

COMP3074 Human-AI Interaction

Lab 1: Modelling and Generating Language

Jérémie Clos

October 2023

Contents

1	Introduction	2
2	Language models	2
2.1	Building a statistical language model	2
2.1.1	Additive smoothing	3
2.1.2	The LM function	4
2.2	Generating language	5
3	Template-Based Natural Language Generation	6
3.1	A simple gap-filling system	6
3.2	More complex templates	6
3.2.1	Content determination	6
3.2.2	Document structuring	9
3.2.3	Aggregation	12
3.2.4	Lexical choice	14
3.2.5	Referring expression generation	16
3.2.6	Realisation	17
3.3	Putting it all together	19
3.3.1	Our starting dataset	19
3.3.2	The NLG pipeline	19
4	Exercises	23
5	Tasks	23


```

9 line_nl_removed = line.replace("\n", " ") # removes newline characters
10 file_nl_removed += line_nl_removed # adds filtered line to the list
11 file_p = "".join([char for char in file_nl_removed if char not in string.
    punctuation]) # joins all the lines in the list in a single string

```

Once the pre-processing is done, we can start building our language model. For simplicity, we will only build unigram, bigram, and trigram language models.

```

1 from nltk.util import ngrams
2 from nltk import word_tokenize, sent_tokenize
3
4 unigram=[]
5 bigram=[]
6 trigram=[]
7 tokenized_text = []
8
9 sents = nltk.sent_tokenize(file_p)
10
11 for sentence in sents:
12     sentence = sentence.lower()
13     sequence = word_tokenize(sentence)
14     for word in sequence:
15         if word == ".":
16             sequence.remove(word)
17         else:
18             unigram.append(word)
19             tokenized_text.append(sequence)
20             bigram.extend(list(ngrams(sequence, 2)))
21 #unigram, bigram, trigram models are created
22             trigram.extend(list(ngrams(sequence, 3)))
23
24 freq_uni = nltk.FreqDist(unigram)
25 freq_bi = nltk.FreqDist(bigram)
26 freq_tri = nltk.FreqDist(trigram)
27
28 print("5 most common unigrams:" + str(freq_uni.most_common(5)))
29 print("5 most common bigrams:" + str(freq_bi.most_common(5)))
30 print("5 most common trigrams:" + str(freq_tri.most_common(5)))

```

Here we used the FreqDist types of NLTK, which is used to represent probability distributions, to build our language models. The advantage of this is that functions from packages tend to be much faster than hand-rolling your own features.

2.1.1 Additive smoothing

Now let us do some smoothing, to avoid ending up with zeroed out probabilities. The following code computes, rather inefficiently, smoothed out language models for unigrams, bigrams, and trigrams. You will notice that it takes ages. You might even decide to take my word for it and cancel the computation half-way through (reminder: you can cancel a program running in the console with Ctrl + C).

```

1 ngrams_all = {1:[], 2:[], 3:[]}
2 for i in range(3):
3     for each in tokenized_text:
4         for j in ngrams(each, i+1):
5             ngrams_all[i+1].append(j)
6
7 ngrams_voc = {1:set([]), 2:set([]), 3:set([])}
8 for i in range(3):
9     for gram in ngrams_all[i+1]:
10         if gram not in ngrams_voc[i+1]:
11             ngrams_voc[i+1].add(gram)
12
13 total_ngrams = {1:-1, 2:-1, 3:-1}

```

```

14 total_voc = {1:-1, 2:-1, 3:-1}
15 for i in range(3):
16     total_ngrams[i+1] = len(ngrams_all[i+1])
17     total_voc[i+1] = len(ngrams_voc[i+1])
18
19 ngrams_prob = {1:[], 2:[], 3:[]}
20 for i in range(3):
21     for ngram in ngrams_voc[i+1]:
22         tlist = [ngram]
23         tlist.append(ngrams_all[i+1].count(ngram))
24         ngrams_prob[i+1].append(tlist)
25
26 for i in range(3):
27     for ngram in ngrams_prob[i+1]:
28         ngram[-1] = (ngram[-1]+1)/(total_ngrams[i+1]+total_voc[i+1])
29
30 for i in range(3):
31     ngrams_prob[i+1] = sorted(ngrams_prob[i+1], key = lambda x:x[1], reverse =
32                                True)
33
34 print("Most common (1,2,3)-grams")
35 print ("Most common unigrams: " + str(ngrams_prob[1][:10]))
36 print ("Most common bigrams: " + str(ngrams_prob[2][:10]))
37 print ("Most common trigrams: " + str(ngrams_prob[3][:10]))

```

2.1.2 The LM function

But why would we waste our time rolling our own implementations when language modelling is such a common task that libraries, such as our beloved NLTK, implement it? Indeed, the `nltk.lm`¹ package contains everything we need to roll our own optimised language models with a fraction of the code (and much faster!). Let us build a nice bigram model.

```

1 import nltk
2
3 n_param = 2
4
5 from nltk import word_tokenize, sent_tokenize
6 from nltk.util import pad_sequence
7 from nltk.lm import MLE, Laplace
8 from collections import Counter
9 from nltk.lm.preprocessing import pad_both_ends, padded_everygram_pipeline
10
11 document = "this is a test sentence! And this is another one. And this, is
12            another one again"
13
14 tokenized = [word_tokenize(sent) for sent in sent_tokenize(document)]
15 corpus, vocab = padded_everygram_pipeline(n_param, tokenized)
16 lm = Laplace(n_param)
17 lm.fit(corpus, vocab)
18
19 print("Size of vocabulary:", str(len(lm.vocab)))
20
21 # Let us print all the unigrams and their scores
22 for vocab_token in lm.vocab:
23     score = lm.score(vocab_token)
24     print("Unigram: [" + vocab_token + "] : " + str(score))
25
26 # Let us print all the bigrams and their scores
27 for vocab_token_prev in lm.vocab:
28     for vocab_token in lm.vocab:
29         score = lm.score(vocab_token, [vocab_token_prev])

```

¹<http://www.nltk.org/api/nltk.lm.html>

```
30 print("Bigram: [" + vocab_token_prev + ", " + vocab_token + "]" : " + str(score
    ))
```

You will notice that this is a lot shorter than what we did before, because NLTK is doing all of the heavy lifting. We define `n_param` to 2. This has an impact on the data preparation by padding it with `<s>` and `</s>` tokens to model the start and end of sentences, and break the text in n-grams through the `everygram_pipeline` function (which does almost everything). We then instantiate a language model with the MLE function, before training it with a specific vocabulary list. Finally, we print the scores of each unigram and bigram in order to check everything is working.

