

Referenzkarte PCL I – Version 5.1 – 28. Dezember 2014

Sequenzielle Datentypen: Listen (lesen/schreiben), Tupel, Strings (nur lesen) <code>l = range(n)</code> Liste von n Zahlen (von 0 bis n-1) <code>l = l1 + l2</code> Verknüpfung von Sequenzen <code>l[i:j]</code> Teilsequenz von l[i] bis l[j-1] <code>l[:k]</code> die ersten k Elemente der Sequenz <code>l[-k:]</code> die letzten k Elemente der Sequenz <code>l[i:j] = ['bla']</code> ersetzen von l[i] bis l[j-1] mit der Liste ['bla'] (nur für Listen) <code>len(l)</code> Anzahl Elemente <code>max(l) ; min(l)</code> max./min. Wert (strings: alph. geordnet) <code>l.count('the')</code> Anzahl auftreten von 'the' in der Liste l <code>l.index('dog')</code> erster Index von 'dog', oder Fehler falls 'dog' nicht in l <code>x in l ; x not in l</code> ist x (nicht) in l enthalten? (Resultat True oder False) <code>l.append('bla')</code> anfügen von 'bla' am Ende von l (nur für Listen) <code>l.insert(i,x)</code> füge x an der Position i in l ein (nur für Listen) <code>l.remove(x)</code> entferne das erste x in l (nur für Listen) <code>l.reverse()</code> Reihenfolge von l umkehren (nur für Listen) <code>l.sort()</code> sortieren, zuerst Zahlen, dann Buchstaben (nur für Listen)	Input/Output <code>import sys</code> ## Systemmodul, u.a. mit Kommandozeil.-Funktionen def main(): <code> filename = sys.argv[1]</code> ## Dateiname auslesen aus Kommandozeile <code> infile = open(filename, 'r')</code> ## die Datei zum Schreiben öffnen <code> my_outfilename = filename + '.letter_a_words'</code> ## Name für Ausgabedatei <code> outfile = open(my_outfilename, 'w')</code> ## Datei zum Schreiben öffnen for line in infile: ## Schleife für jede Zeile in der Datei <code> # <Do something></code> <code> infile.close()</code> ## Schliessen der Dateien <code> outfile.close()</code> if __name__ == '__main__': ## Standard main-Boilerplate mit Funktion-Aufruf <code> main()</code>
Dictionaries / Hashes (d) mit Schlüssel-Wert-Paaren <code>d = {}</code> leeres Dictionary erzeugen <code>d = {'a':34, 'to':23}</code> Dictionary mit zwei Einträgen erzeugen <code>d['the']</code> gibt den Wert von 'the' zurück (in diesem Fall 34) <code>len(d)</code> Anzahl Schlüssel <code>d.copy()</code> erzeugt eine Kopie von d <code>d.items()</code> Liste aller Einträge (Einträge sind 2-Tupel) <code>d.keys()</code> Liste aller Schlüssel <code>d.values()</code> Liste aller Werte <code>x in d</code> True falls der Schlüssel x in d enthalten ist <code>del d['the']</code> entfernt Schlüssel 'the' und zugehörigen Wert aus d	HTML <code>import nltk</code> from urllib import urlopen <code>html = urlopen(given_url).read()</code> ## öffne den Quellcode des Dokuments mit der Adresse given_url zum Lesen <code>text = nltk.clean_html(html)</code> ## nimmt den Quelltext und extrahiert den Text daraus Anderes <code>from __future__ import division</code> nötig, falls man zwei Integer teilt und eine rationale Zahl erhalten will <code>set(list)</code> erzeugt ein Set einer Liste, in dem Duplikate zusammengefasst werden <code>l = nltk.word_tokenize(given_string)</code> tokenisiert given_string in eine Liste <code>foobar".find("bar")</code> sucht das erste Vorkommen von "bar" in "foobar" und liefert den ersten Index der Sequenz zurück
Reguläre Ausdrücke – Syntax <code>^</code> Startposition innerhalb des Strings <code>\$</code> Endposition eines Strings <code>x{m,n}</code> x kommt mindestens m mal vor, aber nicht mehr als n mal gruppiert den Ausdruck x. Kann über i aufgerufen werden (i von 1 bis 9) <code>\d</code> [0-9] <code>\D</code> [^0-9] <code>\w</code> [A-Za-z0-9_]	Reguläre Ausdrücke – Modul re <code>import re</code> <code>l = re.findall(pattern, string)</code> Liste mit allen gefunden Gruppen <code>m = re.search(pattern, string)</code> <code>m.group()</code> gibt alle Übereinstimmungen zurück <code>m.group(1)</code> gibt nur die erste Übereinstimmung zurück <code>s = re.sub(pattern, replacement, given_string)</code> sucht pattern in given_string und ersetzt es durch replacement <code>re.split(pattern, given_string)</code> trennt given_string nach pattern in eine Liste auf <i>Hinweis: Flag (?u) aktiviert Unicode Kategorie für \w and \b → pattern=ur'(?u)(*regex*)'</i> Liste vs. Listenkomprehension <code>word_list = ["foo", "bar", "lorem", "ipsum"]</code> ## Liste mit Stringliterals <code>four_chars_words = [w for w in word_list if len(w)==4]</code> ## Listenkomprehension
Konversion sequenzielle Datentypen: strings (s), lists (l), tuples (t) <code>l = list(t) ; l = list(s) ; t = tuple(l) ; t = tuple(s)</code>	Tastatureingaben <code>var = raw_input()</code> Speichert die Eingabe in einer String-Variable <code>integer = int(s)</code> Konvertiert einen String in einen Integer

