

<p>Permissions</p> <p>Unix/Linux permissions are shown with "ls -l" or "ls -al" (Option a für alle, auch versteckte Dateien):</p> <pre>-rwxr-x--- 1 tom acct 7538 Dec 9 08:48 execfile</pre> <p>The very first part of the permissions tells you whether what you are seeing is an ordinary file, a directory, or something else. Another way to find out what you are looking at is to use the "file" command. For example, "file *"...</p> <p>This file is owned by tom, and has the group ownership of "acct". The owner (tom) has full permissions- he can read this file (display it, read it into an editor, copy it somewhere else), write it (replace its contents with something else), and execute it.</p> <p>Anyone in the group "acct" can read and execute it, but can't write it. Anyone else can't do anything.</p> <p>The x or "execute" permission needs some explanation for those from the Windows world: Unix programs don't use extensions like "bat" or "exe" to be executable- instead, the "x" permissions settings determines this. Of course, the file still has to be a legitimate program- but a legitimate program will not run until it is marked as "x".</p> <p>Erzeuge leere Datei oder überschreibe bestehende:\$ >result.txtA</p>	<p>Directories</p> <p>Permissions on directories can be a little confusing. Read permission ("r") means that you can list the directory- for example, use "ls" on it. But that's all it means: if you only have read permission, you can't cd to that directory, you can't copy files into that directory, and you can't even do an "ls -l" on it. Note that it doesn't help if files or directories under that would be accessible by you: you won't be able to get to them. Write permission ("w") allows you to create new files and to remove files.</p> <p>Normally, write permission is all you need to remove a file- you might not have any permissions on the file itself, but if you have write permission on the directory and the text bit is not set (see below) you will be able to remove the file whether or not you have any permissions on the file itself, and no matter who owns it. However, you will need at least "x" permission also, because "x" is what lets you "search" a directory. You might also think of it as "use its contents". You need "x" permissions on the directory to read a file in that directory, to copy it elsewhere and (as noted above) to delete it.</p> <p>The "id" command (most Unix/Linux) will show you what groups you belong to:</p> <pre>uid=1000(benzro) gid=1000(benzro) groups=1000(benzro),4(adm),24(cdrom),27(sudo),30(dip),46(plugdev),1 15(lpadmin),125(sambashare)</pre> <p>The command "ps auxww" provides complete information about the processes, including all parameters.</p> <p>The <u>locate command</u> is often the simplest and quickest way to find the locations of files and directories on Linux. The command <u>locate file1 dir1</u> would list the absolute paths of all files named <i>file1</i> (or containing <i>file1</i>) and all directories named <i>dir1</i> (or containing <i>dir1</i>) for which the user had access permission. The specificity of locate can be increased by using it together with wildcards or other regular expressions. The command <u>locate *.png less</u> uses the star wildcard to display all files on the system that have the <i>.png file extension</i>.</p>
<p>benzro@benzro-eMachines-E625:~\$ grep --help</p> <p>Usage: grep [OPTION]... PATTERN [FILE]...</p> <p>Search for PATTERN in each FILE or standard input. PATTERN is, by default, a basic regular expression (BRE).</p> <p>Regexp selection and interpretation:</p> <ul style="list-style-type: none"> -E, --extended-regexp PATTERN is an extended regular expression (ERE) -F, --fixed-strings PATTERN is a set of newline-separated fixed strings -G, --basic-regexp PATTERN is a basic regular expression (BRE) -P, --perl-regexp PATTERN is a Perl regular expression -e, --regexp=PATTERN use PATTERN for matching -f, --file=FILE obtain PATTERN from FILE <p>-i, --ignore-case ignore case distinctions</p> <p>-w, --word-regexp force PATTERN to match only whole words</p> <p>-x, --line-regexp force PATTERN to match only whole lines</p> <p>-z, --null-data a data line ends in 0 byte, not newline</p> <p>Miscellaneous:</p> <ul style="list-style-type: none"> -s, --no-messages suppress error messages -v, --invert-match select non-matching lines -V, --version display version information and exit --help display this help text and exit <p>Context control:</p> <ul style="list-style-type: none"> -B, --before-context=NUM print NUM lines of leading context -A, --after-context=NUM print NUM lines of trailing context -C, --context=NUM print NUM lines of output context -NUM same as --context=NUM --color[=WHEN], --colour[=WHEN] use markers to highlight the matching strings; WHEN is 'always', 'never', or 'auto' -U, --binary do not strip CR characters at EOL (MSDOS/Windows) -u, --unix-byte-offsets report offsets as if CRs were not there (MSDOS/Windows) 	<p>Output control:</p> <ul style="list-style-type: none"> -m, --max-count=NUM stop after NUM matches -b, --byte-offset print the byte offset with output lines -n, --line-number print line number with output lines --line-buffered flush output on every line -H, --with-filename print the file name for each match -h, --no-filename suppress the file name prefix on output --label=LABEL use LABEL as the standard input file name prefix -o, --only-matching show only the part of a line matching PATTERN -q, --quiet, --silent suppress all normal output --binary-files=TYPE assume that binary files are TYPE; TYPE is 'binary', 'text', or 'without-match' -a, --text equivalent to --binary-files=text -I equivalent to --binary-files=without-match -d, --directories=ACTION how to handle directories; ACTION is 'read', 'recurse', or 'skip' -D, --devices=ACTION how to handle devices, FIFOs and sockets; ACTION is 'read' or 'skip' -r, --recursive like --directories=recurse -R, --dereference-recursive likewise, but follow all symlinks --include=FILE_PATTERN search only files that match FILE_PATTERN --exclude=FILE_PATTERN skip files and directories matching FILE_PATTERN --exclude-from=FILE skip files matching any file pattern from FILE --exclude-dir=PATTERN directories that match PATTERN will be skipped. -L, --files-without-match print only names of FILES containing no match -l, --files-with-matches print only names of FILES containing matches -c, --count print only a count of matching lines per FILE -T, --initial-tab make tabs line up (if needed) -Z, --null print 0 byte after FILE name <p>grep-Terminalanzeigen:</p> <pre>===== -o (=nur Treffer ausgeben) --colour (=Treffer in Zeile einfärben) -n (=Zeilennummer ausgeben) -v (=inverse match) -h (=keine Dateilangaben)</pre> <p>grep anhalten / wieder / beenden:</p> <pre>===== --> Ctrl+s / Ctrl+q / Ctrl+z oder Ctrl+c</pre>
<p>REGULAR EXPRESSIONS</p> <p>A regular expression is a pattern that describes a set of strings. Regular expressions are constructed analogously to arithmetic expressions, by using</p>	<p>Repetition</p> <p>A regular expression may be followed by one of several repetition operators:</p> <ul style="list-style-type: none"> ? The preceding item is optional and matched at most once.

<p>various operators to combine smaller expressions.</p> <p>grep understands three different versions of regular expression syntax: "basic" (BRE), "extended" (ERE) and "perl" (PRCE). In GNU grep, there is no difference in available functionality between basic and extended syntaxes. In other implementations, basic regular expressions are less powerful. The following description applies to extended regular expressions; differences for basic regular expressions are summarized afterwards. Perl regular expressions give additional functionality, and are documented in pcresyntax(3) and pcrepattern(3), but only work if pcre is available in the system.</p> <p>The fundamental building blocks are the regular expressions that match a single character. Most characters, including all letters and digits, are regular expressions that match themselves. Any meta-character with special meaning may be quoted by preceding it with a backslash.</p>	<ul style="list-style-type: none"> * The preceding item will be matched zero or more times. + The preceding item will be matched one or more times. {n} The preceding item is matched exactly n times. {n,} The preceding item is matched n or more times. {m} The preceding item is matched at most m times. This is a GNU extension. {n,m} The preceding item is matched at least n times, but not more than m times. <p>Concatenation Two regular expressions may be concatenated; the resulting regular expression matches any string formed by concatenating two substrings that respectively match the concatenated expressions.</p>
<p>The period . matches any single character.</p> <p>Character Classes and Bracket Expressions A bracket expression is a list of characters enclosed by [and]. It matches any single character in that list; if the first character of the list is the caret ^ then it matches any character not in the list. For example, the regular expression [0123456789] matches any single digit.</p>	<p>Alternation Two regular expressions may be joined by the infix operator ; the resulting regular expression matches any string matching either alternate expression.</p>
<p>Within a bracket expression, a range expression consists of two characters separated by a hyphen. It matches any single character that sorts between the two characters, inclusive, using the locale's collating sequence and character set. For example, in the default C locale, [a-d] is equivalent to [abcd]. Many locales sort characters in dictionary order, and in these locales [a-d] is typically not equivalent to [abcd]; it might be equivalent to [aBbCcDd], for example. To obtain the traditional interpretation of bracket expressions, you can use the C locale by setting the LC_ALL environment variable to the value C.</p>	<p>Precedence Repetition takes precedence over concatenation, which in turn takes precedence over alternation. A whole expression may be enclosed in parentheses to override these precedence rules and form a subexpression.</p>
<p>Finally, certain named classes of characters are predefined within bracket expressions, as follows. Their names are self explanatory, and they are [:alnum:], [:alpha:], [:cntrl:], [:digit:], [:graph:], [:lower:], [:print:], [:punct:], [:space:], [:upper:], and [:xdigit:]. For example, [:alnum:] means the character class of numbers and letters in the current locale. In the C locale and ASCII character set encoding, this is the same as [0-9A-Za-z].</p> <p>(Note that the brackets in these class names are part of the symbolic names, and must be included in addition to the brackets delimiting the bracket expression.) Most meta-characters lose their special meaning inside bracket expressions. To include a literal] place it first in the list. Similarly, to include a literal ^ place it anywhere but first. Finally, to include a literal - place it last.</p>	<p>Back References and Subexpressions The back-reference \\$n, where n is a single digit, matches the substring previously matched by the nth parenthesized subexpression of the regular expression.</p>
<p>Anchoring The caret ^ and the dollar sign \$ are meta-characters that respectively match the empty string at the beginning and end of a line.</p>	<p>Basic vs Extended Regular Expressions In basic regular expressions the meta-characters ?, +, {, , (, and) lose their special meaning; instead use the backslashed versions \?, \+, \{, \ , \(, and \).</p>
<p>The Backslash Character and Special Expressions The symbols \< and \> respectively match the empty string at the beginning and end of a word. The symbol \b matches the empty string at the edge of a word, and \B matches the empty string provided it's not at the edge of a word. The symbol \w is a synonym for [[:alnum:]] and \W is a synonym for [^[:alnum:]].</p>	<p>Traditional egrep did not support the { meta-character, and some egrep implementations support \{ instead, so portable scripts should avoid { in grep -E patterns and should use {{}} to match a literal {.</p>
<p>Durchsuchen von Datei-namen und Ordner-namen im gesamten Filesystem:</p>	<p>Durchsuchen der Datei-inhalte im lokalen Verzeichnis:</p>
<pre>benzro@benzro-eMachines-E625:~\$ locate -i calendar wc 562 32221 benzro@benzro-eMachines-E625:~\$ locate -i calendar grep -n doc 35:/usr/local/texlive/2015/txmf-dist/doc/generic/pgf/text-en/pgfmanual-en-library-calendar.tex 36:/usr/local/texlive/2015/txmf-dist/doc/generic/pgf/text-en/pgfmanual-en-pgfcalendar.tex 37:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar 38:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/Changes 39:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/README 40:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.ltx 41:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.pdf 42:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.tex 43:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docFR.pdf 44:/usr/local/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docFR.tex 45:/usr/local/texlive/2015/txmf-dist/doc/metapost/bbcard/README.calendar 144:/usr/share/doc/libkronos/calendar4 145:/usr/share/doc/bdmaintools/calendarJudaic.py.gz 146:/usr/share/doc/libkronos/calendar4/changelog.Debian.gz 147:/usr/share/doc/libkronos/calendar4/copyright 148:/usr/share/doc/texlive-doc/generic/pgf/text-en/pgfmanual-en-library-calendar.tex.gz 149:/usr/share/doc/texlive-doc/generic/pgf/text-en/pgfmanual-en-pgfcalendar.tex.gz 150:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar 151:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/README 152:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.ltx.gz 153:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.pdf 154:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docDE.tex.gz 155:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docFR.pdf 156:/usr/share/doc/texlive/2015/txmf-dist/doc/latex/pst-calendar/pst-calendar-docFR.tex.gz benzro@benzro-eMachines-E625:~\$ locate -i calendar grep -n doc wc 24 1734 benzro@benzro-eMachines-E625:~\$ locate -i calendar grep -n -c doc 24</pre>	<p>Suche in allen Dateien (*.*) nach dem Inhalt „Document“ benzro@benzro-eMachines-E625:~\$ grep "Document" *.* callPCL.py:chdir('/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Übungen/PCL_uebung"+intUebung+"_python")</p>
	<p>Gib nur den Dateinamen des Treffers aus (Option -l): benzro@benzro-eMachines-E625:~\$ grep -l "Document" *.* callPCL.py</p>
	<p>Gib den Treffer in der Datei ohne den Dateinamen aus (Option -h): benzro@benzro-eMachines-E625:~\$ grep -h "Document" *.* chdir('/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Übungen/PCL_uebung"+intUebung+"_python/")</p>
	<p>Gib den gesamten Inhalt (Pattern leer matched alles) der versteckten Dateien aus (*): benzro@benzro-eMachines-E625:~\$ grep * less</p>
	<p>Ins richtige Verzeichnis wechseln auf Lubuntu Partition:</p>
<pre>pwd --> /home/benzro cd .. (2 mal) cd media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Übungen/PCL_uebung1 grep</pre>	<p>pwd --> /home/benzro cd .. (2 mal) cd media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Übungen/PCL_uebung1 grep</p>
<pre>benzro@benzro-eMachines-E625:~\$ wc --help Usage: wc [OPTION]... [FILE]... or: wc [OPTION]... --files0-from=F Print newline, word, and byte counts for each FILE, and a total line if more than one FILE is specified. With no FILE, or when FILE is -, read standard input. A word is a non-zero-length sequence of characters delimited by white space.</pre>	<p>benzro@benzro-eMachines-E625:~\$ wc --help Usage: wc [OPTION]... [FILE]... or: wc [OPTION]... --files0-from=F Print newline, word, and byte counts for each FILE, and a total line if more than one FILE is specified. With no FILE, or when FILE is -, read standard input. A word is a non-zero-length sequence of characters delimited by white space.</p>
<p>Dateien Analyse alle Jahrgänge 1930-1939: ----- ->Total Zeilen=Anzahl Artikel + Anzahl Abschnittsanfänge + Anzahl POS-Tags 2741140=923+142101+35+5+24+1+2598051 ->Total an Zeilen: grep -P '^.*' SAC-Jahrbuch_193* wc 2741140 8222497 137232508 ->Anzahl Artikel: grep -P -o '^<article>' SAC-Jahrbuch_193* wc</p>	<p>Bereinigter vertikalisierte Korpus erzeugen: ----- ->(Entferne 923 Artikelzeilen, 142101 Abschnittszeilen und 75 fehlerhafte POS-tags) Option -h (=gibt keine Dateiangaben aus) grep -P -o '^((\w \\$_ \\$ \\$(\)+(\w+\:\w+(\:\w+\?)) \w)' SAC-Jahrbuch_193* Dump2598051_SAC-Jahrbuch_193x wc Dump2598051_SAC-Jahrbuch_193x -> 2598051 7794153 4055402 Die Datei Dump2598051_SAC-Jahrbuch_193x wird für die Aufgaben verwendet</p>

<pre>923 923 39689 -->Anzahl Abschnittsanfänge : grep -P -o '^<s lang=' SAC-Jahrbuch_193* wc 142101 284202 6110343 -->Anzahl erlaubte und unerlaubte POS-Tags : `t((\wW+)+)`t (=W könnte auch ein \t sein, funktioniert nur, weil jede Zeile nur zwei \t hat und keinen dritten.) grep -P -o '^((\wW+)+)`t' SAC-Jahrbuch_193* wc grep -P -o '^((\t+))`t' SAC-Jahrbuch_193* wc 2598116 7794348 128890391 -->Anzahl POS-Tag mit unerlaubtem Zeichen) alleinstehend (Nicht in TAG-Listen): `t(x29)`t (=Hexadezimaler ASCII Code für ")" grep -P '\t(x29)`t' SAC-Jahrbuch_193* wc 35 105 1428 -->Anzahl Pos-Tag mit unerlaubten Zeichen ", &, ` alleinstehend (Nicht in TAG-Listen): `t((x22+) (x26+) (x27+) (x60+))`t (=Hexadezimaler ASCII Code für die Zeichen ", &, `) grep -P '\t((x22+) (x26+) (x27+) (x60+))`t' SAC-Jahrbuch_193* wc 5 15 205 -->Anzahl Pos-Tag mit unerlaubtem Zeichen ":" alleinstehend (Nicht in TAG-Listen): grep -P -h `t(`t(``))`t' SAC-Jahrbuch_193* wc -1 24 -->Anzahl Pos-Tag mit unerlaubtem Zeichen "IN/that"(Nicht in TAG-Listen): grep -P -h `t(IN/that)`t' SAC-Jahrbuch_193* wc -1 1 -->Anzahl erlaubte Pos-Tags (w+:w+(:w+?)) (=Italienische Tags mit einem oder zwei Doppelpunkten) grep -P -h `t((\w\$.\\$,\\$ (\w+:\w+(:w+?)))`t' SAC-Jahrbuch_193* wc -1 2598051</pre>	<p>In wie vielen Wörtern kommt der Buchstabe "r" (klein oder gross) vor? Wörter stehen in der ersten Spalte, --> Suche im Bereich: `^`t Erste Annäherung: grep -P -i --colour '^([\\wöö]*[r][\\wäöü]*)`t' Dump2598051_SAC-Jahrbuch_193x wc -l -->645098 Exakte Lösung: grep -P -colour '^((\\at)*[rR]+[\\at]*)`t' Dump2598051_SAC-Jahrbuch_193x wc -l grep -P -i -colour '^((\\at)*[r]+[\\at]*)`t' Dump2598051_SAC-Jahrbuch_193x wc -l -->708311 In 708311 Wörtern. Regex: [A]t der Hut am Anfang negiert den Klasseninhalt, also kein Tabulator Option -i und [r] hat denselben Effekt wie [Rr]</p> <p>In wie vielen Lemmata kommt die Endung "er" vor? Lemmata stehen in der letzten Spalte, --> Suche im Bereich `t`() Anzahl Zeilen mit "er" in der letzten Spalte: grep -P -i -colour `t((\\at)*er\\at*)\$` Dump2598051_SAC-Jahrbuch_193x wc -l -->284272 Anzahl Zeilen mit "er" am Ende der Zeile grep -P -i -colour `t((\\at)*er)\$` Dump2598051_SAC-Jahrbuch_193x wc -l grep -P -i -colour `er\$` Dump2598051_SAC-Jahrbuch_193x wc -l -->147174 In 147174 Lemmata</p>
<p>Ausgabe von Zeilen einer Datei</p> <p>sed [OPTION]... [script-only-if-no-other-script] [input-file]...</p> <p>Sed is a stream editor. A stream editor is used to perform basic text transformations on an input stream (a file or input from a pipeline). While in some ways similar to an editor which permits scripted edits (such as ed), sed works by making only one pass over the input(s), and is consequently more efficient.</p> <p>To print a specific line (line 5) from a file</p> sed -n 5p <filename> Print all the lines between 10 and 20 of a file sed -n '10,20p' <filename>	<p>Vergleichen von zwei Dateien:</p> <p>This tells you the lines of file1 which do <i>not</i> (-v) appear in file2</p> grep -nFvf testdump1 testdump2 <p>This tells you the lines which differ in file1 and file2</p> diff -y testdump1 testdump2 <p>Die Zeilen, die mit einer spitzen, öffnenden Klammer „<“ beginnen, sind nur in der ersten Datei vorhanden, diejenigen, die mit einer spitzen, schließenden Klammer „>“ beginnen, sind nur in der zweiten Datei vorhanden. Zeilen, die in beiden Dateien gleich sind, werden nicht ausgegeben. Die einzelnen Blöcke werden durch so genannte <i>change commands</i> („Änderungsbefehle“) getrennt, die angeben, welche Aktion (Zeilen hinzufügen -a, ändern -c oder entfernen -d) in welchen Zeilen ausgeführt werden soll.</p>
<pre>grep -P -h `t((\w\$.,(.) ((w+:\w+(:w+?)))`t' SAC-Jahrbuch_193* wc 2598052 7794157 40554052 Die mehrfachen \ und \$ können weggelassen werden: grep -P -h `t((\w\$.\\$,\\$ (\w+:\w+(:w+?)))`t' SAC-Jahrbuch_193* > testdump1 2598052 7794157 40554052 Gruppen-Klammern in Klasse haben keinen Einfluss. Einziger Unterschied zum testdump1 sind 35 zusätzliche Klammern „)“: grep -P -h `t([(\\w)(\\\$.) ((\$.) ((\$0 + (w+:\w+(:w+?)))`t' SAC- Jahrbuch_193* >testdump2 2598087 7794262 40554290 210 alleinstehende „, „, („, PP\$ „ sind hier nicht enthalten grep -P -h `t((\\w (\$.) (\$.) (\$0 + (w+:\w+(:w+?)))`t' SAC- Jahrbuch_193* > testdump3 2597842 7793527 40552666</pre>	Anzeige der Zeilen von testdump2 nicht in testdump1: benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8)PCL-1/Übungen/PCL_uebung1_grep\$ grep -nFvf testdump1 testdump2 707333:))) 877019:))) usw. benzro@benzro-eMachines-E625:~\$ sed -n 707332,707336p testdump2 permitting VVG permit))) . SENT . » (» 2. ADJA @ord@ Anzeige der Zeilen von testdump1 nicht in testdump2: benzro@benzro-eMachines-E625:~\$ grep -nFvf testdump2 testdump1 leere Ausgabe benzro@benzro-eMachines-E625:~\$ sed -n 707332,707336p testdump1 permitting VVG permit))) . SENT . » (» 2. ADJA @ord@ Anzeige unterschiedlichen Zeilen: a benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8)PCL-1/Übungen/PCL_uebung1_grep\$ diff testdump1 testdump2 707332a707333 >))) 877017a877019 >))) usw. benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8)PCL-1/Übungen/PCL_uebung1_grep\$ diff testdump2 testdump1 707333d707332 <))) 877019d877017 <))) usw.
<p>Welches ist das häufigste deutsche Relativpronomen?</p> <p>POS-Tags stehen in der mittleren Spalte, -->Suche im Bereich `t`()</p> <p>Anzahl der deutschen Relativpronomen:</p> grep -P -i -colour `^((PRELS PRELAT))`t' Dump2598051_SAC-Jahrbuch_193x wc -l -->16992 Worte sortieren und aggregieren mit sort uniq -c sort -n grep -P -i -colour `^((PRELS PRELAT))`t' Dump2598051_SAC-Jahrbuch_193x sort uniq -c sort -n	<p>Welche Jahreszahlen kommen im Text vor?</p> <p>Jahreszahlen stehen in der ersten Spalte, --> Suche im Bereich: `^`t</p> <p>Anzahl der deutschen Relativpronomen:</p> grep -P -i -colour `^((PRELS PRELAT))`t' Dump2598051_SAC-Jahrbuch_193x grep -P -i -o `^((\\at)*([0,1,2][1-9]{3}) (^\\at*))`t' Dump2598051_SAC-Jahrbuch_193x grep -P -i -o `^((\\at)*([0,1,2][1-9]{3}) (^\\at*))`t' Dump2598051_SAC-Jahrbuch_193x grep -P -i -o `^((\\at)*([0,1,2][1-9]{3}) (^\\at*))`t`

<pre> 3 Die PRELS d 9 Der PRELS d 40 dessen PRELS d 50 Was PRELS was 78 welchen PRELS welch 111 welchem PRELS PRELS welch 129 welchesPRELS welch 282 deren PRELAT d 313 welcherPRELS welch 314 dessen PRELAT d 346 denen PRELS d 686 welche PRELS welch 749 was PRELS was 1054 dem PRELS d 1267 den PRELS d 1751 das PRELS d 3857 der PRELS d -->5948 die PRELS d "Was" am Satzanfang mit 50 Treffern und "die" allgemein am häufigsten mit 5948 Treffern </pre>	<pre> 20 ADJA 20 NN 24 CD 44 N_C 102 NUM 103 D_card 104 A_card 303 N_P 4319 N_card 4832 CARD </pre> <p>Dieser Befehl zeigt die Jahreszahlen an: Wörter und Pos-Tags stehen in der ersten und zweien Spalte-->Suche im Bereich: <code>\[^t]*()</code> <code>grep -P -i '\^([0,1,2][1-9]{3})([0,1,2][1-9]{3})?\t((CD) (NUM) (N_Card) (CARD))\t'</code> Dump2598051 SAC-Jahrbuch_193x Zusätzlich kommen noch römische Jahreszahlen mit POS-Tag A_ord N_C vor. Z.B. XIIe siècle.</p>																				
<p>Pipe: <code>sort uniq -c sort -n; Erst wird alphabetisch sortiert, dann gleiche Zeilen aggregiert und die Anzahl hinzugefügt (-c), zuletzt wird numerisch sortiert (-n)</code></p> <p>Wieviele Wörter im Text enthalten zwei Vokale (inkl. Umlaute) nacheinander? Wörter stehen in der ersten Spalte, --> Suche im Bereich: <code>\[^t]*()</code> genau zwei: <code>grep -P -i '\^([^\t]*(\^aeiouäöüéââûô)[\^aeiouäöüéââûô]{2}[\^aeiouäöüéââûô])\t*</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->377900 mindestens zwei: <code>grep -P -i '\^([^\t]*(\^aeiouäöüéââûô)[\^aeiouäöüéââûô]{2},[\^aeiouäöüéââûô])\t*</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->394660 Je nachdem, wie man die Frage interpretiert sind es 377900 oder 394660 Wörter</p>	<p>Welche Wörter gibt es, die genau fünf Buchstaben lang sind, wo sowohl der erste und der letzte Buchstabe ein Vokal sind? Wörter stehen in der ersten Spalte, --> Suche im Bereich: <code>\[^t]*()</code> <code>grep -P -i '\^([aeiouäöüéââûô](\^w-\^w)*\\$ \&%\{f..;::>=]\{3\}[aeiouäöüéââûô])\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->8404 Es sind 8404 Worte, die teilweise mehrfach auftreten. Eine Liste mit 645 Wörtern erhält man mit folgendem Befehl: <code>grep -P -i '\^([aeiouäöüéââûô](\^w-\^w)*\\$ \&%\{f..;::>=]\{3\}[aeiouäöüéââûô])\t'</code> Dump2598051 SAC-Jahrbuch_193x sort uniq -c</p>																				
<p>Wie oft ist "soit" eine Verbform? Wie oft eine Konjunktion? Wörter und Pos-Tags stehen in der ersten und zweien Spalte-->Suche im Bereich: <code>\[^t]*()</code> <code>grep -P -i '\^([soit](\{IV\}))\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->147 <code>grep -P -i '\^([soit])\t(\^V+)\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->252 cat Dump2598051 SAC-Jahrbuch_193x grep '^soit_tC_C' -->173</p>	<p>Wieviele Wörter gibt es, die mit einem Vokal beginnen und auf "r", "s", "t" oder "w" enden? Wörter stehen in der ersten Spalte, --> Suche im Bereich: <code>\[^t]*()</code> <code>grep -P -i '\^([aeiouäöüéââûô](\^w*)[rstw])\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->17348 Bemerkung: (\^w)* deckt die Sonderzeichen nicht ab. (\^W)* beinhaltet den t und kann nicht mit \w in Serie geschaltet werden. Es wäre möglich, dass noch mehr Worte vorkommen, die Sonderzeichen enthalten [\w-(\^W)*\&%\{f..;::>=]\{3\}[aeiouäöüéââûô](\^w-\^w)*\\$ \&%\{f..;::>=]\{3\}[aeiouäöüéââûô]\t'</p>																				
<p>Welche Verkleinerungsform ist häufiger, "-chen" oder "-lein"? Wörter und Pos-Tags stehen in der ersten und zweien Spalte-->Suche im Bereich: <code>\[^t]*()</code> Verkleinerungsformen kommen bei Nomen vor. Die POS-Tag für deutsche Nomen sind: NN NE Mit dem folgenden Befehl erhält man eine sehr gute Liste mit Wörtern, in denen -lein die Funktion einer Verkleinerung hat <code>grep -P -i '\^([^\t](3,)(lein)\t(CC NN NE))\t'</code> Dump2598051 SAC-Jahrbuch_193x sort uniq -c <code>grep -P -i '\^([^\t](3,)(lein)\t(CC NN NE))\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->541 Mit -chen ist es etwas schwieriger, da -chen auch bei längeren Wörtern eine andere Funktion als die Verkleinerung einnehmen kann. <code>grep -P -i '\^([^\t](3,)(chen)\t(CC NN NE))\t'</code> Dump2598051 SAC-Jahrbuch_193x sort uniq -c <code>grep -P -i '\^([^\t](3,)(chen)\t(CC NN NE))\t'</code> Dump2598051 SAC-Jahrbuch_193x wc -l -->3420 Die 3420 müssen noch nach unten korrigiert werden, die 541 sind ziemlich exakt. Es ist sehr wahrscheinlich, dass -chen in der Funktion einer Verkleinerung mit Faktor 3 bis 5 häufiger auftritt.</p>	<p>Wie lautet der Befehl, der die Anzahl der Adjektive für jedes Jahrbuch (in deinem Jahrehint) einzeln ausgibt? POS-Tags stehen in der mittleren Spalte, -->Suche im Bereich \t()</p> <p>Pos-Tags für Adjektive: deutsch: ADJA und ADJD französisch: A italienisch: ADJ englisch: JJ und JJR und JJS</p> <p><code>grep -P -i -o "\^(\ADJA/\ADJD/\A/\ADJ/\JJ/\JJR/\JJS)\t"</code> SAC-Jahrbuch_193*</p> <p>-->Listet alle Adjektive mit Dateiangabe auf</p> <table border="0"> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJD</td></tr> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJA</td></tr> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJA</td></tr> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJA</td></tr> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJA</td></tr> <tr><td>SAC-Jahrbuch_1930_mul_columns.txt:</td><td>ADJD</td></tr> </table> <p><code>grep -P -i -o "\^(\w-\^:\^)\t"</code></p> <p>-->Entfernt den POS-Tag, dh. Wählt die erste Spalte mit SAC-Jahrbuch_...: <code> sort uniq -c sort -n</code></p> <p>-->Aggregiert und sortiert</p>	SAC-Jahrbuch_1930_mul_columns.txt:	ADJD	SAC-Jahrbuch_1930_mul_columns.txt:	ADJA	SAC-Jahrbuch_1930_mul_columns.txt:	ADJA	SAC-Jahrbuch_1930_mul_columns.txt:	ADJA	SAC-Jahrbuch_1930_mul_columns.txt:	ADJA	SAC-Jahrbuch_1930_mul_columns.txt:	ADJD								
SAC-Jahrbuch_1930_mul_columns.txt:	ADJD																				
SAC-Jahrbuch_1930_mul_columns.txt:	ADJA																				
SAC-Jahrbuch_1930_mul_columns.txt:	ADJA																				
SAC-Jahrbuch_1930_mul_columns.txt:	ADJA																				
SAC-Jahrbuch_1930_mul_columns.txt:	ADJA																				
SAC-Jahrbuch_1930_mul_columns.txt:	ADJD																				
<p>Aus Musterlösung: cat SAC-Jahrbuch_193*_mul_columns.txt grep -Pc 'NN\t[wäöü]+lein\$' -->296 cat SAC-Jahrbuch_193*_mul_columns.txt grep -Pc 'NN\t[wäöü]+chen\$' -->1749</p>	<p>Die Anzahl der Adjektive über die vier Sprachen D,E,F,I aggregiert und für jedes Jahr gruppiert erhält man mit: <code>grep -P -i -o "\^(\ADJA/\ADJD/\A/\ADJ/\JJ/\JJR/\JJS)\t"</code> SAC-Jahrbuch_193* grep -P -i -o "\^(\w-\^:\^)\t" sort uniq -c sort -n</p> <table border="0"> <tr><td>14245 SAC-Jahrbuch_1931_mul_columns.txt:</td><td></td></tr> <tr><td>14435 SAC-Jahrbuch_1936_mul_columns.txt:</td><td></td></tr> <tr><td>14629 SAC-Jahrbuch_1937_mul_columns.txt:</td><td></td></tr> <tr><td>14983 SAC-Jahrbuch_1933_mul_columns.txt:</td><td></td></tr> <tr><td>14999 SAC-Jahrbuch_1930_mul_columns.txt:</td><td></td></tr> <tr><td>15009 SAC-Jahrbuch_1935_mul_columns.txt:</td><td></td></tr> <tr><td>15222 SAC-Jahrbuch_1934_mul_columns.txt:</td><td></td></tr> <tr><td>15288 SAC-Jahrbuch_1939_mul_columns.txt:</td><td></td></tr> <tr><td>15583 SAC-Jahrbuch_1932_mul_columns.txt:</td><td></td></tr> <tr><td>17043 SAC-Jahrbuch_1938_mul_columns.txt:</td><td></td></tr> </table>	14245 SAC-Jahrbuch_1931_mul_columns.txt:		14435 SAC-Jahrbuch_1936_mul_columns.txt:		14629 SAC-Jahrbuch_1937_mul_columns.txt:		14983 SAC-Jahrbuch_1933_mul_columns.txt:		14999 SAC-Jahrbuch_1930_mul_columns.txt:		15009 SAC-Jahrbuch_1935_mul_columns.txt:		15222 SAC-Jahrbuch_1934_mul_columns.txt:		15288 SAC-Jahrbuch_1939_mul_columns.txt:		15583 SAC-Jahrbuch_1932_mul_columns.txt:		17043 SAC-Jahrbuch_1938_mul_columns.txt:	
14245 SAC-Jahrbuch_1931_mul_columns.txt:																					
14435 SAC-Jahrbuch_1936_mul_columns.txt:																					
14629 SAC-Jahrbuch_1937_mul_columns.txt:																					
14983 SAC-Jahrbuch_1933_mul_columns.txt:																					
14999 SAC-Jahrbuch_1930_mul_columns.txt:																					
15009 SAC-Jahrbuch_1935_mul_columns.txt:																					
15222 SAC-Jahrbuch_1934_mul_columns.txt:																					
15288 SAC-Jahrbuch_1939_mul_columns.txt:																					
15583 SAC-Jahrbuch_1932_mul_columns.txt:																					
17043 SAC-Jahrbuch_1938_mul_columns.txt:																					
<p>Der tr command findet und ersetzt Usage: <code>tr [OPTION]... SET1 [SET2]</code> Translate, squeeze, and/or delete characters from standard input, writing to standard output. -c, -C, --complement use the complement of SET1 -d, --delete delete characters in SET1, do not translate -s, --squeeze-repeats replace each input sequence of a repeated character that is listed in SET1 with a single occurrence of that character -t, --truncate-set1 first truncate SET1 to length of SET2 --help display this help and exit --version output version information and exit</p>	<p>Wie viele Types (unterschiedliche Wörter) enthält der Text? -->Text vertikalisierten mit <code>tr -s '\n' < inputfile > outputfile</code> <code>tr -s '\n' < SecondVariety.txt > DumpVerticalSecondVariety.txt</code> -->Interpunkt entfernen <code>tr -d '.;!-' < DumpVerticalSecondVariety.txt > DumpVerticalSecondVariety_.txt</code> -->Anzahl unterschiedliche Wörter <code>sort DumpVerticalSecondVariety_.txt wc -l</code> 3828 Es gibt 3828 unterschiedliche "Wörter"</p>																				

<p>Wie viele Wörter enthält der Text? benzro@benzro-eMachines-E625:~\$ wc SecondVariety.txt 2973 18500 111571 SecondVariety.txt</p> <p>benzro@benzro-eMachines-E625:~\$ wc DumpVerticalSecondVariety_.txt 19427 18500 104900 DumpVerticalSecondVariety_.txt</p> <p>Es gibt 18500 Wörter auf 19427 Zeilen. Und somit gegen 1000 Leerzeilen.</p>	<p>Was sind die 5 häufigsten Wörter? ->Aggregieren mit sort -f und uniq -i für case insensitive sort -f DumpVerticalSecondVariety_.txt uniq -i -c sort -n</p> <p>330 to 340 and 361 a 375 of 927 1069 the Die fünf häufigsten Wörter case insensitive sind: the, of, a, and, to</p> <p>Aus Musterlösung tr '' '\n' < mein_text.txt sort uniq -c wc sort -nr head -5</p>																																																
<p>Der sort command sortiert Strings und Dateien Usage: sort [OPTION]... [FILE]... Ordering options: -f, --ignore-case fold lower case to upper case characters -g, --general-numeric-sort compare according to general numerical value -i, --ignore-nonprinting consider only printable characters -M, --month-sort compare (unknown) < JAN < ... < DEC' -h, --human-numeric-sort compare human readable numbers (e.g., 2K 1G) -n, --numeric-sort compare according to string numerical value -m, --merge merge already sorted files; do not sort -o, --output=FILE write result to FILE instead of standard output -r, --reverse reverse the result of comparisons -u, --unique with -c, check for strict ordering; without -c, output only the first of an equal run</p>	<p>Der uniq command aggregiert exakt gleiche Zeilen Usage: uniq [OPTION]... [INPUT [OUTPUT]] Filter adjacent matching lines from INPUT (or standard input), writing to OUTPUT (or standard output). With no options, matching lines are merged to the first occurrence. -c, --count prefix lines by the number of occurrences -i, --ignore-case ignore differences in case when comparing</p> <p>sort -u Datei_X.txt ODER sort Datei_X.txt uniq</p>																																																
<p>Welche Wörter folgen auf "nicht"? Gib die fünf häufigsten Wörter an. Aufs Englische nicht folgen: grep -P --colour -o '(not Not)\ (w)+\ ' SecondVariety.txt grep -P -o '\ (w)+\ ' sort -f uniq -c sort -n grep: '(not Not)\ (w)+\ ' liefert not Wort '\ (w)+\ ' liefert Wort</p> <p>grep: \$ funktioniert bei diesem Text nicht als Anker am Ende der Zeile.</p> <p>Aus Musterlösung grep -Pio '\bnicht\b \b \w+\b' mein_text.txt sort uniq -c sort -rn head -5</p>	<p>Wie oft kommt die Wortfolge "Wort1 und Wort2" vor, wo beide Wörter auf "en" enden (z.B. "die grünen und blauen ...")? Im Englischen -y and -y: PCL-1/Übungen/PCL_uebung1\$ grep -P --colour -o '(w)+y\ (and)\ (w)+y\ ' SecondVariety.txt slowly and methodically Es gibt einen Treffer slowly and methodically</p>																																																
<p>Syntax</p> <table border="0"> <tr><td>\</td><td>Escapes the character immediately following it</td></tr> <tr><td>.</td><td>Matches any single character except a newline (unless /s is used)</td></tr> <tr><td>^</td><td>Matches at the beginning of the string (or line, if /m is used)</td></tr> <tr><td>\$</td><td>Matches at the end of the string (or line, if /m is used)</td></tr> <tr><td>*</td><td>Matches the preceding element 0 or more times</td></tr> <tr><td>+</td><td>Matches the preceding element 1 or more times</td></tr> <tr><td>?</td><td>Matches the preceding element 0 or 1 times</td></tr> <tr><td>{...}</td><td>Specifies a range of occurrences for the element preceding it</td></tr> <tr><td>[...]</td><td>Matches any one of the characters contained within the brackets</td></tr> <tr><td>(...)</td><td>Groups subexpressions for capturing to \$1, \$2...</td></tr> <tr><td>(?:...)</td><td>Groups subexpressions without capturing (cluster)</td></tr> <tr><td> </td><td>Matches either the subexpression preceding or following it</td></tr> <tr><td>\1, \2 ...</td><td>The text from the Nth group</td></tr> </table>	\	Escapes the character immediately following it	.	Matches any single character except a newline (unless /s is used)	^	Matches at the beginning of the string (or line, if /m is used)	\$	Matches at the end of the string (or line, if /m is used)	*	Matches the preceding element 0 or more times	+	Matches the preceding element 1 or more times	?	Matches the preceding element 0 or 1 times	{...}	Specifies a range of occurrences for the element preceding it	[...]	Matches any one of the characters contained within the brackets	(...)	Groups subexpressions for capturing to \$1, \$2...	(?:...)	Groups subexpressions without capturing (cluster)		Matches either the subexpression preceding or following it	\1, \2 ...	The text from the Nth group	<p>Escape sequences</p> <p>These work as in normal strings.</p> <table border="0"> <tr><td>\a</td><td>Alarm (beep)</td></tr> <tr><td>\e</td><td>Escape</td></tr> <tr><td>\f</td><td>Formfeed</td></tr> <tr><td>\n</td><td>Newline</td></tr> <tr><td>\r</td><td>Carriage return</td></tr> <tr><td>\t</td><td>Tab</td></tr> <tr><td>\038</td><td>Any octal ASCII value</td></tr> <tr><td>\x7f</td><td>Any hexadecimal ASCII value</td></tr> <tr><td>\x{263a}</td><td>A wide hexadecimal value</td></tr> <tr><td>\cx</td><td>Control-x</td></tr> <tr><td>\N{name}</td><td>A named character</td></tr> </table>	\a	Alarm (beep)	\e	Escape	\f	Formfeed	\n	Newline	\r	Carriage return	\t	Tab	\038	Any octal ASCII value	\x7f	Any hexadecimal ASCII value	\x{263a}	A wide hexadecimal value	\cx	Control-x	\N{name}	A named character
\	Escapes the character immediately following it																																																
.	Matches any single character except a newline (unless /s is used)																																																
^	Matches at the beginning of the string (or line, if /m is used)																																																
\$	Matches at the end of the string (or line, if /m is used)																																																
*	Matches the preceding element 0 or more times																																																
+	Matches the preceding element 1 or more times																																																
?	Matches the preceding element 0 or 1 times																																																
{...}	Specifies a range of occurrences for the element preceding it																																																
[...]	Matches any one of the characters contained within the brackets																																																
(...)	Groups subexpressions for capturing to \$1, \$2...																																																
(?:...)	Groups subexpressions without capturing (cluster)																																																
	Matches either the subexpression preceding or following it																																																
\1, \2 ...	The text from the Nth group																																																
\a	Alarm (beep)																																																
\e	Escape																																																
\f	Formfeed																																																
\n	Newline																																																
\r	Carriage return																																																
\t	Tab																																																
\038	Any octal ASCII value																																																
\x7f	Any hexadecimal ASCII value																																																
\x{263a}	A wide hexadecimal value																																																
\cx	Control-x																																																
\N{name}	A named character																																																
<p>Character classes</p> <table border="0"> <tr><td>[amy]</td><td>Match 'a', 'm' or 'y'</td></tr> <tr><td>[f-j]</td><td>Dash specifies range</td></tr> <tr><td>[f-j-]</td><td>Dash escaped or at start or end means 'dash'</td></tr> <tr><td>[^f-j]</td><td>Caret indicates "match any character <i>except</i> these"</td></tr> </table>	[amy]	Match 'a', 'm' or 'y'	[f-j]	Dash specifies range	[f-j-]	Dash escaped or at start or end means 'dash'	[^f-j]	Caret indicates "match any character <i>except</i> these"	<table border="0"> <tr><td>\d</td><td>A digit</td><td>[0-9]</td></tr> <tr><td>\D</td><td>A nondigit</td><td>[^0-9]</td></tr> <tr><td>\w</td><td>A word character</td><td>[a-zA-Z0-9_]</td></tr> <tr><td>\W</td><td>A non-word character</td><td>[^a-zA-Z0-9_]</td></tr> <tr><td>\s</td><td>A whitespace character</td><td>[\t\n\r\f]</td></tr> <tr><td>\S</td><td>A non-whitespace character</td><td>[^ \t\n\r\f]</td></tr> <tr><td>\C</td><td>Match a byte (with Unicode, '.' matches a character)</td><td></td></tr> </table>	\d	A digit	[0-9]	\D	A nondigit	[^0-9]	\w	A word character	[a-zA-Z0-9_]	\W	A non-word character	[^a-zA-Z0-9_]	\s	A whitespace character	[\t\n\r\f]	\S	A non-whitespace character	[^ \t\n\r\f]	\C	Match a byte (with Unicode, '.' matches a character)																				
[amy]	Match 'a', 'm' or 'y'																																																
[f-j]	Dash specifies range																																																
[f-j-]	Dash escaped or at start or end means 'dash'																																																
[^f-j]	Caret indicates "match any character <i>except</i> these"																																																
\d	A digit	[0-9]																																															
\D	A nondigit	[^0-9]																																															
\w	A word character	[a-zA-Z0-9_]																																															
\W	A non-word character	[^a-zA-Z0-9_]																																															
\s	A whitespace character	[\t\n\r\f]																																															
\S	A non-whitespace character	[^ \t\n\r\f]																																															
\C	Match a byte (with Unicode, '.' matches a character)																																																
<p>Anchors</p> <p>All are zero-width assertions.</p> <table border="0"> <tr><td>^</td><td>Match string start (or line, if /m is used)</td></tr> <tr><td>\$</td><td>Match string end (or line, if /m is used) or before newline</td></tr> <tr><td>\b</td><td>Match word boundary (between \w and \W)</td></tr> <tr><td>\B</td><td>Match except at word boundary (between \w and \w or \W and \W)</td></tr> <tr><td>\A</td><td>Match string start (regardless of /m)</td></tr> <tr><td>\Z</td><td>Match string end (before optional newline)</td></tr> <tr><td>\z</td><td>Match absolute string end</td></tr> <tr><td>\G</td><td>Match where previous m//g left off</td></tr> </table>	^	Match string start (or line, if /m is used)	\$	Match string end (or line, if /m is used) or before newline	\b	Match word boundary (between \w and \W)	\B	Match except at word boundary (between \w and \w or \W and \W)	\A	Match string start (regardless of /m)	\Z	Match string end (before optional newline)	\z	Match absolute string end	\G	Match where previous m//g left off	<p>Quantifiers</p> <p>Quantifiers are greedy by default – match the longest leftmost.</p> <table border="0"> <thead> <tr><th>Maximal</th><th>Minimal</th><th>Allowed range</th></tr> </thead> <tbody> <tr><td>{n,m}</td><td>{n,m}?</td><td>Must occur at least n times but no more than m times</td></tr> <tr><td>{n,}</td><td>{n,}?</td><td>Must occur at least n times</td></tr> <tr><td>{n}</td><td>{n}?</td><td>Must occur exactly n times</td></tr> <tr><td>*</td><td>*?</td><td>0 or more times (same as {0,})</td></tr> <tr><td>+</td><td>+?</td><td>1 or more times (same as {1,})</td></tr> <tr><td>?</td><td>??</td><td>0 or 1 time (same as {0,1})</td></tr> </tbody> </table>	Maximal	Minimal	Allowed range	{n,m}	{n,m}?	Must occur at least n times but no more than m times	{n,}	{n,}?	Must occur at least n times	{n}	{n}?	Must occur exactly n times	*	*?	0 or more times (same as {0,})	+	+?	1 or more times (same as {1,})	?	??	0 or 1 time (same as {0,1})											
^	Match string start (or line, if /m is used)																																																
\$	Match string end (or line, if /m is used) or before newline																																																
\b	Match word boundary (between \w and \W)																																																
\B	Match except at word boundary (between \w and \w or \W and \W)																																																
\A	Match string start (regardless of /m)																																																
\Z	Match string end (before optional newline)																																																
\z	Match absolute string end																																																
\G	Match where previous m//g left off																																																
Maximal	Minimal	Allowed range																																															
{n,m}	{n,m}?	Must occur at least n times but no more than m times																																															
{n,}	{n,}?	Must occur at least n times																																															
{n}	{n}?	Must occur exactly n times																																															
*	*?	0 or more times (same as {0,})																																															
+	+?	1 or more times (same as {1,})																																															
?	??	0 or 1 time (same as {0,1})																																															