What about smoothing? NLTK.lm implements most of the smoothed models, and you just have to instantiate them instead of MLE. For example, you could use a Laplace-smoothed bigram language model as follows:

```
1 from nltk.lm import Laplace
2 lm = Laplace(2)
```

Have a look and investigate the official NLTK page on the lm module² and see the other smoothing models that are provided. Look at the scores in MLE models and smoothed models, and you can observe that MLE has a lot of zeros (unobserved n-grams) while smoothed models do not. This is Cromwell's rule³ in action.

2.2 Generating language

Another advantage of the LM NLTK module is that it makes sampling from a language model very easy. Calling the generate function lets you provide the model with a number of tokens you want to sample, a text seed that you use for sampling, and a random seed for reproducibility (specifying a random seed makes sure that re-running the same program multiple times always uses the same "random" numbers). For example, the following code samples 20 tokens conditioned on the starting token "i".

```
1 print(lm.generate(20, text_seed=['i'], random_seed=1))
```

Testing on my machine, using the Berkeley restaurant corpus, the following result is returned:

```
1 ['dont', 'want', 'to', 'eat', 'lunch', 'id', 'like', 'to', 'drive', '.', 'thai',
  'food', 'only', 'about', 'le', 'bateau', 'ivre', 'tell', 'me', 'a']
```

"I don't want to eat lunch I'd like to drive" is hardly Shakespeare, but it is intelligible English. "Thai food only about le bateau ivre tell me a" however, is a bit more shaky! Now what would happen with a trigram language model? Let us rerun the same sampling process but on a language model trained with `n_param = 3`. The result, again on my machine, was the following:

```
1 ['dont', 'want', 'to', 'eat', 'lunch', 'im', 'willing', 'to', 'go', 'an', '
  ethopian', 'restaurant', 'in', 'albany', 'how', 'about', 'chezpanisse', 'that
  ', 'was', 'cheap']
```

It is a subjective thing, but this looks a lot more like a coherent piece of text than the bigram-generated text. In fact, the higher we put the `n`, the more realistic the text becomes, because we are more likely to sample larger pieces of the original text at once! Set `n_param` large enough, and you will be sampling entire sentences from the original text, making the language model a lot less useful!

²<http://www.nltk.org/api/nltk.lm.html>

³https://en.wikipedia.org/wiki/Cromwell's_rule

3 Template-Based Natural Language Generation

Template-based natural language generation (NLG) refers to the process of producing textual output based on predefined patterns or templates. In this approach, the specific placeholders in the template are filled in with the appropriate information to generate coherent sentences or paragraphs that accomplish a specific informational purpose. It is a rule-driven method where the structure of the output is largely predetermined, as opposed to more dynamic, data-driven approaches. While template-based NLG can be quick and effective for predictable and structured data, it lacks the adaptive flair of machine learning-based NLG techniques.

3.1 A simple gap-filling system

A gap-filling system is the most basic form of template-based NLG. Define the structure of your sentence, define the gaps to be filled, plug in the relevant data and you are done. It is very versatile, provided you fill in the correct data and it allows for generating simple text.

```
1 name = input("What is your name?")
2 country_of_origin = input("Where are you from?")
3
4 output = "Hello {v1}, I see that you are from {v2}.".format(v1=name, v2=
    country_of_origin)
5
6 print(output)
```

It does the job most of the time, but it is not necessarily the most interesting way of doing things, and it leads to a rather robotic-sounding output.

3.2 More complex templates

Let us go back to the essential steps of a template-based NLG system: content determination, document structuring, aggregation, lexical choice, referring expression generation, and realisation.

3.2.1 Content determination

What do we put in the text? If we want to talk about the weather, do we put every piece of weather-related information in it? Obviously not, some selection is required depending on the communicative goal. We provide a few ways you could implement a content determination system, with some toy examples in Python. Some of those examples will make more sense later in the course, when we will have seen text classification, information retrieval, etc.

Rule and Heuristics-Based Approaches For those types of approaches, we set explicit rules and heuristics to decide which pieces of data or facts should be included. This method often relies on if-then-else conditions that are predefined based on the application's requirements.

```
1 def rule_based_description(car):
2     description = ""
3
4     if car['type'] == 'sedan':
5         description += "This is a family-friendly car."
6     elif car['type'] == 'sports':
7         description += "This car is designed for speed and performance."
8
9     return description
10
11 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
12 print(rule_based_description(car_data)) # Output: This car is designed for speed
    and performance.
```

We can also apply heuristic rules derived from domain expertise to make slightly more complex decisions. For example, when generating weather reports, always including temperature and precipitation details could be a heuristic. Another heuristic could be, when describing a car, to a threshold for horsepower at 400 to differentiate decent cars from higher-quality cars.

```
1 def heuristic_description(car):
2     description = "This is a {} car.".format(car['colour'])
3
4     if car['horsepower'] > 400:
5         description += " It has high horsepower, making it powerful."
6     else:
7         description += " It has decent performance for its range."
8
9     return description
10
11 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
12 print(heuristic_description(car_data)) # Output: This is a red car. It has high
    horsepower, making it powerful.
```