Grundlegende Syntax

```

if expr: ## if-Abfrage
    statements
elif expr:
    statements
else:
    statements
while expr: ## while Schleife
    statements

for w in l: ## for Schleife
    statements
## Funktion mit Argument x
## und Rückgabewert x+1
def function_name(x):
    return x+1

```

Import von NLTK: `import nltk`

NLTK – Korpora

`from nltk.corpus import corpus_name`

`corpus_name.fileids()`

`corpus_name.raw(fileid)`

`corpus_name.words(fileid)`

`corpus_name.sents(fileid)`

`corpus_name.paras(fileid)`

importieren des Korpus `corpus_name` aus dem Modul `nltk.corpus`

erzeugt eine Liste mit allen Datei-ID's des Korpus

erzeugt einen unicode-String aus dem Inhalt der Datei mit der ID `fileid`

erzeugt eine „Liste“, bzw. ein listenähnliches Objekt aus dem Text der Datei mit der ID `fileid`.

Kann mit `list()` in eine normale Liste umgewandelt werden. (`List[w]`)

erzeugt eine zweidimensionale „Liste“ aus dem Text, in der jeder Satz eine Element ist,

mit seinen Wörtern wiederum als Liste. (`List[s][w]`)

erzeugt eine dreidimensionale „Liste“ aus dem Text, in der die satzbildenden Listen zusätzlich nach Absätzen zusammengefasst sind (`List[p][s][w]`)

Kategorisierte oder kontext-getaggte Korpora haben zusätzlich die Methode `.categories()`:

`corpus_name.categories()`

erzeugt eine Liste mit allen Kategorien des Korpus.

`corpus_name.categories(fileid)`

erzeugt eine Liste mit den Kategorien, die der Korpus-Datei mit der ID `fileid` zugeordnet sind.

Der Brown-Korpus ist mit Part-of-Speech-Tags (POS) getaggt:

`brown.tagged_words(fileid)`

erzeugt eine Liste von Wort-POS-Tupel (`word, pos_tag`) für jedes Wort

NLTK – Frequency distribution classes

`fdist = nltk.FreqDist(samples)`

erzeugt ein dictionary-ähnliches Objekt: jedes Element (=event) der `samples`-Sequenz

ist ein Schlüssel und die Frequenz des Elements in `sample` ist der Wert

`fdist.N()`

Anzahl Elemente in `samples`

`fdist.max()`

Event mit grösstem Zähler

`fdist.keys()`

Events sortiert nach ihrer Frequenz

`fdist.tabulate()`

tabellarische Anordnung der Frequenz Anordnung

`cfdist = nltk.ConditionalFreqDist(pairs)`

creates frequency distributions of events conditioned on a condition: `pairs` is a list of tuples of the form (`condition, event`)

`cfdist.conditions()`

Alphabetisch sortierte Liste der Bedingungen (`conditions`)

`cfdist[condition]`

Dictionary mit Häufigkeiten von allen events gegeben eine Bedingung (`condition`)

`cfdist[condition][sample]`

Häufigkeit eines bestimmten events gegeben eine Bedingung (`condition`)

`cfdist.tabulate()`

tabellarische Ausgabe von Bedingungen und Events

NLTK – Bigramme erzeugen

`bi_list = nltk.bigrams(list)`

erzeugt ein Generatorobjekt mit Bigramm-Tupeln aus den Elementen der Liste `list`. Umwandlung in eine Liste mit `list()`.

Beispiel:

```
>>> bigramGen = bigrams(["Lorem", "Ipsum", "Dolor", "Sit", "Amet"])
```

```
>>> print list(bigramGen)
```

```
[('Lorem', 'Ipsum'), ('Ipsum', 'Dolor'), ('Dolor', 'Sit'), ('Sit', 'Amet')]
```