<p>Extended constructs</p> <table border="0"> <tr><td>(?#text)</td><td>A comment</td></tr> <tr><td>(?imxs-imsx:...)</td><td>Enable/disable option (as per m// modifiers)</td></tr> <tr><td>(?=...)</td><td>Zero-width positive lookahead assertion</td></tr> <tr><td>(?!...)</td><td>Zero-width negative lookahead assertion</td></tr> <tr><td>(?<=...)</td><td>Zero-width positive lookbehind assertion</td></tr> <tr><td>(?<!...)</td><td>Zero-width negative lookbehind assertion</td></tr> <tr><td>(?>...)</td><td>Grab what we can, prohibit backtracking</td></tr> </table>	(?#text)	A comment	(?imxs-imsx:...)	Enable/disable option (as per m// modifiers)	(?=...)	Zero-width positive lookahead assertion	(?!...)	Zero-width negative lookahead assertion	(?<=...)	Zero-width positive lookbehind assertion	(?<!...)	Zero-width negative lookbehind assertion	(?>...)	Grab what we can, prohibit backtracking	<p>Metacharacters:</p> <p>The following 12 characters need to be escaped with a backslash - “\” - because by default, they mean something special.</p> <table border="0"> <tr><td>\</td><td> </td><td>()</td><td>[{</td><td>^</td><td>\$</td><td>*</td><td>+</td><td>?</td><td>.</td></tr> </table> <table border="0"> <tr><td>.</td><td>Match any one character (except \n)</td></tr> <tr><td> </td><td>Alternation</td></tr> <tr><td>()</td><td>Group and capture</td></tr> <tr><td>[]</td><td>Define character class</td></tr> <tr><td>\</td><td>Modify the meaning of the next char.</td></tr> </table>	\		()	[{	^	\$	*	+	?	.	.	Match any one character (except \n)		Alternation	()	Group and capture	[]	Define character class	\	Modify the meaning of the next char.
(?#text)	A comment																																		
(?imxs-imsx:...)	Enable/disable option (as per m// modifiers)																																		
(?=...)	Zero-width positive lookahead assertion																																		
(?!...)	Zero-width negative lookahead assertion																																		
(?<=...)	Zero-width positive lookbehind assertion																																		
(?<!...)	Zero-width negative lookbehind assertion																																		
(?>...)	Grab what we can, prohibit backtracking																																		
\		()	[{	^	\$	*	+	?	.																										
.	Match any one character (except \n)																																		
	Alternation																																		
()	Group and capture																																		
[]	Define character class																																		
\	Modify the meaning of the next char.																																		
<p>Anchors:</p> <table border="0"> <tr><td>^</td><td>Match at the beginning of a string (or line)</td></tr> <tr><td>\$</td><td>Match at the end of a string (or line)</td></tr> <tr><td>\b</td><td>Match at a ‘word’ boundary</td></tr> <tr><td>\B</td><td>Match at not a ‘word’ boundary</td></tr> </table> <p>These are also known as <i>zero width assertions</i>.</p>	^	Match at the beginning of a string (or line)	\$	Match at the end of a string (or line)	\b	Match at a ‘word’ boundary	\B	Match at not a ‘word’ boundary	<p>Quantifiers:</p> <p>These quantifiers apply to the preceding <i>atom</i>.</p> <table border="0"> <tr><td>*</td><td>Match 0 or more times</td></tr> <tr><td>+</td><td>Match 1 or more times</td></tr> <tr><td>?</td><td>Match 0 or 1 times</td></tr> <tr><td>{N}</td><td>Match exactly N times</td></tr> <tr><td>{N,}</td><td>Match at least N times</td></tr> <tr><td>{N,M}</td><td>Match at least N but not more than M times</td></tr> </table> <p>By default, quantifiers are “greedy”. They attempt to match as many characters as possible. In order to make them match as few characters as possible, follow them with a question mark “?”.</p>	*	Match 0 or more times	+	Match 1 or more times	?	Match 0 or 1 times	{N}	Match exactly N times	{N,}	Match at least N times	{N,M}	Match at least N but not more than M times														
^	Match at the beginning of a string (or line)																																		
\$	Match at the end of a string (or line)																																		
\b	Match at a ‘word’ boundary																																		
\B	Match at not a ‘word’ boundary																																		
*	Match 0 or more times																																		
+	Match 1 or more times																																		
?	Match 0 or 1 times																																		
{N}	Match exactly N times																																		
{N,}	Match at least N times																																		
{N,M}	Match at least N but not more than M times																																		
<p>Character class metacharacters:</p> <table border="0"> <tr><td>^</td><td>If the first character of a class, negates that class</td></tr> <tr><td>-</td><td>Unless first or last character of a class, used for a range</td></tr> </table> <p>Character class shortcuts:</p> <table border="0"> <tr><td>\d</td><td>[0-9]</td><td>A digit</td></tr> <tr><td>\D</td><td>[^0-9]</td><td>A non-digit</td></tr> <tr><td>\s</td><td>[\t\n\r\f]</td><td>A whitespace char.</td></tr> <tr><td>\S</td><td>[^ \t\n\r\f]</td><td>A non-whitespace char.</td></tr> <tr><td>\w</td><td>[a-zA-Z0-9_]</td><td>A ‘word’ char.</td></tr> <tr><td>\W</td><td>[^a-zA-Z0-9_]</td><td>A ‘non-word’ char.</td></tr> </table> <p>These shortcuts can be used either on their own, or within a character class.</p>	^	If the first character of a class, negates that class	-	Unless first or last character of a class, used for a range	\d	[0-9]	A digit	\D	[^0-9]	A non-digit	\s	[\t\n\r\f]	A whitespace char.	\S	[^ \t\n\r\f]	A non-whitespace char.	\w	[a-zA-Z0-9_]	A ‘word’ char.	\W	[^a-zA-Z0-9_]	A ‘non-word’ char.	<p>Modifiers:</p> <p>These modifiers apply to the entire pattern</p> <table border="0"> <tr><td>/i</td><td>Ignore case</td></tr> <tr><td>/g</td><td>Match globally (all)</td></tr> <tr><td>/m</td><td>Let ^ and \$ match next to embedded \n</td></tr> <tr><td>/s</td><td>Let . match \n</td></tr> <tr><td>/x</td><td>Ignore most whitespace and allow comments</td></tr> <tr><td>/e</td><td>Evaluate right hand side of s/// as an expression</td></tr> </table>	/i	Ignore case	/g	Match globally (all)	/m	Let ^ and \$ match next to embedded \n	/s	Let . match \n	/x	Ignore most whitespace and allow comments	/e	Evaluate right hand side of s/// as an expression
^	If the first character of a class, negates that class																																		
-	Unless first or last character of a class, used for a range																																		
\d	[0-9]	A digit																																	
\D	[^0-9]	A non-digit																																	
\s	[\t\n\r\f]	A whitespace char.																																	
\S	[^ \t\n\r\f]	A non-whitespace char.																																	
\w	[a-zA-Z0-9_]	A ‘word’ char.																																	
\W	[^a-zA-Z0-9_]	A ‘non-word’ char.																																	
/i	Ignore case																																		
/g	Match globally (all)																																		
/m	Let ^ and \$ match next to embedded \n																																		
/s	Let . match \n																																		
/x	Ignore most whitespace and allow comments																																		
/e	Evaluate right hand side of s/// as an expression																																		
<pre>recursive grep search grep -ril 'null' /home/al/sarah /var/www # search multiple dirs</pre> <p>Searching for a string in multiple files Our next grep command example searches for all occurrences of the text string joe within all files of the current directory: grep 'joe' *</p> <p>search all files in the current directory that end in the file extension .txt, like this: grep 'joe' *.txt</p> <p>One way to find all the sub-directories in the current directory is to mix the Linux ls and grep commands together in a pipe, like this: ls -al grep '^d'</p> <pre>grep gzip files zgrep foo myfile.gz # all lines containing the pattern 'foo'</pre>	<p>Power file searching with find and grep This is a special way of mixing the Linux find and grep commands together to search every file in every subdirectory of my current location. It searches for the string “foo” in every file below the current directory, in a case-insensitive manner. find . -type f -exec grep -il 'foo' {} \;</p> <ul style="list-style-type: none"> • “.” means “look in the current directory” • “-type f” means “look in files only” • “-exec grep -il foo” means “search for the string ‘foo’ in a case-insensitive manner, and return the matching line and filename when a match is found • “{} \;” is a little bizarre syntax that you need to add to the end of your find command whenever you add the -exec option. I try to think of it as a placeholder for the filenames the find command finds. 																																		

<p>1. Schritt</p> <p>Wo sind wir? pwd</p> <p>Wie kommen wir in das Verzeichnis, in dem wir arbeiten wollen? cd Verzeichnis_X ODER cd ..</p> <p>Wie sehen wir, welche Dateien im Verzeichnis vorliegen? ls -l</p> <p>Wie können wir in eine Datei hineinschauen? more Datei_X.txt (quit with 'q'), [less]</p>	<p>2. Schritt</p> <p>Wie können wir die Wörter in einer Textdatei zählen? wc Datei_X.txt</p> <p>Wie können wir die Anzahl unterschiedlicher Wörter in einem Text zählen?</p> <ol style="list-style-type: none"> 1. Text vertikalisieren 2. Wörter sortieren 3. Duplikate entfernen 4. Wörter zählen <p>Wie können wir einen Text vertikalisieren?</p> <ol style="list-style-type: none"> 1. Suche Leerschlag und ersetze mit \n im Editor oder 2. tr '' '\n' < Datei_X.txt
<p>3. Schritt</p> <p>Wie können wir die Zeilen einer Datei sortieren? sort Datei_X.txt</p> <p>Wie können wir eine Datei sortieren und identische Zeilen entfernen?</p> <p>sort -u Datei_X.txt ODER sort Datei_X.txt uniq</p>	<p>4. Schritt</p> <p>Wie berechnen wir die Häufigkeit für jedes Wort in einer vertikalisierten Textdatei ?</p> <p>sort Datei_X.txt uniq -c</p> <p>Wie sortieren wir nach der Häufigkeit?</p> <p>sort -n (numerische vs. alphabet. Sortierung)</p> <p>Wie sortieren wir rückwärts (= häufigste zuerst)?</p> <p>sort -rn</p>
<p>5. Schritt</p> <p>Wie können wir grosse Textmengen durchsuchen? grep 'Suchwort' Dateimuster</p> <p>Wieviele Sätze kommen im Jahrbuch 1900 vor? grep '<s' SAC_Jahrbuch_1900.xml wc</p> <p>Wieviele Sätze kommen in allen Jahrbüchern des 19. Jhd. vor? grep '<s' SAC_Jahrbuch_18*.xml wc grep -c '<s' SAC_Jahrbuch_18*.xml</p>	<p>6. Schritt</p> <p>Wie oft kommt das Wort 'Berg' vor? grep 'Berg\lt' *18*.xml wc</p> <p>Wie oft kommt das Lemma 'Berg' vor?</p> <p>grep '\tBerg\\$' *18*.xml wc</p> <p>Wieviele Substantive kommen vor?</p> <p>grep '\tNN\lt' *18*.xml wc</p> <p>Wieviele deutsche Adjektive kommen vor?</p> <p>grep '\tADJA\lt' *18*.xml wc</p>
<p>tr = transliterate</p> <p>Vertikalisieren eines Textes: Ersetze jeden Leerschlag mit einem Zeilenumbruch tr '' '\n' < Test1.txt</p> <p>Konvertiere alle Grossbuchstaben in Kleinbuchstaben tr '[A-Z]' '[a-z]' < Test1.txt</p>	<p>Unix-Kommandos</p> <p>Umlenken der Ausgabe mit Überschreiben: '>' Umlenken der Ausgabe mit Anfügen: '>>' Umlenken der Eingabe: '<' Verbinden von zwei Kommandos: Pipe ' '</p>

<p>Wildcards</p> <p>= Stellvertretersymbole</p> <ul style="list-style-type: none"> • Wildcard * == beliebige Anzahl Zeichen • Wildcard ? == ein beliebiges Zeichen <p>Wildcard [] == definiert eine Menge von Zeichen, von denen eines zutreffen muss.</p> <p>Wildcard - == ein Strich in einer Klammer [] definiert einen Bereich von – bis.</p>	<p>Zeichen-Kodierung</p> <p>Latin-1 = iso-8859-1 (Umlaute als 1 byte) Alt! ☷</p> <p>UTF-8 (Umlaute als 2 bytes) Gut! ☺</p> <p>Unix-Kommando zur Bestimmung der Kodierung file -i filename</p> <p>Unix-Kommando zur Konvertierung einer Kodierung in eine andere (from – to) iconv -f -t filename</p>														
<p>grep</p> <p>grep und Zählen (count) der Ergebnisse</p> <p>grep -c pattern file zählt per Datei grep pattern file wc -l zählt die gesamte Ausgabe</p>	<p>grep</p> <p>grep für ein Muster xx und Anzeige des Kontexts</p> <p>grep -B 2 xx file Zeige zwei Zeilen des vorangehenden Kontexts (Before) grep -A 4 xx file Zeige vier Zeilen des folgenden Kontexts (After) grep -C 3 xx file Zeige drei Zeilen des vorangehenden und folgenden Kontexts (Context)</p>														
<p>grep</p> <p>Gegeben ein Text aus dem Text+Berg-Korpus mit Part-of-Speech Tags und Lemmas im 3-Spalten-Format.</p> <table> <tbody> <tr> <td>grep VV</td> <td>findet alle Verben im Deutschen</td> </tr> <tr> <td>grep ADJ</td> <td>findet alle Adjektive im Deutschen</td> </tr> <tr> <td>grep ADV</td> <td>findet alle Adverbien</td> </tr> <tr> <td>grep befreien</td> <td>findet Lemma und Wortform</td> </tr> <tr> <td>grep '^befreien'</td> <td>findet Wortform</td> </tr> <tr> <td>grep '^be.*en\lt'</td> <td>findet alle Wortformen mit <i>be...en</i></td> </tr> <tr> <td>grep '\tbe..en\$'</td> <td>findet alle Lemmas mit <i>be...en</i></td> </tr> </tbody> </table>	grep VV	findet alle Verben im Deutschen	grep ADJ	findet alle Adjektive im Deutschen	grep ADV	findet alle Adverbien	grep befreien	findet Lemma und Wortform	grep '^befreien'	findet Wortform	grep '^be.*en\lt'	findet alle Wortformen mit <i>be...en</i>	grep '\tbe..en\$'	findet alle Lemmas mit <i>be...en</i>	<p>grep</p> <p>grep -o '^be.*en\lt' gibt statt der Zeile nur den Treffer aus → erlaubt Sortieren und Zählen der Hits</p> <p>grep --colour '^be.*en\lt' markiert den Treffer in der Zeile</p>
grep VV	findet alle Verben im Deutschen														
grep ADJ	findet alle Adjektive im Deutschen														
grep ADV	findet alle Adverbien														
grep befreien	findet Lemma und Wortform														
grep '^befreien'	findet Wortform														
grep '^be.*en\lt'	findet alle Wortformen mit <i>be...en</i>														
grep '\tbe..en\$'	findet alle Lemmas mit <i>be...en</i>														
<p>grep und Erweiterungen</p> <p>grep bietet die Grundfunktionalität der Suche mit regulären Ausdrücken.</p> <p>Erweiterungen von grep</p> <p>egrep</p> <p>grep -E ist wie egrep</p> <p>grep -P erlaubt reguläre Ausdrücke wie in Perl. Sehr mächtig! Empfohlen!</p>	<p>grep</p> <table border="1"> <thead> <tr> <th>Befehl</th> <th>Erklärung</th> </tr> </thead> <tbody> <tr> <td>grep '\tbe..en\$'</td> <td>findet alle Lemmas mit 2 Zeichen zwischen be und en</td> </tr> <tr> <td>grep '\tbe.[aeiou].en\$'</td> <td>findet alle Lemmas mit Vokal zwischen 2 Buchstaben</td> </tr> <tr> <td>grep '\tbe.[äöü].en\$'</td> <td>findet alle Lemmas mit Umlaut zwischen 2 Buchstaben</td> </tr> <tr> <td>grep '^[bg]e..en\lt'</td> <td>findet alle Wortformen mit <i>be..en</i> und <i>ge..en</i></td> </tr> </tbody> </table>	Befehl	Erklärung	grep '\tbe..en\$'	findet alle Lemmas mit 2 Zeichen zwischen be und en	grep '\tbe.[aeiou].en\$'	findet alle Lemmas mit Vokal zwischen 2 Buchstaben	grep '\tbe.[äöü].en\$'	findet alle Lemmas mit Umlaut zwischen 2 Buchstaben	grep '^[bg]e..en\lt'	findet alle Wortformen mit <i>be..en</i> und <i>ge..en</i>				
Befehl	Erklärung														
grep '\tbe..en\$'	findet alle Lemmas mit 2 Zeichen zwischen be und en														
grep '\tbe.[aeiou].en\$'	findet alle Lemmas mit Vokal zwischen 2 Buchstaben														
grep '\tbe.[äöü].en\$'	findet alle Lemmas mit Umlaut zwischen 2 Buchstaben														
grep '^[bg]e..en\lt'	findet alle Wortformen mit <i>be..en</i> und <i>ge..en</i>														

<p>Wichtig!</p> <p>Ein regulärer Ausdruck sucht immer den längsten möglichen String (= <i>greedy pattern matching</i>). Unterscheide: grep -o '^be.*en\lt' findet alle Wortformen mit <i>be...en</i> und grep -o '^be.*en' findet alle Zeilen mit <i>be...en</i></p>	<p>grep</p> <p>grep aa grep -i "aa" Case-insensitive! findet "aa" "AA", "aA", "Aa" grep [aeiou][aeiou] findet alle Vokal-Paare! grep [aeiou][aeiou][aeiou] findet alle Vokal-Tripel! grep [aeiou][aeiou][aeiou][aeiou] findet alle Vokal-Quadrupel! Gibt es wirklich? grep '[aeiou]\{4\}' mit Anf-Zeichen!</p>										
<pre>benzro@benzro-eMachines-E625:~\$ grep -P '[aeiou]\{4\}' greptest.txt weeeaioann ienaaeeanenienlödieeeeeakldjije dikeeeeeeeiek dieiiiiieis ieeiadjeiiiiii dikejjffnnneiiiiidd dieiaooeaieiied eieifefioieiiioiod ieowiejiofoajjeejjjeiieieeeeeee dijodljsioejiie dsijeoō ieiei djijoeij ifojfiojf benzro@benzro-eMachines-E625:~\$ grep -P '[aeiou]\{4\}' greptest.txt benzro@benzro-eMachines-E625:~\$ grep -P '[aeiou]\{4\}' greptest.txt weeeaioann ienaaeeanenienlödieeeeeakldjije dikeeeeeeeiek dieiiiiieis ieeiadjeiiiiii dikejjffnnneiiiiidd dieiaooeaieiied eieifefioieiiioiod ieowiejiofoajjeejjjeiieieeeeeee dijodljsioejiie dsijeoō ieiei djijoeij ifojfioif</pre>	<p>grep</p> <p>grep '^[1234567890]' findet Wörter, die mit einer Ziffer beginnen grep '^[1234567890][a-z]' findet Wörter, die mit einer Ziffer + einem Buchstabe beginnen (Beachte grosse und kleine Buchstaben und Umlaute) grep -E '^\\d[a-z]' Alternative Notation!</p>										
<p>Wichtige Symbole</p> <table> <tr> <td>\n</td><td>neue Zeile</td></tr> <tr> <td>\t</td><td>Tabulator</td></tr> <tr> <td>\d</td><td>Ziffern 0-9 ("digits")</td></tr> <tr> <td>\w</td><td>Buchstaben [A-Za-z], Unterstrich und Ziffern</td></tr> <tr> <td>\s</td><td>Leerschlag, Tabulator, neue Zeile</td></tr> </table>	\n	neue Zeile	\t	Tabulator	\d	Ziffern 0-9 ("digits")	\w	Buchstaben [A-Za-z], Unterstrich und Ziffern	\s	Leerschlag, Tabulator, neue Zeile	<p>grep -P</p> <p>grep -P '(ADJ ADV)' mit Alternativen grep -P 'th?al' findet <i>thal</i> und <i>tal</i> grep -P 't(h)al' ist äquivalent zu grep -P 'th?al' !!! grep -P '^schnee*' ist äquivalent zu grep -P '^schne+' !!!</p>
\n	neue Zeile										
\t	Tabulator										
\d	Ziffern 0-9 ("digits")										
\w	Buchstaben [A-Za-z], Unterstrich und Ziffern										
\s	Leerschlag, Tabulator, neue Zeile										

```
benzro@benzro-eMachines-E625:~$ grep -P ee* greptest.txt
weeaioann
ienaeaaenienlenlödieeeeeakldjije dikeieeeeeediek dieiiieis
ieeiadijeiiiiiii dijkfeijjfnnneiiiiid
dieiaöööeiaeiiied
eieifejfioieiioiod ieowiejiofoajjejjjeiieeeeeeee
dijödljsioejiiie
dijdse
dsijeioö eieii
dfsiadjjasije
djijioeij ifojfiojf
diojeoijfndoiaj
benzro@benzro-eMachines-E625:~$ grep -P e+ greptest.txt
weeaioann
ienaeaaenienlenlödieeeeeakldjije dikeieeeeeediek dieiiieis
ieeiadijeiiiiiii dijkfeijjfnnneiiiiid
dieiaöööeiaeiiied
eieifejfioieiioiod ieowiejiofoajjejjjeiieeeeeeee
dijödljsioejiiie
dijdse
dsijeioö eieii
dfsiadjjasije
djijioeij ifojfiojf
dijoeoijfndoiaj
```

```
benzro@benzro-eMachines-E625:~$ grep -P '(\w)\1' greptest.txt
weeaioann
ienaeaaenienlenlödieeeeeakldjije dikeieeeeeediek dieiiieis
ieeiadijeiiiiiii dijkfeijjfnnneiiiiid
dieiaöööeiaeiiied
eieifejfioieiioiod ieowiejiofoajjejjjeiieeeeeeee
dijödljsioejiiie
dsijeioö eieii
dfsiadjjasije
diojeoijfndoiaj
benzro@benzro-eMachines-E625:~$ grep -P '(\w)\1.\1' greptest.txt
ienaeaaenienlenlödieeeeeakldjije dikeieeeeeediek dieiiieis
ieeiadijeiiiiiii dijkfeijjfnnneiiiiid
eieifejfioieiioiod ieowiejiofoajjejjjeiieeeeeeee
```

Reguläre Ausdrücke - Eine rekursive Definition

Jedes Zeichen (auch ein leeres) ist ein regulärer Ausdruck.
Beispiel: "a"

Die **Sequenz** (Folge) von regulären Ausdrücken r_1 und r_2 ist ein regulärer Ausdruck.

Beispiel: " $r_1 r_2$ "

Die **Alternative** von regulären Ausdrücken r_1 und r_2 ist ein regulärer Ausdruck.

Beispiel: " $r_1 | r_2$ "

Ein regulärer Ausdruck r mit **Optionalität, ein- oder mehrmaliger Wiederholung** ist ein regulärer Ausdruck.

Beispiel: " r^* "

Repetition takes precedence over concatenation, which in turn takes precedence over alternation.

Grenzen Regulärer Ausdrücke

aabb ist ein regulärer Ausdruck

aⁿb^{*} ist ein regulärer Ausdruck

Wir können aber **nicht** allgemein formulieren, dass genauso viele **a** wie **b** erkannt werden sollen. Salopp: Reguläre Ausdrücke haben keine Zähler!

Formal:

aⁿbⁿ mit $n \in \mathbb{N}$ ist **keine** reguläre Sprache!

aⁿbⁿ mit $n \in \mathbb{N}$ ist **keine** reguläre Sprache!

Aber

aⁿbⁿ mit $n \in \mathbb{N}, n < 10$ ist eine reguläre Sprache!

grep -P

grep -P '([aeiou])\1' findet alle Vokalpaare

Beachte die **Wiederaufnahme** des ersten Treffers (durch runde Klammern markiert) mit \1

grep -P '^(w)\1' findet alle Buchstaben- oder Ziffernpaare am Wortanfang.

grep -P '(\d)\1(\d)\2' findet zwei gleiche Ziffernpaare.

findet Buchstabentriplegrep -P "(\w)\1(\w)\2"

```
benzro@benzro-eMachines-E625:~$ grep -P '(\w)\1(\w)\2' greptest.txt
weeaioann
ienaeaaenienlenlödieeeeeakldjije dikeieeeeeediek dieiiieis
ieeiadijeiiiiiii dijkfeijjfnnneiiiiid
eieifejfioieiioiod ieowiejiofoajjejjjeiieeeeeeee
```

Reguläre Ausdrücke

a+ == aa* (a einmal oder mehrmals)

a? == a|_ (a oder nichts)

a{3} == aaa

a{2,3} == aa|aaa

[ab] == a|b

\d == 0|1|2|3|4|5|6|7|8|9

Reguläre Sprachen

= Sprachen, die mit regulären Ausdrücken definiert werden können.

Formal:

aⁿbⁿ mit $n \in \mathbb{N}$ ist **keine** reguläre Sprache!

Aber

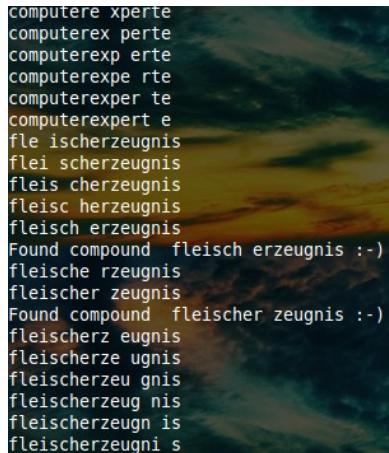
aⁿbⁿ mit $n \in \mathbb{N}, n < 10$ ist eine reguläre Sprache!

aⁿb^m mit $n, m \in \mathbb{N}$ ist eine reguläre Sprache

aⁿb^mc^m mit $n, m \in \mathbb{N}, m < 5$ ist eine reg. Sprache

a^mb⁵c^m mit $m \in \mathbb{N}$ ist **keine** reguläre Sprache!

<h3>Reguläre Ausdrücke vs. Wildcards</h3> <p>sind nicht das selbe!</p> <ul style="list-style-type: none"> – Wildcard * == beliebige Anzahl Zeichen – RegEx * == null oder mehr Wiederholungen eines reg. Ausdrucks – Wildcard ? == ein beliebiges Zeichen – RegEx ? == null oder 1 Instanz eines reg. Ausdrucks 	<h3>First steps in the Python interpreter</h3> <ol style="list-style-type: none"> 1. Start the Python interpreter with 'python' 2. Test Python commands 3. Quit the Python interpreter with 'Ctrl+d' 										
<h3>Write programs in Python</h3> <ol style="list-style-type: none"> 1. Write your Python program in an editor. 2. Save the file to e.g. 'my_program_1.py' 3. Execute the program with 'python my_program_1.py' <p>Note: Interrupt an 'endless' program with 'Ctrl+c'.</p>	<h3>First Python commands</h3> <table border="1" data-bbox="801 718 1428 1100"> <thead> <tr> <th></th><th>Examples</th></tr> </thead> <tbody> <tr> <td>Numbers and Arithmetic</td><td><code>2 + 4</code> <code>21 / 3 * 4**7</code></td></tr> <tr> <td>Variables and Assignments</td><td><code>counter = 21</code> <code>counter = 21 / 7 * 5</code> <code>counter += 2</code> <code>x = y = z = 15</code></td></tr> <tr> <td>Strings</td><td><code>word1 = 'der'</code> <code>word2 = word1 + ' Computer'</code> <code>x = len(word2)</code></td></tr> <tr> <td>Strings with Indexes</td><td><code>letter = word[0]</code> <code>part = word[0:5]</code> <code>suffix = word[-2:]</code></td></tr> </tbody> </table>		Examples	Numbers and Arithmetic	<code>2 + 4</code> <code>21 / 3 * 4**7</code>	Variables and Assignments	<code>counter = 21</code> <code>counter = 21 / 7 * 5</code> <code>counter += 2</code> <code>x = y = z = 15</code>	Strings	<code>word1 = 'der'</code> <code>word2 = word1 + ' Computer'</code> <code>x = len(word2)</code>	Strings with Indexes	<code>letter = word[0]</code> <code>part = word[0:5]</code> <code>suffix = word[-2:]</code>
	Examples										
Numbers and Arithmetic	<code>2 + 4</code> <code>21 / 3 * 4**7</code>										
Variables and Assignments	<code>counter = 21</code> <code>counter = 21 / 7 * 5</code> <code>counter += 2</code> <code>x = y = z = 15</code>										
Strings	<code>word1 = 'der'</code> <code>word2 = word1 + ' Computer'</code> <code>x = len(word2)</code>										
Strings with Indexes	<code>letter = word[0]</code> <code>part = word[0:5]</code> <code>suffix = word[-2:]</code>										
<pre>>>> word="Alibaba" >>> len(word) 7 Element 2 bis aber nicht inklusive 5 >>> istr=word[2:5] >>> istr 'iba' Element 0 bis aber nicht inklusive 5 >>> istr=word[:5] >>> istr 'Aliba' Element 0 bis aber nicht inklusive 7 >>> istr=word[:7] >>> istr 'Alibaba' Element StringLänge-2 bis aber nicht inklusive Stringlänge >>> isuff=word[-2:] >>> isuff 'ba' >>> istr=word[len(word)-2:] >>> istr 'ba'</pre>	<h3>Loops and conditions in Python</h3> <table border="1" data-bbox="833 1190 1412 1538"> <thead> <tr> <th></th><th>Examples</th></tr> </thead> <tbody> <tr> <td>while Loops and Indentation</td><td><code>while counter < 10:</code> <code> print "Hello"</code> <code> counter += 1</code></td></tr> <tr> <td>Comparisons</td><td><code>counter < 50</code> <code>counter == 13</code> <code>word == 'Linguist'</code></td></tr> <tr> <td>if – elif – else Conditions</td><td><code>if x < 10:</code> <code> print "Hello"</code> <code>else:</code> <code> print "Hi"</code></td></tr> <tr> <td>Comments</td><td><code># This is a comment</code></td></tr> </tbody> </table>		Examples	while Loops and Indentation	<code>while counter < 10:</code> <code> print "Hello"</code> <code> counter += 1</code>	Comparisons	<code>counter < 50</code> <code>counter == 13</code> <code>word == 'Linguist'</code>	if – elif – else Conditions	<code>if x < 10:</code> <code> print "Hello"</code> <code>else:</code> <code> print "Hi"</code>	Comments	<code># This is a comment</code>
	Examples										
while Loops and Indentation	<code>while counter < 10:</code> <code> print "Hello"</code> <code> counter += 1</code>										
Comparisons	<code>counter < 50</code> <code>counter == 13</code> <code>word == 'Linguist'</code>										
if – elif – else Conditions	<code>if x < 10:</code> <code> print "Hello"</code> <code>else:</code> <code> print "Hi"</code>										
Comments	<code># This is a comment</code>										
<h3>Variables</h3> <p>When a variable is used for the first time, it must get a value.</p> <p>The first value assignment determines the variable's type (integer, floating point number, string, Boolean).</p> <p>The variable's type restricts the possible operations.</p> <p>Examples:</p> <ul style="list-style-type: none"> • Integer: +, -, *, / • String: +, len() 	<h3>Indentation in Python</h3> <p>Indentation in Python is semantically significant.</p> <p>A block of commands in e.g. loops or conditions must be indented.</p> <p>Note: The output of a Python program (often) changes with indentation. ☺ Be careful!</p>										

	Examples	
for loop	<pre>for my_item in my_list for my_num in range(0, 30) for my_num in range(0, 30, 3) for my_num in range(30)</pre>	Range von 0 bis aber nicht inklusive 5 >>> ilist1=range(0, 5) >>> ilist1 [0, 1, 2, 3, 4] Zweierschritte, Start bei 0, bis aber nicht inklusive 5 >>> ilist2=range(0, 5, 2) >>> ilist2 [0, 2, 4] Liste: + ist nicht dasselbe wie append (siehe unten) >>> ilist=ilist1+ilist2 >>> ilist [0, 1, 2, 3, 4, 0, 2, 4] Set, Menge >>> set(ilist) set([0, 1, 2, 3, 4])
lists	<pre>[2, 4, 6, 8, 10] ['python', 'is', 'great'] ['27.', 'Oktober', 2011]</pre>	
list assignment and concatenation	<pre>my_list1 = [3, 6, 9, 12] my_list2 = [2, 4, 8, 16] my_long_list = my_list1 + my_list2</pre>	Vorsicht mit Methoden wie append, nicht dasselbe wie + >>> ilist=ilist1.append(ilist2) >>> ilist >>> type(ilist) <type 'NoneType'> >>> ilist1 [0, 1, 2, 3, 4, [0, 2, 4]]
accessing and changing items in a list	<pre>my_item = my_list[2] my_part = my_long_list[:5] my_long_list[3] = 64</pre>	
		>>> dir() ['__builtins__', '__doc__', '__name__', '__package__', 'ilist', 'ilist1', 'ilist2', 'istr', 'isuff', 'word'] >>> locals() {'isuff': 'ba', 'word': 'Alibaba', 'ilist2': [0, 2, 4], '__builtins__': <module '__builtin__' (built-in)>, '__package__': None, 'istr': 'ba', 'ilist': None, '__name__': '__main__', 'ilist1': [0, 1, 2, 3, 4, [0, 2, 4]], '__doc__': None} >>> globals() {'isuff': 'ba', 'word': 'Alibaba', 'ilist2': [0, 2, 4], '__builtins__': <module '__builtin__' (built-in)>, '__package__': None, 'istr': 'ba', 'ilist': None, '__name__': '__main__', 'ilist1': [0, 1, 2, 3, 4, [0, 2, 4]], '__doc__': None}
sorting a list - reverse order - according to length of items	<pre>my_list = sorted(my_list) my_list = sorted(my_list, reverse=True) my_list = sorted(my_list, key=len)</pre>	
inserting an item at the end	<pre>my_list.append(42)</pre>	
inserting an item at the beginning - or at any other place	<pre>my_list.insert(0, 'hello') my_list.insert(3, 'world')</pre>	
deleting a list deleting an item from a list	<pre>del my_list del my_list[3]</pre>	
converting a string to a list - by splitting on the blank symbol	<pre>my_text = 'this is magic' my_list = my_text.split(' ')</pre>	
<pre># initialize the variable #text = 'programmieren macht spass' text = '01234567' # while the position is smaller than the length of the text for position in range(3, len(text)+1): # print the trigram print position, ":", text[position-3:position] print "position", len(text)+1, "wird im range nicht angenommen.\n" \ +"(Bis aber nicht inklusive)"</pre>		
<pre># initialize the list variables with lists of words lexicon = ['computer', 'linguistik', 'fleisch', 'fleischer', \ 'erzeugnis', 'zeugnis'] testwords = ['computerlinguistik', 'computerexperte', \ 'fleischerzeugnis'] # for each word in the list testwords for word in testwords: # split the word in all possible ways # the first element is at least 3 characters long for position in range(3, len(word)): ## assign the split parts to two variables part1 = word[:position] part2 = word[position:] ## print the current split for inspection print part1, part2 ## solution with two nested if conditions if part1 in lexicon: if part2 in lexicon: print 'Found compound ', part1, part2, ':-)' ## solution with a complex condition connected with 'and' if (part1 in lexicon) and (part2 in lexicon): print 'Found compound ', part1, part2, ':-)'</pre>		

```

import sys ## necessary for command line arguments
import codecs
## define a variable for the length of the longest word
max_len = 0
## ... and for the word itself
max_word = ""
## get the filename (to be searched) from the command line
filename = sys.argv[1]
filename="/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Beispiele/programme/person_name_and_gender.txt"
## open the file for 'reading'
infile = codecs.open(filename, 'r', 'utf-8')
## for each line in this file do:
for line in infile:
    my_line_list = line.split(' ')
    print "Space delimiter", my_line_list
    my_line_list = line.split()
    print "Default delimiter set", my_line_list
    for word in my_line_list:
        ## if the current word is longer than max_len:
        if len(word) > max_len:
            ## set max_len to the length of the current word
            max_len = len(word)
            ## store the current word in max_word
            max_word = word
            print max_len, max_word
## close the file
infile.close()
## print the length of the longest word and the word itself
print max_len, max_word

```

```

Space delimiter [u'Tom\tm\n']
Default delimiter set [u'Tom', u'm']
3 Tom
Space delimiter [u'Huck\tm\n']
Default delimiter set [u'Huck', u'm']
4 Huck
Space delimiter [u'Becky\tf\n']
Default delimiter set [u'Becky', u'f']
5 Becky
Space delimiter [u'Sid\tm\n']
Default delimiter set [u'Sid', u'm']
Space delimiter [u'Sidney\tm\n']
Default delimiter set [u'Sidney', u'm']
6 Sidney
Space delimiter [u'Jim\tm\n']
Default delimiter set [u'Jim', u'm']
Space delimiter [u'Ben\tm\n']
Default delimiter set [u'Ben', u'm']
Space delimiter [u'Mary\tf\n']
Default delimiter set [u'Mary', u'f']
Space delimiter [u'Peter\tm\n']
Default delimiter set [u'Peter', u'm']
6 Sidney

```

Tom	m
Huck	m
Becky	f
Sid	m
Sidney	m
Jim	m
Ben	m
Mary	f
Peter	m

Finde Median einer Liste mit gerader Länge

```

>>> ilist= [18, 21, 16, 24, 19, 16, 23, 26]
>>> ilist_sorted=sorted(ilist)
>>> ilist_sorted
[16, 16, 18, 19, 21, 23, 24, 26]
>>> ilen=len(ilist)/2
>>> ilen
4
>>> ilist_sorted[ilen]
21

```

Finde Median einer Liste mit ungerader Länge

```

>>> ilist= [18, 21, 16, 24, 19, 16, 23, 26, 14]
>>> ilist_sorted=sorted(ilist)
>>> ilist_sorted
[14, 16, 16, 18, 19, 21, 23, 24, 26]
>>> ilen=len(ilist)/2
>>> ilen
4
>>> ilist_sorted[ilen]
19

```

```

import codecs
## open the file for 'reading'
infile = codecs.open(filename, 'r', 'utf-8')

```

Overview: File handling

	Examples
Command line arguments	<code>import sys my_filename = sys.argv[1] my_arguments = sys.argv</code>
Open file for reading Open file for writing - my_infile / my_outfile are file objects	<code>my_infile = open(my_filename, 'r') my_outfile = open(my_filename, 'w')</code>
Close file	<code>my_infile.close()</code>
Loop over each line in a file	<code>for my_line in my_infile</code>

```

print(""”Dear %(recipient)s,  
... I wish you to leave Sunnydale and never return.  
... Not Quite Love,  
... %(sender)s  
... ”” % {'sender': 'Buffy the Vampire Slayer', 'recipient': 'Spike'}  
Dear Spike,  
  
I wish you to leave Sunnydale and never return.  
  
Not Quite Love,  
Buffy the Vampire Slayer  
"""

```

Hashes (= Python Dictionaries)

	Examples
Define a hash (= dictionary) - with two entries	<code>my_hash = {} my_hash = {'uni':'noun', 'gute':'adj'}</code>
Enter a key – value pair	<code>my_hash['mit'] = 'praep' my_hash.update({"mit":"praep"})</code>
Check if a key is in the hash - and get the value	<code>my_hash.has_key(my_word) my_word in my_hash my_hash.get(my_word)</code>
Loop over a hash	<code>for my_word in my_hash.keys() for my_word in my_hash</code>

```
>>> idict{'a': '2', 'b': 5} Äquivalente Zugriffe=====
```

```
>>> a_tuple = 1, 2, 3, "four"
>>> a_tuple
(1, 2, 3, 'four')
>>> a_tuple = (1, 2, 3, "four")
>>> a_tuple
(1, 2, 3, 'four')
>>> a_tuple = (1, 2, 3, "four", a_tuple)
>>> a_tuple
(1, 2, 3, 'four', (1, 2, 3, 'four'))
>>> a_tuple = (1, 2, 3, "four", a_tuple)
>>> a_tuple
(1, 2, 3, 'four', (1, 2, 3, 'four'), (1, 2, 3, 'four'))
>>> a_tuple[4][4][3]
'four'
>>> 
```

	Examples
get all keys into a list	my_list = my_hash.keys()
get all values into a list	my_list = my_hash.values()
get the number of keys (= entries)	my_length = len(my_hash)
get all key-value tuples into a list	my_list = my_hash.items()

Hashes (= Python Dictionaries): Typical Usages

	Examples
word + PoS tag	my_hash = {'der':'ART', 'in':'APPR', 'Auto':'NN'}
item + frequency	my_hash = {'der':4, 'in':12, 'Auto':2} my_hash = {'ART':7, 'NN':12, 'APPR':8}
word + translation	my_hash = {'der':'the', 'in':'in', 'Auto':'car', 'Haus':'house'}
string + dummy value (for quick access)	my_hash = {'der':1, 'in':1, 'Auto':1}

List

my_list = []
stores a sequence of values $x_1, x_2, x_3 \dots$
preserves order
quick access via number key to value, but **slow** access to a particular value

Hash

my_hash = {}
stores relations $x_1 \rightarrow y_1, x_2 \rightarrow y_2$
random order !
quick access via key to value

Regular Expressions

Search and Substitution

	Examples
Regular expression matching	import re my_match = re.search('abc', my_text) my_match = re.search('^waa..', my_text) print my_match.group()
Get hit from within a match	my_match = re.search('ab(w+)', my_text) print my_match.group(1)
Get multiple matches in a list	my_matches = re.findall('a w+', my_text)
Find and substitute	my_new_text = re.sub('t', 'th', my_text)

Immer prüfen. Wenn none gibts son

```
Textstring
>>> my_text
'dama abc jalj daapa saama jobab jobabc gamaala aslkjabkd jamaboielkdfjabkiie
abumumabizekluhezbaj'
Alle Matches mit.findall()
>>> my_matches = re.findall('a|w+', my_text)
>>> my_matches
['ama', 'abc', 'alj', 'aapa', 'aama', 'ab', 'abc', 'amaala', 'aslkjlabkd', 'amaboielkdfjabkiie',
'abumumabizekluhezbaj']
Erster Match mit search()
>>> my_match = re.search('abc', my_text)
>>> my_match.group()
'abc'
Erster Match mit search()
>>> my_match2 = re.search('w*aa.', my_text)
>>> my_match2
<_sre.SRE_Match object at 0xb7477db0>
>>> my_match2.group()
'daapa'
Erster Match mit search()
>>> my_match3 = re.search('ab(w+)', my_text)
>>> my_match3.group()
'abc'
>>> my_match3.group(1)
'c'
Finden und Ersetzen (Find and Substitute)
>>> my_new_text = re.sub('abc','xyziaik', my_text)
>>> my_new_text
'dama xyziaik jalj daapa saama jobab jobxyziaik gamaala aslkjabkd jamaboielkdfjabkiie
abumumabizekluhezbaj'
```

Miscellaneous

	Examples
Defining functions - the main function	def my_function(): def main():
Global variables (defined outside a function)	
Remove leading (left) "whitespace" characters	my_line.lstrip()
Remove trailing (right) "whitespace" characters	my_line.rstrip()

Python Tuples (are like Lists but immutable)

	Examples
Define a tuple - with three elements	my_tuple = () my_tuple = ('Junge', 'NN', 'ADJA')
Accessing an element	my_tuple[2]
Looping over the elements of a tuple	for x in my_tuple:
Tuples can be sorted	(2, 4, 100) < (2, 5, 1)

```
>>> itupel=2,5,4,7,1>>> sorted
```

'tuple' object does not support item assignment>>

<h3>Using tuples for sorting values from a hash</h3> <pre> my_hash = {'x':3, 'z':10, 'y':4} my_list = [] for key, value in my_hash.items(): ## create a list of tuples my_list.append((value,key)) print sorted(my_list) </pre>	Dictionary mit Key Value Paar >>> my_hash={'y': 4, 'x': 3, 'z': 10} >>> my_list=[] Speichere Key Value Paare in einer Liste aus Tupeln >>> for key, value in my_hash.items(): ... my_list.append((value, key)) ... >>> my_list [(4, 'y'), (3, 'x'), (10, 'z')] Eine Liste aus Tupeln kann sortiert werden >>> print sorted(my_list) [(3, 'x'), (4, 'y'), (10, 'z')]
<h3>Python Functions</h3> <p>We can define and call functions with arguments.</p> <pre> def my_print_function(freq): for x in range(0, freq): print x, 'Hello world' print 'And now we write some lines' my_print_function(40) </pre>	<pre> ## open a file for reading infile = open(filename, 'r') ## create a name for the first output file outfilename_1 = filename + '.a.words.txt' ## open the first file for writing outfile_a = open(outfilename_1, 'w') ## create a name for the second output file outfilename_2 = filename + '.b.words.txt' ## open the second file for writing outfile_b = open(outfilename_2, 'w') ## for each line in the input file for line in infile: # split the line into a list of strings line_list = line.split() # for each word in the line for word in line_list: # remove a comma at the end of the word if word[-1] == ',' or word[-1] == '.': word = word[0:-1] # check if the word starts with 'a' if word[0] == 'a' and len(word)>10: print word # write output to the first file outfile_a.write(word) outfile_a.write('\n') # check if the word starts with 'b' elif word[0] == 'b' and len(word)>10: print word # write output to the second file outfile_b.write(word) outfile_b.write('\n') ## close the files infile.close() </pre>
<pre> >>> istr="asdine" >>> istr[-1] 'e' >>> istr[-1:] 'e' >>> istr[-5:] 'sdine' >>> istr[-5] 's' </pre>	<p>Liste aus drei Tupeln</p> <pre> >>> my_list [(4, 'y'), (3, 'x'), (10, 'z')] Die letzten zwei Tupel in einer Liste >>> my_list[-2:] [(3, 'x'), (10, 'z')] Das zweitletzte Tupel >>> my_list[-2] (3, 'x') Das erste Tupel der zwei letzten Tupel >>> my_list[-2:][1] (10, 'z') ...und davon das erste Element >>> my_list[-2:][1][1] 'z' Das letzte Tupel der letzten zwei Tupel >>> my_list[-2:][-1] (10, 'z') ...als Liste >>> my_list[-2:][-1] [(10, 'z')] ...und Zugriff auf das zweitletzte Element des Ersten Tupels >>> my_list[-2:][-1][0][2] 10 </pre>

```

#!/usr/bin/python

# program to find compound words in a given text
# - assumes the words from the text itself to be the lexicon
# by Martin Volk

import sys ## necessary for command line arguments
def main():
    ## define a list for the lexicon
    lexicon = []
    ## define a type counter
    type_counter = 0
    ## get the filename
    filename = sys.argv[1]
    filename = "/media/benzro/OS/Users/benzro/Documents/LubuntuShared/" + "8) PCL-1/Beispieleprogramme/SecondVariety.txt"
    ## open the input file for reading
    infile = open(filename, 'r')
    ## for each line in the input file
    for line in infile:
        # split the line into a list of strings
        line_list = line.split()
        # for each word in the line
        for word in line_list:
            # if the word is not in the lexicon, then add it
            if word not in lexicon:
                lexicon.append(word)
                # count the words in the lexicon
                type_counter += 1
    ## close the file
    infile.close()
#####

```

```

## assumes the person names to be in the format: name\tgender
##
## by Martin Volk

import sys ## necessary for command line arguments
## define a hash as a global variable
my_person_hash = {}
## define a function that reads the person names from file <person_name_file>
## and stores them in the global hash
def read_and_store_person_names():
    # get the filename of the person name file
    filename = sys.argv[1]
    filename = "/media/benzro/OS/Users/benzro/Documents/LubuntuShared/" + "8) PCL-1/Beispieleprogramme/person_name_and_gender.txt"
    ## open the person name file for reading
    infile = open(filename, 'r')
    ## for each line in the file
    for line in infile:
        # strip away the \n symbol
        clean_line = line.rstrip()
        # split the line into a list of strings
        line_list = clean_line.split('\t')
        name = line_list[0]
        gender = line_list[1]
        # save the name --> gender pair in a hash
        my_person_hash[name] = gender
        # add the genitive form of the name
        # example: Tom --> Toms
        name = name + 's'
        # save the genitive name --> gender pair in the hash
        my_person_hash[name] = gender
#####

```

```

# program to do word-by-word translation based on a bilingual dictionary
# by Martin Volk

# define a hash (= dictionary)
my_lexicon = {'Tante':'aunt',
              'Polly':'Polly',
              'sah':'saw',
              'Mann':'man',
              'mit':'with',
              'Fernglas':'telescope',
              'im':'in the',
              'Garten':'garden',
              'dem':'the',
              'der':'the',
              'den':'the' }

test_sentence = 'Tante Polly sah den Mann mit dem Fernglas'
test_sentence = 'der Mann sah die Tante Polly mit dem Fernglas'
# test_sentence = 'Tante Polly sah den Mann im Garten'
test_sentence = 'Tante Polly sah den Mann im Garten mit dem Fernglas'
print 'Input sentence:', test_sentence
# convert the sentence into a list
test_list = test_sentence.split()
print 'Output sentence:', 
# work through the list and translate each word
for word in test_list:
    ## if the word is in the lexicon as key
    if word in my_lexicon:
        print my_lexicon[word],
    else:
        print '\n>>', word, ' is unknown!!! :-('
print
print '-----'

```

```

#####
## define a compound counter
compound_counter = 0
## check for each word in the lexicon
## if it is composed of two other words from the lexicon
for word in lexicon:
    # split the word in all possible ways
    # the first element must be at least 3 characters long
    # and the last element must be at least 3 characters long
    for position in range(3, len(word)-2):
        ## check if both parts are in the lexicon
        if word[0:position] in lexicon and word[position:] in lexicon:
            print 'Found compound ', word[0:position], word[position:]
            compound_counter += 1
print
print 'I collected', type_counter, 'different words for the lexicon.'
print 'I found', compound_counter, 'compounds ;-)'
## This is the standard boilerplate to call the function main
if __name__ == '__main__':
    main()

```

```

#####
def main():
    # initialize the counters
    person_counter = 0
    gender_counter = 0
    # read the person names into a hash
    read_and_store_person_names()
    ## get the filename of the text file
    filename = sys.argv[2]
    filename = "/media/benzro/OS/Users/benzro/Documents/LubuntuShared/" + "8) PCL-1/Beispieleprogramme/SecondVariety.txt"
    ## open the file for reading
    infile = open(filename, 'r')
    ## for each line in the file
    cnt=0
    for line in infile:
        cnt+=1
        # split the line into a list of strings
        line_list = line.split()
        # for each word in the line
        for word in line_list:
            # if the word is in the person name hash
            if word in my_person_hash:
                gender = my_person_hash[word]
                # print a line of success and joy :-)
                print cnt,'Found name ', word, ' gender ', gender
                # count the number of found name occurrences
                person_counter += 1
                # count the gender occurrences; get(key[, default])
                gender_counter[gender] = gender_counter.get(gender, 0) + 1
    ## close the file
    infile.close()

```

default value fallskey

Encoding einer Datei bestimmen mit file command
benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzr
ed/8) PCL-1/Vorlesung/vl07-vorspann\$ file ***
AE-l1-encoded-as-l1.txt: ISO-8859 text
AE-l1-encoded-as-utf8.txt: UTF-8 Unicode text
AE-l1-encoded-as-utf8.txt=: UTF-8 Unicode text
ae-l1.txt: ISO-8859 text
ae-l1.txt: ISO-8859 text
AE-utf-8-encoded-as-l1.txt: very short file (no magic)
AE-utf-8-encoded-as-utf8.txt: very short file (no magic)

Von L1 nach utf-8 konvertieren
\$ iconv -f iso-8859-1 -t utf8 ae-l1.txt > ae-utf8.txt
Von utf-8 nach L1 konvertieren
\$ iconv -f utf8 -t iso-8859-1 ae-utf8.txt > ae-l1_test.txt
Von utf-8 nach ascii konvertieren
es wird eine Fehlermeldung angezeigt, da das Zeichen ä nicht konvertiert werden kann
\$ iconv -f utf8 -t ascii ae-utf8.txt > ae-ascii.txt
iconv: illegal input sequence at position 1

\$ file *
ae-utf8.txt: UTF-8 Unicode text
ae-ascii.txt: very short file (no magic)
ae-l1_test.txt: ISO-8859 text

Historically a \n was used to move the carriage down, while the \r was used to move the carriage back to the left side of the page.

Two different characters.

\n is used as an end-of-line terminator in Unix text files

\r is used as an end-of-line terminator in Mac text files

\r\n (ie both) are used to terminate lines in Windows and DOS text files.

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

# simon.clematide@uzh.ch
# PCL I

import codecs

#Datei ae-l1.txt Western(ISO-8859-15)
#a@ mit utf-8; a\E4 mit Gedit

#http://unicode-table.com/de/#latin-1-supplement
#http://www.utf8-zeichentabelle.de/
#https://de.wikipedia.org/wiki/ISO_8859-1
print "\nHexadezimale Representation des Speicherinhalts für:"
print "\u00e4: IS08859-1=E4; Unicode=U+00E4(228); utf-8=c3 a4(11000011 10
print "\u00c4: IS08859-1=C4; Unicode=U+00C4(196); utf-8=c3 84(11000011 10

print "\u00d6: IS08859-1=FC; Unicode=U+00FC(252); utf-8=c3 bc"
print "\u00d6: IS08859-1=DC; Unicode=U+00DC(220); utf-8=c3 9c"
print "\u00d6: IS08859-1=F6; Unicode=U+00F6(246); utf-8=c3 b6"
print "\u00d6: IS08859-1=D6; Unicode=U+00D6(214); utf-8=c3 96"

#Encode-Decode
#-----
#Decode: Typ String --> Typ Unicode
#Encode: Typ Unicode --> Typ String
```

```
# Decode from l1 encoded file into unicode strings
f = open("./ae-l1.txt", "r")
f1 = codecs.open("./ae-l1.txt", "r", "l1")
f2 = codecs.open("./ae-l1.txt", "r", "utf-8")

# Encode unicode strings into UTF-8 or l1 encoded file
g = open("./AE-encoded.txt", "w")
g1 = codecs.open("./AE-l1-encoded-as-l1.txt", "w", "l1")
g2 = codecs.open("./AE-l1-encoded-as-utf8.txt", "w", "utf-8")
g3 = codecs.open("./AE-utf-8-encoded-as-l1.txt", "w", "l1")
g4 = codecs.open("./AE-utf-8-encoded-as-utf8.txt", "w", "utf-8")

print "\nDatei ae-l1.txt hat Western(ISO-8859-15)= l1 Kodierung\
und enthält die Zeichen a\E4 (Gedit), also ää; a@ mit utf-u8 \n"

print 'f ist vom Typ String und übernimmt a\E4. E4 wird jedoch' \
'nicht erkannt von der upper() Methode und ignoriert'
print 'g schreibt A\E4 in Datei AE-encoded.txt' \
'(ISO-8859 text) gemäss file *\n' \
'(Es erscheint beim Öffnen: A@ mit utf-8; A\E4 mit Gedit)'
for line in f: #f thought input was utf-8
    g.write(line.upper())
    print "f Type:", type(line) #output as l1
    print "f Canonical repr(line.upper()):", "==>",repr(line.upper())
    print "f Canonical repr(line):", "==>",repr(line), "<=="
    print "f Printed line.upper():", "==>",line.upper(), "<=="
    print "f Printed line:", "==>",line, "<=>\n"


```

```
print 'f1 ist von Typ unicode. f1 wusste, dass Input l1 kodiert' \
'ist, E4 wird von upper() Methode erkannt und auf C4 geändert'
print 'g1 schreibt A\C4 in Datei AE-l1-encoded-as-l1.txt' \
'(ISO-8859 text) gemäss file *\n' \
'(Es erscheint beim Öffnen: A@ mit UTF-8; A\C4 mit Gedit)'
print 'g2 schreibt ÄÄ in Datei AE-l1-encoded-as-utf8.txt' \
'(UTF-8 Unicode text) gemäss file *\n' \
'(Es erscheint beim Öffnen: ÄÄ mit UTF-8; ÄÄ mit Gedit)'
for line in f1: #f1 knew input was l1
    g1.write(line.upper())
    g2.write(line.upper()) #output as utf-8
    print "f1 Type:", type(line) #output as l1
    print "f Canonical repr(line.upper()):", "==>",repr(line.upper())
    print "f Canonical repr(line):", "==>",repr(line), "<=="
    print "f Printed line.upper():", "==>",line.upper(), "<=="
    print "f Printed line:", "==>",line, "<=>\n"
```

```
print 'f2 ist von Typ unicode. f2 dachte Input sei utf-8 und' \
'ignorierte das unbekannte ä'
print 'g3 schreibt A in Datei AE-utf-8-encoded-as-l1.txt' \
'(very short file (no magic)) gemäss file *\n' \
'(Es erscheint beim Öffnen: A mit UTF-8; A mit Gedit)'
print 'g4 schreibt Ä in Datei AE-utf-8-encoded-as-utf8.txt' \
'(very short file (no magic)) gemäss file *\n' \
'(Es erscheint beim Öffnen: A mit UTF-8; A mit Gedit)'
for line in f2: #f2 thought input was utf-8
    g3.write(line.upper())
    g4.write(line.upper()) #output as utf-8
    print "f2 Type:", type(line) #output as l1
    print "f2 Canonical repr(line.upper()):", "==>",repr(line.upper())
    print "f2 Canonical repr(line):", "==>",repr(line), "<=="
    print "f2 Printed line.upper():", "==>",line.upper(), "<=="
    print "f2 Printed line:", "==>",line, "<=>\n"

print type(f), type(f1), type(f2),

f1.close()
g1.close()
f2.close()
g2.close()
```

```
Hexadezimale Representation des Speicherinhalts für:
ä: IS08859-1=E4; Unicode=U+00E4(228); utf-8=c3 a4(11000011 10
10100100)
Ä: IS08859-1=C4; Unicode=U+00C4(196); utf-8=c3 84(11000011 10
10000100)
ü: IS08859-1=FC; Unicode=U+00FC(252); utf-8=c3 bc
Ü: IS08859-1=DC; Unicode=U+00DC(220); utf-8=c3 9c
ö: IS08859-1=F6; Unicode=U+00F6(246); utf-8=c3 b6
Ö: IS08859-1=D6; Unicode=U+00D6(214); utf-8=c3 96

Datei ae-l1.txt hat Western(ISO-8859-15)= l1 Kodierung und
enthält die Zeichen a\E4 (Gedit), also ää; a@ mit utf-u8

f ist vom Typ String und übernimmt a\E4. E4 wird jedoch nicht
erkannt von der upper() Methode und ignoriert
g schreibt A\E4 in Datei AE-encoded.txt (ISO-8859 text)
gemäß file *
(Es erscheint beim Öffnen: A@ mit utf-8; A\E4 mit Gedit)
f Type: <type 'str'>
f Canonical repr(line.upper()): ==> 'A\xe4\n' <==
f Canonical repr(line): ==> 'a\xe4\n' <==
f Printed line.upper(): ==> A@<==>
<==>
f Printed line: ==> a@<==>
<==>
```

```
f1 ist von Typ unicode. f1 wusste, dass Input l1 kodiert ist,
E4 wird von upper()Methode erkannt und auf C4 geändert
g1 schreibt A\c4 in Datei AE-l1-encoded-as-l1.txt (ISO-8859
text) gemäss file *
(Es erscheint beim Öffnen: Ä mit UTF-8; A\c4 mit Gedit)
g2 schreibt ÄÄ in Datei AE-l1-encoded-as-utf8.txt (UTF-8
Unicode text) gemäss file *
(Es erscheint beim Öffnen: ÄÄ mit UTF-8; ÄÄ mit Gedit)
f1 Type: <type 'unicode'>
f Canonical repr(line.upper()): ==> u'A\xc4\n' <==
f Canonical repr(line): ==> u'a\xe4\n' <==
f Printed line.upper(): ==> ÄÄ
<==
f Printed line: ==> äÄ
<==

f2 ist von Typ unicode. f2 dachte Input sei utf-8 und
ignorierte das unbekannte ä
g3 schreibt A in Datei AE-utf-8-encoded-as-l1.txt (very short
file (no magic)) gemäss file *
(Es erscheint beim Öffnen: A mit UTF-8; A mit Gedit)
g4 schreibt A in Datei AE-utf-8-encoded-as-utf8.txt (very
short file (no magic)) gemäss file *
(Es erscheint beim Öffnen: A mit UTF-8; A mit Gedit)
f2 Type: <type 'unicode'>
f Canonical repr(line.upper()): ==> u'A' <==
f Canonical repr(line): ==> u'a' <==
f Printed line.upper(): ==> A <==
f Printed line: ==> a <==

<type 'file'> <type 'instance'> <type 'instance'>
```

```
#Versuch den String als utf-8 interpretiert(encoded) darzustellen
print "2:", repr(v), str(v)
#'\xc4pple' @apple
#aber im Interpreter:
#>>> str('\xc4pple')
#'\xc4pple'
#>>> print str('\xC4pple')
#@apple

print "3:", repr(u)
#u'\xc4pple'

#Versuch aus einem Typ unicode einen String darzustellen
#scheitert daran, dass ascii ein Problem hat.
#print "3:",str(u)
# UnicodeEncodeError: 'ascii' codec can't encode character
# u'\xc4' in position 0: ordinal not in range(128)

#Versuch den String als utf-8 interpretiert(encoded) darzustellen
print "4:", v
#@apple

#Der Typ unicode wird korrekt verarbeitet
print "5:", u
#@Äpple

print "11:", tell_me_about(uv2)
#(<type 'unicode'>, u'\xc3\x84pple')

#decode String (in utf-8) to unicode with utf-8 table
#Es wird ein Typ unicode erzeugt mit unicode (und L1) Kodierung
uv3 = "Äpple".decode("utf-8")
print "12:", repr(uv3)
#u'\xc4pple'      #c4 ist Ä in unicode (und L1)

print "13:", tell_me_about(uv3)
#(<type 'unicode'>, u'\xc4pple')

#A little more illustration - with "Ä"
#-----
#Typ unicode ist gleich Typ unicode
print "14:", u"Ä" == u"\xc4", repr(u"Ä"), repr(u"\xc4")
#True u'\xc4' u'\xc4'

#Typ String in utf-8 ist nicht gleich Typ unicode
print "15:", "Ä" == u"\xc4", repr("Ä"), repr(u"\xc4")
#False 'u'\xc3\x84' u'\xc4'

#Typ String in utf-8 ist gleich Typ unicode im Ascii Bereich
print "16:", "a" == u"a", repr("a"), repr(u"a")

#Typ String (in utf-8) in Typ unicode dekodiert ist gleich Typ unic
print "17:", "Ä".decode('utf8') == u"\xc4", \
      repr("Ä".decode('utf-8')), repr(u"\xc4")
#True u'\xc4' u'\xc4'
```

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

#ohne die coding Zeile gibt es eine Fehlermeldung, wenn diese
#Datei einen Umlaut enthält (wie in enthält)
#(Line 4 SyntaxError: Non-ASCII character '\xc3' in file)
#Roland Benz

#Die unicode, L1 und utf-8 Tabellen enthalten die Ascii Tabelle
#Die unicode Tabelle enthält die L1 Tabelle (oder zumindest Teile d
#Die unicode und utf-8 Tabellen ähneln zwischen den Zeichen 128-191
#utf-8 hat in diesem Bereich lediglich eine c2 oder c3 davor.

def tell_me_about(s): return (type(s), s)

#A plain string
#-----
v = "\xc4pple" # type string
u = u'\xc4pple' # type unicode

print "0:", tell_me_about(v)
#(<type 'str'>, '\xc4pple')

print "1:", tell_me_about(u)
#(<type 'unicode'>, u'\xc4pple')
```

```
#Decoding a iso8859-1 string - convert plain string to unicode
#-----
#decode string to unicode with L1 table
uv = v.decode('iso-8859-1') #benutzt L1 Tabelle
print "6:", repr(uv)
#u'\xc4pple'      #c4 ist Ä in unicode (und L1)

print "7:", tell_me_about(uv)
#(<type 'unicode'>, u'\xc4pple')

#Der Typ unicode wird korrekt verarbeitet
print "8:", uv, "\xC4pple".decode("iso-8859-1")
#Äpple Äpple

print "9:", v.decode('iso-8859-1') == u'\xc4pple'
#True

#Versuch den String mit utf-8 Tabelle zu decoden schlägt fehl
#uv1 = v.decode("utf-8")
#print "9:", repr(uv)
# return codecs.utf_8_decode(input, errors, True)
# UnicodeDecodeError: 'utf8' codec can't decode byte 0xc4
# in position 0: invalid continuation byte

#decode String (in utf-8) to unicode with (useless?) L1 table
#Es wird ein Typ unicode erzeugt mit utf-8 Kodierung
uv2 = "Äpple".decode("iso-8859-1")
print "10:", repr(uv2), repr("Äpple")
#u'\xc3\x84pple'      #c3 84 ist utf-8 für Ä
```

```
#Typ String (in utf-8) ist ungleich Typ String (in L1, unicode)
print "18:", "Ä" == "\xc4"
#False

#Encode-Decode
#-----
#Decode: Typ String --> Typ Unicode
#Encode: Typ Unicode --> Typ String

#Encoding to UTF
#-----
#v = "\xc4pple" # type string
# convert Type String (in iso-8859-1) to Type unicode
# and to Type String (in utf-8)
u8 = v.decode('iso-8859-1').encode('utf-8')
print "19:", repr(u8), str(u8), u8
# 'u'\xc3\x84pple' Äpple Äpple

print "20:", tell_me_about(u8)
#(<type 'str'>, '\xc3\x84pple')

#Convert to String in utf-16
u16 = v.decode('iso-8859-1').encode('utf-16')
print "21:", tell_me_about(u16), str(u16), u16
#(<type 'str'>, '\ufffe\xc4\x00p\x00p\x001\x00e\x00')

# 000 pple 000 pple

#Convert to Type unicode
print "22:", tell_me_about(u8.decode('utf8'))
#(<type 'unicode'>, u'\xc4pple')
```

<pre>#Convert to Type unicode print "23:", tell_me_about(u16.decode('utf16')) #(<type 'unicode'>, u'\xc4pple') #Relationship between unicode and UTF and latin1 #----- #Typ String in utf-8 print "24:", u8 #Apple #Typ unicode print "25:", u8.decode('utf-8') # printing unicode #Apple #Typ String in utf-16, print statement interprets u16 as utf-8 print "26:", u16 # printing 'bytes' of u16 #00c4pple #Typ unicode print "27:", u16.decode('utf16') #Apple # printing unicode #Typ String in L1 vs Typ String in utf-8 print "28:", v == u8 #False # v is a iso8859-1 string; u8 is a utf-8 string #Typ unicode vs Typ String in utf-8 print "29:", v.decode('iso8859-1') == u8 #False # v.decode(...) returns unicode</pre>	<pre>#Strings decode to the same representation in Typ unicode print "30:", u8.decode('utf-8') == v.decode('latin1') == u16.decode(#True # all decode to the same unicode memory representation # (latin1 is iso-8859-1) #Unicode Exceptions #----- #Versuch Typ String in utf-8 direkt nach Typ String in l1 zu kodieren print "31:", u8.encode('iso8859-1') #Traceback (most recent call last): # File "<stdin>", line 1, in <module> #UnicodeDecodeError: 'ascii' codec can't decode byte 0xc3 in position # ordinal not in range(128) print "32:", u16.encode('iso8859-1') #Traceback (most recent call last): # File "<stdin>", line 1, in <module> #UnicodeDecodeError: 'ascii' codec can't decode byte 0xff in position # ordinal not in range(128) v.encode('iso8859-1') #Traceback (most recent call last): # File "<stdin>", line 1, in <module> #UnicodeDecodeError: 'ascii' codec can't decode byte 0xc4 in position # ordinal not in range(128)</pre>
<pre>0: (<type 'str'>, '\xc4pple') 1: (<type 'unicode'>, u'\xc4pple') 2: '\xc4pple' Äpple 3: u'\xc4pple' 4: Äpple 5: Äpple 6: u'\xc4pple' 7: (<type 'unicode'>, u'\xc4pple') 8: Äpple Äpple 9: True 10: u'\xc3\x84pple' '\xc3\x84pple' 11: (<type 'unicode'>, u'\xc3\x84pple') 12: u'\xc3\x84pple' 13: (<type 'unicode'>, u'\xc4pple') 14: True u'\xc4' u'\xc4' 15: False '\xc3\x84' u'\xc4' 16: True 'a' u'a' 17: True u'\xc4' u'\xc4' 18: False 19: '\xc3\x84pple' Äpple Äpple 20: (<type 'str'>, '\xc3\x84pple') 21: (<type 'str'>, '\xff\xfe\x40\x00p\x00p\x00l\x00e\x00') 000 p p l e 000 p p l e 22: (<type 'unicode'>, u'\xc4pple') 23: (<type 'unicode'>, u'\xc4pple') 24: Äpple 25: Äpple 26: 000 p p l e 27: Äpple 28: False 29: False 30: True 31:</pre>	<pre>#!/usr/bin/python # -*- coding: utf-8 -*- # simon.clematide@uzh.ch # PCL I # Representation of strings in UTF-8 encoded files print "Length of 'a':", len('a'), "Canonical:", repr('a') print "Length of 'ä':", len('ä'), "Canonical:", repr('ä') print "Length of u'a':", len(u'a'), "Canonical:", repr(u'a') print "Length of u'ä':", len(u'ä'), "Canonical:", repr(u'ä') </pre> <p style="background-color: #e0e0ff; padding: 5px;">/str_representation_utf8.py"</p> <pre>Length of 'a': 1 Canonical: 'a' Length of 'ä': 2 Canonical: '\xc3\x84' Length of u'a': 1 Canonical: u'a' Length of u'ä': 1 Canonical: u'\xe4' Process finished with exit code 0</pre>
<pre>#!/usr/bin/python # -*- coding: iso-8859-1 -*- # simon.clematide@uzh.ch # PCL I # Representation of in latin-1 encoded files print "length of 'a':", len('a'), "Canonical", repr('a') print "length of 'ä':", len('ä'), "Canonical", repr('ä') print "length of u'a':", len(u'a'), "Canonical", repr(u'a') print "length of u'ä':", len(u'ä'), "Canonical", repr(u'ä')</pre>	<pre>length of 'a': 1 Canonical 'a' length of 'ä': 1 Canonical '\xe4' length of u'a': 1 Canonical u'a' length of u'ä': 1 Canonical u'\xe4' Process finished with exit code 0</pre>

```

#!/usr/bin/python
# -*- coding: utf-8 -*-

# simon.clematide@uzh.ch
# PCL I

# Data type unicode

# Einzelige (!) Zeichenkette
# mit Escape-Sequenzen
s1 = u'\u20aca\n\xe4'
print "Canonical s1: ",repr(s1)
print "Printed s1: ",s1

# Rohe Sequenz ur"..."
# \uNNNN werden aufgelöst!
s2 = ur'\u20aca\n\xe4'
print "Canonical s2: ",repr(s2)
print "Printed s2: ",s2

# Longstring
s3 = u"""\u20aca
a"""
print "Canonical s3: ",repr(s3)
print "Printed s3: ",s3

# Roher Longstring
s4 = ur"""€a
\u00e4"""
print "Canonical s4: ",repr(s4)
print "Printed s4: ",s4

print
print "Type of s1:", type(s1)
print "Type of s2:", type(s2)
print "Type of s3:", type(s3)
print "Type of s4:", type(s4)

print
print "Test: s1 == s2 Result:", s1 == s2
print "Test: s1 == s3 Result:", s1 == s3
print "Test: s2 == s4 Result:", s2 == s4
print "Test: s3 == s4 Result:", s3 == s4

```

```

/unicode_literals.py"
Canonical s1: u'\u20aca\n\xe4'
Printed s1: €a
ä
Canonical s2: u'\u20aca\n\xe4'
Printed s2: €a\n\xe4
Canonical s3: u'\u20aca\n\xe4'
Printed s3: €a
ä
Canonical s4: u'\u20aca\n\xe4'
Printed s4: €a
ä

Type of s1: <type 'unicode'>
Type of s2: <type 'unicode'>
Type of s3: <type 'unicode'>
Type of s4: <type 'unicode'>

Test: s1 == s2 Result: False
Test: s1 == s3 Result: True
Test: s2 == s4 Result: False
Test: s3 == s4 Result: True

```

```

#!/usr/bin/python
# -*- coding: utf-8 -*-

# simon.clematide@uzh.ch
# PCL I

# Data type str

# Einzelige (!) Zeichenkette
# mit Escape-Sequenzen
s1 = "a\n\x61"
print "Canonical s1: ",repr(s1)
print "Printed s1: ",s1

# Rohe Sequenz r"..."
# ohne Escapes
s2 = r"a\n\x61"
print "Canonical s2: ",repr(s2)
print "Printed s2: ",s2

# Longstring
s3 = """a
a"""
print "Canonical s3: ",repr(s3)
print "Printed s3: ",s3

```

```

/LubuntuShared/8) PCL-1/Vorlesung/vl07-vorspann/str_literals
.py"
Canonical s1: 'a\na'
Printed s1: a
a
Canonical s2: 'a\n\x61'
Printed s2: a\n\x61
Canonical s3: 'a\na'
Printed s3: a
a
Canonical s4: 'a\na'
Printed s4: a
a

Type of s1: <type 'str'>
Type of s2: <type 'str'>
Type of s3: <type 'str'>
Type of s4: <type 'str'>

Test: s1 == s2 Result: False
Test: s1 == s3 Result: True
Test: s2 == s4 Result: False
Test: s1 == s4 Result: True

Process finished with exit code 0

```

```

# Rohe Longstring
s4 = r""""a
a"""
print "Canonical s4: ",repr(s4)
print "Printed s4: ",s4

print
print "Type of s1:", type(s1)
print "Type of s2:", type(s2)
print "Type of s3:", type(s3)
print "Type of s4:", type(s4)

print
print "Test: s1 == s2 Result:", s1 == s2
print "Test: s1 == s3 Result:", s1 == s3
print "Test: s2 == s4 Result:", s2 == s4
print "Test: s1 == s4 Result:", s1 == s4

```

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

import re

text = u"Viele Köche verderben den Brei."
pattern_ur = ur"(\w+)"
pattern_u = u"(\w+)"
pattern_r = r"(\w+)"
pattern = "(\\w+)"
print "\n", pattern_ur, pattern_u, pattern_r, pattern, "\n"
print repr(pattern_ur), repr(pattern_u), repr(pattern_r), \
    repr(pattern), "\n"
# Alle Matches finden
m_ur = re.findall(pattern_ur, text)
m_u = re.findall(pattern_u, text)
m_r = re.findall(pattern_r, text)
m = re.findall(pattern, text)
for g in m_ur:
    print g
print ""
for g in m_u:
    print g
print ""
for g in m_r:
    print g
print ""
for g in m:
    print g
```

```
(\w+) (\w+) (\w+) (\w+)
u'(\w+') u'(\w+') '(\w+') '(\w+')

Viele
K
che
verderben
den
Brei

Viele
K
che
verderben
den
Brei

Viele
K
che
verderben
den
Brei

Viele
K
che
verderben
den
Brei
```

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

# simon.clematide@uzh.ch
# PCL I

import re

text = u"Viele Köche verderben den Brei."
pattern = ur"(\w+)" # Das Flag (?u) aktiviert Unicode-Kategorien fuer \w und \b

# Resultat ist eine Liste
m = re.findall(pattern, text)

for s in m:
    print s
```

```
/re.findall_flag_u.py  
Viele  
Köche  
verderben  
den  
Brei
```

```
#!/usr/bin/python
# -*- coding: utf-8 -*-
import re

text = "Hässliche Köche verdürben das Gebräu" #utf8
text2 = u"Hässliche Köche verdürben das Gebräu" #unicode

pattern = r"([aeioäöü]+)" #utf-8
pattern2 = ur"([aeioäöü]+)" #unicode

print "\n",repr(text),repr(pattern),repr(text2),repr(pattern2)

# Im Ersetzungstext können gematchte Gruppen eingefügt werden.
# # N (N ist die N-te gruppierende Klammer im Pattern)
replacement = r"\[1]"
replacement2 = ur"[\\1]"

print "\n",repr(replacement),repr(replacement2),"\n"

#pattern und text müssen beide Typ String oder beide Typ unicode
#sein, um die Umlaute zu finden
print re.sub(pattern, replacement, text) #alle
print re.sub(pattern, replacement, text2) #ohne äöü
print re.sub(pattern, replacement2, text) #alle
print re.sub(pattern2, replacement, text) #ohne äöü
print re.sub(pattern2, replacement2, text2) #ohne äöü
print re.sub(pattern2, replacement, text2) #alle
print re.sub(pattern2, replacement2, text2) #alle
```

```
'H\xc3\xaa4ssliche K\xc3\xb6ch[e] verd\xc3\xbcrb[e]n d[a]s G[e]br[\xe4]u  
Gebr\xc3\xaa4u' '([aeio\xc3\xa4\xc3\xb6\xc3\xbcl]+)'  
u'H\xe4sslische K\xf6ch[e] verd\xfcrb[e]n das Gebr\xe4u'  
u'([aeio\xe4\xf6\xfc]+)'  
  
'[\x11]' u'[\x11]'  
  
H[\xe4]ssl[i]ch[e] K[\x80]ch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]br[\xe4]u  
Hässl[i]ch[e] Köch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]bräu  
H[\xe4]ssl[i]ch[e] K[\x80]ch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]br[\xe4]u  
Hässl[i]ch[e] Köch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]bräu  
H[\xe4]ssl[i]ch[e] Köch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]bräu  
H[\xe4]ssl[i]ch[e] Köch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]bräu  
H[\xe4]ssl[i]ch[e] K[\x80]ch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]br[\xe4]u  
H[\xe4]ssl[i]ch[e] K[\x80]ch[e] v[e]rd\xfcl[ü]rb[e]n d[a]s G[e]br[\xe4]u
```

text und pattern müssen dieselbe Kodierung

Pattern immer mitr"" für text ""oderur"" für te-

```
#!/usr/bin/python
# -*- coding: utf-8 -*-

# simon.clematide@uzh.ch
# PCL I

import re

#Vorsicht: Innerhalb einer Klasse [.] hat der Punkt die
#Bedeutung eines Punktes und nicht jedes Zeichens außer \n
text = 'my email is simon.clematide@uzh.ch and yours?'
pattern = r'([\w.]+)@([\w.]+)'

m = re.search(pattern, text)

# get whole match
print "Whole match:", m.group()

# get the first matched group
print "First group:", m.group(1)

# get the second matched group
print "Second group:", m.group(2)

# What is m?
print m
```

I need to extract all letters after the + sign or at the beginning of a string like this:

```
formula = "X+BC+DAF"
```

I tried so, and I do not want to see the + sign in the result. I wish see only ['X', 'B', 'D'] .

```
>>> re.findall("^[A-Z]+|[+][A-Z]", formula)
['X', 'B', '+D']
```

When I grouped with parenthesis, I got this strange result:

```
re.findall("^(?:[A-Z])|([+])([A-Z])", formula)
[('X', ''), ('+', 'B'), ('+', 'D')]
```

Why it created tuples when I try to group ? How to write the regexp directly such that it returns ['X', 'B', 'D'] ?

```
/re_search_group.py"
Whole match: simon.clematide@uzh.ch
First group: simon.clematide
Second group: uzh.ch
<sre.SRE_Match object at 0xb74a2f08>
```

If there are any capturing groups in the regular expression then `re.findall` returns only the values captured by the groups. If there are no groups the entire matched string is returned.

```
re.findall(pattern, string, flags=0)
```

Return all non-overlapping matches of pattern in string, as a list of strings. The string is scanned left-to-right, and matches are returned in the order found. If one or more groups are present in the pattern, return a list of groups; this will be a list of tuples if the pattern has more than one group. Empty matches are included in the result unless they touch the beginning of another match.

How to write the regexp directly such that it returns ['X', 'B', 'D'] ?

Instead of using a capturing group you can use a non-capturing group:

```
>>> re.findall(r"(?:^|[^+])([A-Z])", formula)
['X', 'B', 'D']
```

Or for this specific case you could try a simpler solution using a word boundary:

```
>>> re.findall(r"\b[A-Z]", formula)
['X', 'B', 'D']
```

Or a solution using `str.split` that doesn't use regular expressions:

```
>>> [s[0] for s in formula.split('+')]
['X', 'B', 'D']
```

```
# -*- coding: utf-8 -*-
# simon.clematide@uzh.ch
# PCL I

import re
#Achtung beim Bilden von Gruppen (). Gibt es capturing groups ()
#dann werden nur die Werte der capturing groups zurückgegeben.
#Sollen alle Matches zurückgegeben werden müssen alle Gruppen
#non capturing groups (?: ) sein.
text = "That U.S.A. poster-print costs $12.40... per m^2."
pattern = r'''(?x)
            # set flag (?x) to allow verbose regexes
           (?:[A-Z]\.)*          # abbreviations, e.g. U.S.A.
            | \$?\d+(?:[,]\d+)*%? # currency/percentages, $12.40, 82%
            | \w+(?:-\w+)*        # words with optional internal hyphens
            | \.\.\.               # ellipsis
            | [.,;:]*              # punctuation
            | \S+                  # catch-all for non-layout characters
            ...'''

m = re.findall(pattern, text)
print m
```

```
/re.findall_tokenizer.py"
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...', 'per', 'm', '^2.1']
```

Interne Hilfe

```
>>> help(len)
Help on built-in function len in module __builtin__:

len(...)
    len(object) -> integer

    Return the number of items of a sequence or collection.
```

Erklärung der Erklärung

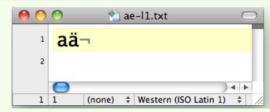
- ▶ Die Signatur (*signature*) beschreibt die Typen der Funktionsargumente und des resultierenden Funktionswerts
- ▶ Bedeutung wird informell erklärt.
- ▶ Mehr Info in der Online-Doku

<https://docs.python.org/2/library/functions.html?#len>.

Hilfe Modus: Methoden/Funktionen
>>> help(dict)

Help on class dict in module __builtin__:

<p>Hilfe Modus: help()</p> <pre>>>>help() modules help>__builtin__ help>q >>></pre>	<p>Hilfe Modus: Topics/Module</p> <pre>>>>help("topics") >>>help("DICTIONARIES") Mapping Types --- "dict" >>>help("modules") >>>help("__builtin__") >>></pre>
<pre>>>> " sep ".join(["1","2","3","4"]) '1 sep 2 sep 3 sep 4'</pre>	<pre>>>> dict([('a',1),('2','b')]) {'a': 1, '2': 'b'} >>> dict([('a',1),('b','2')]) {'a': 1, 'b': '2'} >>> f'a: {1,b: 2}' {'a': 1, 'b': '2'} >>> dict(a=1,b="2") {'a': 1, 'b': '2'} >>> idict=dict() >>> idict["a"] = 1 >>> idict.update({'b': "2"}) {'a': 1, 'b': '2'}</pre>
<p>Interne Hilfe</p> <pre>>>> help(re.split) Help on function split in module re: split(pattern, string, maxsplit=0, flags=0) Split the source string by the occurrences of the pattern returning a list containing the resulting substrings.</pre> <p>Erklärung der Erklärung</p> <ul style="list-style-type: none"> ▶ Name=Wert zeigt Standardwert optionaler Argumente. ▶ Ergebnis- und Argumenttypen werden im Fließtext informell erklärt. ▶ Mehr Info in der Online-Doku http://docs.python.org/2/. 	<p>Erst modul importieren:</p> <pre>>>>import re >>>help(re.split)</pre> <p>Oder via modul help</p> <pre>>>>help("re")</pre>
<p>Python ist eine objektorientierte Programmiersprache.</p> <p>Alle Daten (Werte, Datenstrukturen) in Python sind Objekte.</p> <p>Python ist eine dynamisch getypte Programmiersprache.</p> <p>Alle Objekte haben einen Typ.</p> <ul style="list-style-type: none"> ▶ Warum "dynamisch"? Der Typ einer Variablen wird nicht statisch im Quelltext deklariert, er wird dynamisch zur Laufzeit bestimmt. 	
<p>Alle Objekte haben einen Typ</p> <p>Eingebaute Funktion type()</p> <p>Sie bestimmt den Typ von jedem Objekt.</p> <p>Typen bestimmen</p> <pre>>>> type(1) >>> x = 1+3*42 >>> type(x) >>> type("ABBA") >>> type("AB"+'BA') >>> type("A" == "a") >>> type(['a', 'b']) >>> type(['a', 'b'][0]) >>> type({}) >>> type(re.search('X', 'aaa')) >>> type(None) >>> type(re.search('a', 'aaa'))</pre>	<pre>>>> import re print(type(1)) x=1+3*42 print("1",type(x)) print("2",type("ABBA")) print("3",type("AB"+'BA')) print("4",type("A"=="a")) print("5",type(['a', 'b'])) print("6",type(['a', 'b'][0])) print("7",type({})) print("8",type(re.search('X', 'aaa'))) print("9",type(None)) print("10",type(re.search('a', 'aaa'))) <type 'int'> ('1', <type 'int'>) ('2', <type 'str'>) ('3', <type 'str'>) ('4', <type 'bool'>) ('5', <type 'list'>) ('6', <type 'str'>) ('7', <type 'dict'>) ('8', <type 'NoneType'>) ('9', <type 'NoneType'>) ('10', <type '_sre.SRE_Match'>)</pre>
<p><u>Libre Office Finden und Ersetzen mit REGEX</u></p> <p>\n will match a newline (Shift-Enter) \$ on its own will match a paragraph mark (Enter) [] will match a space</p> <p>\<([^\]+)\[\]+\\1 finds duplicate words separated by spaces</p>	

<p>Die eingebaute Funktion <code>repr()</code> Erzeugt eine kanonische Zeichenkette aus jedem Objekt.</p> <pre>>>> repr("a") "'a'" >>> repr(['a', "b"]) "[‘a’, ‘b’]"</pre> <p>Read-Eval-Print-Loop (REPL) ² Im interaktiven Gebrauch wird der zuletzt eingegebene Ausdruck (<i>expression</i>) evaluiert und das Resultat als kanonischer String mit <code>print</code> ausgegeben.</p> <pre>>>> ['a', "b"] ["a", 'b'] >>> print repr(['a', "b"]) ['a', 'b']</pre>	<pre>>>> repr("a") ...a" >>> repr("a") ...a" >>> repr('a') ...a" >>> repr(['a', "b"]) "[‘a’, ‘b’]" >>> ['a', "b"] ['a', 'b'] >>> print repr(['a', "b"]) ['a', 'b']</pre>																									
<p>Objekte haben Methoden. Methoden sind Funktionen, die von einem Objekt aus aufgerufen werden.</p> <p>Methoden aufrufen Methoden-Aufrufe (<i>invocation</i>) evaluieren zu einem Objekt. Mit entsprechenden Argumenten!</p> <pre>>>> "ABBA".count("B") 2 >>> "ABBA".lower() 'abba'</pre>	<p>Methoden sind Objekte Methoden-Namen evaluieren zu einem abstrakten Objekt.</p> <pre>>>> "ABBA".count <built-in method count of str object at 0x10f743cc0></pre> <p>Methoden = aufrufbare Attribute Eine Methode ist ein aufrufbares Attribut eines Objekts. Eine Funktion ist ein aufrufbares Objekt eines Moduls.</p> <p>Dokumentationsstrings von Funktionen als Attribut <code>__doc__</code></p> <pre>>>> len.__doc__ 'len(object) -> integer\n\nReturn the number of items of a sequence or collection.'[...] >>> help(len) Help on built-in function len in module __builtin__: len(...) len(object) -> integer Return the number of items of a sequence or collection.</pre>																									
<h3>Zeichenkodierungen und Zeichensätze</h3> <p>Einschränkungen von ASCII Die weitverbreitete Zeichenkodierung mit 128 Kodes (7-Bit) unterstützt keine nicht-englischen Buchstaben.</p> <p>Verschiedene Erweiterungen mit 8-Bit Kodierungen Zeichenkodierungen mit 256 Kodes (1 Byte) für verschiedene Alphabete oder Betriebssysteme: ISO-8859-1, ISO-8859-9, Mac-Roman, Windows-1252</p> <p>Universale Lösung Unicode Unicode weist jedem Zeichen einen eindeutigen Zahlen-Kode und eine Kategorie zu. Unicode 6.0 definiert 109'449 graphische Zeichen</p>	<h3>Speicher- und Transportformat UTF-8</h3> <p>Persistente Speicherung und Datenübertragung mit Unicode UTF (Abk. für <i>Unicode Transformation Format</i>) beschreibt Methoden, einen Unicode-Wert auf eine Folge von Bytes abzubilden.</p> <p>Beispiele für UTF-8-Kodierungen</p> <table border="1" data-bbox="810 1170 1421 1298"> <thead> <tr> <th>Zeichen</th> <th>Unicode</th> <th>Unicode binär</th> <th>UTF-8 binär</th> <th>UTF-8 hexadezimal</th> </tr> </thead> <tbody> <tr> <td>Buchstabe y</td> <td>U+0079</td> <td>00000000 11111001</td> <td>01111001</td> <td>0x79</td> </tr> <tr> <td>Buchstabe ä</td> <td>U+00E4</td> <td>00000000 11100100</td> <td>11000011 10100100</td> <td>0xC3 0xA4</td> </tr> <tr> <td>Eingetragene Marke ®</td> <td>U+00AE</td> <td>00000000 10101110</td> <td>11000010 10101110</td> <td>0xC2 0xAE</td> </tr> <tr> <td>Eurozeichen €</td> <td>U+20AC</td> <td>00100000 10101100</td> <td>11100010 10000010 10101100</td> <td>0xE2 0x82 0xAC</td> </tr> </tbody> </table> <p>Quelle: http://de.wikipedia.org/wiki=UTF-8</p>	Zeichen	Unicode	Unicode binär	UTF-8 binär	UTF-8 hexadezimal	Buchstabe y	U+0079	00000000 11111001	01111001	0x79	Buchstabe ä	U+00E4	00000000 11100100	11000011 10100100	0xC3 0xA4	Eingetragene Marke ®	U+00AE	00000000 10101110	11000010 10101110	0xC2 0xAE	Eurozeichen €	U+20AC	00100000 10101100	11100010 10000010 10101100	0xE2 0x82 0xAC
Zeichen	Unicode	Unicode binär	UTF-8 binär	UTF-8 hexadezimal																						
Buchstabe y	U+0079	00000000 11111001	01111001	0x79																						
Buchstabe ä	U+00E4	00000000 11100100	11000011 10100100	0xC3 0xA4																						
Eingetragene Marke ®	U+00AE	00000000 10101110	11000010 10101110	0xC2 0xAE																						
Eurozeichen €	U+20AC	00100000 10101100	11100010 10000010 10101100	0xE2 0x82 0xAC																						
<h3>Textdatei als Bytefolge</h3> <p>Die Repräsentation der Zeichen mit Kodes > 127 sind unterschiedlich.</p> <p>Datei in UTF-8-Kodierung</p>  <pre>ä = 2 Bytes = C3 A4 \$ hexdump ae-utf8.txt 000000 61 c3 a4 0a 000004</pre> <p>Datei mit 4 Bytes</p>	<p>\$ grep . AE-11-encoded-as-l1.txt A♦ \$ sed -n 1,2p AE-11-encoded-as-l1.txt A♦ \$ hexdump -x -c AE-11-encoded-as-l1.txt 000000 c441 000a 000000 A ♦ \n 000000 41 c4 0a A... </p> <p>Datei in Latin-1-Kodierung</p>  <pre>ä = 1 Byte = E4 \$ hexdump ae-l1.txt 000000 61 e4 0a 000003</pre> <p>Datei mit 3 Bytes</p>																									
<pre>\$ echo abc hexdump -e '/1 "% ax)" "-e '/1 "%02X" "\n" 0) 61 1) 62 2) 63 3) 0A</pre>	<p>U+0061 a 61 U+00E4 ä C3 a4 U+0041 A 41 U+00C4 Ä C3 84</p> <p>\$ echo aäÅÄ hexdump -x -c 000000 c361 41a4 84c3 000a 000000 a ♦ ♦ A ♦ 204 \n 000000 61 c3 a4 41 c3 84 0a a... </p>																									

<pre> U+0031 1 31 U+0032 2 32 U+0038 8 38 \$ echo 128 xxd -b 0000000: 00110001 00110010 00110000 00001010 128. \$ echo 128 xxd 0000000: 3132 380a 128. \$ echo 128 hexdump -C 00000000 31 32 38 0a 128. </pre>	
<p>Zeichenkode berechnen aus Zeichen</p> <pre> >>> ord("A") 65 >>> ord('a') 97 >>> ord('\n') 10 >>> ord("\t") 9 >>> ord('\x20') # Hexadezimal 32 >>> ord("\'")</pre>	<p>Zeichen berechnen aus Zeichenkode</p> <pre> >>> chr(65) 'A' >>> chr(97) 'a' >>> chr(10) '\n' >>> chr(9) '\t' >>> chr(32) ' ' >>> chr(39) '"'</pre>
	<p>Buchstabe A</p> <pre> >>> ord("A") 65 >>> chr(65) 'A' >>> hex(65) '0x41' >>> chr(0x41) 'A'</pre> <p>Buchstabe ä als Type unicode</p> <pre> >>> ord(u"ä") 228 >>> unichr(228) u'\xe4' >>> u"\xe4" u'\xe4' >>> print u"\xe4" ä >>> u"\xe4".encode("utf-8") '\xc3\xaa' >>> print '\xc3\xaa' ä</pre> <p>Buchstabe ä als Type String</p> <pre> >>> ord("ä") Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: ord() expected a character, but string of length 2 found >>> chr(228) '\xe4' >>> chr(0xE4) "\xe4" >>> print chr(0xE4) >>> print '\xe4' ❶</pre>
<p>Zahlen-Konversion in Terminal</p> <pre> \$ echo \$((0xE4)) 228 \$ echo 'ibase=16; E4'bc 228 \$ echo 'ibase=10; obase=16;228'bc E4</pre>	<p>Zahlenkonversion in Python</p> <pre> >>> int('010110', 2) 22 >>> hex(22) '0x16' >>> hex(int('010110', 2)) '0x16' >>> int('0x16', 16) 22</pre>
<p>Datentyp bestimmen und testen</p> <pre> >>> type('ABBA') <type 'str'> >>> type("ABBA") == str True</pre>	<p>String-Literale notieren ▶3</p> <pre> # Einzelige (!) Zeichenkette # mit Escape-Sequenzen s1 = "a\n\x61" # Rohe Sequenz r"..." # ohne Escapes s2 = r"a\n\x61"</pre> <p>Unicode Zeichenkodes</p> <pre> >>> ord(u'€') 8364 >>> unichr(8364) u'\u20aca'</pre> <p>Datentyp bestimmen und testen</p> <pre> >>> type(u'A') <type 'unicode'> >>> type(u'ab') == unicode True >>> isinstance('ab',unicode) False</pre> <p>String-Literale notieren ▶4</p> <pre> # Einzelige (!) Zeichenkette # mit Escape-Sequenzen s1 = u"\u20aca\n\xe4" # Rohe Sequenz ur"..." # \uNNNN werden aufgelöst! s2 = ur'\u20aca\n\xe4' # Longstring s3 = u"""\u20aca # Roher Longstring s4 = ur"""€a \x00e4""" >>> isinstance(u"ab",unicode)True</pre>

<p>Kodierung der Python-Quellkodes deklarieren</p> <p>Kodierungskommentar für UTF-8-kodierte Quelltexte Deklariere Kodierung immer mit Kodierungskommentar, wenn Nicht-ASCII-Zeichen vorkommen!</p> <p>Datei in UTF-8-Kodierung</p> <pre>#!/usr/bin/python # -*- coding: utf-8 -*- print "Length of 'a':", len('a'), "Canonical:", repr('a') print "Length of 'ä':", len('ä'), "Canonical:", repr('ä') print "Length of u'a':", len(u'a'), "Canonical:", repr(u'a') print "Length of u'ä':", len(u'ä'), "Canonical:", repr(u'ä') Für Latin-1: # -*- coding: iso-8859-1 -*- \$ iso-8859-1 ist in Python 2 Standard.</pre>	<p>Enkodieren und Dekodieren von Zeichenketten</p> <p>Explizites Dekodieren von UTF-8-Repräsentation</p> <pre>>>> text = 'B\xc3\xa4h' # UTF-8-Repräsentation von Bäh >>> unicodetext = text.decode('utf-8') Es entsteht ein Unicode-Objekt!</pre> <p>Explizites Enkodieren von Unicode-Zeichen als UTF-8-Bytefolge</p> <pre>>>> unicode_text = u'Bäh' >>> utf8_text = unicode_text.encode('utf-8') Es entsteht eine Byte-String!</pre>
<p>Das Modul codecs</p> <p>codecs: Kodieren und Dekodieren Funktionen für Lesen und Schreiben von Unicode-Strings</p> <p>Einlesen von Latin-1 und schreiben von UTF-8</p> <pre>import codecs # Decode from l1 encoded file into unicode strings f = codecs.open("./ae-l1.txt", "r", "l1") # Encode unicode strings into UTF-8 encoded file g = codecs.open("./AE-l1-encoded-as-utf8.txt", "w", "utf-8") for line in f: g.write(line.upper()) </pre> <p>\$ Beim Einlesen entstehen Zeichenketten vom Typ unicode.</p>	<p>Funktion re.sub(): Ersetzen mit regulären Ausdrücken</p> <p>Globales Ersetzen mit Rückreferenz</p> <pre>import re text = u"Ässliche Köche verdürben das Gebräu" pattern = ur"([aeioäöü]+)" # Im Ersetzungstext können gematchte Gruppen eingefügt werden. # \N ist die N-te gruppierende Klammer im Pattern) replacement = ur"\[1\]" print re.sub(pattern, replacement, text)</pre> <p>\$ replacement ist eine Zeichenkette, kein regulärer Ausdruck! Falls nichts gemacht wird, bleibt die Zeichenkette unverändert!</p>
<p>Warum raw strings für Reguläre Ausdrücke in Python?</p> <ul style="list-style-type: none"> r'REGEX' oder ur'REGEX' ► Generell empfohlen in http://docs.python.org/2/library/re.html. ► Für viele Escape-Sequenzen macht es zwar keinen Unterschied, weil Python-Strings und Reguläre Ausdrücke letztlich dieselben Zeichen bedeuten: \a, \f, \n, \r, \t, \v ► Andere Escape-Sequenzen existieren nur in der Regulären Notation und Python lässt den Backslash stehen: \A, \B, \w, \W, \d, \D, \s, \S ► Aber andere Reguläre Notationen würden beim Einlesen von Nicht-Raw-Strings missinterpretiert: Um einen einzelnen Backslash zu matchen, müssten wir schreiben: re.match("\\\\\", "\\\") ► \b: Bell-Zeichen (ASCII-Code 8) im String; aber Grenze zwischen Wortzeichen und Nicht-Wortzeichen in Regex. re.sub("\bthe\b", "THE", "Other uses of the") ► Numerische Rückreferenzen \1 	<p>Gruppierung mit/ohne Rückreferenzen</p> <p>Runde Klammern ergeben referenzierbare Gruppierungen</p> <pre>>>> text = 'Blick-Leser, A-Post-Fans, andere Bindestrich-Komposita' >>> re.sub(r'(\w+)+(w+)', r'\2', text) Nicht alle Gruppen müssen referenzierbar sein! (Effizienzgründel)</pre> <p>Nichtreferenzierbare Gruppierung: (?::REGEX)</p> <pre>>>> text = 'Blick-Leser, A-Post-Fans, andere Bindestrich-Komposita' >>> re.sub(r'(?::\w+)+(w+)', r'\1', text)</pre>
<pre>>>> import re >>> text = 'Blick-Leser, A-Post-Fans, andere Bindestrich-Komposita' >>> re.sub(r'(\w+)+(w+)', r'\2', text) 'Leser, Fans, andere Komposita' >>> re.sub(r'(?::\w+)+(w+)', r'\1', text) 'Blick-, Post-, andere Bindestrich-'. >>> re.sub(r'(?::\w+)+(w+)', 'jummy', text) 'jummy, jummy, andere jummy' >>> import re text = "Viele Köche verderben den Brei." pattern = ur"(\w+)" # Alle Matches finden m = re.findall(pattern, text) for g in m: print g Viele K öche verderben den Brei</pre>	<p>Funktion re.findall(): Globale Suche</p> <p>Alle nicht-überlappenden Matches extrahieren</p> <pre>import re text = u"Viele Köche verderben den Brei." pattern = ur"(\w+)" # Alle Matches finden m = re.findall(pattern, text) for g in m: print g</pre> <p>\$ Pattern und Text müssen immer vom gleichen String-Typ sein!</p>

<p>Regex-Alternative ist nicht kommutativ!</p>	<pre>>>> import re print re.findall(r'a aa',"Saal") print re.findall(r'aa a',"Saal") ['a', 'a'] ['aa']</pre>
<p>Die Reihenfolge in einer Regex-Alternative ist nicht beliebig!</p> <pre>import re print re.findall(r'a aa',"Saal") print re.findall(r'aa a',"Saal")</pre>	
<p>Gruppierung mit/ohne Rückreferenzen: <code>re.findall()</code></p> <p><code>re.findall()</code> und gruppierende Klammern</p> <p>⚠ Unterschiedliche Funktionalität, falls gruppierende Klammern im regulären Ausdruck sind oder nicht!</p> <ul style="list-style-type: none"> ▶ Ohne: Liste der Matches ▶ Mit: Liste von Tupeln, wobei jedes Tupel-Element den gematchten Inhalt der entsprechenden gruppierenden Klammer enthält. <pre>>>> re.findall(r'(a h) a(a)', "kahler Saal") [('h', ''), ('', 'a')] # 1. möglicher Match [('', 'a')] # 2. möglicher Match</pre> <p>Was bedeutet das?</p> <pre>kahler Saal ('h', '') # 1. möglicher Match (' ', 'a') # 2. möglicher Match</pre>	<pre>>>> re.findall(r'a(h) a(a)', "kahler Saal aataah brama") [('h', ''), ('', 'a'), ('', 'a'), ('', 'a')] >>> re.findall(r'a(h) a(a)', "kahler Saal ahtaah brama") [('h', ''), ('', 'a'), ('h', ''), ('', 'a')] >>> re.findall(r'a(h).+a(a)', "kahler Saal ahtaah brama") [('h', 'a')] >>> re.findall(r'a(?h).+a(?a)', "kahler Saal ahtaah brama") ['ahler Saal ahtaah'] >>> re.findall(r'a(?h) a(?a)', "kahler Saal ahtaah brama") ['ah', 'aa', 'ah', 'aa']</pre> <p>Letzes Beispiel: letztes aa wird gefunden, aber überschneidendes ah nicht.</p>
<p>Unicode Flag: Was ist ein Wortzeichen?</p> <p>Unicode-Kategorien aktivieren</p> <pre>import re text = u"Viele Köche verderben den Brei." pattern = ur"(\w+)" # Das Flag (?u) aktiviert Unicode-Kategorien fuer \w und \b pattern = ur"(\?u)(\w+)" # Resultat ist eine Liste m = re.findall(pattern, text) for s in m: print s</pre> <p>⚠ Das Unicode-Flag (?u) zählt nicht als Gruppe wie alle (?....).</p>	<pre>>>> import re text = u"Viele Köche verderben den Brei." pattern = ur"(\w+)" # Das Flag (?u) aktiviert Unicode-Kategorien fuer \w und \b pattern = ur"(\?u)(\w+)" # Resultat ist eine Liste m = re.findall(pattern, text) for s in m: print s</pre> <pre>Viele Köche verderben den Brei</pre>
<p>Lesbare und kommentierte reguläre Ausdrücke</p> <p>Was matcht dieser Ausdruck?</p> <p>(?:[A-Z]\.)+ \w+(?:-\w+)* \$\?d+[.\d]*% \.\\. [,;?]+ \\$+</p> <p>Lesbare und kommentierbare Ausdrücke dank Flag (?x)</p> <pre>import re text = "That U.S.A. poster-print costs \$12.40..." pattern = r'''(?x) (?:[A-Z]\.)+ # set flag (?x) to allow verbose regexps \\$?\?d+(?:[.,]\d+)*% # currency/percentages, \$12.40, 82% \w+(?:-\w+)* # words with optional internal hyphens \.\\. # ellipsis [.,;?]+ # punctuation \\$+ # catch-all for non-layout characters ''' m = re.findall(pattern, text) print m</pre>	<p>Single and double quoted strings in Python are identical. The only difference is that single-quoted strings can contain unescaped double quote characters, and vice versa. For example:</p> <pre>'a "quoted" word' "another 'quoted' word"</pre> <p>Module vs Package=====A Python module is simply a Python source file.</p>

Keywords [edit]	
Python has the following keywords or <i>reserved words</i> ; they cannot be used as identifiers. ^{[3][4]}	
• and	• False [note 2]
• as	• finally
• assert	• for
• break	• from
• class	• global
• continue	• if
• def	• import
• del	• in
• elif	• is
• else	• lambda
• except	• None
• exec [note 1]	• nonlocal [note 2]
Notes	
1. ^ a b Starting from Python 3, <code>exec</code> and <code>print</code> are functions, so they are not keywords anymore.	
2. ^ a b c Starting from Python 3, keywords <code>True</code> , <code>False</code> and <code>nonlocal</code> were introduced.	

To illustrate, the following examples all return a dictionary equal to {"one": 1, "two": 2, "three": 3}:
>>> a = dict(one=1, two=2, three=3) >>> b = {'one': 1, 'two': 2, 'three': 3} >>> c = dict(zip(['one', 'two', 'three'], [1, 2, 3])) >>> d = dict([('two', 2), ('one', 1), ('three', 3)]) >>> e = dict({'three': 3, 'one': 1, 'two': 2}) >>> a == b == c == d == e True

foo function in Python:

```
def foo(x):
    if x == 0:
        bar()
        baz()
    else:
        qux(x)
        foo(x - 1)
```

Anweisung (engl. Statement)(Programmierung)

Als **Anweisung** (engl. *statement*) bezeichnet man in der Informatik ein zentrales Element vieler imperativer Programmiersprachen. Die Programme derartiger Sprachen setzen sich primär aus einer oder mehreren Anweisungen zusammen. Eine Anweisung stellt eine in der Syntax einer Programmiersprache formulierte einzelne Vorschrift dar, die im Rahmen der Abarbeitung des Programms auszuführen ist. Wie eine Anweisung syntaktisch auszusehen hat, wird durch die jeweilige Programmiersprache bzw. deren Spezifikation festgelegt und ist innerhalb eines Programms nicht änderbar. Anweisungen sind üblicherweise Zuweisungen, Kontrollanweisungen (wie Sprünge, Schleifen und bedingte Anweisungen) und Prozeduraufrufe. Abhängig von der Programmiersprache sind teilweise auch Zusicherungen, Deklarationen, Klassen- und Funktionsdefinitionen Anweisungen.

Im Gegensatz zu Ausdrücken haben Anweisungen nicht immer einen Wert. Einige Anweisungen können jedoch auch Ausdrücke sein, so z. B. Zuweisungen, Inkrementoperatoren oder Funktionsaufrufe.

In computer programming, a **statement** is the smallest standalone element of an imperative programming language that expresses some action to be carried out. It is an instruction written in a high-level language that commands the computer to perform a specified action.^[1] A program written in such a language is formed by a sequence of one or more statements. A statement may have internal components (e.g., expressions).

Many languages (e.g. C) make a distinction between statements and definitions, with a statement only containing executable code and a definition declaring an identifier, while an expression evaluates to a value only.^[2] A distinction can also be made between simple and compound statements; the latter may contain statements as components.

Assembler-Anweisung	MOV AX,BX
Initialisierung	TYPE SALARY = INTEGER
Deklaration	VAR A:INTEGER
Zuweisung	A := A + 1
Block	begin WRITE('Number? '); READLN(NUMBER); end
Bedingte Anweisung	if A > 3 then WRITELN(A) else WRITELN("NOT YET") end
Switch-Case-Anweisung	switch (c) { case 'a': alert(); break; case 'q': quit(); break; }
While-do-Schleife	while NOT EOF DO begin READLN end
Do-while-Schleife	do {computation(&i); } while (i < 10);
For-Schleife	for A:=1 to 10 do WRITELN(A) end
Unterprogrammaufruf	CLEARSCREEN()
Return-Anweisung	return 5;
Goto-Anweisung	goto 1
Assertion	assert(ptr != NULL);

Expression (computer science)

An **expression** in a programming language is a combination of one or more explicit values, constants, variables, operators, and functions that the programming language interprets (according to its particular rules of precedence and of association) and computes to produce ("to return", in a stateful environment) another value. This process, as for mathematical expressions, is called evaluation. The returned value can be of various types, such as numerical, string, and logical.

In many programming languages a function, and hence an expression containing a function, may have side effects. An expression with side effects does not normally have the property of referential transparency. In many languages (e.g. C++), expressions may be ended with a semicolon (;) to turn the expression into an expression statement. This asks the implementation to evaluate the expression for its side-effects only, and to disregard the result of the expression.

Statement/Anweisung (Instruction/Command/Befehl)

- Nur Seiteneffekte, returniert keinen Wert
- Kann Ausdrücke enthalten
- Definition: Von Variablen, von Funktionen (inkl. Deklaration des Funktionsprototypen/der Signatur)
- Funktionsaufruf

Expression/Ausdruck:

- Evaluiert und teilt einen Wert
- Hat immer einen Wert

<pre>#!/usr/bin/python #Encoding for letter Ä #http://www.utf8-zeichentabelle.de/ print """\nEncodings for letter Ä: iso-8859-1 (Latin 1): C4 (=196 or 11000100); utf-8 PC snippets_l1.py: Reload or Convert to ISO-8859-1 The encoding you've chosen ('ISO-8859-1') may change the contents of 'snippets_l1.py'. Do you want to reload the file from disk or convert the text and save in the new encoding? Convert print print print print int(5) print bin(5) Incompatible Encoding: ISO-8859-1 Please do not convert to 'ISO-8859-1'. Encoding 'ISO-8859-1' does not support some characters from the text. apple_type_ print type(print apple #String type str with utf-8 encoding</pre>	<pre>#!/usr/bin/python # -*- coding: utf-8 -*- # import sys # reload(sys) # sys.setdefaultencoding('utf-8') import sys, locale, os print(sys.stdout.encoding) print(sys.stdout.isatty()) print(locale.getpreferredencoding()) print(sys.getfilesystemencoding()) print(os.environ["PYTHONIOENCODING"]) print(chr(246), unichr(9786), unichr(9787)), "\n" print """File is encoded with utf-8 File Coding needs to specified (see 2. line) so that python is able to read letter Ä""" #Encoding for letter Ä #http://www.utf8-zeichentabelle.de/ print """\nEncodings for letter Ä: iso-8859-1 (Latin 1): C4 (=196 or 11000100); unicode : U+00C4 (=196 or 11000100); utf-8 : c3 84 (=50052 or 11000011 10000100)"""</pre>
<pre>#Conversion between hex, decimal and binary print "\nConversion between hex, decimal and binary" print int('0xc4',16) #196 print bin(196) #0b11000100 print int('0xc384',16) #50052 print bin(50052) #0b11000011 10000100 #String type str with iso-8859-1 encoding print "\nString type str with iso8859-1_encoding" apple_type_str_l1 = "\xc4apple" #Apple with iso-8859-1 encoding print type(apple_type_str_l1) #<type 'str'> print apple_type_str_l1 #apple #String type str with utf-8 encoding print "\nString type str with utf-8 encoding" apple_type_str_utf8 = "\xc3\x84apple" #Apple with utf-u encoding print type(apple_type_str_utf8) #<type 'str'> print apple_type_str_utf8 #Apple #String type unicode print "\nString type unicode" apple_type_unicode = u"\xc4apple" #Apple with unicode encoding print type(apple_type_unicode) #<type 'unicode'> print apple_type_unicode #Apple</pre>	<pre>#String type unicode print "\nstring type unicode" apple_type_unicode = u"\xc4apple" #Apple with unicode encoding print type(apple_type_unicode) #<type 'unicode'> print apple_type_unicode #Apple #Conversion between iso8859-1, unicode and utf-8 print "\nConversion between iso8859-1, unicode and utf-8" apple_type_unicode = apple_type_str_l1.decode('iso-8859-1') print type(apple_type_unicode) #<type 'unicode'> print apple_type_unicode #Apple with unicode encoding apple_type_str_utf8 = apple_type_unicode.encode('utf-8') print type(apple_type_str_utf8) #<type 'str'> print apple_type_str_utf8 #Apple with utf-u encoding print "\nConversion between utf-8, unicode and iso8859-1" apple_type_unicode = apple_type_str_utf8.decode('utf-8') print type(apple_type_unicode) #<type 'unicode'> print apple_type_unicode #Apple with unicode encoding apple_type_str_l1 = apple_type_unicode.encode('iso-8859-1') print type(apple_type_str_l1) #<type 'str'> print apple_type_str_l1 #apple with iso-8859-1 encoding</pre>
<pre>def describe_number(n): if n > 1000000: return "LARGE" elif n > 1000: return "Medium" else: return "small" print "Never printed!"</pre>	<pre>from nltk.book import * words = set(text1) longwords = [w for w in words if len(w)>15] print longwords</pre>
<pre>V = [x**2 for x in range(10)] print V</pre>	<pre># simon.clematide@uzh.ch # PCL I # Datei: def_block.py def wc1(textfile): c = 0 for line in textfile: for word in line.split(): c += 1 return c #count words of whole file def wc2(textfile): c = 0 for line in textfile: for word in line.split(): c += 1 return c #count words of first line def wc3(textfile): c = 0 for line in textfile: for word in line.split(): c += 1 return c #no count, return 1 #ruft sich selber auf print "wc1:", wc1(open('def_block.py','r')) print "wc2:", wc2(open('def_block.py','r')) print "wc3:", wc3(open('def_block.py','r'))</pre>

<pre> def foo(a): """ Return foo of a """ result = 0 for item in a: result += item return result c = foo([5,10,23]) print c </pre>	<pre> s = 'Python is great!' def f(): s = "But that's strange." print s print s f() print s # Output # Python is great! # But that's strange. # Python is great! </pre>
<pre> s = 'Python is great!' def f(): global s #das sollte man nicht tun s = "But that's strange." print s print s f() print s #Output # Python is great! # But that's strange. # But that's strange. </pre>	<pre> # Importiere Modul book aus Package nltk import nltk.book # Objekte und Funktionen aus nltk.book können nur in # vollqualifizierter Punktnotation bezeichnet werden. print "Zweites Wort aus text1:", nltk.book.text1[1] # Objekte und Funktionen können nicht direkt bezeichnet werden: print text1[1] #Traceback (most recent call last): </pre>
<pre> # Lade Modul book aus Package nltk und # importiere alle Objekte und Funktionen ins aktuelle Modul from nltk.book import * # Objekte und Funktionen aus nltk.book können ohne # Modulpräfixe bezeichnet werden. print "Zweites Token aus text1:", text1[1] # Die vollqualifizierter Punktnotation geht dann nicht print "Zweites Wort aus text1:", \ nltk.book.text1[1] #Traceback (most recent call last): </pre>	
Locate mit Regex kombiniert: \$ locate DateiNamenMuster Beispiele: \$ locate -i --regex nltk\.+corpora less \$ locate -i --regex nltk\.+corpora\+.txt less	Locate mit Regex und Grep kombiniert \$ locate DateiNamenMuster xargs grep DateiInhaltMuster Beispiel: Durchsucht die mit locate gefundenen Dateien nach dem Muster „Moby Dick“ \$ locate -i --regex nltk\.+corpora\+.txt xargs grep -l "Moby Dick" /home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt
In NANO Shell Befehle ausführen. Resultate werden im Shell Modus angezeigt. <pre> !grep -n --color "\bman\b" "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" !grep -c "\bman\b" "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" !grep "\bman\b" "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" wc </pre> <pre> 22559: Starbuck. I am old;--shake hands with me, man. 22564: a brave man that weeps; how great the agony of the persuasion then! 22612: main-truck--"Ha! he soars away with it!--Where's the old man now? 22739: inboard again; the third man helplessly dropping astern, but still 22818: whole aspect, and spite of all that mortal man could do, the solid !done (press RETURN) 499 6040 33307 !done (press RETURN) </pre>	xargs - build and execute command lines from standard input xargs reads items from the standard input, delimited by blanks or newlines, and executes the command (default is /bin/echo) one or more times with any initial-arguments followed by items read from standard input. xargs [options] [command [initial-arguments]] Examples: Durchsucht die mit locate gefundenen Dateien nach dem Muster „Moby Dick“ \$ locate -i --regex nltk\.+corpora\+.txt xargs grep -l "Moby Dick" Generates a compact listing of all the users on the system. cut -d: -f1 </etc/passwd sort xargs echo
Anzahl Zeichen \$ cat "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" wc 22924 212030 1242990 Anzahl Token (vereinfacht) \$ grep "\b.*\b" "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" wc 19650 212029 1236439 Anzahl Leerzeichen \$ grep -v "\b.*\b" "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" wc 3274 1 6551	Text vertikalisieren mit tr -s. grep anwenden um Sätzzeichen und Leerzeichen zu entfernen. alphabetisch sortieren und aggregieren mit uniq, numerisch sortieren \$ tr -s '\n' '' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" grep -P -o "\b.*\b" sort uniq -c sort -n 2866 that 3889 in 4477 a 4496 to 5958 and 6507 of 13624 the

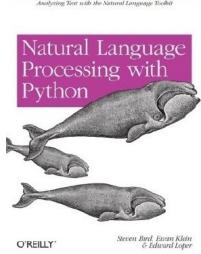
<pre>ohne grep : 12177 vs 13624 the vs (13721 bei NLTK) Begründungsversuche mit nachfolgenden commands \$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" sort uniq -c sort -n 2478 that 3273 3531 in 3979 a 4117 to 5298 and 5953 of 12177 the mit grep aber ohne Option -o : 12177 vs 13624 the (vs 13721 bei NLTK) \$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" grep -P "\b.*b\b" sort uniq -c sort -n 2478 that 3531 in 3979 a 4117 to 5298 and 5953 of 12177 the</pre>	<p>Nach der vereinfachten Vertikalisierung gibt es noch einige Komposita, welche „the“ enthalten</p> <pre>\$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" grep -P -o "\b.*b\b" sort uniq -c grep --color ".*the.*" sort -n 1 aesthetically 1 aesthetics 1 again--there 1 again!--there 1 ago!--then 1 Ahab--there's</pre> <p>In nur drei Fällen folgt ein nichtwortzeichen</p> <pre>tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" grep -P -o "\b.*\b\b" sort uniq -c sort -n grep -P ".*\bthe\b.*" 1 Jonas-in-the-Whale 1 the"--but 2 Season-on-the-Line</pre>
<pre>tr gefolgt von sort und uniq -c liefert zwei verschiedene Zeilen mit z.B. the \$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" sort uniq -c sort -n 522 of 583 and 1427 the 2478 that 3273 3979 a 4117 to 5298 and 5953 of 12177 the</pre>	<pre>\$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" sort uniq -c grep --color ".*the.*" sort -n 320 other 391 there 522 they 529 their 1427 the 12177 the</pre>
<p>Des Rätsels Lösung:</p> <p>In diesen Fällen stand die nicht am Ende der Zeile</p> <pre>\$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" sort uniq -c grep -P --color " the\\$" sort -n 12177 the</pre> <p>In diesen Fällen stand die am Ende der Zeile und der Carriage Return \r ist noch vorhanden!</p> <pre>\$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" sort uniq -c grep -P --color " the\\$" sort -n 1427 the</pre> <p>Mit Echo lässt sich das Piping gut testen</p> <pre>\$ echo 'Morgen bringen wir Cäsar um' tr [a-zA-Z] [n-za-mN-ZA-M] Zbetra oevatra jve Päfne hz \$ echo 'Zbetra oevatra jve Päfne hz' tr [a-zA-Z] [n-za-mN-ZA-M] Morgen bringen wir Cäsar um</pre>	<p>Tr muss mehrfach angewendet werden. Es gibt zwei Varianten</p> <p>Wie bis anhin tr -s '\n' und zusätzlich tr -d '\r'</p> <pre>\$ tr -s '\n' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" tr -d '\r' sort uniq -c grep -P --color " the\\$" sort -n 13604 the</pre> <p>Erst den Carriage Return in Space umwandeln und dann Space und Spaces in Newline</p> <pre>Das funktioniert da tr -s '\n' mehrere Spaces als einen betrachtet \$ tr -s '\r' < "/home/benzro/nltk_data/corpora/gutenberg/melville-moby_dick.txt" tr -s '\n' sort uniq -c grep -P --color " the\\$" sort -n 13604 the</pre>
<pre>from nltk.book import * print "from 0 :" texts() #Texte sind Sequenzen von Tokens # (Textanzeige) print "from 1 :" print text1[0:10] print "from 2 :" print text2[0:10] #Konkordanzen erstellen #KWIC (Keyword in Context) # (Stichwort mit etwas Text davor -B6 oder danach -A6) print "from 3 :" text1.concordance("man") #Keyword=man print "from 4 :" text2.concordance("man") print "from 5 :" text1.concordance("woman") #Keyword=woman print "from 6 :" text2.concordance("woman") #Ähnlichkeit Similarity # (andere Wörter, welche in ähnlichen Kontexten auftauchen) print "from 7 :" text1.similar("woman") print "from 8 :" text2.similar("woman") print "from 9 :" text1.similar("love") print "from 10 :" text2.similar("love")</pre>	<p>NLTK_DemoTour.py"</p> <p>*** Introductory Examples for the NLTK Book ***</p> <p>Loading text1, ..., text9 and sent1, ..., sent9</p> <p>Type the name of the text or sentence to view it.</p> <p>Type: 'texts()' or 'sents()' to list the materials.</p> <p>...gekürzt</p> <p>from 0 :</p> <p>text1: Moby Dick by Herman Melville 1851</p> <p>text2: Sense and Sensibility by Jane Austen 1811</p> <p>text3: The Book of Genesis</p> <p>...gekürzt</p> <p>from 1 :</p> <p>[u'I', u'Moby', u'Dick', u'by', u'Herman', u'Melville', u'1851', u''], u'ETYMOLOGY', u'']</p> <p>from 2 :</p> <p>[u'I', u'Sense', u'and', u'Sensibility', u'by', u'Jane', u'Austen', u'1811', u''], u'CHAPTER']</p> <p>from 3 :</p> <p>Displaying 25 of 527 matches:</p> <p>Civitas) which is but an artificial man . -- OPENING SENTENCE OF HOBSES 'S y of that sort that was killed by any man , such is his fierceness and swiftnes</p> <p>...gekürzt</p> <p>from 4 :</p> <p>Displaying 25 of 121 matches:</p> <p>ate owner of this estate was a single man , who lived to a very advanced age ,</p> <p>The son , a steady respectable young man , was amply provided for by the fortu</p> <p>...gekürzt</p> <p>from 5 :</p> <p>Displaying 10 of 10 matches:</p> <p>include that , like the dyspeptic old woman , he must have " broken his digester</p> <p>flexions by the sight of a freckled woman with yellow hair and a yellow gown ,</p> <p>...gekürzt</p> <p>from 6 :</p> <p>Displaying 25 of 68 matches:</p> <p>ties . Had he married a more amiable woman , he might have been made still more</p> <p>t was so much the greater , and to a woman in Mrs . Dashwood 's situation , wi</p> <p>...gekürzt</p> <p>from 7 :</p> <p>man king shark serpent camel job father hare bull fleece fellow</p> <p>surveyor as cobbler monks laugh kings ship times fiddler</p> <p>from 8 :</p> <p>man men word year friend letter moment way part gentleman situation</p> <p>person living creature style day little thing kind state</p> <p>from 9 :</p> <p>man sea it ship captain me boats them bone matter sharks nature life</p> <p>air hand way head water blubber place</p> <p>from 10 :</p> <p>affection town heart see mother sister time me word elinor life it</p> <p>marianne dear family him bed do regard her</p>

<pre># Häufigkeit in einem Korpus # (word count \$wc oder grep -c) print "from 11 :" print text1.count('love') print "from 12 :" print text2.count('love') # Häufigkeitsdistributionen # Alle Häufigkeiten aller verschiedenen Tokens berechnen # print "from 13_1 :" fd1 = FreqDist(text1) print type(fd1) vocabulary1= sorted(fd1, key=fd1.get, reverse=True) for w in vocabulary1[:50]: print w, "\t\t", fd1[w] print "from 13_2 :" fd2 = FreqDist(text2) print type(fd2) vocabulary2= sorted(fd2, key=fd2.get, reverse=True) for w in vocabulary2[:20]: print w, "\t\t", fd2[w] print "from 14_1 :" fd1.plot(50,cumulative=False) print "from 14_2 :" fd2.plot(20,cumulative=False)</pre>	from 11 : 24 from 12 : 77 from 13_1: <class 'nltk.probability.FreqDist'> , 18713 the 13721 . 6862 of 6536 and 6024 a 4569 ...gekürzt from 13_2: <class 'nltk.probability.FreqDist'> , 9397 to 4063 . 3975 the 3861 of 3565 and 3350 ...gekürzt from 14_1: from 14_2:
<pre># Statistische Kollokationen # Welche Wörter kommen unerwartet häufig nebeneinander vor? print "from 15 :" text1.collocations() print "from 16 :" text2.collocations() # Dispersions Plots # Wo tauchen auf einer Zeitachse Wörter wie oft auf? # Amerikanische Ansprachen print "from 17 :" text4.distribution_plot(["freedom", "war"]) print "from 18 :" text4.distribution_plot(["economy", "war"])</pre>	from 15 : Sperm Whale; Moby Dick; White Whale; old man; Captain Ahab; sperm whale; Right Whale; Captain Peleg; New Bedford; Cape Horn; cried Ahab; years ago; lower jaw; never mind; Father Mapple; cried Stubb; chief mate; white whale; ivory leg; one hand from 16 : Colonel Brandon; Sir John, Lady Middleton; Miss Dashwood; every thing; thousand pounds; dare say; Miss Steeles; said Elinor; Miss Steele; every body; John Dashwood; great deal; Harley Street; Berkeley Street; Miss Dashwoods; young man; Combe Magna; every day; next morning from 17 : from 18 :
<pre># Textgenerierung # Statistische Generierung von Texten aus Trigramm-Statistiken von Wörtern # Typischer Wortverbindungen aus einem Korpus. # (In NLTK 3 momentan kaputt.) Politische Rede? # Sich inspirieren lassen von präsidenialen # Ansprachen der vergangenen Präsidenten der USA. print "from 19 :" #text4.generate() '''Fellow - Countrymen : At this second gathering , our commerce and share them . Instead of undertaking particular recommendations on this occasion , feeling the emotions which the Executive , in seeking practical plans for mediation , conciliation , and evil is real , and defense ; if a preference , upon a party , and we have piled deficit upon deficit , mortgaging our future is this which gives inestimable value to those who came here believed they could to merit it can deliver ; from inflated rhetoric that postures instead of others , and by our own''' # Listenkomprehension print "from 20 :" V = set(text1) long_words = [w for w in V if len(w)>15] print long_words</pre>	from 19 : from 20 : [u'hermaphroditical', u'subterraneousness', u'uninterpenetratingly', u'irresistibleness', u'responsibilities', u'comprehensiveness', u'uncompromisedness', u'superstitiousness', u'uncomfortableness', u'supernaturalness', u'circumnavigating', u'characteristically', u'cannibalistically', u'circumnavigations', u'indispensableness', u'preternaturalness', u'apprehensiveness', u'CIRCUMNAVIGATION', u'simultaneousness', u'undiscriminating', u'Physiognomically', u'physiognomically', u'circumnavigation', u'indiscriminately'] Process finished with exit code 0
<pre>import nltk text = """At eight o'clock on Thursday morning Arthur didn't feel very good.""" #Tokenizer tokens = nltk.word_tokenize(text) print "\nTokenizer\n", tokens #Tagger tagged = nltk.pos_tag(tokens) print "\nTagger\n", tagged #Chunker entities = nltk.chunk.ne_chunk(tagged) print "\nChunker\n", entities #Stemmer stemmer = nltk.stem.PorterStemmer() print "\nStemmer" for tok in nltk.word_tokenize('He believes in stemming.'): print tok, "\t", stemmer.stem(tok) #Lemmatizer print "\nLemmatizer" lemmatizer = nltk.stem.WordNetLemmatizer() for tok in nltk.word_tokenize('He believes in stemming.'): print tok, "\t", lemmatizer.lemmatize(tok)</pre>	Tokenizer (make_words) ['At', 'eight', 'o\'clock', 'on', 'Thursday', 'morning', 'Arthur', 'did', "n't", 'feel', 'very', 'good', '.'] Tagger (add pos-tags) [('At', 'IN'), ('eight', 'CD'), ('o\'clock', 'NN'), ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN'), ('Arthur', 'NNP'), ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'), ('very', 'RB'), ('good', 'JJ'), ('.', '.') Chunker (add constituents, noun groups, verb groups etc) <u>Shallow parsing</u> (also <u>chunking</u> , <u>light parsing</u>) is an analysis of a sentence which identifies the constituents (noun groups, verbs, verb groups, etc.), but does not specify their internal structure, nor their role in the main sentence. (S At/IN eight/CD o'clock/NN on/IN Thursday/NNP morning/NN (PERSON Arthur/NNP) did/VBD n't/RB feel/VB very/RB good/JJ .) Stemmer (token as input and adds a truncated form called stem) He He believes believ in in stemming stem .

	<p>Lemmatizer (token as input and adds the dictionary or base form called lemma)</p> <pre>He He believes belief in in stemming stemming .</pre>
<pre>#!/usr/bin/env python # -*- coding: utf-8 -*- import nltk text = """At eight o'clock on Thursday morning Arthur didn't feel very good.""" print "\nRaw Text\n",text text_ = 'He believes in stemming.' print "\nRaw Text_\n",text_ #Tokenizer tokens = nltk.word_tokenize(text) print "\nTokenizer\n",tokens itokens=text_.split() print "\niSimpleSplitTokenizer\n",itokens #Text itext=nltk.Text(itokens) print "\niText\n ", itext #to print the first line for i in itext: print type(i), i print " Text Länge und Typ\n ", len(itext), type(itext) #Tagger tagged = nltk.pos_tag(tokens) print "\nTagger\n",tagged itagged = nltk.pos_tag(itokens) print "\niTagger\n",itagged</pre>	<pre>/usr/bin/python2.7 "/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Vorlesung/vl08/NLTK_Einf.py"</pre> <p>Raw Text At eight o'clock on Thursday morning Arthur didn't feel very good.</p> <p>Raw Text_</p> <p>He believes in stemming.</p> <p>Tokenizer</p> <pre>[At', 'eight', "o'clock", 'on', 'Thursday', 'morning', 'Arthur', 'did', "n't", 'feel', 'very', 'good', '.']</pre> <p>iSimpleSplitTokenizer</p> <pre>['He', 'believes', 'in', 'stemming.']}</pre> <p>iText</p> <pre><Text: He believes in stemming....> <type 'str'> He <type 'str'> believes <type 'str'> in <type 'str'> stemming. Text Länge und Typ 4 <class 'nltk.text.Text'></pre> <p>Tagger</p> <pre>[('At', 'IN'), ('eight', 'CD'), ("o'clock", 'NN'), ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN'), ('Arthur', 'NNP'), ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'), ('very', 'RB'), ('good', 'JJ'), ('.', '.')]</pre> <p>iTagger</p> <pre>[('He', 'PRP'), ('believes', 'VBZ'), ('in', 'IN'), ('stemming.', 'NN')]</pre>
<pre>#Chunker entities = nltk.chunk.ne_chunk(tagged) print "\nChunker\n",entities ientities = nltk.chunk.ne_chunk(itagged) print "\niChunker\n",ientities #Stemmer stemmer = nltk.stem.PorterStemmer() print "\nStemmer" for tok in nltk.word_tokenize('He believes in stemming.'): print tok,"\\t", stemmer.stem(tok) print "\niStemmer" for tok in itext: print tok,"\\t", stemmer.stem(tok) #Lemmatizer print "\nLemmatizer" lemmatizer = nltk.stem.WordNetLemmatizer() for tok in nltk.word_tokenize('He believes in stemming.'): print tok,"\\t", lemmatizer.lemmatize(tok) print "\niLemmatizer" for tok in itext: print tok,"\\t", lemmatizer.lemmatize(tok)</pre>	<p>Chunker</p> <pre>(S At/IN eight/CD o'clock/NN on/IN Thursday/NNP morning/NN (PERSON Arthur/NNP) did/VBD n't/RB feel/VB very/RB good/JJ .)</pre> <p>iChunker</p> <pre>(S He/PRP believes/VBZ in/IN stemming./NN)</pre> <p>Stemmer</p> <pre>He He believes belief in in stemming stem .</pre> <p>iStemmer</p> <pre>He He believes belief in in stemming. stemming.</pre> <p>Lemmatizer</p> <pre>He He believes belief in in stemming stemming .</pre> <p>iLemmatizer</p> <pre>He He believes belief in in stemming. stemming.</pre>

<p>Short and dense: http://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html</p> <p>The goal of both stemming and lemmatization is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form.</p> <p>However, the two words differ in their flavor. Stemming usually refers to a crude heuristic process that chops off the ends of words in the hope of achieving this goal correctly most of the time, and often includes the removal of derivational affixes. Lemmatization usually refers to doing things properly with the use of a vocabulary and morphological analysis of words, normally aiming to remove inflectional endings only and to return the base or dictionary form of a word, which is known as the lemma.</p> <p>From the NLTK docs:</p> <p>Lemmatization and stemming are special cases of normalization. They identify a canonical representative for a set of related word forms.</p>	<p>Lemmatisation is closely related to stemming. The difference is that a stemmer operates on a single word without knowledge of the context, and therefore cannot discriminate between words which have different meanings depending on part of speech. However, stemmers are typically easier to implement and run faster, and the reduced accuracy may not matter for some applications.</p> <p>For instance:</p> <ol style="list-style-type: none"> 1. The word "better" has "good" as its lemma. This link is missed by stemming, as it requires a dictionary look-up. 2. The word "walk" is the base form for word "walking", and hence this is matched in both stemming and lemmatization. 3. The word "meeting" can be either the base form of a noun or a form of a verb ("to meet") depending on the context, e.g., "in our last meeting" or "We are meeting again tomorrow". Unlike stemming, lemmatisation can in principle select the appropriate lemma depending on the context.
<pre>tuples2 = [('a',1),('b',2),('c',3),('d',4)] #tuples2= "klklk" #funktioniert nicht da kein Tupel l1 = [(c2,c1) for (c1,c2) in tuples2] print l1 #[(1, 'a'), (2, 'b'), (3, 'c'), (4, 'd')]</pre>	<pre># Einfaches Statement print "Something to print" # If-Statement mit komplexen Ausdrücken drin if len("A "+"String") > 5: print "A "+"String".lower() # Prinzip: # Statements sind keine Expressions # Frage: # Aber warum kann ich interaktiv schreiben? "A Statement" # oder 3*4+8</pre>
<pre>#Demo von Mengen- und Dictionarykomprehension #Alphanumeric is a combination of alphabetic and numeric # characters, and is used to describe the collection of # Latin letters and Arabic digits or a text constructed # from this collection. mc = {x.lower() for x in "Das Alphabet. 12" if x.isalnum()} print mc #set(['a', 'b', 'e', 'd', 'h', 'l', 'i', 'p', # 's', '2', 't']) der Punkt fällt weg ms = set() for x in "Das Alphabet. 12": if x.isalnum(): ms.add(x.lower()) print ms #set(['a', 'b', 'e', 'd', 'h', 'l', 'i', 'p', # 's', '2', 't']) Der Punkt fällt weg text = "abarakadabra. 12" dc = {c:text.count(c) for c in set(text)} print dc #{'a': 5, ' ': 1, 'b': 2, 'd': 1, 'k': 1, '.': 1, # 'l': 1, 'r': 2, '2': 1} # Wie würde man das als For-Schleife implementieren? ds = dict() for c in set(text): ds.update({c:text.count(c)}) print ds #{'a': 5, ' ': 1, 'b': 2, 'l': 1, 'd': 1, '2': 1, # 'k': 1, 'r': 2, '.': 1}</pre>	<pre>/usr/bin/python2.7 "/media/benzro/OS/Users/benzro/Documents /LubuntuShared/8) PCL-1/Vorlesung/vl09/other_comprehensions.py" set(['a', 'b', 'e', 'd', 'h', 'l', 'i', 'p', 's', '2', 't']) set(['a', 'b', 'e', 'd', 'h', 'l', 'i', 'p', 's', '2', 't']) {'a': 5, ' ': 1, 'b': 2, 'd': 1, 'k': 1, '.': 1, 'l': 1, 'r': 2, '2': 1} {'a': 5, ' ': 1, 'b': 2, 'l': 1, 'd': 1, '2': 1, 'k': 1, 'r': 2, '.': 1} Process finished with exit code 0</pre>
<pre>from nltk.corpus import brown # Korpus als Liste von 2-Tupeln (Wort, POS-Tag) brown_tagged_words = brown.tagged_words() # Das balancierte Korpus umfasst Texte aus 15 Kategorien print "brown.categories ", brown.categories(),"\n" # Es werden alle 500 Einzeldateien verarbeitet! cnt=0 for f in brown.fileids()[1:10]: cnt+=1 print cnt, ": ", f print brown.abspath(f) # brown print 60*' ', \ '\nKorpus als Liste von Paaren von Wörtern und POS-Tags:' print brown_tagged_words, "\n", \ len(brown_tagged_words), "\n" #1 161 192 Zeichen print brown_tagged_words[0:20], "\n"</pre>	<pre>brown.categories [u'adventure', u'belles_lettres', u'editorial', u'fiction', u'government', u'hobbies', u'humor', u'learned', u'lore', u'mystery', u'news', u'religion', u'reviews', u'romance', u'science_fiction'] 1 : ca02 /home/benzro/nltk_data/corpora/brown/ca02 2 : ca03 /home/benzro/nltk_data/corpora/brown/ca03 3 : ca04 /home/benzro/nltk_data/corpora/brown/ca04 4 : ca05 /home/benzro/nltk_data/corpora/brown/ca05 5 : ca06 /home/benzro/nltk_data/corpora/brown/ca06 6 : ca07 /home/benzro/nltk_data/corpora/brown/ca07 7 : ca08 /home/benzro/nltk_data/corpora/brown/ca08 8 : ca09 /home/benzro/nltk_data/corpora/brown/ca09 9 : ca10 /home/benzro/nltk_data/corpora/brown/ca10</pre>

<pre># Wo liegen die Dateien dieses Korpus? #Python print "brown.root path: ", brown.root #Unix ## locate -i --regex nltk\.corpora\.+brown ## ls -l "/home/benzro/nltk_data/corpora/brown/README" # -rw-rw-r-- 1 benzro benzro 387 Nov 13 14:46 ## ls -l "/home/benzro/nltk_data/corpora/brown/CONTENTS" # -rw-rw-r-- 1 benzro benzro 145896 Nov 13 14:46 ## less "/home/benzro/nltk_data/corpora/brown/CONTENTS" ## less "/home/benzro/nltk_data/corpora/brown/README"</pre>	<p>Korpus als Liste von Paaren von Wörtern und POS-Tags: [(u'The', u'AT'), (u'Fulton', u'NP-TL'), ...] 1161192</p> <p>[(u'The', u'AT'), (u'Fulton', u'NP-TL'), (u'County', u'NN-TL'), (u'Grand', u'JJ-TL'), (u'Jury', u'NN-TL'), (u'said', u'VBD'), (u'Friday', u'NR'), (u'an', u'AT'), (u'investigation', u'NN'), (u'of', u'IN'), (u'Atlanta's', u'NP\$'), (u'recent', u'JJ'), (u'primary', u>NN'), (u'election', u'NN'), (u'produced', u'VBD'), (u'no', u'AT'), (u'evidence', u'NN'), (u'"', u'"'), (u''), (u'that', u'CS')]</p> <p>brown.root path /home/benzro/nltk_data/corpora/brown</p>
<pre>from nltk.corpus import gutenberg # oder absoluter Pfad einer Textdatei filename = 'austen-emma.txt' # Korpus als eine lange Zeichenkette emma_chars = gutenberg.raw(filename) print 60*'*', '\nKorpus als eine lange Zeichenkette: ' print emma_chars[-324:] #die letzten 324 Zeichen # Korpus als Liste von Wörtern # (Wort ist Zeichenkette) emma_words = gutenberg.words(filename) print 60*'*', '\nKorpus als Liste von Wörtern: ' print emma_words[-30:] #die letzten 30 Wörter # Korpus als Liste von Sätzen # (Satz ist Liste von Wörtern) emma_sents = gutenberg.sents(filename) print 60*'*', '\nKorpus als Liste von Saetzen: ' print emma_sents[-5:] #die letzten 5 Sätze # Korpus als Liste von Paragraphen # (Paragraph ist Liste von Sätzen) emma_paras = gutenberg.paras(filename) print 60*'*', '\nKorpus als Liste von Paragraphen: ' print emma_paras[-2:] #die letzten 2 Paragraphen</pre>	<p>***** Korpus als eine lange Zeichenkette: white satin, very few lace veils; a most pitiful business!--Selina would stare when she heard of it.--But, in spite of these deficiencies, the wishes, the hopes, the confidence, the predictions of the small band of true friends who witnessed the ceremony, were fully answered in the perfect happiness of the union. FINIS ***** Korpus als Liste von Wörtern: [u't', u'he', u'confidence', u'!', u'the', u'predictions', u'of', u'the', u'small', u'band', u'of', u'true', u'friends', u'who', u'witnessed', u'the', u'ceremony', u'!', u'were', u'fully', u'answered', u'in', u'the', u'perfect', u'happiness', u'of', u'the', u'union', u'!', u'FINIS'] ***** Korpus als Liste von Saetzen: [gekürzt..., [u'The', u'wedding', u'was', u'very', u'much', u'like', u'other', u'weddings', u'!', u'where', u'the', u'parties', u'have', u'no', u'taste', u'for', u'finery', u'or', u'parade', u'!', u'and', u'Mrs', u'!', u'Elton', u'!', u'from', u'the', u'particulars', u'detailed', u'by', u'her', u'husband', u'!', u'thought', u'it', u'all', u'extremely', u'shabby', u'!', u'and', u'very', u'inferior', u'to', u'her', u'own', u'...', u'Very', u'little', u'white', u'satin', u'!', u'very', u'few', u'lace', u'veils', u'!', u'a', u'most', u'pitiful', u'business', u'!', u'Selina', u'would', u'stare', u'when', u'she', u'heard', u'of', u'it', u'"], [u'!', u'But', u'!', u'in', u'spite', u'of', u'the', u'deficiencies', u'!', u'the', u'wishes', u'!', u'the', u'hopes', u'!', u'the', u'confidence', u'!', u'the', u'predictions', u'of', u'the', u'small', u'band', u'of', u'true', u'friends', u'who', u'witnessed', u'the', u'ceremony', u'!', u'were', u'fully', u'answered', u'in', u'the', u'perfect', u'happiness', u'of', u'the', u'union', u''], [u'FINIS']] ***** Korpus als Liste von Paragraphen: [[[u'The', u'wedding', u'was', u'very', u'much', u'like', u'other', u'weddings', u'!', u'where', u'the', u'parties', u'have', u'no', u'taste', u'for', u'finery', u'or', u'parade', u'!', u'and', u'Mrs', u'!', u'Elton', u'!', u'from', u'the', u'particulars', u'detailed', u'by', u'her', u'husband', u'!', u'thought', u'it', u'all', u'extremely', u'shabby', u'!', u'and', u'very', u'inferior', u'to', u'her', u'own', u'...', u'Very', u'little', u'white', u'satin', u'!', u'very', u'few', u'lace', u'veils', u'!', u'a', u'most', u'pitiful', u'business', u'!', u'Selina', u'would', u'stare', u'when', u'she', u'heard', u'of', u'it', u'"], [u'!', u'But', u'!', u'in', u'spite', u'of', u'the', u'deficiencies', u'!', u'the', u'wishes', u'!', u'the', u'hopes', u'!', u'the', u'confidence', u'!', u'the', u'predictions', u'of', u'the', u'small', u'band', u'of', u'true', u'friends', u'who', u'witnessed', u'the', u'ceremony', u'!', u'were', u'fully', u'answered', u'in', u'the', u'perfect', u'happiness', u'of', u'the', u'union', u''], [u'FINIS']]]</p>
<pre># Wieviele Paragraphen? number_of_paras=len(gutenberg.paras(filename)) print "\n number of paras: ",number_of_paras,"\\n" # Wieviele Sätze? number_of_sents=len(gutenberg.sents(filename)) print "\n number of sents: ",number_of_sents,"\\n" # Wieviele Sätze pro Paragraph im Schnitt? emma_paras_sents_mean=number_of_sents/number_of_paras print "\n number sents per para: ",emma_paras_sents_mean,"\\n" # Wieviele Wörter pro Satz im Schnitt? number_of_words=len(gutenberg.words(filename)) emma_sents_words_mean=number_of_words/number_of_sents print "\n number of words per sent: ",emma_sents_words_mean,"\\n" # Wie kann man die durchschnittliche Wortlänge berechnen? number_of_chars=len(gutenberg.raw(filename)) emma_words_chars_mean=number_of_chars/number_of_words print "\n number of chars per word: ",emma_words_chars_mean,"\\n" # Wo liegen die Dateien dieses Korpus? print "\n gutenberg.root directory: ", gutenberg.root,"\\n"</pre>	<p>number of paras: 2371</p> <p>number of sents: 7752</p> <p>number sents per para: 3</p> <p>number of words per sent: 24</p> <p>number of chars per word: 4</p> <p>gutenberg.root directory: /home/benzro/nltk_data/corpora/gutenberg</p>

<h3>NLTK (Natural Language Toolkit)</h3> <p>NLTK-Framework http://www.nltk.org</p> <ul style="list-style-type: none"> ▶ Sammlung von Open-Source-Python-Modulen für die Sprachverarbeitung (Natural Language Processing, NLP) ▶ Ressourcen: frei verfügbare Korpora, Treebanks, Lexika ... ▶ Applikationen: Tokenizer, Stemmer, Tagger, Chunker, Parser, Semantik, Alignierung... (oft eher Toy-Implementationen für Lehrzwecke; andere Toolkits sind industrial-strength http://spacy.io) ▶ Module für Evaluation, Klassifikation, Clustering, Maschinelles Lernen (Schnittstellen zu State-of-the-Art-Bibliotheken) ▶ API (Application Programming Interface) für WordNet und Lexika <p>Installationsanleitungen für Win, Mac, Linux http://www.nltk.org/install.html</p>	<h3>Bird et. al (2009): Natural Language Processing in Python¹</h3>  <ul style="list-style-type: none"> ▶ Praktische Einführung in NLP mit Hilfe von NLTK 3 ▶ Anwendungsorientiert, keine vertiefte Einführung in Python-Konzepte ▶ Kursbuch für PCL I (und auch II) ▶ 2. Edition (nur online!) benutzt Syntax von Python 3 (kompatibel mit futurisiertem Python 2) ▶ 1. Edition (NLTK 2!) als PDF abgelegt im Materialordner ▶ Weiteres NLTK-basiertes Buch mit vielen NLP-Rezepten: [PERKINS 2010] 																																							
<h3>Modulübersicht</h3> <table border="1"> <thead> <tr> <th>Language processing task</th> <th>NLTK modules</th> <th>Functionality</th> </tr> </thead> <tbody> <tr> <td>Accessing corpora</td> <td>nltk.corpus</td> <td>standardized interfaces to corpora and lexicons</td> </tr> <tr> <td>String processing</td> <td>nltk.tokenize, nltk.stem</td> <td>tokenizers, sentence tokenizers, stemmers</td> </tr> <tr> <td>Collocation discovery</td> <td>nltk.collocations</td> <td>t-test, chi-squared, point-wise mutual information</td> </tr> <tr> <td>Part-of-speech tagging</td> <td>nltk.tag</td> <td>n-gram, backoff, Brill, HMM, Trigram</td> </tr> <tr> <td>Classification</td> <td>nltk.classify, nltk.cluster</td> <td>decision tree, maximum entropy, naïve Bayes, EM, k-means</td> </tr> <tr> <td>Chunking</td> <td>nltk.chunk</td> <td>regular expression, n-gram, named-entity</td> </tr> <tr> <td>Parsing</td> <td>nltk.parse</td> <td>chart, feature-based, unification, probabilistic, dependency</td> </tr> <tr> <td>Semantic interpretation</td> <td>nltk.sem, nltk.inference</td> <td>lambda calculus, first-order logic, model checking</td> </tr> <tr> <td>Evaluation metrics</td> <td>nltk.metrics</td> <td>precision, recall, agreement coefficients</td> </tr> <tr> <td>Probability and estimation</td> <td>nltk.probability</td> <td>frequency distributions, smoothed probability distributions</td> </tr> <tr> <td>Applications</td> <td>nltk.app, nltk.chat</td> <td>graphical concordancer, parsers, WordNet browser, chatbots</td> </tr> <tr> <td>Linguistic fieldwork</td> <td>nltk.toolbox</td> <td>manipulate data in SIL Toolbox format</td> </tr> </tbody> </table>	Language processing task	NLTK modules	Functionality	Accessing corpora	nltk.corpus	standardized interfaces to corpora and lexicons	String processing	nltk.tokenize, nltk.stem	tokenizers, sentence tokenizers, stemmers	Collocation discovery	nltk.collocations	t-test, chi-squared, point-wise mutual information	Part-of-speech tagging	nltk.tag	n-gram, backoff, Brill, HMM, Trigram	Classification	nltk.classify, nltk.cluster	decision tree, maximum entropy, naïve Bayes, EM, k-means	Chunking	nltk.chunk	regular expression, n-gram, named-entity	Parsing	nltk.parse	chart, feature-based, unification, probabilistic, dependency	Semantic interpretation	nltk.sem, nltk.inference	lambda calculus, first-order logic, model checking	Evaluation metrics	nltk.metrics	precision, recall, agreement coefficients	Probability and estimation	nltk.probability	frequency distributions, smoothed probability distributions	Applications	nltk.app, nltk.chat	graphical concordancer, parsers, WordNet browser, chatbots	Linguistic fieldwork	nltk.toolbox	manipulate data in SIL Toolbox format	<h3>Verzeichnisstruktur vom NLTK 3</h3> <pre>\$ tree -F /Library/Python/2.7/site-packages/nltk /Library/Python/2.7/site-packages/nltk ├── __init__.py ├── _util.py ├── _util.pyc └── util.py ├── __init__.py ├── __init__.pyc ├── chartparser.py ├── chartparser_app.py ├── chartparser_app.pyc ├── chunkparser.py └── chunkparser_app.pyc ... 24 directories, 648 files</pre> <ul style="list-style-type: none"> ▶ Module: Dateien mit Python-Quellkode: <code>util.py</code> ▶ Maschinenunabhängig kompilierter Bytekode: <code>chunkparser_app.pyc</code> ▶ Packages: Verzeichnisse wie <code>nltk</code> oder <code>app</code> mit <code>__init__.py</code>
Language processing task	NLTK modules	Functionality																																						
Accessing corpora	nltk.corpus	standardized interfaces to corpora and lexicons																																						
String processing	nltk.tokenize, nltk.stem	tokenizers, sentence tokenizers, stemmers																																						
Collocation discovery	nltk.collocations	t-test, chi-squared, point-wise mutual information																																						
Part-of-speech tagging	nltk.tag	n-gram, backoff, Brill, HMM, Trigram																																						
Classification	nltk.classify, nltk.cluster	decision tree, maximum entropy, naïve Bayes, EM, k-means																																						
Chunking	nltk.chunk	regular expression, n-gram, named-entity																																						
Parsing	nltk.parse	chart, feature-based, unification, probabilistic, dependency																																						
Semantic interpretation	nltk.sem, nltk.inference	lambda calculus, first-order logic, model checking																																						
Evaluation metrics	nltk.metrics	precision, recall, agreement coefficients																																						
Probability and estimation	nltk.probability	frequency distributions, smoothed probability distributions																																						
Applications	nltk.app, nltk.chat	graphical concordancer, parsers, WordNet browser, chatbots																																						
Linguistic fieldwork	nltk.toolbox	manipulate data in SIL Toolbox format																																						
<h3>Module importieren</h3> <p>Anweisung: import Module ▶1</p> <pre># Importiere Modul book aus Package nltk import nltk.book # Objekte und Funktionen aus nltk.book können nur in vollqualifizierter Punktnotation bezeichnet werden. print "Zweites Wort aus text1:", nltk.book.text1[1] # Objekte und Funktionen können nicht direkt bezeichnet werden: print text1[1]</pre>	<h3>Alle Objekte und Funktionen aus Modulen importieren</h3> <p>Anweisung: from Module import * ▶2</p> <pre># Lade Modul book aus Package nltk und importiere alle Objekte und Funktionen ins aktuelle Modul from nltk.book import * # Objekte und Funktionen aus nltk.book können ohne Modulpräfix bezeichnet werden. print "Zweites Token aus text1:", text1[1] # Die vollqualifizierter Punktnotation geht dann nicht print "Zweites Wort aus text1:", nltk.book.text1[1]</pre>																																							
<h3>Python 2 futurisieren</h3> <p>Wichtige Direktiven zur Kompatibilität von Python 2 mit Python 3</p> <pre>from __future__ import print_function, unicode_literals, division</pre> <table border="1"> <tr> <td data-bbox="171 1724 457 1776">Python 2: print-Statement</td> <td data-bbox="472 1724 758 1776">Python 3: print-Funktion ▶</td> </tr> <tr> <td data-bbox="171 1776 457 1805">>>> print u'sähe:', count</td> <td data-bbox="472 1776 758 1805">>>> print('sähe:', count)</td> </tr> <tr> <td data-bbox="171 1805 457 1857">Python 2: Ganzzahldivision</td> <td data-bbox="472 1805 758 1857">Python 3: Division</td> </tr> <tr> <td data-bbox="171 1857 457 1886">>>> print 1/3</td> <td data-bbox="472 1857 758 1886">>>> print(1/3)</td> </tr> <tr> <td data-bbox="171 1886 457 1915">0</td> <td data-bbox="472 1886 758 1915">0.3333333333333333</td> </tr> </table>	Python 2: print-Statement	Python 3: print-Funktion ▶	>>> print u'sähe:', count	>>> print('sähe:', count)	Python 2: Ganzzahldivision	Python 3: Division	>>> print 1/3	>>> print(1/3)	0	0.3333333333333333	<h3>Eine Tour durch Kapitel 1</h3> <ul style="list-style-type: none"> ▶ Repräsentation von Text-Korpora als Objekt vom Typ <code>Text</code> (im Wesentlichen als Liste von String-Token) ▶ KWIC (<i>Keyword in context</i>): Konkordanzen erstellen und anzeigen ▶ Vorkommensähnlichkeit (<i>similarity</i>): Welche unterschiedlichen Wörter erscheinen häufig in ähnlichen Kontexten? ▶ Häufigkeitsverteilungen (<i>frequency distribution</i>) berechnen: Wie oft kommt welche Wortform vor? ▶ Statistische Kollokationen (<i>collocations</i>): Welche Wortpaare kommen viel häufiger zusammen vor als zufällig zu erwarten wäre? ▶ Dispersion-Plot (Korpuslinguistik): An welchen Stellen in einem Korpus kommt ein Wort vor? [BAKER et al. 2006] <p>Visualisierungen mit Plotting benötigen separate Diagramm-Bibliothek <code>matplotlib</code> (Download via matplotlib.org).</p>																													
Python 2: print-Statement	Python 3: print-Funktion ▶																																							
>>> print u'sähe:', count	>>> print('sähe:', count)																																							
Python 2: Ganzzahldivision	Python 3: Division																																							
>>> print 1/3	>>> print(1/3)																																							
0	0.3333333333333333																																							

<h3>Listenkomprehension (list comprehension)</h3>	<h3>Listenkomprehension mit Bedingungen ►</h3>
<p>Mathematische Mengenkomprehension</p> <p>Die Menge aller Quadratzahlen aus der Grundmenge der Zahlen von 0 bis 9.</p> $\{x^2 \mid x \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}\}$ <p>In Python mit Listen als Mengen:</p> <pre>>>> [x**2 for x in range(10)] [0, 1, 4, 9, 16, 25, 36, 49, 64, 81]</pre>	<p>Syntaxschema (aus [BIRD et al. 2009, 19])</p> $\{w \mid w \in V \& P(w)\} \quad \text{Die Liste aller Elemente } w \text{ aus } V, \text{ für die die Eigenschaft } P(w) \text{ wahr ist.}$ <p>[w for w in V if p(w)]</p>
<p>Syntax</p> <div style="display: flex; align-items: center;"> <div style="border: 1px solid orange; padding: 2px 10px; margin-right: 10px;">[x**2]</div> <div style="border: 1px solid blue; padding: 2px 10px; margin-right: 10px;">for x in range(10)]</div> <div style="display: flex; justify-content: space-around; width: 100%;"> Bildungsvorschrift Grundmenge, deren Werte x durchläuft </div> </div>	<p>Filtern von Vokabularlisten ►3</p> <pre>from nltk.book import * words = set(text1) longwords = [w for w in words if len(w)>15] set(text1) erzeugt Menge aller Listenelemente aus text1.</pre>
<pre>>>> from nltk.book import * *** Introductory Examples for the NLTK Book *** Loading text1, ..., text9 and sent1, ..., sent9 Type the name of the text or sentence to view it. Type: 'texts()' or 'sents()' to list the materials. text1: Moby Dick by Herman Melville 1851 text2: Sense and Sensibility by Jane Austen 1811 text3: The Book of Genesis text4: Inaugural Address Corpus text5: Chat Corpus text6: Monty Python and the Holy Grail text7: Wall Street Journal text8: Personals Corpus text9: The Man Who Was Thursday by G . K . Chesterton 1908 >>> </pre>	<p>Package nltk</p> <pre>>>> from nltk import * >>> help <module 'nltk.help' from '/usr/local/lib/python2.7/dist-packages/nltk/help.pyc'></pre> <p>Für das NLTK BUCH gibt es eine book collection welche mit dem Befehl</p> <pre>>>> nltk.download() runtergeladen werden kann</pre> <p>Modul book aus Package nltk</p> <pre>>>> from nltk.book import * *** Introductory Examples for the NLTK Book *** Loading text1, ..., text9 and sent1, ..., sent9</pre>
<p>6. Now you're ready to load your own corpus, using the following code:</p> <pre>>>> >>> >>> f=open('ccc.txt','rU') >>> text=f.read() >>> text1=text.split() >>> abstracts=nltk.Text(text1) >>> </pre>	<p>Modul corpus aus Package nltk</p> <p>Corpus Gutenberg</p> <p>Pfad des Ordners gutenberg im Ordner corpora</p> <pre>>>> from nltk.corpus import gutenberg >>> print gutenberg.root /home/benzro/nltk_data/corpora/gutenberg Dateien des Ordners gutenberg >>> gutenberg.fileids() ['austen-emma.txt', 'austen-persuasion.txt', 'austen-sense.txt', 'bible-kjv.txt', 'blake-poems.txt', 'bryant-stories.txt', 'burgess-busterbrown.txt', 'carroll-alice.txt', 'chesterton-ball.txt', 'chesterton-brown.txt', 'chesterton-thursday.txt', 'edgeworth-parents.txt', 'melville-moby_dick.txt', 'milton-paradise.txt', 'shakespeare-caesar.txt', 'shakespeare-hamlet.txt', 'shakespeare-macbeth.txt', 'whitman-leaves.txt']</pre> <p>Corpus Brown</p> <pre>>>> from nltk.corpus import brown >>> print brown.root /home/benzro/nltk_data/corpora/brown >>> print brown.fileids() [u'ca01', u'ca02', u'ca03', u'ca04', u'ca05', u'ca06', u'ca07', u'ca08', u'ca09...']</pre> <p>Keine Ordner sondern Te</p>
<h3>Funktionen definieren und aufrufen</h3>	<h3>Mehrere return-Anweisungen</h3>
<p>Definition der einstellige Funktion foo()</p> <pre>def foo(a): """ Return foo of a """ result = 0 for item in a: result += item return result</pre>	<p>Effekt der return-Anweisung</p> <pre>def describe_number(n): if n > 1000000: return "Large" elif n > 1000: return "Medium" else: return "small" print "Never printed!"</pre>
<p>return-Statement ergibt Rückgabewert, d.h. den Funktionswert. Die Parameter der Funktion (Platzhalter für Argumente) stehen in Klammern. Die erste Zeichenkette im Rumpf ist der Doc-String.</p> <p>Funktionsaufruf (call)</p> <pre>c = foo([5,10,23])</pre>	<ul style="list-style-type: none"> Verarbeitung der return-Anweisung beendet die Ausführung der Funktion. Beliebige Objekte können als Funktionswert zurückgegeben werden, auch Listen.

Wann und wozu sind Funktionen gut?	Skopus (Erreichbarkeit) von Variablennamen
<pre>def foo(a): """ Return foo of a """ result = 0 for item in a: result += item return result</pre>	<p>Modulweit erreichbare globale Variablen (<i>globals</i>)</p> <ul style="list-style-type: none"> ▶ (Variablen-)Namen, die in einem Modul zugewiesen werden, sind danach im ganzen Modul erreichbar. ▶ Modul <i>foo</i> ist Python-Quellcode aus Datei <i>foo.py</i>.
<p>Heilmittel gegen Spaghettikode</p> <ul style="list-style-type: none"> ▶ Abstraktion: Eine Funktion kann einige Zeilen Kode bezeichnen, welche oft gebraucht werden. ▶ Schnittstelle: Die Parameter einer Funktion machen den Kode an ganz bestimmten Stellen variabel (= Parametrisierung). ▶ Klarheit: Eine gute Funktion hat eine klar beschreibbare (=spezifizierbare) Funktionalität. 	<p>Funktionsweit erreichbare lokale Variablen (<i>locals</i>)</p> <ul style="list-style-type: none"> ▶ Parameter <i>a</i> und <i>b</i> einer Funktion <i>foo(a,b)</i>, sind nur innerhalb der Funktion <i>foo()</i> erreichbar. ▶ (Variablen)-Namen, die in einer Funktion definiert werden, sind nur innerhalb der Funktion erreichbar.
<pre># Listenbildung mit iterativen Statements sl = list() for c in "St. Moritz-Str. 23": if c.isalnum(): sl.append(c.lower()) print sl # Listenbildung mit Listenkomprehensionsausdruck el = [c.lower() for c in "St. Moritz-Str. 23" if c.isalnum()] print el</pre>	<pre>PCL-1/Vorlesung/vl09/iter_statement_expression.py" ['S', 't', ' ', 'M', 'o', 'r', 'i', 'z', ' ', 'S', 't', ' ', 'r', 'i', 'z', ' ', 's', 't', ' ', 'r', 'i', 'z', ' ', '2', ' ', '3'] ['S', 't', ' ', 'M', 'o', 'r', 'i', 'z', ' ', 'S', 't', ' ', 'r', 'i', 'z', ' ', 's', 't', ' ', 'r', 'i', 'z', ' ', '2', ' ', '3'] Process finished with exit code 0</pre>
<pre># Als Statement, Anweisung sl = [] for c in "St. Moritz-Str. 23": if c.isalnum(): sl.append(c) else: sl.append(' ') print sl # Als Expression, Ausdruck el = [c if c.isalnum() else ' ' for c in "St. Moritz-Str. 23"] print el</pre>	<pre>PCL-1/Vorlesung/vl09/if_then_else_statement_vs_expression.py" ['S', 't', ' ', 'M', 'o', 'r', 'i', 'z', ' ', 'S', 't', ' ', 'r', 'i', 'z', ' ', 's', 't', ' ', 'r', 'i', 'z', ' ', '2', ' ', '3'] ['S', 't', ' ', 'M', 'o', 'r', 'i', 'z', ' ', 'S', 't', ' ', 'r', 'i', 'z', ' ', 's', 't', ' ', 'r', 'i', 'z', ' ', '2', ' ', '3'] Process finished with exit code 0</pre>
<pre># Demonstration von Wichtigem from nltk.book import * #Set words = set(text1) print "\n", len(words), "\n" #print words[:20] #set hat kein Attribut __getitem__ #String print ', '.join(words)[:20], "\n" #ersten 20 chars #List print list(words)[:20], "\n" #ersten 20 list elements #List longwords = [w for w in words if len(w)>15] print longwords #List l = [] l.append(3) # ein Element anhängen l.extend((4, 'x', 5)) # eine ganze Sequenz anhängen print l del l[3] # ein Element löschen l[2] = 2 # eine Element austauschen l.sort(reverse=True) # in-place umgekehrt sortieren print l</pre>	<p>19317 funereal, unscientific [u'funereal', u'unscientific', u'divinely', u'foul', u'four', u'gag', u'prefix', u'woods', u'clotted', u'Duck', u'hanging', u'plaudits', u'woody', u'Until', u'marching', u'disobeying', u'canes', u'granting', u'advantage', u'Westers'] [u'hermaphroditical', u'subterraneousness', u'uninterpenetratingly', u'irresistibleness', u'responsibilities', u'comprehensiveness', u'uncompromisedness', u'superstitiousness', u'uncomfortableness', u'supernaturalness', u'circumnavigating', u'characteristically', u'cannibalistically', u'circumnavigations', u'indispensableness', u'preternaturalness', u'apprehensiveness', u'CIRCUMNAVIGATION', u'simultaneousness', u'undiscriminating', u'Physiognomically', u'physiognomically', u'circumnavigation', u'indiscriminately'] [3, 4, 'x', 5] [4, 3, 2]</p>

```

import re

teststr = '''This
is a
Test...'''

# Traditionelle Funktionsdefinition
def sf(s):
    # \s is any white space character
    # and equivalent to [ \t \n \r \f]
    return re.sub(r'\s+', '', s)
print sf(teststr)

def isf(s):
    return re.sub(r'[ \t \n \r \f]+', '', s)
print isf(teststr)

# Funktionsdefinition mit Lambda-Ausdruck
ef = lambda s: re.sub(r'\s+', '', s)
print ef(teststr)

# Was passiert, wenn kein return Statement existiert?
def sf2(s):
    re.sub(r'\s+', '', s)
print sf2(teststr)

# Was passiert, wenn ich schreibe:
# SyntaxError: invalid syntax -> return
# ef = lambda s: return re.sub(r'\s+', '', s)
# print sf(teststr)

```

```

import nltk
from nltk.corpus import gutenberg

#Load text
emma_words = gutenberg.words('austen-emma.txt')
print "\n", type(emma_words), "\n", len(emma_words)

#Calculate frequency distribution of words
emma_fd = nltk.FreqDist(emma_words)
print "\n", type(emma_fd), "\n", emma_fd

#Print frequency distribution of words
emma_fd_sorted=sorted(emma_fd, key=emma_fd.get, reverse=True )
print "\n", type(emma_fd_sorted)," \n", emma_fd_sorted[:20], "\n"
for w in emma_fd_sorted[:5]:
    print w, "\t", emma_fd[w]

#Tabuliere die häufigsten 20 Wörter
print "\n", type(emma_fd.tabulate)
emma_fd.tabulate(20)

#Plot frequency distribution of words
emma_fd.plot(10, cumulative=False)

# Finde alle Wörter für die gilt:
# - mehr als 10 Buchstaben und
# - kommen mindestens 7 mal vor
# - printe nur die ersten 20
wl = sorted([w for w in emma_fd.keys()
            if len(w)>10 and emma_fd[w]> 7])
print "\n", wl[:20], "\n"

```

```

import nltk
from nltk.corpus import brown

#Conditional Frequency Distribution
cfд = nltk.ConditionalFreqDist([
    (genre, word)
    for genre in brown.categories()
    for word in brown.words(categories=genre)])
print "\n", type(cfд), "\n", cfд, "\n"

#Wie oft kommen die modals (samples)
# conditional in den genres (conditions) vor
genres = ['news', 'religion', 'hobbies',
          'science_fiction', 'romance', 'humor']
modals = ['can', 'could', 'may', 'might', 'must', 'will']
cfд.tabulate(conditions=genres, samples=modals)

```

```

'-----'
PCL-1/Vorlesung/vl09/functions_statement_vs_expression.py"
ThisisaTest.
ThisisaTest.
ThisisaTest.
None

```

```
Process finished with exit code 0
```

```

/LubuntuShared/8) PCL-1/Vorlesung/vl09/freqdist_emma.py"
<class 'nltk.corpus.reader.util.StreamBackedCorpusView'>
192427

<class 'nltk.probability.FreqDist'>
<FreqDist with 7811 samples and 192427 outcomes>

<type 'list'>
[u',', u'.', u'to', u'the', u'and', u'of', u'I', u'a', u'was',
 u'her', u';', u'it', u'in', u'not', u'''', u've', u'she',
 u'that', u'you', u'had']

,           11454
.           6928
to          5183
the         4844
and         4672

<type 'instancemethod'>
    . to the and of I a was her ; it
    in not " be she that you had
11454 6928 5183 4844 4672 4279 3178 3004 2385 2381 2199 2128
2118 2101 2004 1970 1778 1730 1677 1606

[u'accomplished', u'acknowledge', u'acknowledged',
 u'acquaintance', u'affectionate', u'approbation',
 u'arrangement', u'circumstance', u'circumstances',
 u'comfortable', u'comfortably', u'communicated',
 u'communication', u'compliments', u'concealment',
 u'consciousness', u'consequence', u'considerable',
 u'consideration', u'considering']

```

```

/LubuntuShared/8) PCL-1/Vorlesung/vl09/CondFreqDist_brown.py"
<class 'nltk.probability.ConditionalFreqDist'>
<ConditionalFreqDist with 15 conditions>

      can could may might must will
news   93   86   66   38   50   389
religion 82   59   78   12   54   71
hobbies 268   58  131   22   83   264
science_fiction 16   49   4   12   8   16
romance 74  193   11   51   45   43
humor   16   30   8   8   9   13

Process finished with exit code 0

```

Example	Description
<code>fdist = FreqDist(samples)</code>	Create a frequency distribution containing the given samples
<code>fdist.inc(sample)</code>	Increment the count for this sample
<code>fdist['monstrous']</code>	Count of the number of times a given sample occurred
<code>fdist.freq('monstrous')</code>	Frequency of a given sample
<code>fdist.N()</code>	Total number of samples
<code>fdist.keys()</code>	The samples sorted in order of decreasing frequency
<code>for sample in fdist:</code>	Iterate over the samples, in order of decreasing frequency
<code>fdist.max()</code>	Sample with the greatest count
<code>fdist.tabulate()</code>	Tabulate the frequency distribution
<code>fdist.plot()</code>	Graphical plot of the frequency distribution
<code>fdist.plot(cumulative=True)</code>	Cumulative plot of the frequency distribution
<code>fdist1 < fdist2</code>	Test if samples in fdist1 occur less frequently than in fdist2

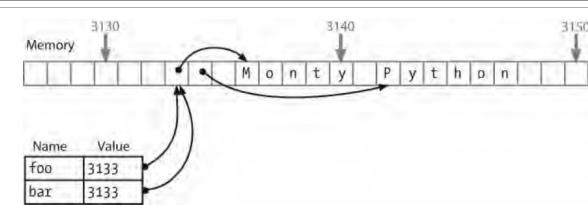


Figure 4-1. List assignment and computer memory: Two list objects `foo` and `bar` reference the same location in the computer's memory; updating `foo` will also modify `bar`, and vice versa.

Assignment always copies the value of an expression, but a value is not always what you might expect it to be. In particular, the “value” of a structured object such as a list is actually just a *reference* to the object. In the following example, ❶ assigns the reference of `foo` to the new variable `bar`. Now when we modify something inside `foo` on line ❷, we can see that the contents of `bar` have also been changed.

```
>>> foo = ['Monty', 'Python']
>>> bar = foo ❶
>>> foo[1] = 'Bodkin' ❷
>>> bar
['Monty', 'Bodkin']
```

The line `bar = foo` ❶ does not copy the contents of the variable, only its “object reference.” To understand what is going on here, we need to know how lists are stored in the computer’s memory. In Figure 4-1, we see that a list `foo` is a reference to an object stored at location 3133 (which is itself a series of pointers to other locations holding strings). When we assign `bar = foo`, it is just the object reference 3133 that gets copied.

Let’s experiment some more, by creating a variable `empty` holding the empty list, then using it three times on the next line.

```
>>> empty = []
>>> nested = [empty, empty, empty]
>>> nested
[[], [], []]
>>> nested[1].append('Python')
>>> nested
[[[], ['Python'], []], ['Python'], ['Python']]
>>> nested[1] = ['Monty']
>>> nested
[[['Python'], ['Monty'], ['Python']]
```

Observe that changing one of the items inside our nested list of lists changed them all. This is because each of the three elements is actually just a reference to one and the same list in memory.

We began with a list containing three references to a single empty list object. Then we modified that object by appending ‘Python’ to it, resulting in a list containing three references to a single list object ['Python']. Next, we **overwrote one of those references with a reference to a new object ['Monty']**. This last step modified one of the three object references inside the nested list. However, the ['Python'] object wasn’t changed, and is still referenced from two places in our nested list of lists. It is crucial to appreciate this difference between modifying an object via an object reference (`.append`) and overwriting an object reference (`=`).

Important: To copy the items from a list `foo` to a new list `bar`, you can write `bar = foo[:]`. This copies the object references inside the list. To copy a structure without copying any object references, use `copy.deepcopy()`.

```
>>> ilist[1, 2, 3]>>> ilis 2=ilist>>> ilist3=ilist[:]>>> ilist4=
from copy import deepcopy
ilist1 = ['a', 'b', ['ab', 'ba']]
```

```
>>> nested = [[] * 3
>>> nested[1].append('Pythons')
>>> nested[1].append('Pythons')
>>> nested[1].append('Pythons')
>>> nested
[['Pythons', 'Pythons', 'Pythons'], ['Pythons', 'Pythons', 'Pythons'],
['Pythons', 'Pythons', 'Pythons']]
>>> nested[1] = ['Monty']
>>> nested[1].append('Pythons')
>>> nested
[['Pythons', 'Pythons', 'Pythons'], ['Monty', 'Pythons', ['Pythons',
'Pythons']],
['Pythons', 'Pythons']]
>>> nested[0].append('Monty')
>>> nested[0].append('Montys')
>>> nested[0].append('Montys')
>>> nested
[['Pythons', 'Pythons', 'Pythons', 'Montys', 'Montys', ['Monty',
'Pythons']], ['Pythons', 'Pythons', 'Pythons', 'Montys', 'Montys'],
['Pythons', 'Pythons']]
```

Python provides two ways to check that a pair of items are the same. The `is` operator tests for object identity. We can use it to verify our earlier observations about objects. First, we create a list containing several copies of the same object, and demonstrate that they are not only identical according to `==`, but also that they are one and the same object:

```
>>> size = 5
>>> python = ['Python']
>>> snake_nest = [python] * size
>>> snake_nest[0] == snake_nest[1] == snake_nest[2] == snake_nest[3] == snake_nest[4]
True
>>> snake_nest[0] is snake_nest[1] is snake_nest[2] is snake_nest[3] is snake_nest[4]
True
```

```
>>> snake_nest = [python] * 5
>>> snake_nest
[['Python'], ['Python'], ['Python'], ['Python'], ['Python']]
>>> python = ['Python']
>>> snake_nest = [python] * 5
>>> snake_nest
[['Python'], ['Python'], ['Python'], ['Python'], ['Python']]
>>> snake_nest[3]=[Cobra']
>>> snake_nest
[['Python'], ['Python'], ['Python'], ['Cobra'], ['Python']]
>>> snake_nest[1].append('Boa')
>>> snake_nest
[['Python', 'Boa'], ['Python', 'Boa'], ['Python', 'Boa'], ['Python', 'Boa'],
['Python', 'Boa']]
>>> [id(snake) for snake in snake_nest]
[2942339148L, 2942339148L, 2942339148L, 2941265932L, 2942339148L]
```

<p>Conditionals</p> <p>In the condition part of an if statement, a non-empty string or list is evaluated as true, while an empty string or list evaluates as false.</p> <pre>>>> mixed = ['cat', '', ['dog'], []] >>> for element in mixed: ... if element: ... print element ... cat ['dog']</pre> <p>That is, we don't need to say if <code>len(element) > 0</code>: in the condition.</p>	<pre>>>> sent = ['No', 'good', 'fish', 'goes', 'anywhere', 'without', 'a', 'porpoise', '.'] >>> all(len(w) > 4 for w in sent) False >>> any(len(w) > 4 for w in sent) True >>> [(len(w) > 4 for w in sent)] [<generator object <genexpr> at 0xb653e43c>] >>> [len(w) > 4 for w in sent] [False, False, False, True, True, False, True]</pre>														
<p>Caution!</p> <p>Tuples are constructed using the comma operator. Parentheses are a more general feature of Python syntax, designed for grouping. A tuple containing the single element 'snark' is defined by adding a trailing comma, like this: 'snark'. The empty tuple is a special case, and is defined using empty parentheses () .</p>	<pre>>>> first_list = [1,2,3,4] my_set=set(first_list) my_list = list(my_set) >>> first_list [1, 2, 3, 4] >>> my_set set([1, 2, 3, 4]) >>> my_list [1, 2, 3, 4]</pre>														
<p><i>Table 4-1. Various ways to iterate over sequences</i></p> <table border="1" data-bbox="160 743 768 893"> <thead> <tr> <th>Python expression</th> <th>Comment</th> </tr> </thead> <tbody> <tr> <td><code>for item in s</code></td> <td>Iterate over the items of s</td> </tr> <tr> <td><code>for item in sorted(s)</code></td> <td>Iterate over the items of s in order</td> </tr> <tr> <td><code>for item in set(s)</code></td> <td>Iterate over unique elements of s</td> </tr> </tbody> </table>	Python expression	Comment	<code>for item in s</code>	Iterate over the items of s	<code>for item in sorted(s)</code>	Iterate over the items of s in order	<code>for item in set(s)</code>	Iterate over unique elements of s	<p>Python expression Comment</p> <hr/> <table border="1" data-bbox="810 743 1275 810"> <tbody> <tr> <td><code>for item in reversed(s)</code></td> <td>Iterate over elements of s in reverse</td> </tr> <tr> <td><code>for item in set(s).difference(t)</code></td> <td>Iterate over elements of s not in t</td> </tr> <tr> <td><code>for item in random.shuffle(s)</code></td> <td>Iterate over elements of s in random order</td> </tr> </tbody> </table> <p>The sequence functions illustrated in Table 4-1 can be combined in various ways; for example, to get unique elements of s sorted in reverse, use <code>reversed(sorted(set(s)))</code>.</p>	<code>for item in reversed(s)</code>	Iterate over elements of s in reverse	<code>for item in set(s).difference(t)</code>	Iterate over elements of s not in t	<code>for item in random.shuffle(s)</code>	Iterate over elements of s in random order
Python expression	Comment														
<code>for item in s</code>	Iterate over the items of s														
<code>for item in sorted(s)</code>	Iterate over the items of s in order														
<code>for item in set(s)</code>	Iterate over unique elements of s														
<code>for item in reversed(s)</code>	Iterate over elements of s in reverse														
<code>for item in set(s).difference(t)</code>	Iterate over elements of s not in t														
<code>for item in random.shuffle(s)</code>	Iterate over elements of s in random order														
<p>We can convert between these sequence types. For example, <code>tuple(s)</code> converts any kind of sequence into a tuple, and <code>list(s)</code> converts any kind of sequence into a list. We can convert a list of strings to a single string using the <code>join()</code> function, e.g., <code>'.'.join(words)</code>.</p> <p>Some other objects, such as a <code>FreqDist</code>, can be converted into a sequence (using <code>list()</code>) and support iteration:</p> <pre>>>> raw = 'Red lorry, yellow lorry, red lorry, yellow lorry.' >>> text = nltk.word_tokenize(raw) >>> fdist = nltk.FreqDist(text) >>> list(fdist) ['lorry', ',', 'yellow', '.', 'Red', 'red'] >>> for key in fdist: ... print fdist[key], ... 4 3 2 1 1 1</pre>	<p>In the next example, we use tuples to re-arrange the contents of our list. (We can omit the parentheses because the comma has higher precedence than assignment.)</p> <pre>>>> words = ['I', 'turned', 'off', 'the', 'spectroroute'] >>> words[2], words[3], words[4] = words[3], words[4], words[2] >>> words ['I', 'turned', 'the', 'spectroroute', 'off']</pre>														
<p>There are also functions that modify the structure of a sequence, which can be handy for language processing. Thus, <code>zip()</code> takes the items of two or more sequences and "zips" them together into a single list of pairs. Given a sequence s, <code>enumerate(s)</code> returns pairs consisting of an index and the item at that index.</p> <pre>>>> words = ['I', 'turned', 'off', 'the', 'spectroroute'] >>> tags = ['noun', 'verb', 'prep', 'det', 'noun'] >>> zip(words, tags) [(I, 'noun'), ('turned', 'verb'), ('off', 'prep'), ('the', 'det'), ('spectroroute', 'noun')] >>> list(enumerate(words)) [(0, I), (1, 'turned'), (2, 'off'), (3, 'the'), (4, 'spectroroute')]</pre>	<pre>>>> text = nltk.corpus.nps_chat.words() >>> cut = int(0.9 * len(text)) ❶ >>> training_data, test_data = text[:cut], text[cut:] ❷ >>> text == training_data + test_data ❸ True >>> len(training_data) / len(test_data) ❹ 9</pre>														
<pre>>>> words = 'I turned off the spectroroute'.split() >>> wordlens = [(len(word), word) for word in words] >>> wordlens [(1, 'I'), (6, 'turned'), (3, 'off'), (3, 'the'), (12, 'spectroroute')] >>> wordlens.sort() >>> wordlens [(1, 'I'), (3, 'off'), (3, 'the'), (6, 'turned'), (12, 'spectroroute')] >>> ' '.join(w for _, w in wordlens) 'I off the turned spectroroute' >>> ' '.join({w:k} for (k, w) in wordlens) Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: sequence item 0: expected string, dict found >>> ' '.join(w*k for (k, w) in wordlens) 'I offoffoff the the turnedturnturnturnturnturnturnturn spectroroute spectroroute spectroroute spectroroute spectroroute spectroroute spectroroute spect'</pre>	<p>Here, a lexicon is represented as a list because it is a collection of objects of a single type—lexical entries—of no predetermined length. An individual entry is represented as a tuple because it is a collection of objects with different interpretations, such as the orthographic form, the part-of-speech, and the pronunciations (represented in the SAMPA computer-readable phonetic alphabet; see http://www.phon.ucl.ac.uk/home/sampa/). Note that these pronunciations are stored using a list. (Why?)</p> <pre>>>> lexicon = [... ('the', 'det', ['Di:', 'D@']), ... ('off', 'prep', ['Of', 'O:f']) ...]</pre>														

```
>>> max([w.lower() for w in nltk.word_tokenize(text)]) ❶
'word'
>>> max(w.lower() for w in nltk.word_tokenize(text)) ❷
'word'
```

The second line uses a **generator expression**. This is more than a notational convenience: in many language processing situations, generator expressions will be more efficient. In ❶, storage for the list object must be allocated before the value of `max()` is computed. If the text is very large, this could be slow. In ❷, the data is streamed to the calling function. Since the calling function simply has to find the maximum value—the word that comes latest in lexicographic sort order—it can process the stream of data without having to store anything more than the maximum value seen so far.

Finally, let's combine the idioms we've been exploring. First, we create a list of *cie* and *cei* words, then we loop over each item and print it. Notice the comma at the end of the print statement, which tells Python to produce its output on a single line.

```
>>> tricky = sorted([w for w in set(text2) if 'cie' in w or 'cei' in w])
>>> for word in tricky:
...     print word,
ancient ceiling conceit conceited conceive conscience
conscientious conscientiously deceitful deceive ...
```

There is a common pattern to all of these examples: `[w for w in text if condition]`, where `condition` is a Python “test” that yields either true or false. In the cases shown in the previous code example, the condition is always a numerical comparison. However, we can also test various properties of words, using the functions listed in [Table 1-4](#).

Table 1-4. Some word comparison operators

Function	Meaning
<code>s.startswith(t)</code>	Test if s starts with t
<code>s.endswith(t)</code>	Test if s ends with t
<code>t in s</code>	Test if t is contained inside s
<code>s.islower()</code>	Test if all cased characters in s are lowercase
<code>s.isupper()</code>	Test if all cased characters in s are uppercase
<code>s.isalpha()</code>	Test if all characters in s are alphabetic
<code>s.isalnum()</code>	Test if all characters in s are alphanumeric
<code>s.isdigit()</code>	Test if all characters in s are digits
<code>s.istitle()</code>	Test if s is titlecased (all words in s have initial capitals)

Gutenberg-Projekte: Elektronische Edition älterer Texte

Definition (Korpus (sächlich, im Plural Korpora))

Ein Korpus ist eine Sammlung von Texten.

Sammlung vorwiegend englischsprachiger Bücher

 Sammlung von über 50'000 frei verfügbaren Büchern, deren Copyright abgelaufen ist in den USA.
<http://www.gutenberg.org>

Sammlung deutschsprachiger Bücher

 Sammlung von über 7'000 frei verfügbaren Büchern, deren Copyright abgelaufen ist in Deutschland. D.h. 70 Jahre nach dem Tod des Autors oder Übersetzers.
<http://gutenberg.spiegel.de>

Das Package `nltk.corpus`

Enthält Packages und Module, mit denen Korpora in verschiedenen Formaten eingelesen werden können.

Das Objekt `nltk.corpus.gutenberg`

Stellt eine Auswahl von 18 englischsprachigen Gutenberg-Texten (ASCII) als Teil der NLTK-Korpusdaten zum Einlesen zur Verfügung.

Funktionen des Objekts `nltk.corpus.gutenberg`

Repräsentationen für reine Text-Korpora (ASCII oder iso-latin-1) ▶1

```
from nltk.corpus import gutenberg

filename = 'austen-emma.txt' # oder absoluter Pfad einer Textdatei

# Korpus als eine lange Zeichenkette
emma_chars = gutenberg.raw(filename)

# Korpus als Liste von Wörtern (Wort ist Zeichenkette)
emma_words = gutenberg.words(filename)

# Korpus als Liste von Sätzen (Satz ist Liste von Wörtern)
emma_sents = gutenberg.sents(filename)

# Korpus als Liste von Paragraphen (Paragraph ist Liste von Sätzen)
emma_paras = gutenberg.paras(filename)
```

Methoden des Objekts `nltk.corpus.brown`²

Zusätzlich zu den Funktionen von Textkorpora, gibts Listen mit Paaren aus einem Token und seinem POS-Tag.

Repräsentationen für getaggte Korpora

```
from nltk.corpus import brown

# Korpus als Liste von 2-Tupeln (Wort, POS-Tag)
brown_tagged_words = brown.tagged_words()
```

Eigenheiten des Brownkorpus: Unterschiedliche Textsorten

```
# Das balancierte Korpus umfasst Texte aus 15 Kategorien
brown.categories()
```

Arten von Korpora: Korpus-Typologie

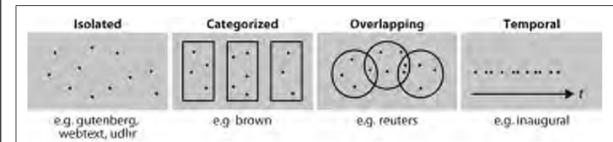
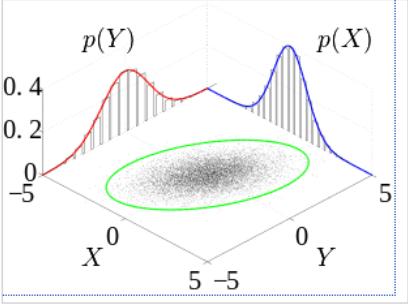


Figure 2-3. Common structures for text corpora: The simplest kind of corpus is a collection of isolated texts with no particular organization; some corpora are structured into categories, such as genre (Brown Corpus); some categorizations overlap, such as topic categories (Reuters Corpus); other corpora represent language use over time (Inaugural Address Corpus).

Quelle: [Bain et al., 2009, 30]

- Die Texte in einem Korpus (Textsammlung, d.h. mehrere Texte) können in unterschiedlicher Ordnung zueinander stehen.
- Ein Korpus kann **balanciert** (repräsentativ zusammengestellt) oder **opportunistisch** (nimm, was du kannst!) sein.

<p style="text-align: center;">Word Tally</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>the</td><td> - - - </td></tr> <tr><td>been</td><td> - - </td></tr> <tr><td>message</td><td> </td></tr> <tr><td>persevere</td><td> </td></tr> <tr><td>nation</td><td> - </td></tr> </table> <ul style="list-style-type: none"> ▶ Deskriptive Statistiken: Fundamentale Funktionalität in korpuslinguistischen Auswertungen ▶ Buchstaben, Wörter, Eigennamen, Sätze, Paragraphen zählen ▶ Minima, Maxima und Mittelwerte ermitteln ▶ Verteilung der Häufigkeiten darstellen (Tabelle, Plots) ▶ Letztlich: Verteilungen vergleichen 	the	- - -	been	- -	message		persevere		nation	-	<h3>Häufigkeitsverteilungen als allgemeine Datenstruktur</h3> <ul style="list-style-type: none"> ▶ Allgemein: Häufigkeit der Items einer variierenden Grösse (eine statistische Variable) auszählen ▶ NLTK-Klasse <code>nltk.FreqDist</code> ist eine Datenstruktur zum einfachen Erstellen von Häufigkeitsverteilungen (<i>frequency distribution</i>) <p>(Abstrakte) Datenstrukturen</p> <p>"In der Informatik und Softwaretechnik ist eine Datenstruktur ein Objekt zur Speicherung und Organisation von Daten. Es handelt sich um eine Struktur, weil die Daten in einer bestimmten Art und Weise angeordnet und verknüpft werden, um den Zugriff auf sie und ihre Verwaltung effizient zu ermöglichen. Datenstrukturen sind nicht nur durch die enthaltenen Daten charakterisiert, sondern vor allem durch die Operationen auf diesen Daten, die Zugriff und Verwaltung ermöglichen und realisieren." (http://de.wikipedia.org/wiki/Datenstruktur)</p>												
the	- - -																						
been	- -																						
message																							
persevere																							
nation	-																						
<p>www.nltk.org/api/nltk.html#nltk.probability.FreqDist.N</p> <pre>>>> from nltk.tokenize import word_tokenize >>> from nltk.probability import FreqDist >>> sent = 'This is an example sentence' >>> fdist = FreqDist() >>> for word in word_tokenize(sent): ... fdist[word.lower()] += 1</pre> <p>An equivalent way to do this is with the initializer:</p> <pre>>>> fdist = FreqDist(word.lower() for word in word_tokenize(sent))</pre> <p>Formally, a frequency distribution can be defined as a function mapping from each sample to the number of times that sample occurred as an outcome. Frequency distributions are generally constructed by running a number of experiments, and incrementing the count for a sample every time it is an outcome of an experiment. For example, the following code will produce a frequency distribution that encodes how often each word occurs in a text:</p>	<h3>Berechnen der häufigsten längsten Wörter</h3> <pre>import nltk from nltk.corpus import gutenberg emma_words = gutenberg.words('austen-emma.txt') emma_fd = nltk.FreqDist(emma_words) # Finde alle Wörter für die gilt: # - mehr als 10 Buchstaben und # - kommen mindestens 7 mal vor wl = sorted([w for w in emma_fd.keys() if len(w)>10 and emma_fd[w]> 7])</pre>																						
<h3>Bivariate (bedingte) Häufigkeitsverteilungen</h3> <table border="1" style="display: inline-table; vertical-align: top;"> <tr><td>Condition: News</td><td>Condition: Romance</td></tr> <tr><td>the</td><td> - - - </td></tr> <tr><td>cute</td><td> - - </td></tr> <tr><td>Monday</td><td> - - </td></tr> <tr><td>could</td><td> </td></tr> <tr><td>will</td><td> - </td></tr> </table> <table border="1" style="display: inline-table; vertical-align: top;"> <tr><td>the</td><td> - - </td></tr> <tr><td>cute</td><td> </td></tr> <tr><td>Monday</td><td> </td></tr> <tr><td>could</td><td> - - </td></tr> <tr><td>will</td><td> </td></tr> </table> <p>Gemeinsame Häufigkeit der Items von 2 variierenden Größen (zweier statistischer Variable) auszählen</p> <ul style="list-style-type: none"> ▶ Sprechweise: Die eine Variable heisst in NLTK Bedingung (<i>condition</i>), die andere Ereignis (<i>event, sample</i>) ▶ Eine bedingte Häufigkeitsverteilung besteht pro Bedingung aus einer einfachen Häufigkeitsverteilung. ▶ NLTK-Klasse <code>nltk.ConditionalFreqDist</code> umfasst geeignete Methoden für Frequenzdistributionen von Paaren (=2er-Tupel): (<i>condition, sample</i>) ▶ Beispiel: Mit den 15 Kategorien im Brownkorpus ergeben sich 15 Bedingungen mit insgesamt 1'161'192 Events (Wörtern). 	Condition: News	Condition: Romance	the	- - -	cute	- -	Monday	- -	could		will	-	the	- -	cute		Monday		could	- -	will		<p>Joint probability distribution</p> <p>In the study of probability, given at least two random variables X, Y, \dots, that are defined on a probability space, the joint probability distribution for X, Y, \dots is a probability distribution that gives the probability that each of X, Y, \dots falls in any particular range or discrete set of values specified for that variable. In the case of only two random variables, this is called a bivariate distribution, but the concept generalizes to any number of random variables, giving a multivariate distribution.</p> <p>The joint probability distribution can be expressed either in terms of a joint cumulative distribution function or in terms of a joint probability density function (in the case of continuous variables) or joint probability mass function (in the case of discrete variables). These in turn can be used to find two other types of distributions: the marginal distribution giving the probabilities for any one of the variables with no reference to any specific ranges of values for the other variables, and the conditional probability distribution giving the probabilities for any subset of the variables conditional on particular values of the remaining variables.</p>
Condition: News	Condition: Romance																						
the	- - -																						
cute	- -																						
Monday	- -																						
could																							
will	-																						
the	- -																						
cute																							
Monday																							
could	- -																						
will																							

 <p>Many sample observations (black) are shown from a joint probability distribution. The marginal densities are shown as well.</p>	<h3>Bedingte Häufigkeiten berechnen</h3> <p>Modalverben in Abhängigkeit von Textkategorien</p> <pre>import nltk from nltk.corpus import brown cf = nltk.ConditionalFreqDist([(genre, word) for genre in brown.categories() for word in brown.words(categories=genre)])</pre> <pre>genres = ['news', 'religion', 'hobbies', 'science_fiction', 'romance', 'humor'] modals = ['can', 'could', 'may', 'might', 'must', 'will'] cf.tabulate(conditions=genres, samples=modals)</pre>																																				
<h3>Funktionen der Klasse nltk.ConditionalFreqDist</h3> <p>Table 2-4. NLTK's conditional frequency distributions: Commonly used methods and idioms for defining, accessing, and visualizing a conditional frequency distribution of counters</p> <table border="1"> <thead> <tr> <th>Example</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td><code>cfdist = ConditionalFreqDist(pairs)</code></td> <td>Create a conditional frequency distribution from a list of pairs</td> </tr> <tr> <td><code>cfdist.conditions()</code></td> <td>Alphabetically sorted list of conditions</td> </tr> <tr> <td><code>cfdist[condition]</code></td> <td>The frequency distribution for this condition</td> </tr> <tr> <td><code>cfdist[condition][sample]</code></td> <td>Frequency for the given sample for this condition</td> </tr> <tr> <td><code>cfdist.tabulate()</code></td> <td>Tabulate the conditional frequency distribution</td> </tr> <tr> <td><code>cfdist.tabulate(samples, conditions)</code></td> <td>Tabulation limited to the specified samples and conditions</td> </tr> <tr> <td><code>cfdist.plot()</code></td> <td>Graphical plot of the conditional frequency distribution</td> </tr> <tr> <td><code>cfdist.plot(samples, conditions)</code></td> <td>Graphical plot limited to the specified samples and conditions</td> </tr> <tr> <td><code>cfdist1 < cfdist2</code></td> <td>Test if samples in <code>cfdist1</code> occur less frequently than in <code>cfdist2</code></td> </tr> </tbody> </table>	Example	Description	<code>cfdist = ConditionalFreqDist(pairs)</code>	Create a conditional frequency distribution from a list of pairs	<code>cfdist.conditions()</code>	Alphabetically sorted list of conditions	<code>cfdist[condition]</code>	The frequency distribution for this condition	<code>cfdist[condition][sample]</code>	Frequency for the given sample for this condition	<code>cfdist.tabulate()</code>	Tabulate the conditional frequency distribution	<code>cfdist.tabulate(samples, conditions)</code>	Tabulation limited to the specified samples and conditions	<code>cfdist.plot()</code>	Graphical plot of the conditional frequency distribution	<code>cfdist.plot(samples, conditions)</code>	Graphical plot limited to the specified samples and conditions	<code>cfdist1 < cfdist2</code>	Test if samples in <code>cfdist1</code> occur less frequently than in <code>cfdist2</code>	<h3>Sequenz-Datentypen: list, str, unicode, tuple</h3> <p>Definition (Sequenz = Endliche Folge von Objekten)</p> <ul style="list-style-type: none"> ▶ Zugriff auf Elemente einer Sequenz mittels ganzzahligem Index: <code>s[i]</code> ▶ Zugriff auf Abschnitte (slice) mittels Angabe von Start- und exklusiver Endposition: <code>s[start:end]</code> ▶ Bestimmen der Anzahl Element mittels <code>len(s)</code> 																
Example	Description																																				
<code>cfdist = ConditionalFreqDist(pairs)</code>	Create a conditional frequency distribution from a list of pairs																																				
<code>cfdist.conditions()</code>	Alphabetically sorted list of conditions																																				
<code>cfdist[condition]</code>	The frequency distribution for this condition																																				
<code>cfdist[condition][sample]</code>	Frequency for the given sample for this condition																																				
<code>cfdist.tabulate()</code>	Tabulate the conditional frequency distribution																																				
<code>cfdist.tabulate(samples, conditions)</code>	Tabulation limited to the specified samples and conditions																																				
<code>cfdist.plot()</code>	Graphical plot of the conditional frequency distribution																																				
<code>cfdist.plot(samples, conditions)</code>	Graphical plot limited to the specified samples and conditions																																				
<code>cfdist1 < cfdist2</code>	Test if samples in <code>cfdist1</code> occur less frequently than in <code>cfdist2</code>																																				
<h3>Listen: Veränderliche (<i>mutable</i>) Sequenzen ►</h3> <p>Typische Modifikationen für Listen</p> <pre>l = [] # Zuweisung ist keine Listenmodifikation! l.append(1) # ein Element anhängen l.extend((4, 'x', 5)) # eine ganze Sequenz anhängen del l[3] # ein Element löschen l[2] = 3 # eine Element austauschen l.sort(reverse=True) # in-place rückwärts sortieren print l</pre> <ul style="list-style-type: none"> ▶ Nur bei Listen, d.h. veränderlichen Sequenzen, können Elemente (oder Abschnitte (<i>slices</i>)) gelöscht, ersetzt oder ergänzt werden. ▶ ⚡: Listen-Methoden, welche in-place-Modifikationen durchführen, liefern als Rückgabewert <code>None</code> zurück und nicht die Liste! 	<p>Set aus Liste mit unveränderlichen Elementen generieren->ok</p> <pre>>>> iset=set("m",1,"b","m")) >>> iset set(['m', 'b', 'm'])</pre> <p>Set aus Tupel (enthält selbst nur unveränderliche Elemente) generieren->ok</p> <pre>>>> iset=set(("m",1,"b","m")) >>> iset set((1, 'm', ('b', 'm')))</pre> <p>Set ohne umschließende Klammern um Elemente generieren->error</p> <pre>>>> iset=set("m",1,"b","m")) Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: set expected at most 1 arguments, got 3 Set mit Element Liste oder anderen veränderlichen Datentypen->error</pre> <pre>>>> iset=set(("m",1,"b","m")) Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: unhashable type: 'list' Set Index abfragen->error</pre> <pre>>>> iset[3] Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: 'set' object does not support indexing</pre>																																				
<h3>Syntaktischer Zucker für Methoden von Sequenzen</h3> <p>Python bietet für wichtige Methoden von Sequenzen Spezialnotation an. Ob die Spezialnotation funktioniert, hängt nur davon ab, ob mein Objekt die entsprechende Methode kann!</p> <table border="1"> <tr> <td data-bbox="171 1731 468 1866"> Enthalten (Membership) <pre>>>> 3 in [1,2,3] True >>> [1,2,3].__contains__(3) True</pre> </td> <td data-bbox="468 1731 785 1866"> i-tes Element <pre>>>> ('a','c')[1] 'c' >>> ('a','c').__getitem__(1) 'c'</pre> </td> </tr> <tr> <td data-bbox="171 1866 468 2034"> Abschnitt (Slicing) <pre>>>> "ABBA"[1:3] 'BB' >>> "ABBA".__getslice__(1,3) 'BB'</pre> </td> <td data-bbox="468 1866 785 2034"> <pre>>>> help(str.__getitem__) Help on wrapper_descriptor: __getitem__(...) x.__getitem__(y) <=> x[y]</pre> </td> </tr> </table>	Enthalten (Membership) <pre>>>> 3 in [1,2,3] True >>> [1,2,3].__contains__(3) True</pre>	i-tes Element <pre>>>> ('a','c')[1] 'c' >>> ('a','c').__getitem__(1) 'c'</pre>	Abschnitt (Slicing) <pre>>>> "ABBA"[1:3] 'BB' >>> "ABBA".__getslice__(1,3) 'BB'</pre>	<pre>>>> help(str.__getitem__) Help on wrapper_descriptor: __getitem__(...) x.__getitem__(y) <=> x[y]</pre>	<h3>Objektorientierte Modellierung</h3> <p>Kontoinhaber</p> <table border="1"> <tr> <td>Name</td> <td>Eigenschaften</td> </tr> <tr> <td>Vorname</td> <td></td> </tr> <tr> <td>Anschrift</td> <td></td> </tr> <tr> <td>Beruf</td> <td></td> </tr> <tr> <td>Geburtsdatum</td> <td></td> </tr> <tr> <td>Wohnort wechseln</td> <td>Methoden</td> </tr> <tr> <td>Beruf wechseln</td> <td></td> </tr> </table> <p>Konto</p> <table border="1"> <tr> <td>Kontoinhaber</td> <td></td> </tr> <tr> <td>Kontonummer</td> <td></td> </tr> <tr> <td>Kreditrahmen</td> <td></td> </tr> <tr> <td>Kontostand</td> <td></td> </tr> <tr> <td>Verfügungsberecht.</td> <td></td> </tr> <tr> <td>Einzahlen</td> <td></td> </tr> <tr> <td>Auszahlen</td> <td></td> </tr> <tr> <td>Überweisungen</td> <td></td> </tr> <tr> <td>Daueraufträge</td> <td></td> </tr> </table> <p>Mittels Klassen können eigene (Daten-)Typen geschaffen werden.</p> <p>Quelle: http://www Python-kurs.eu/Klassen.php</p>	Name	Eigenschaften	Vorname		Anschrift		Beruf		Geburtsdatum		Wohnort wechseln	Methoden	Beruf wechseln		Kontoinhaber		Kontonummer		Kreditrahmen		Kontostand		Verfügungsberecht.		Einzahlen		Auszahlen		Überweisungen		Daueraufträge	
Enthalten (Membership) <pre>>>> 3 in [1,2,3] True >>> [1,2,3].__contains__(3) True</pre>	i-tes Element <pre>>>> ('a','c')[1] 'c' >>> ('a','c').__getitem__(1) 'c'</pre>																																				
Abschnitt (Slicing) <pre>>>> "ABBA"[1:3] 'BB' >>> "ABBA".__getslice__(1,3) 'BB'</pre>	<pre>>>> help(str.__getitem__) Help on wrapper_descriptor: __getitem__(...) x.__getitem__(y) <=> x[y]</pre>																																				
Name	Eigenschaften																																				
Vorname																																					
Anschrift																																					
Beruf																																					
Geburtsdatum																																					
Wohnort wechseln	Methoden																																				
Beruf wechseln																																					
Kontoinhaber																																					
Kontonummer																																					
Kreditrahmen																																					
Kontostand																																					
Verfügungsberecht.																																					
Einzahlen																																					
Auszahlen																																					
Überweisungen																																					
Daueraufträge																																					

Typen/Klassen und Objekte		Typ-Aufrufe als Objekt-Konstruktoren								
Gattungen und Individuen in der Welt		Konstruiere Objekte von einem bestimmten Typ, indem du den Typ wie eine Funktion aufrufst!								
<table border="1"> <thead> <tr> <th>Gattung</th><th>Individuum</th></tr> </thead> <tbody> <tr> <td>Mensch</td><td>Elvis Presley</td></tr> <tr> <td>Hauptstadt</td><td>Paris</td></tr> </tbody> </table>		Gattung	Individuum	Mensch	Elvis Presley	Hauptstadt	Paris	Viele Konstruktor-Funktionen erlauben Argumente. Erklärungen gibt <code>help(type)</code> .		
Gattung	Individuum									
Mensch	Elvis Presley									
Hauptstadt	Paris									
Typen/Klassen und Objekte/Klassen-Instanzen in Python		Default-Objekte								
<table border="1"> <thead> <tr> <th>Typ/Klasse</th><th>Objekt/Instanz</th></tr> </thead> <tbody> <tr> <td>int</td><td>3</td></tr> <tr> <td>str</td><td>'abc'</td></tr> <tr> <td>list</td><td>[1,2,3]</td></tr> </tbody> </table>		Typ/Klasse	Objekt/Instanz	int	3	str	'abc'	list	[1,2,3]	<pre>>>> unicode() u'' >>> int() 0 >>> list() [] >>> dict() {} >>> set() set()</pre>
Typ/Klasse	Objekt/Instanz									
int	3									
str	'abc'									
list	[1,2,3]									
Wichtig: Objekte sind Instanzen eines Typs oder einer Klasse.		Konstruktoren mit Parametern								
Objekte konstruieren mittels Klassen/Typ-Konstruktor		<pre>>>> str(123) '123' >>> int('10110',2) 22 >>> set([3,3,2,2,'a',1.1,'a']) set(['a', 2, 3, 1.1]) >>> list(set([2,1,'a'])) ['a', 1, 2] >>> dict(a='DET',do='VB') {'a': 'DET', 'do': 'VB'}</pre>								
Konstruktoren (<i>constructors</i>) vs Literale (<i>literals, displays</i>)		8.3.1. Counter objects								
<ul style="list-style-type: none"> 💡 Nur die Objekte der wichtigsten eingebauten Datentypen können als Literale oder mit Spezialnotation konstruiert werden. ▶ Konstrukturen der Form <code>TYPE()</code> sind immer möglich! 		A counter tool is provided to support convenient and rapid tallies. For example:								
<table border="1"> <thead> <tr> <th>Typ/Klasse</th><th>Objekt/Instanz herstellen</th></tr> </thead> <tbody> <tr> <td>nltk.probability.FreqDist</td><td>nltk.FreqDist("abarakadabra")</td></tr> <tr> <td>collections.Counter</td><td>collections.Counter('abarakadabra')</td></tr> </tbody> </table>		Typ/Klasse	Objekt/Instanz herstellen	nltk.probability.FreqDist	nltk.FreqDist("abarakadabra")	collections.Counter	collections.Counter('abarakadabra')	<pre>>>> # Tally occurrences of words in a list >>> cnt = Counter() >>> for word in ['red', 'blue', 'red', 'green', 'blue', 'blue']: ... cnt[word] += 1 >>> cnt Counter({'blue': 3, 'red': 2, 'green': 1}) >>> # Find the ten most common words in Hamlet >>> import re >>> words = re.findall(r'\w+', open('hamlet.txt').read().lower()) >>> Counter(words).most_common(10) [('the', 1143), ('and', 966), ('to', 762), ('of', 669), ('i', 631), ('you', 554), ('a', 546), ('my', 514), ('hamlet', 471), ('in', 451)]</pre>		
Typ/Klasse	Objekt/Instanz herstellen									
nltk.probability.FreqDist	nltk.FreqDist("abarakadabra")									
collections.Counter	collections.Counter('abarakadabra')									
Hinweis: Die Klasse <code>nltk.probability.FreqDist</code> erweitert die eingebaute Python-Klasse <code>collections.Counter</code> ▶!		Ein Blick hinter die Kulisse: Methodenaufrufe								
Klassenbezeichner sind wichtig! Keine Objekte ohne Typen/Klassen!		Schein: Objekt ruft seine Methode auf OBJECT.METHOD(ARG1)								
Common patterns for working with <code>Counter</code> objects:		Sein: Klasse/Type ruft Methode mit Objekt als 1. Argument auf CLASS.METHOD(OBJECT, ARG1)								
<pre>sum(c.values()) # total of all counts c.clear() # reset all counts list(c) # list unique elements set(c) # convert to a set dict(c) # convert to a regular dictionary c.items() # convert to a list of (elem, cnt) pairs Counter(dict(list_of_pairs)) # convert from a list of (elem, cnt) pairs c.most_common()[:n-1:-1] # n least common elements c += Counter() # remove zero and negative counts</pre>		<table border="1"> <thead> <tr> <th>Schein</th><th>Sein</th></tr> </thead> <tbody> <tr> <td>"A Test".lower()</td><td>str.lower("A Test")</td></tr> <tr> <td>"ABBA".count('A')</td><td>str.count("ABBA", 'A')</td></tr> </tbody> </table>	Schein	Sein	"A Test".lower()	str.lower("A Test")	"ABBA".count('A')	str.count("ABBA", 'A')		
Schein	Sein									
"A Test".lower()	str.lower("A Test")									
"ABBA".count('A')	str.count("ABBA", 'A')									
Unterschied zwischen <i>Statements</i> und <i>Expressions</i>		Das macht VIEL Sinn! Die Addition definiert man auch auf der Klasse der Ganzzahlen, nicht für jede Zahl einzeln! Objekte rufen Methoden auf, welche auf Klassenebene definiert sind!								
Anweisungen (<i>statements</i>) ▶5		Listenbildung via Anweisungen und Ausdruck								
werden vom Python-Interpreter ausgeführt und evaluieren zu keinem Wert .		Listenbildung mit iterativen Statements								
<pre>print Statement ▶6 print "Something to print"</pre>		<pre>s1 = [] for c in "St. Moritz-Str. 23": if c.isalnum(): s1.append(c.lower())</pre>								
Ausdrücke (<i>expressions</i>) ▶7		Listenbildung mit einem Ausdruck: Listenkomprehension								
werden zu einem Wert (Objekt) evaluiert und enthalten keine Statements .		<pre>el = [c.lower() for c in "St. Moritz-Str. 23" if c.isalnum()]</pre>								
Boole'sche und andere Ausdrücke innerhalb von Statements										
<pre># If-Statement mit komplexen Ausdrücken drin if len("A "+"String") > 5: print "A "+"String".lower()</pre>										

If-then-else als Anweisung und If-Else als Ausdruck	Funktionsdefinition via Anweisungen und Ausdruck
<p>Listenbildung mit iterativen Statements</p> <pre>sl = [] for c in "St. Moritz-Str. 23": if c.isalnum(): sl.append(c) else: sl.append(' ')</pre> <p>Default-if-else Ausdruck</p> <pre>el = [c if c.isalnum() else ' ' for c in "St. Moritz-Str. 23"]</pre> <p>⌚: Abweichende Reihenfolge von if-then-else-Bestandteilen, da typischerweise der Then-Ausdruck der Standardwert ist.</p>	<p>Funktionsdefinition mit iterativem Statement ▶8</p> <pre>def sf(s): return re.sub(r'\s+', '', s)</pre> <p>Funktionsdefinition via Lambda-Ausdruck</p> <pre>ef = lambda s: re.sub(r'\s+', '', s)</pre> <p>Lambda-Ausdrücke (Lambda expression))</p> <p>Mathematische Notation zur Definition von anonymen Funktionen:</p> <ul style="list-style-type: none"> ▶ Funktionsdefinition (Rechenvorschrift): $(\lambda x : x + 1)$ ▶ Funktionsevaluation: $(\lambda x : x + 1)(3) = 4$ ▶ Lambda bindet/abstrahiert die Funktionsparameter im Funktionsrumpf ▶ Kurz: Parametrisierte Ausdrücke
Komprehension von Mengen und Dictionaries	<p>Tupel vs Set</p> <pre>#Tupel CAN contain mutable objects: >>> itupel = ([1, 2, 3], {3, 2, 1}) >>> itupel ([1, 2, 3], {3, 2, 1}) >>> itupel = ("m", 1, ["b", "m"]) >>> itupel ('m', 1, ['b', 'm']) #Set CANNOT contain mutable objects: >>> iset = set("m", 1, "b", "m")) Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: 'set' object is not iterable #Tupel ARE indexable >>> itupel[2][1] 'm' # Set ARE NOT indexable >>> iset[3] Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: 'set' object does not support indexing #Tupel CAN change mutable objects: >>> itupel[2][1] = "aha" >>> itupel ('m', 1, ['b', 'aha']) #Tupel CANNOT change immutable objects: >>> itupel[1][1] = "aha" Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: 'int' object does not support item assignment >>> del(itupel[1][1]) Traceback (most recent call last): File "<input>", line 1, in <module> TypeError: 'int' object does not support item deletion</pre>
<p>Mengenkomprehension und iterative Lösung</p> <pre>mc = {x.lower() for x in "Das Alphabet" if x.isalnum()}</pre> <pre>ms = set() for x in "Das Alphabet": if x.isalnum(): ms.add(x.lower())</pre> <p>Komprehension von Dictionaries</p> <pre>text = "abradabra" dc = {c:text.count(c) for c in set(text)}</pre> <p>Wie würde man das iterativ programmieren?</p>	
<pre>#List lc = [x.lower() for x in "Das Alphabet. 12" if x.isalnum()] print lc lc = list(x.lower() for x in "Das Alphabet. 12" if x.isalnum()) print lc #Sets mc = {x.lower() for x in "Das Alphabet. 12" if x.isalnum()} print mc mc = set(x.lower() for x in "Das Alphabet. 12" if x.isalnum()) print mc #Dictionary text = "abradabra. 12" dc = {c:text.count(c) for c in set(text)} print dc dc = dict((c, text.count(c)) for c in set(text)) print dc</pre>	<pre>#List ['d', 'a', 's', 'a', ' ', 'p', 'h', 'a', 'b', 'e', 't', ' ', ' ', '2'] ['d', 'a', 's', 'a', ' ', 'p', 'h', 'a', 'b', 'e', 't', ' ', ' ', '2'] #Sets set(['a', 'b', 'e', 'd', 'h', 't', ' ', 'p', 's', ' ', '2', 't']) set(['a', 'b', 'e', 'd', 'h', 't', ' ', 'p', 's', ' ', '2', 't']) #Dictionary {'a': 5, ' ': 1, 'b': 2, 'd': 1, 'h': 1, 'e': 1, 't': 1, 'p': 2, 's': 2, ' ': 1} {'a': 5, ' ': 1, 'b': 2, 'd': 1, 'h': 1, 'e': 1, 't': 1, 'p': 2, 's': 2, ' ': 1}</pre>
<p>Sequence Types -- str, unicode, list, tuple, buffer, xrange</p> <p>An iterable which supports efficient element access using integer indices via the <code>__getitem__()</code> special method and defines a <code>len()</code> method that returns the length of the sequence. Some built-in-sequence types are list, str, tuple, and unicode. Note that dict also supports <code>__getitem__()</code> and <code>__len__()</code>, but is considered a mapping rather than a sequence because the lookups use arbitrary immutable keys rather than integers.</p> <p>Mapping Types -- dict</p> <p>A mapping object maps immutable values (keys) to arbitrary objects (values). Mappings (dict) are mutable objects. There is currently only one standard mapping type, the <i>dictionary</i>. A dictionary's keys are almost arbitrary values. Only values containing lists, dictionaries or other mutable types (that are compared by value rather than by object identity) may not be used as keys.</p> <p>Iterable Types -- str, unicode, list, tuple, dict, file ...</p> <p>Examples of iterables include all sequence types (such as list, str, and tuple) and some non-sequence types like dict and file and objects of any classes you define with an <code>__iter__()</code> or <code>__getitem__()</code> method. Iterables can be used in a for loop and in many other places where a sequence is needed (zip(), map(), ...).</p>	<p>Wortlisten als Lexika</p> <p>Definition (Wortlisten)</p> <p>Die einfachste Form von Lexika sind Wortlisten. Als Rohtext-Datei typischerweise 1 Wort pro Zeile und sortiert.</p> <p>Stopwortlisten (stopwords) in NLTK ▶1</p> <pre>stopwords_en = nltk.corpus.stopwords.words('english') print len(stopwords_en), stopwords_en[:20] # >>> 127 ['i', 'herself', 'was', 'because', 'from', 'any', 't']</pre> <p>Hinweis: Spezialsyntax <code>[:20]</code> gibt jedes 20. Element zurück.</p>

```

import nltk

stopwords_en = nltk.corpus.stopwords.words('english')
print len(stopwords_en), stopwords_en[::20] # 127
    # [u'i', u'herself', u'was', u'because', u'from',
    # u'any', u't']

### Using German stopword list this way suffers from
### encoding problems
stopwords_de = nltk.corpus.stopwords.words('german')
print len(stopwords_de), sorted(stopwords_de)[::20] # 231
    # [u'aber', u'auch', u'dem', u'dieser', u'er', u'hin',
    # u'jeder', u'manchem', u'nur', u'solches', u'warst',
    # u'wollen']

```

```

from __future__ import division
import nltk

#Brown words
brown_words=nltk.corpus.brown.words()
print "\nfrom 0.0: ", len(brown_words), type(brown_words)
print brown_words[::20]

#List: Inffizienterer Lookup für die Stopwörter mit Listen
stopwords_en = nltk.corpus.stopwords.words('english')
print "\nfrom 0.1: ", len(stopwords_en), stopwords_en[::20]

# Was berechnet foo()?
def foo(text):
    """ Hier fehlt Dokumentation... """
    bar = [w for w in text
           if w.lower() not in stopwords_en]
    print "\nfrom 1.1: ", len(bar), len(text)
    return len(bar)/len(text)*100.
print "\nfrom 1.2: ", foo(brown_words)

# Set: Effizienterer Lookup für die Stopwörter mit Mengen
stopwords_en_set = {w for w in stopwords_en}
print "\nfrom 2: ", len(stopwords_en_set), stopwords_en_set

# Optimierte Funktion
def non_stopwords_percentage(text):
    """ Return the percentage of non-stopwords
        in a list of English tokens. """
    non_stopwords = [w for w in text
                     if w.lower() not in stopwords_en_set]
    print "\nfrom 3.1: ", len(non_stopwords), len(text)
    return len(non_stopwords)/len(text)*100.
print "\nfrom 3.2: ", non_stopwords_percentage\
    (brown_words)
print "\nfrom 3.3: "
help(non_stopwords_percentage)

# Wie misst man die Effizienz?
=====
import cProfile

# Vorberechnen der Liste mit 1 Mio. Token
brown_words_list = list(nltk.corpus.brown.words())
print "\nfrom 4: ", len(brown_words_list), type(brown_words_list)
print brown_words_list[::20]

# Ineffizientere Verarbeitung mit Stopwortlisten
print "\nfrom 5.1: "
cProfile.run("foo(brown_words)")
print "\nfrom 5.2: "
cProfile.run("foo(brown_words_list)")

# Effizientere Verarbeitung mit Stopwortmenge
# (oder Dictionary mit Dummy-Werten)
print "\nfrom 6.1: "
cProfile.run("non_stopwords_percentage(brown_words)")
print "\nfrom 6.2: "
cProfile.run("non_stopwords_percentage(brown_words_list)")

# Wie kann man die Interpunktionsstoken ausfiltern?
=====
import re

def delete_punctuation(s):
    """ Return string with all punctuation symbols
        of iso-latin 1 deleted. """
    p = r'![!#%&\x27`()*,.-/:;?@[\]_{}]\xa1\xab\xb7\xbb\xbf]'
    return re.sub(p, '', s)

def content_word_percentage(text):
    """ Return the percentage of content words in
        a list of English tokens. """

```

Anteil von echten Inhaltswörtern bestimmen

Wie kann man die Interpunktionsstoken eliminieren?

►3

```

import re

def delete_punctuation(s):
    """ Return string with all punctuation symbols of iso-latin 1 deleted.
        p = r'![!#%&\x27`()*,.-/:;?@[\]_{}]\xa1\xab\xb7\xbb\xbf'
        return re.sub(p, '', s)

def content_word_percentage(text):
    """ Return the percentage of content words in a list of English tokens.
        content_words = [w for w in text
                         if delete_punctuation(w) != ''
                            and w.lower() not in stopwords_en_set]

```

foo_fraction_en.py"

Die Ausgabe sieht aus wie eine Liste ist aber ein ineffizientes Objekt

from 0.0: 1161192 <class 'nltk.corpus.reader.util.ConcatenatedCorpusView'>
[u'The', u'Fulton', u'County', u'Grand', u'Jury', u'said', u'Friday', u'an', u'investigation', u'of',
u'Atlanta's', u'recent', u'primary', u'elect', u'produced', u'~, u'no', u'evidence', u'~~', u'that']

from 0.1: 127 <type 'list'>
[u'i', u'herself', u'was', u'because', u'from', u'any', u't']

Neue Divisionsregeln importieren sonst berechnet die Funktion immer 0

from 1.2:
from 1.1: 688946 1161192
59.3309289075

Mit set Datenstruktur geht's schneller

from 2: 127 set([u'all', u'just', u'being', u'over', u'both', u'through', u'yourselves', u'its', u'before',
u'herself', u'had', u'should', u'to', u'only', u'under', u'ours', u'has', u'do', u'them', u'his', u'very',
gekürzt...])

Neue Divisionsregeln importieren sonst berechnet die Funktion immer 0

from 3.2:
from 3.1: 688946 1161192
59.3309289075

from 3.3:

Help on function non_stopwords_percentage in module __main__:

non_stopwords_percentage(text)
 Return the percentage of non-stopwords
 in a list of English tokens.

Mit List Datenstruktur geht's schneller

from 4: 1161192 <type 'list'>
[u'The', u'Fulton', u'County', u'Grand', u'Jury', u'said', u'Friday', u'an', u'investigation', u'of',
u'Atlanta's', u'recent', u'primary', u'elect', u'produced', u'~, u'no', u'evidence', u'~~', u'that']

Effizienzmessungen

Inhaltswörter: class 'nltk.corpus.reader.util.ConcatenatedCorpusView'
Stopwörter: List
from 5.1:
from 1.1: 688946 1161192
10488536 function calls (10488532 primitive calls) in 48.990 seconds

Inhaltswörter: List
Stopwörter: List
from 5.2:
from 1.1: 688946 1161192
1161199 function calls in 10.615 seconds

Inhaltswörter: class 'nltk.corpus.reader.util.ConcatenatedCorpusView'
Stopwörter: Set
from 6.1:
from 3.1: 688946 1161192
10488536 function calls (10488532 primitive calls) in 41.013 seconds

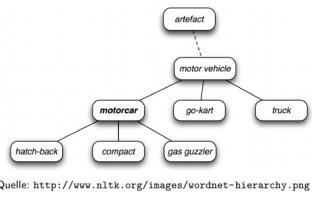
Inhaltswörter: List
Stopwörter: Set
from 6.2:
from 3.1: 688946 1161192
1161199 function calls in 4.211 seconds

Die Inhaltswörter machen nicht 59% sondern nur 46% aus
Es gibt nicht mehr 688946 sondern nur noch 541074 Inhaltswörter

from 7.2:
from 7.1: 541074 1161192
46.5964284976

Delete_Punctuation benötigt um die 20 zusätzliche Sekunden
from 8:

from 7.1: 541074 1161192
5658095 function calls in 26.922 seconds

<pre> content_words = [w for w in text if delete_punctuation(w) != '' and w.lower() not in stopwords_en_set] print "\nfrom 7.1: ", len(content_words), len(text) return len(content_words)/len(text)*100. print "\nfrom 7.2: ", \ content_word_percentage(brown_words_list) print "\nfrom 8: " cProfile.run("content_word_percentage(brown_words_list)") </pre>	<p>Tokens welche die Funktion verändert. Alle Satzzeichen, welche vorher als Inhaltswort mitgezählt wurden, sind nun aus der Berechnung entfernt worden.</p> <pre> . creativity-oriented creativityoriented , ; people-oriented peopleoriented M.D. MD father's fathers ; </pre>										
<p>CMU (Carnegie Mellon University) Pronouncing Dictionary</p> <p>File Format: Each line consists of an uppercased word, a counter (for alternative pronunciations), and a transcription. Vowels are marked for stress (1=primary, 2=secondary, 0=no stress). E.g.: NATURAL 1 N AE1 CH ERO AHO L</p> <p>The dictionary contains 127069 entries. Of these, 119400 words are assigned a unique pronunciation, 6830 words have two pronunciations, and 839 words have three or more pronunciations. Many of these are fast-speech variants.</p> <p>Phonemes: There are 39 phonemes, as shown below:</p> <table border="1"> <thead> <tr> <th>Phoneme Example Translation</th> <th>Phoneme Example Translation</th> </tr> </thead> <tbody> <tr> <td>AA odd</td> <td>AA D</td> </tr> <tr> <td>AH hut</td> <td>HH AH T</td> </tr> <tr> <td>AE at</td> <td>AO ought</td> </tr> <tr> <td>AE T</td> <td>AO T</td> </tr> </tbody> </table> <pre> \$ grep -w RESEARCH /Users/siclemat/nltk_data/corpora/cmudict/cmudict RESEARCH 1 R IY0 S ER1 CH RESEARCH 2 R IY1 S ERO CH </pre> <p>Wie soll man solche Information in Python als Daten repräsentieren?</p>	Phoneme Example Translation	Phoneme Example Translation	AA odd	AA D	AH hut	HH AH T	AE at	AO ought	AE T	AO T	<p>CMU (Carnegie Mellon University) Pronouncing Dictionary</p> <p>Strukturierte Lexikoneinträge CMU besteht aus Paaren von Lemma und Listen von phonetischen Kodes.</p> <p>Filtern von Lexikoneinträgen</p> <pre> import nltk entries = nltk.corpus.cmudict.entries() print entries[71607] # ('love', ['L', 'AH1', 'V']) # Finde alle Wörter auf -n, welche als -M ausgesprochen werden. print [word for (word,pron) in entries if pron[-1] == 'M' and word[-1] == 'n'] </pre>
Phoneme Example Translation	Phoneme Example Translation										
AA odd	AA D										
AH hut	HH AH T										
AE at	AO ought										
AE T	AO T										
<p>WordNet: Ein Netz von Bedeutungsbeziehungen¹</p> <p>Wie lässt sich die Bedeutung eines Worts angeben?</p> <ul style="list-style-type: none"> ► Klassische Charakterisierung: Umschreibung, Definition ► Relationale lexikalische Semantik = Bedeutungsbeziehungen ► Angabe von Synonymen, Hypernymen, Hyponymen, Antonymen usw., welche ein Netz (Hierarchie) von verknüpften Bedeutungen ergeben  <p>Quelle: http://www.nltk.org/images/wordnet-hierarchy.png</p>	<p>WordNet: Komplexe lexikalische Datenstruktur²</p> <p>Speziell zugeschnittene Datenstruktur benötigt</p> <ul style="list-style-type: none"> ► Zugriff auf Bedeutungen (<i>synsets</i>) und Wörter (<i>lemmas</i>) ► Navigation im Wortnetz entlang der semantischen Relationen (Oberbegriffe, Unterbegriffe, Gegenbegriffe) ► Berechnen von semantischer Verwandtschaft (Nähe, Bezüge) im Netz <p>WordNet in NLTK</p> <pre> import nltk from nltk.corpus import wordnet as wn # Welche Bedeutungen hat das Wort "car"? print wn.synsets('car') # Definition einer Bedeutung print wn.synset('car.n.01').definition() # Alle hyponymen Bedeutungen eines Lemmas berechnen print wn.synset('car.n.01').hyponyms() </pre>										
<pre> import nltk # Mit "as" kann man Kurznamen einführen. from nltk.corpus import wordnet as wn # <h2>Synsets (unterschiedliche Bedeutungen) eines Lemmas berechnen</h2> # Welche Bedeutungen hat das Wort "car"? print "\nfrom 1: ", wn.synsets('car') print "\nfrom 2: ", wn.synsets('motorcar') # Definition einer Bedeutung print "\nfrom 3: ", wn.synset('car.n.01').definition() # <h2>Alle synonymen Lemmas eines Synsets berechnen</h2> print "\nfrom 4: ", wn.synset('car.n.01').lemma_names() print "\nfrom 5: ", wn.synset('car.n.01').lemmas() # Was ist eine Lemma-Datenstruktur? print "\nfrom 6: ", type(wn.synset('car.n.01').lemmas())[0] print "\nfrom 7: " help(nltk.corpus.reader.wordnet.Lemma) # Alle hyponymen Bedeutungen eines Lemmas berechnen print "\nfrom 8: ", wn.synset('car.n.01').hyponyms() </pre>	<pre> /LubuntuShared/8) PCL-1/Vorlesung/vl10/nltk_wordnet.py" from 1: [Synset('car.n.01'), Synset('car.n.02'), Synset('car.n.03'), Synset('car.n.04'), Synset('cable_car.n.01')] from 2 [Synset('car.n.01')] from 3: a motor vehicle with four wheels; usually propelled by an internal combustion engine from 4: [u'car', u'auto', u'automobile', u'machine', u'motorcar'] from 5: [Lemma('car.n.01.car'), Lemma('car.n.01.auto'), Lemma('car.n.01.automobile'), Lemma('car.n.01.machine'), Lemma('car.n.01.motorcar')] from 6: <class 'nltk.corpus.reader.wordnet.Lemma'> from 7: Help on class Lemma in module nltk.corpus.reader.wordnet: class Lemma(_WordNetObject) The lexical entry for a single morphological form of a sense-disambiguated word. from 8: [Synset('ambulance.n.01'), Synset('beach_wagon.n.01'), Synset('bus.n.04'), Synset('cab.n.03'), Synset('compact.n.03'), Synset('convertible.n.01'), Synset('coupe.n.01'), Synset('cruiser.n.01'), Synset('electric.n.01'), Synset('gas_guzzler.n.01'), Synset('hardtop.n.01'), Synset('hatchback.n.01'), Synset('horseless_carriage.n.01'), Synset('hot_rod.n.01'), Synset('jeep.n.01'), Synset('limousine.n.01'), Synset('loaner.n.02'), Synset('minicar.n.01'), Synset('minivan.n.01'), Synset('model_t.n.01'), Synset('pace_car.n.01'), Synset('racer.n.02'), Synset('roadster.n.01'), Synset('sedan.n.01'), Synset('sport_utility.n.01'), Synset('sports_car.n.01'), Synset('stanley_steamer.n.01'), Synset('stock_car.n.01'), Synset('subcompact.n.01'), Synset('touring_car.n.01'), Synset('used-car.n.01')] </pre>										

<pre>import nltk entries = nltk.corpus.cmudict.entries() print entries[71607] # ('love', ['L', 'AH1', 'V']) # Finde alle Wörter auf -n, welche als -M ausgesprochen werden. print [word for (word,pron) in entries if pron[-1] == 'M' and word[-1] == 'n']</pre>	(u'love', [u'L', u'AH1', u'V']) [u'autumn', u'column', u'condemn', u'damn', u'goddamn', u'hymn', u'solemn']
<pre>global_list = [('der',1200),('die',1000),('das',900)] print "\nfrom 0: ", global_list def del_first(l): print "\nfrom 1.1: ", 'global_list is parameter l:',\ global_list is l del l[0] print "\nfrom 1.2: ", 'global_list is parameter l:',\ global_list is l del_first(global_list) print "\nfrom 2: ", global_list del_first(global_list) print "\nfrom 3: ", global_list</pre>	<pre>from 0: [('der', 1200), ('die', 1000), ('das', 900)] from 1.1: global_list is parameter l: True from 1.2: global_list is parameter l: True from 2: [('die', 1000), ('das', 900)] from 1.1: global_list is parameter l: True from 1.2: global_list is parameter l: True from 3: [('das', 900)]</pre> <p>Process finished with exit code 0</p>
<h3>Elemente hinzuzählen</h3> <p>Inkrementieren</p> <pre>In [18]: fdist.update('a') fdist</pre> <p>Out[18]: Counter({'a': 5, 'b': 2, 'd': 1, 'k': 1, 'r': 2})</p> <h3>Wieviele Elemente gibt es insgesamt?</h3> <p>Die Methode N():</p> <pre>In [19]: fdist.N()</pre> <p>Out[19]: 11</p> <p>Was ist die absolute Häufigkeit eines Elements?</p> <pre>In [20]: fdist['r']</pre> <p>Out[20]: 2</p>	<p>Was ist die relative Häufigkeit eines Elements?</p> <pre>In [21]: fdist.freq('r')</pre> <p>Out[21]: 0.18181818181818182</p> <p>Welches Element kommt am häufigsten vor?</p> <pre>In [22]: fdist.max()</pre> <p>Out[22]: 'a'</p>
<h3>Distributionen addieren</h3> <pre>In [28]: fdist2 = nltk.FreqDist('Hokuspokus') fdist + fdist2</pre> <p>Out[28]: Counter({'H': 1, 'a': 5, 'b': 2, 'd': 1, 'k': 3, 'o': 2, 'p': 1, 'r': 2, 's': 2, 'u': 2})</p>	<h3>Iteriere durch alle Elemente nach Häufigkeit</h3> <pre>In [9]: for item in fdist: print item, fdist[item]</pre> <pre>In [10]: for item,freq in fdist.most_common(): print item, freq</pre> <pre>In [26]: for item in sorted(fdist,key=fdist.get,reverse=True): print item, fdist[item]</pre> <p>a 5 r 2 b 2 k 1 d 1</p>

<pre>Zahlen runden round(0.05,1) ==> 0.1 up round(0.15,1) ==> 0.1 round(0.25,1) ==> 0.3 up round(0.35,1) ==> 0.3 round(0.45,1) ==> 0.5 up round(0.55,1) ==> 0.6 round(0.65,1) ==> 0.7 up round(0.75,1) ==> 0.8 up round(0.85,1) ==> 0.8 round(0.95,1) ==> 0.9 round(1.05,1) ==> 1.1 up round(1.15,1) ==> 1.1 round(1.25,1) ==> 1.3 up round(1.35,1) ==> 1.4 up round(1.45,1) ==> 1.4 round(1.55,1) ==> 1.6 up round(1.65,1) ==> 1.6 round(1.75,1) ==> 1.8 up round(1.85,1) ==> 1.9 up round(1.95,1) ==> 1.9 'Runden' durch die Formatierung '%lf' % 0.05 ==> 0.1 up '%lf' % 0.15 ==> 0.1 '%lf' % 0.25 ==> 0.2 '%lf' % 0.35 ==> 0.3 '%lf' % 0.45 ==> 0.5 up '%lf' % 0.55 ==> 0.6 up '%lf' % 0.65 ==> 0.7 up '%lf' % 0.75 ==> 0.8 up '%lf' % 0.85 ==> 0.8</pre>	<pre># Demonstration von Funktionen mit yield #===== def quadriere(iterierbar): for i in iterierbar: yield i*i quadrat = quadriere([10,11]) print "from 1: ", quadrat print "from 2: ", type(quadrat) # next() liefert nächstes Element zurück. print "from 3: ", quadrat.next() print "from 4: ", quadrat.next() #Falls keines mehr erhältlich ist, entsteht eine Ausnahme. # In Iterationskontexten führt diese Ausnahme dazu, # dass die Iteration beendet wird. # Die Ausnahme wird "entsprechend" behandelt. #print "from 5: ", quadrat.next() # Generatoren sind erschöpft nach einem Durchgang print "from 7: ", sum(quadrat) #Aggregationsfunktion auf Generator anwenden print "from 6: ", sum(quadriere([10,11]))</pre> <p>Output</p> <pre>from 1: <generator object quadriere at 0xb7253edc> from 2: <type 'generator'> from 3: 100 from 4: 121 from 7: 0 from 6: 221</pre>
<pre># Demonstration von Generatorausdrücken quadrat = (i*i for i in [10,11]) print quadrat print type(quadrat) # generator hat keine len() # print len(quadrat) # next() liefert nächstes Element zurück. print quadrat.next() print quadrat.next() # Durchiterieren quadrat = (i*i for i in xrange(10,13)) for q in quadrat: print q</pre>	<pre>/LubuntuShared/8) PCL-1/Vorlesung/vl11/generators_next.py" <generator object <genexpr> at 0xb7218edc> <type 'generator'> 100 121 100 121 144</pre>
<pre>while True: x = raw_input('Please type in a number: ') try: float(x) print "from 1: " break # nur hier kann Schleife beendet werden print "from 2: " except ValueError: print "from 3: " pass print "from 4: " finally: print "from 5: " print "The number was:", x</pre>	<pre># Please type in a number: null # from 3: # from 4: # from 5: # Please type in a number: eins # from 3: # from 4: # from 5: # Please type in a number: 7,78 # from 3: # from 4: # from 5: # Please type in a number: 90.25 # from 1: # from 5: # The number was: 90.25</pre>
<pre>import sys #Den Inhalt dieser Datei # ohne Leerzeilen printen filename = "with_open.py" with open(filename,'r') as f: for l in f: if l.rstrip() != '': sys.stdout.write(l)</pre>	<pre># Modul zur Zeitmessung von Python-Statements #===== import timeit setup = 'import random' # Effekt von xrange vs. range bei random.sample # Konstruiere 2 Timer-Objekte tr = timeit.Timer('random.sample(range(1000000),100)', setup) tx = timeit.Timer('random.sample(xrange(1000000),100)', setup) # Führe Timings je einmal durch und speichere Anzahl Sekunden trsecs = tr.timeit(1) txsecs = tx.timeit(1) print "Aufgabe: Sample 100 Zahlen aus dem Bereich 0 bis 999999." print "Zeit mit xrange: ", txsecs, "Sekunden" print "Zeit mit range: ", trsecs, "Sekunden" print "xrange ist etwa", trsecs/txsecs, "Mal schneller!" Output Aufgabe: Sample 100 Zahlen aus dem Bereich 0 bis 999999. Zeit mit xrange: 0.000220060348511 Sekunden Zeit mit range: 0.13285045624 Sekunden xrange ist etwa 605.674972914 Mal schneller!</pre>
<pre># Rechenzeit und Speichereffizienz messen #===== # HOWTO und Erklärungen # (1) Programm via Terminal im Hintergrund aufstarten (&): # \$ python timeit_listcomprehension.py & # \$ python timeit_generator.py & # Das Programm gibt jeweils die Prozess-ID aus, unter der es läuft. # Damit können die Programme auseinander gehalten werden. # Im Folgenden steht PID1 bzw. PID2 für die jeweilige nummerische ID.</pre>	<p>Im Terminal:</p> <p>mit „&“ wird der Prozess im Hintergrund gestartet</p> <pre>\$ python timeit_listcomprehension.py & [1] 3720 benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Vorlesung/vl11\$ process id: 3720 python timeit_generator.py & [2] 3722 benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8) PCL-1/Vorlesung/vl11\$ process id: 3722</pre>

```
# (2) Speicherverbrauch messen im Terminal mit ps:
#   $ ps -vp PID1
#   $ ps -vp PID2
#   Die entscheidende Information steckt im Attribut RSS
(resident set size),
#   das den physikalisch belegten Speicher in Kilobytes (im
Gegensatz zum
#   virtuellen Speicher) misst.
# (3) Speicherverbrauch messen mit Task Manager (WIN) oder
Aktivitätsanzeige (MACOS)
#   Am besten Filtern auf python, danach den Prozess über seine
Prozess-ID identifizieren.

timeit_listcomprehension.py
import nltk, timeit, time, os
print 'process id:', os.getpid()
words = nltk.corpus.brown.words()
def test_listcomprehension():
    return set([w.lower() for w in words])
# Initialisiere Timer-Objekt
tl = timeit.Timer(test_listcomprehension)
# Timing von Listenkomprehension
print 'Timed list comprehension (seconds):', tl.timeit(1)
# Speicherverbrauch muss von außerhalb gemessen werden.
# Es gibt keinen einfachen Weg, das innerhalb von Python zu
machen.
time.sleep(600) # 600 Sekunden

timeit_generator.py
import nltk, timeit, time, os
print 'process id:', os.getpid()
words = nltk.corpus.brown.words()
def test_generator():
    return set(w.lower() for w in words)
# Initialisiere Timer-Objekt
tg = timeit.Timer(test_generator)
# Timing von Generatorausdruck
print 'Timed generator (seconds):', tg.timeit(1)
# Speicherverbrauch muss von außerhalb gemessen werden.
# Es gibt keinen einfachen Weg, das innerhalb von Python zu
machen.
time.sleep(600) # 600 Sekunden
```

Instanzen (`isinstance(Object,Class)`)

Relation zwischen einem Objekt und seiner Klasse (Typ)!

```
>>> isinstance([1,2,3], list)
True
>>> isinstance([1,2,3], dict)
```

Unterklassen (`issubclass(lowerclass,upperclass)`)

Relation zwischen 2 Klassen/Typen!

```
>>> issubclass(nltk.FreqDist, dict)
True
>>> issubclass(nltk.FreqDist, collections.Counter)
True
>>> issubclass(collections.Counter, dict)
True
```

benzro@benzro-eMachines E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8)
PCL-1/Vorlesung/vl11\$
Timed list comprehension (seconds): 18.6455891132
Timed generator (seconds): 17.7329790592
ps-vp 3720
PID TTY STAT TIME MAJFL TRS DRS RSS %MEM COMMAND
3720 pts/1 RI 0:08 0 3385 152566 84860 2.1 python timeit_li
ps-vp 3722
PID TTY STAT TIME MAJFL TRS DRS RSS %MEM COMMAND
3722 pts/1 SI 0:10 0 3385 108694 41244 1.0 python timeit_ge
benzro@benzro-eMachines-E625:/media/benzro/OS/Users/benzro/Documents/LubuntuShared/8)
PCL-1/Vorlesung/vl11\$

Im Task Manager:

Command	User	CPU%	RSS	VM-Size	PID ▲	State	Prio	PPID
python	benzro	0%	85.3 MB	154.3 MB	3693	S	0	1
python	benzro	0%	40.3 MB	109.5 MB	3696	S	0	1
bash	benzro	0%	4.9 MB	7.5 MB	3703	S	0	10579
python2.7	benzro	0%	85.1 MB	154.3 MB	3718	S	0	10579
python	benzro	0%	85.0 MB	154.3 MB	3720	S	0	3703
python	benzro	0%	40.3 MB	109.5 MB	3722	S	0	3703

Zusammenfassung:

Die Listenkomprehension (ID3720) braucht 154 MB RAM und 18.6 Sekunden
Der Generator (ID3722) braucht 109.5 MB RAM und 17.3 Sekunden

Listenkomprehension: mit eckigen Klammern**return set([w.lower() for w in words])****Testgenerator: ohne eckige Klammern****return set(w.lower() for w in words)**

>>> help(issubclass)

Help on built-in function issubclass in module __builtin__:

```
issubclass(...)
issubclass(C, B) -> bool
```

Return whether class C is a subclass (i.e., a derived class) of class B.

- ▶ Klassen spezifizieren und implementieren die **Eigenschaften** (=Attribute) und **Funktionen** (=Methoden) von Objekten.
- ▶ Klassen **abstrahieren** gemeinsame Eigenschaften und Funktionalitäten.
- ▶ Klassen sind in Unter-/Oberklassen (*superclass/subclass*) organisiert (**Vererbung**).
- ▶ Vererbung heisst, dass Eigenschaften/Methoden einer Oberklasse **defaultmäßig** auch in der Unterklassie zur Verfügung stehen.
- ▶ Die Methoden können in der Unterklassie aber auch **umdefiniert** werden (**Flexibilität**).
- ▶ Jede selbstdefinierte Klasse muss **eine Oberklasse** haben.
- ▶ Eine selbstdefinierte Klasse kann auch mehrere Oberklassen haben, d.h. **Mehrfachvererbung** ist möglich.

<h3>Die Eigentümlichkeiten der obersten Klasse <code>object</code></h3> <p>Die Oberklasse aller Klassen in Python heisst <code>object</code>.</p> <p>Die Klasse <code>object</code> ist trotz ihres Namens eine Klasse!</p> <pre>>>> help(object) Help on class object in module __builtin__: class object The most base type</pre> <p>Objekte (Instanzen) der Klasse <code>object</code></p> <pre>>>> o = object() >>> type(o) <type 'object'></pre> <p>Gibt es eine Oberklasse der Klasse <code>object</code>?</p> <pre>>>> issubclass(object, object) True</pre>	<h3>Klassen definieren: Case-insensitive Strings</h3> <p>Motivation: Konsistenter Umgang mit Zeichenketten, wo Gross-/Kleinschreibung keine Rolle spielt.</p> <p>Definition der Klasse, der Konstruktorfunktion und einer Methode</p> <pre>class Istr(object): """Case-insensitive string class""" def __init__(self, s): self._is = s.lower() # self ist Instanzparameter def endswith(self, s): # Methode endswith(s) return self._is.endswith(s.lower())</pre> <p>Instantiierung eines Objekts und Methodenaufruf</p> <pre>s = Istr('ABC') # Konstruktion eines Objekt der Klasse Istr s.endswith('c') # Methoden-Aufruf</pre>
<h3>Zusammenhang von Definition und Verwendung</h3> <p>Das Diagramm zeigt den Ablauf von links nach rechts:</p> <ul style="list-style-type: none"> Klassendefinition: Ein Kasten enthält den Code: <code>class Istr(object):</code> Konstruktordefinition: Ein Kasten enthält den Code: <code>def __init__(self,s): self._is = s.lower()</code> Methodendefinition: Ein Kasten enthält den Code: <code>def find(self,s): ls = s.lower() return self._is.find(ls)</code> Objektinstantiierung: Ein orangefarbener Kasten enthält den Code: <code>s = Istr('ABC')</code> Methodenaufruf: Ein orangefarbener Kasten enthält den Code: <code>s.find('bC')</code> 	<h3>Instanzvariablen</h3> <p>Jede Objekt-Instanz kann Attribute mit individuellen Werten in sich tragen. Normalerweise werden diese Instanzvariablen beim Konstruieren des Objekts mit einem Wert belegt.</p> <p>Instanzvariablen</p> <pre>>>> s = Istr('ABC') >>> print s._is abc >>> s2 = Istr('XYZ') >>> print s2._is xyz</pre> <p>Grade der Öffentlichkeit: Namenskonvention</p> <ul style="list-style-type: none"> Öffentliche Instanzvariablen ohne Unterstrich am Anfang: Überall frei benutzbar! Private Instanzvariablen beginnen mit Unterstrich: Sollen nur innerhalb der Klassendefinition verwendet werden! Grund: Datenabstraktion: Interne Implementation kann ändern, ohne das Klassenbenutzung sich ändert muss
<h3>Instanzvariablen</h3> <p>Jede Objekt-Instanz kann Attribute mit individuellen Werten in sich tragen. Normalerweise werden diese Instanzvariablen beim Konstruieren des Objekts mit einem Wert belegt.</p> <p>Instanzvariablen</p> <pre>>>> s = Istr('ABC') >>> print s._is abc >>> s2 = Istr('XYZ') >>> print s2._is xyz</pre> <p>Grade der Öffentlichkeit: Namenskonvention</p> <ul style="list-style-type: none"> Öffentliche Instanzvariablen ohne Unterstrich am Anfang: Überall frei benutzbar! Private Instanzvariablen beginnen mit Unterstrich: Sollen nur innerhalb der Klassendefinition verwendet werden! Grund: Datenabstraktion: Interne Implementation kann ändern, ohne das Klassenbenutzung sich ändert muss 	<h3>Fazit Objektorientierte Programmierung (OOP)</h3> <p>Kernkonzepte nach http://en.wikipedia.org/wiki/Object-oriented_programming</p> <ul style="list-style-type: none"> Datenkapselung I: Bündeln von Datenstrukturen und zugehöriger Funktionalität unter einer Adresse (=Objekt) Datenkapselung (Abstraktion) II: Klare Schnittstelle, welche Attribute und Methoden für öffentliche und welche für private (objektinterne) Zwecke nutzbar sind Klassenzugehörigkeit: Objekte sind Instanzen einer Klasse Vererbung: Unterklassen können Attribute/Methoden von ihren Oberklassen erben Dynamische Bindung: Welche Methode (d.h. Methode von welcher (Ober-)Klasse) ein Objekt benutzt, wird erst beim Aufruf der Methode festgelegt anhand der <i>method resolution order</i>. Selbst-Parameter (<code>self</code>): Platzhalter für das Instanzobjekt in der Definition einer Klasse

```
>>> class iclass(object):... def __init__(self,s):...
>>> class iclass3(type):... def __init__(self,s):...
```

<h3>Namen, Zuweisung, Bindung und Objekte</h3> <p>Zuweisung (Assignment)</p> <pre>a = 5*8</pre> <p>Was passiert beim Verarbeiten der Zuweisungsanweisung?</p> <ol style="list-style-type: none"> Evaluiere (evaluate) RHS-Ausdruck $5*8$ zu einem Ganzzahl-Objekt. (RHS=right-hand side) Bind (binding) das evaluierte Ganzzahl-Objekt an den Namen a. <p>Was passiert beim Statement $c = b$?</p>	<p>Namen referieren auf Objekte</p> <ol style="list-style-type: none"> Rechteck = Objekte Kreis = Referenz auf Objekt <p>Quelle: [SUMMERFIELD 2008, 14]</p> <p>Zuweisungsverhalten bei einfachen Datentypen</p> <pre>b erhält die Referenz aufs Objekt 8 >>> b=8 c erhält die Referenz aufs Objekt 8 >>> c=b b erhält die Referenz aufs Objekt 7 >>> b=7 c behält die Referenz aufs Objekt 8 >>> c 8</pre> <p>versus Zuweisungsverhalten von Referenzdatentypen</p>				
<h3>Mehrfaches Binden eines Namens (rebinding)</h3> <p>Was passiert, wenn derselbe Name mehrfach zugewiesen wird?</p>	<p>Unreferenzierte Objekte und Müllsammlung</p> <pre>s = "Ein String" # s macht "Ein String" zugänglich # im nachfolgenden Programm. s += " wird zusammengesetzt!" # Nach dieser Anweisung ist # "Ein String" nicht mehr # zugänglich via Name s. print s # s referenziert ein neues Objekt. Ein String wird zusammengesetzt!</pre> <ul style="list-style-type: none"> Unbenannte Objekte sind ausserhalb des Ausdrucks, in dem sie vorkommen, nicht mehr zugänglich und benutzbar: Sie sind Datenmüll (garbage). Nicht mehr zugängliche Objekte können periodisch gesammelt und entsorgt werden (garbage collection). Dadurch wird Speicherplatz frei. <p>Python weiss für jedes Objekt, wie viele Referenzen (Namen) darauf existieren. Das Modul gc ist eine Schnittstelle zur garbage collection.</p>				
<h3>Unreferenzierte Objekte und Müllsammlung</h3> <pre>s = "Ein String" # s macht "Ein String" zugänglich # im nachfolgenden Programm. s += " wird zusammengesetzt!" # Nach dieser Anweisung ist # "Ein String" nicht mehr # zugänglich via Name s. print s # s referenziert ein neues Objekt. Ein String wird zusammengesetzt!</pre> <ul style="list-style-type: none"> Unbenannte Objekte sind ausserhalb des Ausdrucks, in dem sie vorkommen, nicht mehr zugänglich und benutzbar: Sie sind Datenmüll (garbage). Nicht mehr zugängliche Objekte können periodisch gesammelt und entsorgt werden (garbage collection). Dadurch wird Speicherplatz frei. <p>Python weiss für jedes Objekt, wie viele Referenzen (Namen) darauf existieren. Das Modul gc ist eine Schnittstelle zur garbage collection.</p>	<p>Identität eines Objekts (<code>id()</code>) und mutable data</p> <p>Identität bei veränderlichen Datenstrukturen</p> <p>Verschiedene Namen können auf dasselbe Objekt referenzieren. Die eingebaute Funktion <code>id()</code> identifiziert jedes Objekt über eine Ganzzahl (entspricht ungefähr seiner Speicheradresse). Python garantiert, dass 2 verschiedene Objekte gleichzeitig nie dieselbe ID haben.</p> <table border="1"> <tr> <td data-bbox="793 1381 1087 1560"> <p>(Re-)Binding einer Variable</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l = ['a'] >>> id(l) 4299634664</pre> </td> <td data-bbox="1087 1381 1431 1560"> <p>Veränderliche Datenstrukturen</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l[0] = 'b' >>> id(l) 4300400112</pre> </td> </tr> <tr> <td data-bbox="793 1560 1087 1630"> <p>Weshalb?</p> </td> <td data-bbox="1087 1560 1431 1630"> <p>Weshalb?</p> </td> </tr> </table>	<p>(Re-)Binding einer Variable</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l = ['a'] >>> id(l) 4299634664</pre>	<p>Veränderliche Datenstrukturen</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l[0] = 'b' >>> id(l) 4300400112</pre>	<p>Weshalb?</p>	<p>Weshalb?</p>
<p>(Re-)Binding einer Variable</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l = ['a'] >>> id(l) 4299634664</pre>	<p>Veränderliche Datenstrukturen</p> <pre>>>> l = ['a'] >>> id(l) 4300400112 >>> l[0] = 'b' >>> id(l) 4300400112</pre>				
<p>Weshalb?</p>	<p>Weshalb?</p>				
<h3>Binding in for-Konstrukten</h3> <p>Identität vs. Wertgleichheit (equality, Äquivalenz)</p> <ul style="list-style-type: none"> <code>o1 == o2</code> testet, ob 2 Objekte/Variablen denselben Wert haben <code>o1 is o2</code> testet, ob 2 Objekte/Variablen dieselbe Identität haben, d.h. identisch sind, d.h. <code>id(o1) == id(o2)</code> <p>Was wird hier ausgegeben? ▶2</p> <pre>>>> l = [('der',1200),('die',1000),('das',900)] >>> for i,e in enumerate(l): ... print e is l[i] True True True</pre> <p>§ In for-Konstrukten werden bestehende Objekte an neue Namen gebunden!</p>	<p>Kopieren von Listen</p> <p>Eine spannende Verbindung: Binding und Listen ▶2</p> <p>Wie spielen Zuweisung von Listen-Namen und Veränderbarkeit zusammen?</p> <p>Kopieren oder Binding?</p> <pre>l1 = [('der',1200),('die',1000)] # Binding l2 = l1 # Kopieren via Slicing l3 = l1[:] # Welche Listen werden modifiziert? l1[0] = ('der',1201)</pre> <p>Kopieren via allgemeinem Modul zum Kopieren von Objekten</p> <pre>import copy l4 = copy.copy(l1)</pre>				

<pre>#!/usr/bin/env python # -*- coding: utf-8 -*- l1 = [('der',1200),('die',1000)] l2=[] for i,e in enumerate(l1): print "\nfrom 1: ", i #iterator print "\nfrom 2: ",e #element #l2[i]=(i,e) #funktioniert nicht, index out of range l2.append((i,e)) print "\nfrom 3: ",l1[i] print "\nfrom 4: ",e is l1[i] print "\nfrom 5: ",enumerate(l1) print "\nfrom 6: ",l2 help(enumerate) # class enumerate(object) # enumerate(iterable[, start]) -> # iterator for index, value of iterable # enumerate is useful for obtaining an indexed list: # (0, seq[0]), (1, seq[1]), (2, seq[2]), ... </pre>	<pre>/LubuntuShared/8) PCL-1/Vorlesung/vl11/snippets.py" from 1: 0 from 2: ('der', 1200) from 3: ('der', 1200) from 4: True from 1: 1 from 2: ('die', 1000) from 3: ('die', 1000) from 4: True from 5: <enumerate object at 0xb71dce3c> from 6: [(0, ('der', 1200)), (1, ('die', 1000))]</pre>
<pre>global_list = [('der',1200),('die',1000),('das',900)] print "from 1: ", global_list def del_first(l): print "from 2: " print 'global_list is parameter l:', global_list is l del l[0] del_first(global_list) print "from 3: ",global_list</pre>	<pre>from 1: [('der', 1200), ('die', 1000), ('das', 900)] from 2: global_list is parameter l: True from 3: [('die', 1000), ('das', 900)]</pre>
<pre>l1 = [('der',1200),('die',1000)] print "from 1: ", l1 l2 = l1 print "from 2: ", l2 # Kopieren via Slicing l3 = l1[:] print "from 3: ", l3 # Kopieren via copy modul import copy l4 = copy.copy(l1) print "from 4: ", l4 # Welche Listen werden modifiziert? l1[0] = ('der',99999999) print "from 5.1: ", l1 print "from 5.2: ", l2 print "from 5.3: ", l3 print "from 5.4: ", l4 print "\n",40*"-","\n"</pre>	<pre>from 1: [('der', 1200), ('die', 1000)] from 2: [('der', 1200), ('die', 1000)] from 3: [('der', 1200), ('die', 1000)] from 4: [('der', 1200), ('die', 1000)] from 5.1: [('der', 99999999), ('die', 1000)] from 5.2: [('der', 99999999), ('die', 1000)] from 5.3: [('der', 1200), ('die', 1000)] from 5.4: [('der', 1200), ('die', 1000)] ===== =====</pre>
<p>Sortieren und maximieren</p> <p>Ordnung erzeugen bei Dictionaries</p> <ul style="list-style-type: none"> ▶ min(), max(), in, sorted() etc. operieren über Schlüsseln. ▶ dict.values() ist Liste aller Werte. ▶ Höchster Schlüssel: max(d) ▶ Höchster Wert: max(d.values()) ▶ Schlüssel mit höchstem Wert: max(d,key=d.get) ▶ Nach Schlüssel sortieren: sorted(d) ▶ Nach Werten sortieren: sorted(d,key=d.get) ▶ Umgekehrt nach Werten sortieren: <code>sorted(d, key=d.get, reverse=True)</code> 	<pre>Dict mit (Key:Value) pairs >>> d={"neun":9,"fünf":5, "a":2,"z":99,"gg":"h"} Größter Key >>> max(d.keys()) 'z' >>> max(d) 'z' Größter Value >>> max(d.values()) 'h' Key des größten Values >>> max(d.keys(), key=d.get) 'gg' >>> max(d,key=d.get) 'gg' Value des größten Keys >>> d[max(d.keys())] 99 Undefinierbar?? >>> max(d.values(),key=d.get) 2</pre>

<p>Motivation</p> <p>key word in context</p> <p>Ziel: KWIC in Python</p> <p>Eine Klasse programmieren, welche eine Konkordanz über einem gestemmten Index anzeigt.</p> <p>Beispiel-Output</p> <pre>>>> text.concordance('die', width=30) BLACK KNIGHT : Then you shall die . ARTHUR : I command you Camelot . He was not afraid to die , O brave Sir Robin . Concorde , you shall not have died in vain ! CONCORDE : Uh 2 : Oh , he ' s died ! FATHER : And I want that ? MAYNARD : He must have died while carving it . LAUNCE ARTHUR : Look , if he was dying , he wouldn ' t bother</pre>	<p>Beispiel: Konkordanzprogramm über gestemmten Wörtern</p> <p>KWIC als Klasse: Datenstrukturen und Funktionalitäten</p> <ul style="list-style-type: none"> ▶ Text: Folge von Wörtern ▶ Index: Abbildung von (gestemmtem) Wort zu allen Vorkommenspositionen im Text ▶ Stemmer: Stemming von Wortformen ▶ KWIC-Anzeige: Formatierung der Treffer im KWIC-Stil <p>Benötigte Kompetenzen</p> <ul style="list-style-type: none"> ▶ Wie lassen sich (einfache) Klassen definieren? ▶ Definition eigener Regex-Stemmer oder Benutzung von NLTK-Stemmern ▶ Formatierung von zentriertem textuellem Output mit Format-Ausdrücken
<pre># Ziehe eine zufällige Auswahl von Vorkommen # von Token aus einem Korpus ===== import nltk, random corpus = nltk.corpus.nps_chat.words() #xrange_=0,1,2,3...len(corpus) xrange_=xrange(len(corpus)) print "\nfrom 1: ", xrange_, len(corpus) #ziehe Anzahl (anz) Zufallszahlen aus xrange_ anz=20 randsamp=random.sample(xrange_,anz) print "\nfrom 2: ", randsamp #print die Zufallswörter for i in randsamp: print "\nfrom 3: ", corpus[i] # as a reusable function with a generator return value def sample_corpus1(text,size): return (text[i] for i in random.sample(xrange(len(text)),size)) # as a reusable function with a list return value def sample_corpus2(text,size): return [text[i] for i in random.sample(xrange(len(text)),size)] print "\nfrom 4: ", sample_corpus2(corpus,20)</pre>	<pre>random_sample_xrange.py" from 1: xrange(45010) 45010 from 2: [18928, 29965, 26596, 33993, 21320, 11471, 35897, 33148, 3189, 2246, 22251, 43138, 38653, 2010, 4934, 3503, 13063, 1062, 18675, 38180] from 3: adds from 3: o from 3: one gekürzt... from 4: [u'Pat', u'thats', u'..', u'U121', u'....', u'....', u'for', u':)', u'!', u'pfft', u'she', u'didnt', u'the', u'missed', u'!', u'welome', u'U35', u'song', u'welcome', u'i'] Process finished with exit code 0</pre>
<pre># Example from NLTK book p. 107-108 import nltk, re class RegexStemmer(object): #ObjektInitialisierung def __init__(self, r= r'^(.?)(ing ly ed ious ies ive es s ment)?\$'): #Eigener Regex zur Wortstammbildung self._r = r print "\nfrom 0.0: ", type(self._r), self._r #Hier wird der WortStamm bestimmt def stem(self,word): m = re.match(self._r, word) #Stamm ^(.?? try: return m.group(1) except: return "" #falls keine Gruppe 1 besteht class IndexedText(object): #ObjektInitialisierung def __init__(self, stemmer, text): #Text ist Wortliste self._text = text print "\nfrom 0.1: ", type(self._text), self._text[:10] #Stemmer zur Stammbildung self._stemmer = stemmer print "\nfrom 0.2: ", type(self._stemmer), self._stemmer #Textindex mit WortStamm, dict, {WortStamm, AlleIndices} self._index = nltk.Index((self._stem(word), i) for (i, word) in enumerate(text)) print "\nfrom 0.3: ", type(self._index) print "\nfrom 0.3: is dict? ",issubclass(nltk.Index,dict) #Stamm in Kleinbuchstaben def _stem(self, word): return self._stemmer.stem(word).lower() #Konkordanzprint def concordance(self, word, width=20): #Aufruf der Methode stem key = self._stem(word) # stemmed keyword #width: chars vor und nach word wc = width/4 # words of context #Iteriere über AlleIndices des key=Wortstamm for i in self._index[key]: #linker Kontext ohne word lcontext = ' '.join(self._text[i-wc:i])</pre>	<pre>stemmed_kwic.py" from 1.1: from 1.2: from 0.1: <class 'nltk.corpus.reader.util.StreamBackedCorpusView'> [u'SCENE', u'1', u':', u'[', u'wind', u']', u'[', u'clop', u'clop', u'clop'] from 0.2: <class 'nltk.stem.porter.PorterStemmer'> <PorterStemmer> from 0.3: <class 'nltk.util.Index'> from 0.3: is dict? True from 1.3: KWIC mit Porter-Stemmer Defeat at the castle seems to have utterl afraid our life must seem very dull and q PIGLET : Well , what seems to be the trou . You know , it seemed a bit daft me , who , when he seemed about to reco thur and his knights seemed hopeless , wh ===== from 2.1: from 0.0: <type 'str'> ^(.?)(ing ly ed ious ies ive es s ment)?\$</pre> <p>from 2.test: seem</p> <p>from 2.2:</p> <pre>from 0.1: <class 'nltk.corpus.reader.util.StreamBackedCorpusView'> [u'SCENE', u'1', u':', u'[', u'wind', u']', u'[', u'clop', u'clop', u'clop'] from 0.2: <class '__main__.RegexStemmer'> <__main__.RegexStemmer object at 0xb1b8b5ac> from 0.3: <class 'nltk.util.Index'> from 0.3: is dict? True</pre> <p>from 2.3:</p> <p>KWIC mit simpel Regex-Stemmer</p> <pre>Defeat at the castle seems to have utterl afraid our life must seem very dull and q PIGLET : Well , what seems to be the trou . You know , it seemed a bit daft me , who , when he seemed about to reco thur and his knights seemed hopeless , wh =====</pre> <p>from 3.0:</p> <pre>GHT : Then you shall die . ARTHUR : I He was not afraid to die , O brave Sir</pre>

<pre> #word (index i) und rechter Kontext rcontext = ' '.join(self._text[i:i+wc]) ldisplay = '%*s' % (width, lcontext[-width:]) rdisplay = '%-*s' % (width, rcontext[:width]) print ldisplay, rdisplay #Text laden text = nltk.corpus.webtext.words('grail.txt') #Neues nltk stemmer objekt ===== print "\nfrom 1.1:" porter_stemmer = nltk.PorterStemmer() #Neues Text Objekt, mit nltk Stemmer und Index print "\nfrom 1.2:" porter_index = IndexedText(porter_stemmer, text) #Konkordanz von gestemmtem Wort seem print "\nfrom 1.3: ", "\nKWIC mit Porter-Stemmer\n" porter_index.concordance('seem') print "\n", "40*=", "\n" #Neues Stemmer Objekt mit eigener Regex ===== print "\nfrom 2.1:" regex_stemmer = RegexStemmer() #Test: Stemmer anwenden auf seeming print "\nfrom 2. test: ", regex_stemmer.stem('seeming') #Neues Text Objekt mit eigenem Stemmer und Index print "\nfrom 2.2:" regex_index = IndexedText(regex_stemmer, text) #Konkordanz von gestemmtem Wort seem print "\nfrom 2.3: ", "\nKWIC mit simplem Regex-Stemmer\n" regex_index.concordance('seem') print "\n", "40*=", "\n" ## Wie gut ist unser simpler Stemmer? ===== print "\nfrom 3.0:" porter_index.concordance('dies') print "\nfrom 3.1:" regex_index.concordance('dies') print "\n", "40*=", "\n" print "\nfrom 3. test: ", regex_stemmer.stem('dies') print "\nfrom 3. test: ", regex_stemmer.stem('dying') print "\nfrom 3. test: ", regex_stemmer.stem('is') print "\n", "40*=", "\n" ## Wie gut ist unser better Stemmer? ===== regex = r'^({2,}?)({ing ly ed ious ies ive es s ment})\$' print "\nfrom 3.2: ", regex better_regex_index = IndexedText(RegexStemmer(r=regex), text) print "\nfrom 3.3: ", "\nKWIC mit better Regex-Stemmer\n" better_regex_index.concordance('dies') print "\n", "40*=", "\n" </pre>	<pre> , you shall not have died in vain ! CONCO Oh , he ' s died ! FATHER : And YNARD : He must have died while carving i Look , if he was dying , he wouldn ' t from 3.1: # 1 : Where ' d you get the coconu : No , they ' d have to have it , simple ! They ' d just use a strand very nice . And how d ' you get that here . Oh ! How d ' you do ? at me , they ' d put me away ! of that . I ' d rather -- FATHER : ?! HERBERT : I ' d rather ... [music Lucky , so you ' d better get used to ARTHUR : Then we ' d best leave them he Three . And we ' d better not risk an ' aaggggh ' . He ' d just say it ! ===== from 3.test: d from 3.test: dy from 3.test: i ===== from 3.2: ^({2,}?)({ing ly ed ious ies ive es s ment})\$ from 0.0: <type 'str'> ^({2,}?)({ing ly ed ious ies ive es s ment})\$ from 0.1: <class 'nltk.corpus.reader.util.StreamBackedCorpusView'> [u'SCENE', u'1', u'!', u'', u'wind', u'!', u'[', u'clap', u'clap', u'clap'] from 0.2: <class '_main___.RegexStemmer'> <_main__.RegexStemmer object at 0xb1b96eac> from 0.3: <class 'nltk.util.Index'> from 0.3: is dict? True from 3.3: KWIC mit better Regex-Stemmer , you shall not have died in vain ! CONCO Oh , he ' s died ! FATHER : And carry on on foot . Dis - mount ! TIM YNARD : He must have died while carving i ===== Process finished with exit code 0 </pre>
--	--

<h3>Formatierung mit Hilfe von Format-Ausdrücken</h3> <ul style="list-style-type: none"> ► Flexiblere Kontrolle für Ausgabe von Zahlen und Strings ist erwünscht ► Formatierungsausdruck: 'STRINGTEMPLATE WITH FORMATS' % TUPLE ► Ein Formatierungsausdruck (<i>string formatting expression</i>) trennt Layout (Platzhalter %d, %f für Zahlen, %s für Strings) von den variablen Daten (Tupel) ► Anzahl Nachkommastellen ('%.2f'), Padding mit Leerzeichen ('% 4.2f'), linksbündig ('%-7s'), rechtsbündig ('%7s') <pre> >>> 'a string:%s and an integer:% 4d' % ('abc',3) 'a string:abc and an integer: 3' >>> 'Padding a string:%-6s and a float:% 8.2f' % ('abc',3.175) 'Padding a string:abc and a float: 3.17' </pre>	<h3>Formatierungsausdrücke</h3> <table border="1"> <tr> <td data-bbox="810 1320 1091 1545"> Überraschung <pre> >>> '%.1f' % 0.05 '0.1' >>> '%.1f' % 0.15 '0.1' >>> round(0.05,1) 0.1 >>> round(0.15,1) 0.1 </pre> </td><td data-bbox="1091 1320 1421 1545"> Prozentzeichen schützen mit % <pre> >>> '%.1f%%' % 0.15 '0.1%' </pre> </td></tr> <tr> <td data-bbox="810 1545 1091 1621"> Schulregel mit aufzurundendem .5 verzerrt systematisch (bias) </td><td data-bbox="1091 1545 1421 1621"> Variables Padding mit * <pre> >>> width = 8 >>> '%*s' % (width, 'abc') ' abc' >>> '%-*s' % (width, 'abc') 'abc ' </pre> </td></tr> </table>	Überraschung <pre> >>> '%.1f' % 0.05 '0.1' >>> '%.1f' % 0.15 '0.1' >>> round(0.05,1) 0.1 >>> round(0.15,1) 0.1 </pre>	Prozentzeichen schützen mit % <pre> >>> '%.1f%%' % 0.15 '0.1%' </pre>	Schulregel mit aufzurundendem .5 verzerrt systematisch (bias)	Variables Padding mit * <pre> >>> width = 8 >>> '%*s' % (width, 'abc') ' abc' >>> '%-*s' % (width, 'abc') 'abc ' </pre>
Überraschung <pre> >>> '%.1f' % 0.05 '0.1' >>> '%.1f' % 0.15 '0.1' >>> round(0.05,1) 0.1 >>> round(0.15,1) 0.1 </pre>	Prozentzeichen schützen mit % <pre> >>> '%.1f%%' % 0.15 '0.1%' </pre>				
Schulregel mit aufzurundendem .5 verzerrt systematisch (bias)	Variables Padding mit * <pre> >>> width = 8 >>> '%*s' % (width, 'abc') ' abc' >>> '%-*s' % (width, 'abc') 'abc ' </pre>				

<h3>Regex-Stemmer-Klasse definieren</h3> <p>Klasse mit optionalem Konstruktor-Argument</p> <pre> class RegexStemmer(object): def __init__(self, r=r'^({.*?}(ing ly ed ious ies ive es s ment))\$'): self._r = r def stem(self, word): m = re.match(self._r, word) return m.group(1) </pre> <p>Initialisierung und Verwendung</p> <pre> regex_stemmer = RegexStemmer() regex_stemmer.stem('seeming') </pre>	<h3>Textindex als Klasse IndexedText definieren</h3> <p>IndexedText</p> <pre> class IndexedText(object): def __init__(self, stemmer, text): self._text = text self._stemmer = stemmer self._index = nltk.Index((self._stem(word), i) for (i, word) in enumerate(text)) </pre> <p>enumerate(s) generiert Paare (Position,Element) aus Sequenz</p> <pre> l = ['wenn', 'fliegen', 'hinter', 'fliegen', 'fliegen'] >>> list(enumerate(l)) [(0, 'wenn'), (1, 'fliegen'), (2, 'hinter'), (3, 'fliegen'), (4, 'fliegen')] </pre> <p>nltk.Index(Pairs) erzeugt invertierten Index aus (Element,Position)</p> <pre> >>> index = nltk.Index((w,i) for (i,w) in enumerate(l)) >>> index['fliegen'] [1, 3, 4] </pre>
---	--

<h3>Private und öffentliche Methode von IndexedText</h3> <p>Unterstrich markiert Privatheit: Nur für Benutzung in der Klasse</p> <pre><code>def _stem(self, word): return self._stemmer.stem(word).lower()</code></pre> <h3>Öffentliche Methode für Formatierung</h3> <pre><code>def concordance(self, word, width=40): key = self._stem(word) # stemmed keyword wc = width/4 # words of context for i in self._index[key]: lcontext = ' '.join(self._text[i-wc:i]) rcontext = ' '.join(self._text[i:i+wc]) ldisplay = '%-*s' % (width, lcontext[-width:]) rdisplay = '%-*s' % (width, rcontext[:width]) print ldisplay, rdisplay</code></pre> <p>⌚ Noch privater sind Attribute der Form <code>o._NAME</code>.</p>	<h3>Generatoren ausdrücke (<i>generator expressions</i>)</h3> <p>Listenkomprehension: Prinzip "Liste aller Dinge, die ..."</p> <p>Bau die Liste aller kleingeschriebenen Wörter aus dem Brown-Korpus und erzeuge danach aus der Liste eine Menge!</p> <pre><code>set([w.lower() for w in nltk.corpus.brown.words()])</code></pre> <p>Generatoren ausdrücke: Prinzip "Der Nächste, bitte!"</p> <p>Nimm ein kleingeschriebenes Wort nach dem andern und mache es zum Element der Menge!</p> <pre><code>set(w.lower() for w in nltk.corpus.brown.words())</code></pre>
<h3>Listenkomprehension vs. Generatoren ausdrücke</h3> <p>Generatoren ausdrücke statt Listenkomprehension</p> <p>Im NLTK-Buch wird aus Effizienzgründen <code>set(w.lower() for w in text)</code> statt <code>set([w.lower() for w in text])</code> notiert.</p> <ul style="list-style-type: none"> ▶ Listenkomprehension erzeugt im Arbeitsspeicher immer eine Liste aller Elemente. ▶ Generatoren ausdrücke sind spechereffizient. Sie übergeben ihre Element auf Verlangen einzeln der auswertenden Funktion (intern <code>g.next()</code>). ▶ Generatoren ausdrücke unterstützen darum Längenmethode <code>len()</code> nicht. ▶ Generatoren ausdrücke unterstützen kein Slicing: <code>1[:10]</code>. ▶ Mit <code>list(generator)</code> wird jeder Generator zur Liste. ▶ Spechereffizienz ist bei allen Funktionen optimiert, welche Daten vom Typ <code>iterable</code> verarbeiten: <code>max()</code>, <code>sum()</code>, <code>set()</code> usw. ▶ Generatoren sind nach 1 Durchgang erschöpft, d.h. aufgebraucht! 	<h3>Generatoren ausdrücke und die Iteratorfunktion <code>next()</code> ▶</h3> <pre><code>>>> quadrat = (i*i for i in [10,11]) >>> quadrat <generator object <genexpr> at 0x16a6f80> >>> type(quadrat) <type 'generator'> >>> quadrat.next() 100 >>> quadrat.next() 121 >>> quadrat.next() Traceback (most recent call last): File "<stdin>", line 1, in <module> StopIteration</code></pre> <p>⌚ Die Ausnahme (<i>exception</i>) <code>StopIteration</code> erscheint, wenn der Generator erschöpft ist. Konsumenten von Generatoren müssen Ausnahme korrekt behandeln.</p>
<h3>Generatorfunktionen mit <code>yield</code></h3> <pre><code>def quadriere(iterierbar): for i in iterierbar: yield i*i quadrat = quadriere([10,11]) >>> quadrat <generator object quadriere at 0x103945500> >>> type(quadrat) <type 'generator'> >>> quadrat.next() 100 >>> sum(quadrat) 121 >>> sum(quadriere([10,11])) 221</code></pre>	<h3>Rechenzeit und Speicherverbrauch messen</h3> <p>Programm mit Generatoren ausdrücken</p> <pre><code>import nltk, timeit, time, os words = nltk.corpus.brown.words() def test_generator(): return set(w.lower() for w in words) # Initialisiere Timer-Objekt tg = timeit.Timer(test_generator) # Timing von Generatoren ausdruck print 'Timed generator (seconds):', tg.timeit(1)</code></pre> <p>⌚ Der Speicherverbrauch muss extern gemessen werden.</p>

<h3>Rechenzeit und Speicherverbrauch messen</h3> <p>Programm mit Listencomprehension</p> <pre>import nltk, timeit, time, os words = nltk.corpus.brown.words() def test_listcomprehension(): return set([w.lower() for w in words]) # Initialisiere Timer-Objekt tl = timeit.Timer(test_listcomprehension) # Timing von Listencomprehension print 'Timed list comprehension (seconds):', tl.timeit(1) ❸ Der Speicherverbrauch muss extern gemessen werden.</pre>	<h3>Effizienz in Rechenzeit</h3> <p>Zufällige Auswahl von Elementen aus einem Bereich</p> <pre># Modul zur Zeitmessung von Python-Statements import timeit setup = 'import random' # Konstruiere 2 Timer-Objekte tr = timeit.Timer('random.sample(range(1000000),100)', setup) tx = timeit.Timer('random.sample(xrange(1000000),100)', setup) # Führe Timings je einmal durch und speichere Anzahl Sekunden trsecs = tr.timeit(1) txsecs = tx.timeit(1) print "Aufgabe: Sample 100 Zahlen aus dem Bereich 0 bis 999999." print "Zeit mit xrange:", txsecs, "Sekunden" print "Zeit mit range:", trsecs, "Sekunden" print "xrange ist etwa", txsecs/trsecs, "Mal schneller!"</pre>								
<h3>Zufälliges Auswählen von Wörtern</h3> <p>Das Ziehen einer zufälligen Stichprobe (<code>sample</code>) aus einem Korpus. ▶9</p> <pre>import nltk, random corpus = nltk.corpus.nps_chat.words() # for demonstration for i in random.sample(xrange(len(corpus)),20): print corpus[i] # as a reusable function with a generator return value def sample_corpus1(text,size): return (text[i] for i in random.sample(xrange(len(text)),size)) # as a reusable function with a list return value def sample_corpus2(text,size): return [text[i] for i in random.sample(xrange(len(text)),size)]</pre>	<h3>Häufige Exceptions</h3> <p>Ausnahmen können in jeder Stufe der Programmausführung auftreten!</p> <table border="0"> <tr> <td style="vertical-align: top;"> SyntaxError <pre>print 1 2</pre> </td><td style="vertical-align: top;"> KeyError <pre>a = {} a["test"]</pre> </td></tr> <tr> <td style="vertical-align: top;"> NameError <pre>print a</pre> </td><td style="vertical-align: top;"> RuntimeError <pre>def x(): return x() x()</pre> </td></tr> <tr> <td style="vertical-align: top;"> ZeroDivisionError <pre>1 / 0</pre> </td><td style="vertical-align: top;"> TypeError <pre>sum(["1", "2", "3"])</pre> </td></tr> <tr> <td style="vertical-align: top;"> IndexError <pre>a = [1, 2, 3] a[3]</pre> </td><td></td></tr> </table>	SyntaxError <pre>print 1 2</pre>	KeyError <pre>a = {} a["test"]</pre>	NameError <pre>print a</pre>	RuntimeError <pre>def x(): return x() x()</pre>	ZeroDivisionError <pre>1 / 0</pre>	TypeError <pre>sum(["1", "2", "3"])</pre>	IndexError <pre>a = [1, 2, 3] a[3]</pre>	
SyntaxError <pre>print 1 2</pre>	KeyError <pre>a = {} a["test"]</pre>								
NameError <pre>print a</pre>	RuntimeError <pre>def x(): return x() x()</pre>								
ZeroDivisionError <pre>1 / 0</pre>	TypeError <pre>sum(["1", "2", "3"])</pre>								
IndexError <pre>a = [1, 2, 3] a[3]</pre>									
<h3>Wie gehe ich mit Fehlern um?</h3> <p>x = raw_input()</p> <p>Robuste Programmierung</p> <ul style="list-style-type: none"> ► Wir wollen x in eine Zahl umwandeln, bei ungültiger Eingabe eine neue Eingabe verlangen. ► <code>float(x)</code> führt zu Programmterminierung ► <code>x.isdigit()</code> akzeptiert nur Teilmenge aller Zahlen ► Verkettung von Regeln möglich, aber umständlich <p>Ausnahmen (Exceptions)</p> <ul style="list-style-type: none"> ► Ausnahmen können im Programm abgefangen werden, anstatt dass sie zur Terminierung führen. ► Oft eleganter, als Ausnahmen zu vermeiden. 	<h3>Ausnahmen auffangen: try-Konstrukt</h3> <table border="0"> <tr> <td style="vertical-align: top;"> Syntax-Schema <pre>try: block1 except E: block2</pre> </td><td style="vertical-align: top;"> Syntax-Schema mit finally <pre>try: block1 except E: block2 finally: block3</pre> </td></tr> <tr> <td colspan="2"> Erklärung <ul style="list-style-type: none"> ► Führe block1 aus. ► Wenn währenddessen eine Ausnahme vom Typ E auftritt, führe block2 aus ► Führe block3 auf jeden Fall am Schluss aus. </td></tr> <tr> <td colspan="2"> ❸ Ausnahmen sind ebenfalls Objekte und haben infolgedessen einen Typ! </td></tr> </table>	Syntax-Schema <pre>try: block1 except E: block2</pre>	Syntax-Schema mit finally <pre>try: block1 except E: block2 finally: block3</pre>	Erklärung <ul style="list-style-type: none"> ► Führe block1 aus. ► Wenn währenddessen eine Ausnahme vom Typ E auftritt, führe block2 aus ► Führe block3 auf jeden Fall am Schluss aus. 		❸ Ausnahmen sind ebenfalls Objekte und haben infolgedessen einen Typ!			
Syntax-Schema <pre>try: block1 except E: block2</pre>	Syntax-Schema mit finally <pre>try: block1 except E: block2 finally: block3</pre>								
Erklärung <ul style="list-style-type: none"> ► Führe block1 aus. ► Wenn währenddessen eine Ausnahme vom Typ E auftritt, führe block2 aus ► Führe block3 auf jeden Fall am Schluss aus. 									
❸ Ausnahmen sind ebenfalls Objekte und haben infolgedessen einen Typ!									
<h3>Ausnahmen ignorieren</h3> <p>▶10</p> <pre>while True: x = raw_input('Please type in a number: ') try: float(x) break except ValueError: pass</pre> <p>Leere Blöcke</p> <ul style="list-style-type: none"> ► Blöcke müssen immer mindestens eine Anweisung enthalten ► <code>pass</code> für leere Blöcke (<i>no operation, no(o)p</i>) 	<h3>Philosophien der Fehlerbehandlung: LBYL vs EAFP</h3> <table border="0"> <tr> <td style="vertical-align: top;"> LBYL <pre>if w in freqs: freqs[w] += 1 else: freqs[w] = 1</pre> </td><td style="vertical-align: top;"> EAFP <pre>try: freqs[w] += 1 except KeyError: freqs[w] = 1</pre> </td></tr> <tr> <td colspan="2"> Look before you leap. </td></tr> <tr> <td colspan="2">  It's easier to ask for forgiveness than for permission. </td></tr> </table>	LBYL <pre>if w in freqs: freqs[w] += 1 else: freqs[w] = 1</pre>	EAFP <pre>try: freqs[w] += 1 except KeyError: freqs[w] = 1</pre>	Look before you leap.		 It's easier to ask for forgiveness than for permission.			
LBYL <pre>if w in freqs: freqs[w] += 1 else: freqs[w] = 1</pre>	EAFP <pre>try: freqs[w] += 1 except KeyError: freqs[w] = 1</pre>								
Look before you leap.									
 It's easier to ask for forgiveness than for permission.									

<h3>Fehler auffangen: Wie spezifisch?</h3> <pre>try: ... (ganz viel Code) ... except: pass</pre> <p>Welche Exceptions soll man abfangen?</p> <ul style="list-style-type: none"> ▶ Zu allgemeine except-Klauseln erschweren das Bemerken und Finden von Programmierfehlern. ▶ Setze try/except-Klauseln gezielt ein. ▶ Bestimme den Ausnahmentyp, der abgefangen werden soll. 	<h3>with-Konstrukt für Datei-Handling ▶11</h3> <p>with-Konstrukt für Datei-Handling</p> <ul style="list-style-type: none"> ▶ Das Betriebssystem erlaubt nicht, dass Hunderte von Dateien von einem Prozess geöffnet sind. ▶ Bei Prozessen, welche vielen Dateien lesen/schreiben, müssen die Dateien geschlossen werden. ▶ Das with-Konstrukt mit Datei-Objekten macht dies automatisch (was auch immer für Ausnahmesituationen beim Dateiverarbeiten entstehen). <pre>filename = "with_open.py" with open(filename,'r') as f: for l in f: if l.rstrip() != '': sys.stdout.write(l)</pre>