## Regex

| | | |
|---|---|---|
| (*) indicates that the preceding character can occur 0 or more times. | meo*w | mew, meow, meooow, and meoooo-oooooooow |
| ? - character can appear either 0 or 1 time | humou?r | humour humor |
| . and it can match any single character (letter, number, symbol or whitespace) in a piece of text | ......... | any 9-char-acter text |
| [] will match any of the characters included within the brackets | con[sc]en[sc]us | consensus, concensus, consencus, and concencus |
| {} contains the exact quantity | roa{3}r | roaaar |
| {}n. the quantity range of characters to be matched | roa{3,6}r | roaaar, roaaaar, roaaaaar, or roaaaaaar |
| \|, allows for the matching of either of two subexpres-sions. | baboon-s\|g-orillas | will match the text baboons as well as the text gorillas. |

## Regex (cont)

| | | |
|---|---|---|
| Anchors (hat ^ and dollar sign $) are used in regular expressions to match text at the start and end of a string, respectively. | ^Monkeys: my mortal enemy$ | will completely match the text Monkeys: my mortal enemy but not match Spider Monkeys: my mortal enemy or Monkeys: my mortal enemy in the wild |
| [letter-letter] or [n-n] | a range of characters that can be matched | [A-Z]. : match any uppercase letter [a-z]. : match any lowercase letter [0-9]. : match any digit [A-Za-z] : match any uppercase or lowercase letter |
| Shorthand character classes simplify writing regular expressions | \w represents the regex range [A-Za-z0-9_], \d represents [0-9], | \W represents [A-Za-z0-9_] matching any character not included by \w, \D represents [0-9] matching any character not included by \d |

## Regex (cont)

| | | |
|---|---|---|
| Negated character set | [^cdh]are | will match the m in mare. |
| + ndicates that the preceding character can occur 1 or more times | meo+w | will match meow, meooow, and meoooo-ooo-ooooow, but not match mew |

## Text Preprocessing

| | | |
|---|---|---|
| Noise removal | import re result = re.sub(r'[\.\?-\!\,\:\;\"]', '', text) | Removes Punctu-ation |
| Tokenization is the text prepro-cessing task of breaking up text into smaller components of text | from nltk.t-okenize import word_t-okenize text = "This is a text to tokenize" tokenized = word_tokeniz-e(text) | print(tok-enized) # ["This", "-is", "a", "-text", "to", "tokeni-ze"] |
| In natural language processing, normalization encompasses many text prepro-cessing tasks including | stemming, lemmatiza-tion, | upper or lowerc-asing, and stopwords removal. |

## Text Preprocessing (cont)

| | | |
|---|---|---|
| Stemming In natural language processing, stemming is the text preprocessing normalization task concerned with bluntly removing word affixes (prefixes and suffixes). | from nltk.stem import Porter-Stemmer tokenized = ["So", "many", "-squids", "are", "-jumping"] stemmer = PorterStemmer() stemmed = [stemmer.stem(token) for token in tokenized] | # ['So', 'mani', 'squid', 'are', 'jump'] |
| Lemmatization In natural language processing, lemmatization is the text preprocessing normalization task concerned with bringing words down to their root forms. | from nltk.stem import WordNetLemmatizer tokenized = ["So", "many", "-squids", "are", "-jumping"] lemmatizer = WordNetLemmatizer() lemmatized = [lemmatizer.lemmatize(token) for token in tokenized] | ['So', 'many', 'squid', 'be', 'jump'] |

## Text Preprocessing (cont)

| | | |
|---|---|---|
| stopword removal is the process of removing words from a string that don't provide any information about the tone of a statement. | from nltk.corpus import stopwords # define set of English stopwords stop_words = set(stopwords.words('english')) | # remove stopwords from tokens in dataset statement_no_stop = [word for word in word_tokens if word not in stop_words] |
| parser.chunk.RegexpParser | Uses a set of regular expression patterns to specify the behavior of the parser | {<DT\|JJ>} # chunk determiners and adjectives |
| Token = Smaller Component of Text Stem = Remove prefix and suffix Lemmatization = Bring down to root Stopword = Remove meaningless | | |

## Lists and Strings

| | | |
|---|---|---|
| z = 'Natural Language Processing' | z.replace(' ', '\n') | 'Natural\nLanguage\nProcessing' |
| | list(z) | Split text into character tokens |
| | set(z) | Unique tokens |
| x = ['Natural', 'Language', 'Toolkit'] | x.insert(0, 'Python') | ['Language', 'Natural', 'Python', 'Toolkit'] |