

Python für Linguisten

Dozentin: Wiebke Petersen & Co-Dozentin: Esther Seyffarth

Weiterführende Packages in Python

Pakete in Python

- Bisher haben wir im Unterricht nur Standardfunktionen und mitgelieferte Pakete von Python verwendet.
- Standardfunktionen sind ohne weiteres verfügbar, wie z.B. `sort()` oder `print()`.
- Mitgelieferte Pakete sind Sammlungen von Modulen, die bestimmte Aufgabenbereiche abdecken. Wir haben schon das Paket für reguläre Ausdrücke kennengelernt:

```
1 >>> import re
2 >>> pattern = re.compile("...")
```

- Man kann auch selbstgeschriebene Module importieren, wenn man ein Projekt bearbeitet, das sich in verschiedene Teilbereiche aufteilen lässt (wie z.B. in einigen Projektgruppen im Seminar).
- Wenn man versucht, ein Modul zu importieren, wird zunächst das Verzeichnis mit den Standardpaketen durchsucht (`<Pythonpath>/Lib/` und `<Pythonpath>/Lib/site-packages/`). Erst danach wird das aktuelle Arbeitsverzeichnis durchsucht!

Pakete installieren

- Selbstgeschriebene Module, die wiederkehrende Aufgaben erfüllen, kann man auch als Paket im Pythonpfad ablegen, dann kann man sie jederzeit in anderen Projekten importieren.
- Viele kluge Köpfe entwickeln permanent geniale Pythonpakete, die man sich einzeln je nach Bedarf installieren kann.
- Zum Installieren verwendet man in der Kommandozeile/Terminal folgenden Befehl: `pip install [packagename]`
- Einige Packages kann man sich auch von den GitHub-Seiten der Entwickler/innen herunterladen.
- Nach dem Installieren kann man sofort die Pakete importieren und verwenden. Zur Benutzung am besten immer die Dokumentation des jeweiligen Pakets lesen!

sqlite3 (mitgeliefert)

- Gutes Modul zum Erstellen und Bearbeiten von lokalen Datenbanken im Sqlite-Format.

```
1 >>> # Modul importieren
2 >>> import sqlite3
3 >>> # Verbindung zur Datenbankdatei herstellen
4 >>> conn = sqlite3.connect("german.sqlite")
5 >>> # Cursor zum Lesen der DB erzeugen
6 >>> c = conn.cursor()
7 >>> # mithilfe des Cursors eine SQL-Abfrage ausführen
8 >>> row = c.execute('''SELECT Wort, POS FROM wort WHERE
9                       Wort LIKE "%python%";''')
10 >>> # fetchall() holt alle Ergebnisse der Anfrage, fetchone()
11     # holt ein einziges Ergebnis
12 >>> print(row.fetchall())
13 [('Baumpython', 'NN'), ('Felsenpython', 'NN'), ('Fleckenpython',
14 'NN'), ('Netzpython', 'NN'), ('Python', 'NN'), ('Pythons', 'NN'),
15 ('Ringpython', 'NN'), ('Schwarzkopfpython', 'NN'),
16 ('Teppichpython', 'NN'), ('Wasserpython', 'NN')]
```

csv (mitgeliefert)

- Modul zum Erstellen und Bearbeiten von CSV-Dateien.

```
1 >>> # Modul importieren
2 >>> import csv
3 >>> with open('corpus.csv', 'w', newline='') as csvfile:
4         csvwriter = csv.writer(csvfile, delimiter=';', \
5                                 quotechar='|', quoting=csv.QUOTE_MINIMAL)
6         csvwriter.writerow(["1. Spalte", "Zweite", "Dritte"])
```

urllib (mitgeliefert)

- Gutes Modul zum Öffnen von Webseiten. Auch in der Lage, Cookies zu speichern und Formulare auszufüllen.

```
1 >>> # Modul importieren
2 >>> import urllib.request
3 >>> url = "https://en.wikipedia.org/wiki/Special:Random"
4
5 >>> # Webseite anfragen und in r speichern
6 >>> r = urllib.request.urlopen(url)
7
8 >>> # Achtung: Der Quellcode wird nicht als normaler String
9 # geliefert, sondern als Bytestring. Er kann aber einfach
10 # umgewandelt werden.
11 >>> r = r.read()
12 >>> r = r.decode("utf-8")
13
14 >>> # Der Webseitentext kann nun wie ein String verarbeitet
15 # werden:
16 >>> print(r.split("<head>", 1)[1].split("<body>", 1)[0])
```

Weitere Pakete

- simplejson: Zum Schreiben und Lesen von JSON-Dateien oder zum Verarbeiten von API-Rückgaben.
- xml: Zum Schreiben und Lesen von XML-Dateien.
- collections: Ermöglicht zusätzliche Funktionen bei Dictionaries, Tupeln, Listen und Mengen.
- time: Zum Messen der Dauer eines Programms, aber auch zum Erzeugen von Timestamps usw., z.B. beim Schreiben von Logdateien
- tweepy: Zum Abfragen und Posten von Tweets, extrem gut geeignet zum Erstellen einer Datenbasis (die man z.B. als Sqlite-Datenbank oder als CSV-Datei speichert)
- pep8: Ein Tool, das prüft, ob der Code eines selbstgeschriebenen Moduls den Vorgaben nach PEP8 entspricht.
- traceback: Ermöglicht eine weiterführende Behandlung von Fehlermeldungen, z.B. nützlich beim Schreiben von Logdateien.

Übung: Erstellen eines Wikipedia-Korpus

- Laden Sie sich die Datei `wikipedia_corpus.py` von der Kurswebseite herunter.
- Machen Sie sich mit der Funktionalität von `get_website_content(url)` vertraut. Öffnen Sie mehrere Webseiten mithilfe dieser Funktion und lassen Sie sie anzeigen.
- Welche Funktion hat die Zeile `page.decode("utf-8")`?
- Welche Funktion hat die Zeile `html.unescape(...)`?
- Wir wollen ein Korpus aus zufälligen Wikipedia-Artikeln erstellen. Die Funktion `make_corpus()` ist der Einstiegspunkt in das Programm.
- Ergänzen Sie den notwendigen Code, damit jeweils der erste Absatz jedes Wikipedia-Artikels in die CSV-Datei geschrieben wird.