Data-Driven Methods Using statistical or machine learning models that have been trained on large datasets to learn the significance of the pieces of content. These models can prioritise and select content based on patterns in the training data. Note that the following example will not run; it is just for illustration purposes.

```
1 from sklearn.dummy import DummyClassifier
2
3 # Example dataset: type, colour, horsepower mapped to a description label
4 X_train = [['sedan', 'blue', 100], ['sports', 'red', 500], ['SUV', 'black', 300]]
5 y_train = ['family', 'performance', 'utility']
6
7 # Train a dummy classifier for demonstration
8 model = DummyClassifier(strategy="most_frequent")
9 model.fit(X_train, y_train)
10
11 def data_driven_description(car):
12     prediction = model.predict([[car['type'], car['colour'], car['horsepower']]])
13     if prediction[0] == 'family':
14         return "This is a family-friendly car."
15     elif prediction[0] == 'performance':
16         return "This car is designed for speed and performance."
17     else:
18         return "This is a versatile vehicle."
19
20 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
21 print(data_driven_description(car_data)) # Might output: This car is designed
    for speed and performance.
```

User preferences In this case, we adapt content determination based on user profiles or feedback. Over time, the system can learn to present information that is more closely aligned with the interests or preferences of the individual user.

```
1 user_preferences = {
2     'colour_interest': 'red',
3     'type_preference': 'sports'
4 }
5
6 def user_pref_description(car, preferences):
7     description = ""
8     if car['colour'] == preferences['colour_interest']:
9         description += "This car's colour matches your preference!"
10    if car['type'] == preferences['type_preference']:
11        description += " It's the type of car you like!"
```

```

12     return description
13
14 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
15 print(user_pref_description(car_data, user_preferences)) # Output: This car's
    colour matches your preference! It's the type of car you like!

```

Ontologies In the context of content determination and natural language generation, ontologies are structured frameworks that represent knowledge about a specific domain, categorising entities, and defining relationships between them. They facilitate the understanding and generation of content by providing a semantic hierarchy and context, ensuring that the generated language aligns with the structured knowledge of a particular subject or field. For example, in the medical domain, if discussing a particular disease, an ontology might suggest including symptoms, causes, and treatments.

```

1 # Let us assume a simplified ontology represented as a Python dictionary
2 car_ontology = {
3     'sports': {'related_terms': ['speed', 'race', 'performance']},
4     'SUV': {'related_terms': ['family', 'travel', 'space']}
5 }
6
7 def ontology_description(car, ontology):
8     car_type = car['type']
9     if car_type in ontology:
10         related_terms = ', '.join(ontology[car_type]['related_terms'])
11         return f"This {car_type} car is associated with {related_terms}."
12     return "Unknown car type."
13
14 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
15 print(ontology_description(car_data, car_ontology)) # Output: This sports car is
    associated with speed, race, performance.

```

Query-Driven Determination When the NLG system is tied to a query system (e.g., a search engine, a personal assistant, or a chatbot), content determination can be directly influenced by the user's query. Only information directly responding or related to the query would be selected in order to make your text as concise as possible while still meeting the communicative goal.

```

1 user_query = "red powerful cars"
2
3 def query_driven_description(car, query):
4     description = ""
5     if 'red' in query and car['colour'] == 'red':
6         description += "The car is red,"
7     if 'powerful' in query and car['horsepower'] > 400:
8         description += " and it's indeed powerful."
9     return description
10
11 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
12 print(query_driven_description(car_data, user_query)) # Output: The car is red,
    and it's indeed powerful.

```

Interactive Systems Allow users to interactively guide the content determination process. This can be done by allowing users to specify constraints, ask follow-up questions, or fine-tune the scope of information.

```

1 def interactive_description(car):
2     print("What do you want to know about the car?")
3     attribute = input("Choose: type, colour, or horsepower: ")
4     if attribute in car:

```



```

5         return f"The car's {attribute} is {car[attribute]}."
6         return "That's not a feature I know about."
7
8 car_data = {'type': 'sports', 'colour': 'red', 'horsepower': 500}
9 print(interactive_description(car_data)) # The user input will guide the
      description.

```

3.2.2 Document structuring

Document structuring focusses on how to organise the output of our system. Like the other steps, there are several ways to do so.

Rule, heuristics and templates We can use predefined rules or schemas as well as domain-specific heuristics that dictate the sequence or hierarchy of information presentation, based on the nature and relationship of the content. For example, in a research paper, starting with an introduction, followed by methods, results, and then discussion.

```

1 def rule_based_structure(city_data):
2     report = ""
3
4     if "historical_landmarks" in city_data:
5         report += f"Historical Landmarks:\n{city_data['historical_landmarks']}\n\n"
6
7     if "tourist_attractions" in city_data:
8         report += f"Tourist Attractions:\n{city_data['tourist_attractions']}\n\n"
9
10    if "restaurants" in city_data:
11        report += f"Recommended Restaurants:\n{city_data['restaurants']}"
12
13    return report
14
15 city_info = {
16     'tourist_attractions': 'Central Park, City Museum',
17     'restaurants': 'Bistro Cafe, Riverside Grill',
18     'historical_landmarks': 'Old Clock Tower, Heritage Building'
19 }
20 print(rule_based_structure(city_info))

```

Alternatively, focusing more on domain heuristics.

```

1 def heuristic_structure(city_data):
2     report = ""
3
4     # Always start with an introduction if the city name is provided
5     if "name" in city_data:
6         report += f"Discover the wonders of {city_data['name']}!\n\n"
7
8     # Always list tourist attractions before restaurants
9     if "tourist_attractions" in city_data:
10        report += f"Tourist Attractions:\n{city_data['tourist_attractions']}\n\n"
11
12    if "restaurants" in city_data:
13        report += f"Dine at:\n{city_data['restaurants']}. "
14
15    return report
16
17 print(heuristic_structure({'name': 'SampleCity', 'tourist_attractions': 'City
      Museum', 'restaurants': 'Bistro Cafe'}))

```

A similar approach would be to use preset templates with designated slots for different types of information. By filling these slots, the generated content adheres to a predictable structure.

```

1 TEMPLATE = """
2 Historical Landmarks:
3 {historical_landmarks}
4
5 Tourist Attractions:
6 {tourist_attractions}
7
8 Recommended Restaurants:
9 {restaurants}
10 """
11
12 def template_structure(city_data):
13     return TEMPLATE.format(**city_data)
14
15 print(template_structure(city_info))

```

Machine Learning We can also use machine learning and train models on structured documents to learn the typical patterns and sequences of this kind of content. Once trained, the model can predict the best structure for new information.

```

1 from sklearn.tree import DecisionTreeClassifier
2 import numpy as np
3
4 # Mock dataset: tourist attractions, restaurants, historical landmarks mapped to
5 # a sequence
6 X_train = np.array([
7     [1, 0, 1], # 1 if present, 0 if absent
8     [1, 1, 0],
9     [0, 1, 1]
10 ])
11 y_train = [1, 2, 3] # 1: historical-tourist-restaurant, 2: tourist-restaurant-
12 # historical, 3: restaurant-historical-tourist
13
14 model = DecisionTreeClassifier()
15 model.fit(X_train, y_train)
16
17 def ml_structure(city_data):
18     sequence = model.predict([[
19         1 if 'tourist_attractions' in city_data else 0,
20         1 if 'restaurants' in city_data else 0,
21         1 if 'historical_landmarks' in city_data else 0
22     ]])[0]
23
24     sections = {
25         1: f"Historical Landmarks:\n{city_data.get('historical_landmarks', '')}",
26         2: f"Tourist Attractions:\n{city_data.get('tourist_attractions', '')}",
27         3: f"Recommended Restaurants:\n{city_data.get('restaurants', '')}"
28     }
29
30     return "\n\n".join([sections[i] for i in range(1, 4) if sequence == i])
31
32 print(ml_structure(city_info))

```

Discourse Models We can make use of models that understand the conventions of textual coherence, ensuring that the document follows recognised patterns of discourse, such as introducing a topic before delving into details.

```

1 # For simplicity, let us assume a basic discourse model that always starts with
2 # an introduction, followed by details, and ends with recommendations.
3 def discourse_structure(city_data):
4     report = f"Welcome to the report on {city_data['name']}. \n\n"

```

```

5     if "historical_landmarks" in city_data:
6         report += f"Historical Landmarks:\n{city_data['historical_landmarks']}\n\n"
7
8     if "tourist_attractions" in city_data:
9         report += f"Tourist Attractions:\n{city_data['tourist_attractions']}\n\n"
10
11    report += f"Lastly, when you're hungry, consider visiting:\n{city_data['restaurants']}\n"
12
13    return report
14
15 print(discourse_structure({'name': 'SampleCity', 'historical_landmarks': '
    Heritage Building', 'tourist_attractions': 'City Museum', 'restaurants': '
    Bistro Cafe'}))

```

User-guided Structuring We could also allow users to define or tweak the structure according to their preferences, making the generated content customisable (and potentially better suited to user needs). This makes sense if we have a reason to think that users will have large differences in how they consume the information you provide.

```

1 def user_guided_structure(city_data):
2     order = input("What order do you prefer? (options: historical_landmarks,
3         tourist_attractions, restaurants, separated by commas): ").split(',')
4
5     report = ""
6     for section in order:
7         if section in city_data:
8             report += f"{section.replace('_', ' ').title()}\n{city_data[section]}\n\n"
9
10    return report
11
12 # Note: This function is interactive and waits for user input.
13 # print(user_guided_structure(city_info))

```

Graph-based Methods We can represent content as nodes in a graph, where edges indicate relationships. Algorithms can then traverse the graph to determine an optimal structure to present the information.

```

1 # Using a simple adjacency list to represent a graph, where each section points
2   to the next section.
3 graph_structure = {
4     'intro': 'historical_landmarks',
5     'historical_landmarks': 'tourist_attractions',
6     'tourist_attractions': 'restaurants',
7     'restaurants': None # end
8 }
9
10 def graph_based_structure(city_data, graph):
11     report = "Welcome to the city report.\n\n"
12     section = 'intro'
13
14     while graph[section]:
15         report += f"{graph[section].replace('_', ' ').title()}\n{city_data[graph[section]]}\n\n"
16         section = graph[section]
17
18     return report
19
20 print(graph_based_structure(city_info, graph_structure))

```

3.2.3 Aggregation

Aggregation in NLG involves combining separate pieces of information into coherent, grammatical, and often more concise sentences. This step helps ensure that the generated content is not overly verbose and is easy to understand.

Syntactic Aggregation This method merges pieces of information based on the syntactic structures of sentences, without necessarily considering deeper semantics or conceptual meanings. It focusses mainly on grammar and syntax to form coherent sentences. For example:

- "John plays football."
- "John plays basketball."

Those two statements would yield the following syntactic aggregation: "John plays football and basketball." In the above example, the aggregation is performed based on the common subject "John" and the verb "plays", combining the two objects "football" and "basketball" with the conjunction "and".

```
1 def syntactic_aggregate(sentence1, sentence2):
2     # Split sentences to extract the subject, verb, and object
3     subj1, verb1, obj1 = sentence1.split()
4     subj2, verb2, obj2 = sentence2.split()
5
6     # Check for common subject and verb
7     if subj1 == subj2 and verb1 == verb2:
8         return f"{subj1} {verb1} {obj1} and {obj2}."
9     else:
10        return None
11
12 sentence1 = "John plays football."
13 sentence2 = "John plays basketball."
14
15 print(syntactic_aggregate(sentence1[:-1], sentence2[:-1])) # Removing trailing
    periods for simplicity
```

Conceptual Aggregation This method merges information on the basis of their underlying concepts or meanings. It takes into account the semantics of the sentences and can potentially reformulate the sentences to convey combined concepts more effectively. For example:

- "John is a software developer."
- "John creates mobile applications."

Those two statements would yield the following conceptual aggregation: "John is a software developer specialising in mobile applications." Here, the aggregation goes beyond just the syntactic structure. It recognises that "creating mobile applications" is a specialisation of being a "software developer", and it merges the two ideas into a single, more descriptive sentence. For this, let us use a simple heuristic approach where we have predefined patterns for conceptual relationships.

```
1 def conceptual_aggregate(sentence1, sentence2):
2     # Recognise patterns for software development specialisations if "software
3     # developer" in sentence1 and "creates mobile applications" in sentence2:
4     return "John is a software developer specialising in mobile applications."
5
6     return None
7
8 sentence1 = "John is a software developer."
9 sentence2 = "John creates mobile applications."
10
11 print(conceptual_aggregate(sentence1, sentence2))
```

Of course, this is an overly simplistic example. In practice, special data structures such as knowledge graphs can be used to keep track of conceptual relationships and aggregate them appropriately.

```

1 class KnowledgeGraph:
2     def __init__(self):
3         # Basic graph structure: {entity: {relation: value}}
4         self.graph = {}
5
6     def add_relation(self, entity, relation, value):
7         if entity not in self.graph:
8             self.graph[entity] = {}
9             self.graph[entity][relation] = value
10
11    def get_relations(self, entity):
12        return self.graph.get(entity, {})
13
14 def conceptual_aggregate(sentences):
15     kg = KnowledgeGraph()
16
17     # Extract relations and populate the knowledge graph
18     for sentence in sentences:
19         words = sentence.split()
20         entity = words[0]
21         if "software developer" in sentence:
22             kg.add_relation(entity, 'profession', 'software developer')
23         elif "specialises in" in sentence:
24             specialisation = words[-2] + " " + words[-1].strip('.')
25             kg.add_relation(entity, 'specialises', specialisation)
26
27     # Aggregation Logic
28     aggregated_statements = []
29     for entity, relations in kg.graph.items():
30         if relations.get('profession') == 'software developer' and 'mobile
31         applications' in relations.get('specialises', ''):
32             aggregated_statements.append(f"{entity} is a software developer
33             specialising in mobile applications.")
34             # ... add more rules as needed
35
36     return aggregated_statements
37
38 sentences = ["John is a software developer.", "John specialises in mobile
39             applications."]
40 print(conceptual_aggregate(sentences))

```

A more complex version of that would use an ontology to keep track of different software development specialisations (e.g., games, websites, AI, etc.) and produce appropriate aggregates. In practice, effective aggregation often combines both syntactic and conceptual methods to ensure grammatical correctness while preserving or even improving the conveyed meaning.

Approaches for aggregation There are different methods and strategies for implementing aggregation, depending on the nature of the data and the desired result.

- **Rule-based Aggregation** Define explicit rules to determine how and when to aggregate information. For example, if two sentences have the same subject, they might be combined using a conjunction.
- **Statistical Methods** Use statistical measures or patterns derived from large datasets to decide on aggregation. For example, certain sentence structures might be more common and can be preferred during aggregation.
- **Template-based Aggregation** Define templates where multiple pieces of information can be inserted into predefined slots in the template, creating a single aggregated sentence.

- **Machine Learning Approaches** Train models on examples of aggregated content, allowing the model to learn and predict how new content should be aggregated.
- **Semantic Methods** Understand the semantic relationships between pieces of information and use this understanding to meaningfully aggregate content.
- **User-driven Aggregation** Provide interfaces for users to manually or semiautomatically guide the aggregation process according to their preferences.

3.2.4 Lexical choice

Lexical Choice is one of the steps in natural language generation (NLG), where appropriate words or phrases are chosen to convey a specific meaning, sentiment, style, or tone. It deals with selecting the best words or word sequences to use in a generated sentence, considering the context and the intended message.

Factors influencing lexical choice:

- **Specificity:** Some information requires a more specific term, while others can be more generic. For instance, "canine" vs. "dog".
- **Formality:** The level of formality required can influence the choice of words. For example, "purchase" (formal) vs. "buy" (casual).
- **Tone and Sentiment:** Depending on whether the content is meant to be positive, negative, or neutral, the word choice can vary, e.g., "challenging" vs. "problematic".
- **Variability:** To avoid repetition or to engage the reader, different lexical items might be chosen for variety, even if they mean the same thing.
- **Audience Understanding:** Consideration of the audience's familiarity with terms is vital. For example, technical jargon might be appropriate for an expert audience, but not for laypeople.

Much like the other steps, there are several possible ways of building lexical choice into an NLG system.

Rule-based Methods Predefined rules or templates dictate which terms or phrases to use based on specific conditions or contexts.

```

1 def rule_based_choice(context):
2     rules = {
3         "greeting": "Hello",
4         "farewell": "Goodbye",
5         "thank": "Thank you"
6     }
7     return rules.get(context, "Unknown")
8
9 print(rule_based_choice("greeting")) # Outputs: Hello

```

Thesaurus-based Methods Use a thesaurus or lexicon to find synonyms and select them according to context. In this example, we use WordNet⁴ (2), an ontology accessible from NLTK.

⁴<https://wordnet.princeton.edu/>

```

1 import nltk
2 from nltk.corpus import wordnet
3
4 # Ensure WordNet data is downloaded - this typically does not need to be run more
  than once as the content is cached for further use
5 nltk.download('wordnet')
6
7 def thesaurus_based_choice(word):
8     synonyms = []
9
10    # Fetch synonyms from WordNet
11    for syn in wordnet.synsets(word):
12        for lemma in syn.lemmas():
13            synonyms.append(lemma.name())
14
15    # Remove duplicates and the original word
16    synonyms = list(set(synonyms))
17    if word in synonyms:
18        synonyms.remove(word)
19
20    return synonyms[0] if synonyms else word
21
22 print(thesaurus_based_choice("happy")) # Might output: joyful

```

Statistical Methods Statistical language models can be used to determine the likelihood of the appropriateness of a word or phrase in a particular context. Feel free to go back to the previous section of this lab to have a more in-depth treatment of language models.

```

1 corpus = "I love programming. I love Python. I hate errors."
2
3 def bigram_choice(word):
4     words = corpus.split()
5     bigrams = [(words[i], words[i+1]) for i in range(len(words)-1)]
6     choices = [b[1] for b in bigrams if b[0] == word]
7     return max(set(choices), key=choices.count) if choices else None
8
9 print(bigram_choice("I")) # Outputs: love

```

Machine Learning Approaches Deep learning models, like transformers, can be trained on vast amounts of text to make contextually appropriate lexical choices. For example, using the transformers library from Huggingface (warning: this involves installing a rather large library with `pip install transformers`, which might take a long time or not work at all on some computers). You might recognise the name of this language model, GPT-2, as being the ancestor of GPT-3, which is the basis of a famous chatbot.

```

1 # You might need to start by running pip install transformers
2 from transformers import GPT2LMHeadModel, GPT2Tokenizer
3
4 tokenizer = GPT2Tokenizer.from_pretrained("gpt2-medium")
5 model = GPT2LMHeadModel.from_pretrained("gpt2-medium")
6
7 def ml_choice(partial_text):
8     inputs = tokenizer.encode(partial_text, return_tensors="pt")
9     outputs = model.generate(inputs, max_length=inputs.shape[1]+1, do_sample=
10    False)
11    decoded = tokenizer.decode(outputs[0])
12    return decoded.split()[-1] # Return the last word
13
14 print(ml_choice("Machine learning is a")) # Might output: process

```

User Feedback We can also incorporate user feedback to refine lexical choices over time, especially in interactive NLG systems. This is particularly useful if you do not know enough about the user and want to leave them the choice to personalise their experience based on their preferences and level of expertise.

```

1 suggestions = {
2     "food": ["apple", "banana", "cherry"],
3     "vehicle": ["car", "bike", "bus"]
4 }
5
6 def feedback_choice(topic):
7     print(f"Suggested word for {topic}: {suggestions[topic][0]}")
8     feedback = input("Is this okay? (yes/no) ")
9     if feedback == "yes":
10        return suggestions[topic][0]
11    else:
12        # Just moving to the next suggestion for simplicity
13        return suggestions[topic][1]
14
15 print(feedback_choice("food")) # Outputs apple or banana based on feedback

```

3.2.5 Referring expression generation

Referring Expression Generation (REG) is a subtask of Natural Language Generation (NLG) that focusses on determining how to refer to entities in a way that they can be unambiguously recognised by the listener or reader. When generating natural language, it is crucial to decide how to refer to objects or entities so that the intended audience can easily identify and distinguish them. REG concerns itself with selecting the most appropriate noun phrases to accomplish this.

For instance, if there are two apples on a table, one green and one red, instead of just saying "the apple," a proper referring expression would specify "the green apple" or "the red apple" to disambiguate between the two.

Here are some techniques and considerations for implementing REG:

- **Attribute Selection** Determine which attributes (like color, size, type, etc.) are essential to distinguish the entity from others in its context.
- **Definiteness** Decide whether to use definite ('the cat') or indefinite ('a cat') articles based on the familiarity or uniqueness of the entity.
- **Pronominal Reference** Decide whether it is appropriate to use a pronoun ("he," "she," "it") based on prior mentions and clarity.
- **Salience** Entities that have been recently mentioned or play a significant role in the discourse are more salient and might be referred to differently.

One of the most popular approaches to REG is the Incremental Algorithm (1) which attempts to find the shortest candidate description that makes an entity distinguishable from the others. Revisiting our previous example with software engineers would yield the following toy example.

```

1 developers = [
2     {"name": "Alice", "specialization": "mobile apps", "years_of_experience": 5,
3      "known_languages": ["Java", "Kotlin"], "company": "TechCorp"},
4     {"name": "Bob", "specialization": "web development", "years_of_experience":
5      3, "known_languages": ["JavaScript", "React"], "company": "WebSoft"},
6     {"name": "Charlie", "specialization": "mobile apps", "years_of_experience":
7      7, "known_languages": ["Swift"], "company": "TechCorp"},
8     {"name": "Dana", "specialization": "embedded systems", "years_of_experience":
9      10, "known_languages": ["C", "C++"], "company": "HardwareTech"},
10 ]

```



```

8 def discriminate(entity, candidates, attribute):
9     """Return true if the attribute discriminates the entity from the candidates.
10    """
11     value = entity[attribute]
12     return sum(1 for e in candidates if e[attribute] == value) == 1
13
14 def incremental_reg(entity, developers):
15     # Initially, all developers are candidates except the target entity
16     candidates = [e for e in developers if e != entity]
17
18     # Order of attributes to try (can be optimised based on context)
19     attributes_order = ["specialisation", "years_of_experience", "known_languages",
20                        "company"]
21
22     description = []
23
24     for attribute in attributes_order:
25         if discriminate(entity, candidates, attribute):
26             description.append(str(entity[attribute]))
27             break # The entity is now uniquely identified
28         else:
29             # If not discriminative, add to the description and reduce the
30             # candidate list
31             description.append(str(entity[attribute]))
32             candidates = [e for e in candidates if e[attribute] == entity[
33                        attribute]]
34
35     return ' '.join(description)
36
37 developer_to_refer = {"name": "Charlie", "specialisation": "mobile apps", "
38                      years_of_experience": 7, "known_languages": ["Swift"], "company": "TechCorp"}
39 print(incremental_reg(developer_to_refer, developers))
40 # Might output: mobile apps 7 Swift, depending on the order of attributes and
41                their distinctiveness.

```

3.2.6 Realisation

The realisation stage is the final step in the Natural Language Generation (NLG) process, where structured intermediate representations or abstract specifications are transformed into coherent, grammatically correct, and fluent natural language text.

Here are the main activities in the realisation stage:

- **Morphological Inflection** At this stage, the correct morphological forms of words are chosen based on the context. For instance, ensuring verbs agree with their subjects in number and tense.
- **Sentence Aggregation** This is about combining smaller sentences or clauses into more complex and fluent ones. For example, "John is a teacher." and "John teaches maths." might be aggregated as "John is a teacher who teaches maths."
- **Word Order Determination** In languages with flexible word order, the best order is chosen based on fluency and emphasis. For example, in English, we prefer "green big ball" over "big green ball."
- **Insertion of Function Words** Words like 'of,' 'is,' 'the,' etc., which might not be present in the abstract representation, are inserted for grammatical correctness.
- **Punctuation and Formatting** Proper punctuation marks are inserted, and the text is correctly formatted.

Now if you have paid attention, you might be thinking that Sentence Aggregation sounds a lot like the Aggregation stage we saw earlier. And that would not be completely incorrect. Although there is a subtle overlap, particularly around sentence aggregation, the two stages have different primary objectives: **Aggregation** is about **content synthesis and elimination of redundancy**, while **Realisation** is about **linguistic correctness and fluency**. In a well-structured NLG system, aggregation occurs before realisation. By the time you reach the realisation stage, the content should be well-organised, nonredundant, and ready to be transformed into fluent text.

In practice, how distinct these stages are might vary depending on the specific NLG system or framework. Some systems might blend these stages, especially when working with simpler tasks or when the intermediate representations are close to natural language. But in more sophisticated NLG systems, especially those dealing with complex data-to-text tasks, the distinction is clearer and helps in producing high-quality text.

Let us return to our running example by building an abstract representation of our software developers and their various attributes, and write a small program to generate a short bio for each of them.

```

1 developers = [
2     {"subject": "Alice", "verb": "work", "duration": 5, "field": "mobile apps", "
3     languages": ["Java", "Kotlin"], "company": "TechCorp", "location": "at"},
4     {"subject": "Bob", "verb": "specialise", "duration": 3, "field": "web
5     development", "languages": ["JavaScript", "React"], "company": "WebSoft", "
6     location": "with"},
7     {"subject": "Charlie", "verb": "work", "duration": 7, "field": "mobile apps",
8     "languages": ["Swift"], "company": "TechCorp", "location": "at"},
9 ]
10
11 def realise(developer):
12     # Morphological Inflection
13     if developer["verb"] == "work" and developer["duration"] > 1:
14         verb_form = "has worked"
15     elif developer["verb"] == "specialise" and developer["duration"] > 1:
16         verb_form = "has specialised"
17     else:
18         verb_form = developer["verb"]
19
20     # Handling list of languages
21     if len(developer["languages"]) > 1:
22         languages = ', '.join(developer["languages"][:-1]) + ' and ' + developer[
23         "languages"][-1]
24     else:
25         languages = developer["languages"][0]
26
27     # Constructing the sentence based on word choice and structure
28     if developer["location"] == "at":
29         sentence = f"{developer['subject']} {verb_form} for {developer['duration']
30         } years in {developer['field']} at {developer['company']} and is proficient
31         in {languages}."
32     else:
33         sentence = f"{developer['subject']} {verb_form} for {developer['duration']
34         } years in {developer['field']} with {developer['company']} and knows {
35         languages}."
36
37     return sentence
38
39 for developer in developers:
40     bio = realise(developer)
41     print(bio)

```

The output of this program should be the following:

```

1 Alice has worked for 5 years in mobile apps at TechCorp and is proficient in Java
  and Kotlin.

```

```

2 Bob has specialised for 3 years in web development with WebSoft and knows
  JavaScript and React.
3 Charlie has worked for 7 years in mobile apps at TechCorp and is proficient in
  Swift.

```

As you can see, the realisation process ensures that verbs are inflected correctly based on duration (work to has worked, specialise to has specialised), and chooses between "proficient in" and "knows" based on the context.

3.3 Putting it all together

Looking at the previous sections in isolation might seem daunting, so let us try to put everything together with a simple NLG system.

3.3.1 Our starting dataset

We will be working with the details of two of our star developers: Alice and Bob.

```

1 developers = [
2     {
3         "name": "Alice",
4         "role": "Senior Developer",
5         "experience": 5,
6         "education": {"degree": "MSc in Computer Science", "university": "
Imperial College London"},
7         "expertise": "mobile apps",
8         "tech_stack": [{"language": "Java", "proficiency": "expert"}, {"language"
: "Kotlin", "proficiency": "intermediate"}],
9         "past_employers": ["MobiTech", "AppStart"],
10        "current_project": {"name": "MobileBank", "description": "A next-gen
mobile banking solution."},
11        "accolades": ["Best Mobile App Developer 2022", "Tech Innovator Award
2021"],
12        "hobbies": ["photography", "cycling"]
13    },
14    {
15        "name": "Bob",
16        "role": "Frontend Developer",
17        "experience": 3,
18        "education": {"degree": "BSc in Software Engineering", "university": "
University of Manchester"},
19        "expertise": "web development",
20        "tech_stack": [{"language": "JavaScript", "proficiency": "expert"}, {"
language": "React", "proficiency": "advanced"}],
21        "past_employers": ["WebGenius", "SiteMakers"],
22        "current_project": {"name": "WebStore", "description": "A comprehensive e
-commerce platform."},
23        "accolades": ["Web Developer of the Year 2020"],
24        "hobbies": ["painting", "chess"]
25    }
26 ]

```

3.3.2 The NLG pipeline

First, we start by implementing content determination. To keep it simple, we will just take the entire dataset. If we wanted to be more subtle, we could focus on a specific subset of the data. For example, if we wanted to craft the bios of the developers for a grant proposal, we would focus on specific achievements and technical skills.

```

1 def content_determination(dev):
2     # For simplicity, we will select all attributes. But this step may involve
  more complex logic in real-world applications.

```

```
3     return dev
```

We then move on to document structuring. Once again, we do not need to get too complicated: we simply organise the information we have from our previous stage in a way that makes it easier for the following stages to transform into usable natural language text.

```
1 def document_structuring(dev):
2     return {
3         "introduction": [dev['name'], dev['role'], dev['experience'], dev['education']],
4         "skills": [dev['expertise'], dev['tech_stack']],
5         "experience": [dev['past_employers'], dev['accolades']],
6         "current_work": [dev['current_project']],
7         "personal": [dev['hobbies']]
8     }
```

The next stage is aggregation. We simply merge the information into concise and complete sentences. If we wanted to improve this, we could add alternative templates and some form of randomness to pick one of several alternative ways of stating the same sentences.

```
1 def aggregation(structured_dev):
2     aggregated_data = {}
3
4     aggregated_data["introduction"] = f"{structured_dev['introduction'][0]}, a {structured_dev['introduction'][1]} with {structured_dev['introduction'][2]} years of experience, graduated with a {structured_dev['introduction'][3]} degree" from {structured_dev['introduction'][3]} university".
5
6     tech_stack = ', '.join([f"{t['language']} ({t['proficiency']})" for t in structured_dev["skills"][1]])
7     aggregated_data["skills"] = f"Specialising in {structured_dev['skills'][0]}, they are proficient in {tech_stack}."
8
9     past_employers = ', '.join(structured_dev['experience'][0])
10    accolades = ', '.join(structured_dev['experience'][1])
11    aggregated_data["experience"] = f"Having previously worked at {past_employers}, they have been awarded with {accolades}."
12
13    project = structured_dev['current_work'][0]
14    aggregated_data["current_work"] = f"Currently, they're involved in the project '{project['name']}', which is described as: {project['description']}."
15
16    hobbies = ', '.join(structured_dev['personal'][0][:-1]) + ' and ' + structured_dev['personal'][0][-1]
17    aggregated_data["personal"] = f"Outside of professional life, their interests include {hobbies}."
18
19    return aggregated_data
```

After this, we move on to the lexical choice. For the sake of illustration, we will simply add a synonym. You can refer to the previous section to see alternative ways of dealing with lexical choice to vary parameters such as the tone or complexity of the result.

```
1 def lexical_choice(aggregated_dev):
2     text = aggregated_dev["skills"].replace("proficient", "skilled")
3     aggregated_dev["skills"] = text
4     return aggregated_dev
```

We then need to manage references within the text. Let us have some fun and deal with things like gender (defaulting to gender neutral if not provided), anaphoric references (ensuring that the developer's attributes are referred back to when relevant), and sensitivity to contextual details (past employers or specific hobbies).

```
1 def referring_expression(aggregated_dev):
```

```

2 # Initial reference with full name
3 intro = aggregated_dev["introduction"].replace(aggregated_dev['name'], f"{
aggregated_dev['name']} (hereinafter referred to as the developer)")
4
5 # Determine the pronoun based on gender
6 pronoun = 'they'
7 possessive_pronoun = 'their'
8 if 'gender' in aggregated_dev:
9     if aggregated_dev['gender'].lower() == 'male':
10         pronoun = 'he'
11         possessive_pronoun = 'his'
12     elif aggregated_dev['gender'].lower() == 'female':
13         pronoun = 'she'
14         possessive_pronoun = 'her'
15
16 # Replace subsequent direct name mentions with the pronoun
17 skills = aggregated_dev["skills"].replace(aggregated_dev['name'], pronoun)
18 experience = aggregated_dev["experience"].replace(aggregated_dev['name'],
pronoun)
19 current_work = aggregated_dev["current_work"].replace(aggregated_dev['name'],
pronoun)
20
21 # Adjust for verb conjugation with 'they'
22 if pronoun == 'they':
23     current_work = current_work.replace('has', 'have')
24
25 # Referring back to past employers when discussing current projects, if
relevant
26 if any(employer in aggregated_dev["current_work"] for employer in
aggregated_dev["experience"][0]):
27     current_work = current_work.replace("Currently", f"Following {
possessive_pronoun} tenure at {' '.join(aggregated_dev['experience'][0])}")
28
29 # Handling anaphoric references in personal info (e.g., hobbies)
30 personal = aggregated_dev["personal"]
31 if 'coding' in personal:
32     personal = personal.replace('coding', 'this activity') # Assuming a
mention of coding as a hobby after introducing profession
33
34 referred_content = {
35     "introduction": intro,
36     "skills": skills,
37     "experience": experience,
38     "current_work": current_work,
39     "personal": personal
40 }
41
42 return referred_content

```

For the final stage, we simply join all the resulting words together and add a bit of flair to the descriptions. This is of course a lot simpler than the methods described in the previous section on realisation, but we are mostly trying to illustrate the different phases. Nowadays, complex NLG frameworks automate much of the details.

```

1 def realisation(lexical_dev):
2     intro = lexical_dev["introduction"]
3
4     # