Wettercode
2	2012-01-01	8	9.825	14	58
3	2012-01-02	7	7.4375	12	
4	2012-01-03	8	5.5375	18	63
5	2012-01-04	4	5.6875	19	80
6	2012-01-05	6	5.3	23	80
7	2012-01-06	3	2.625	10	
8	2012-01-07	7	6.528571	14	61

# LEARNING RESOURCES

- Work through [\*\*this\*\*](#) introduction to working with Pandas (only Lesson 1).
- Watch [\*\*this\*\*](#) video on importing data as a Pandas DataFrame (18 minutes).
- Work through [\*\*this\*\*](#) introduction to creating visualizations with Matplotlib (only Lesson 1).
- Watch [\*\*this\*\*](#) video (4 minutes) to understand the relevance of confidence intervals.

# TASKS

- Create a GitHub Codespace and save the files 'kiwo.csv', 'umsatzdaten\_gekuerzt.csv', and 'wetter.csv' from this GitHub repository:  
<https://github.com/opencampus-sh/einfuehrung-in-data-science-und-ml>
- Create a bar chart showing the average sales per weekday.
- In a second step, add confidence intervals for the sales per weekday.
- In a further step, sort the weekdays from Monday to Sunday.

### Durchschnittliche Umsätze je Wochentag mit Konfidenzintervallen

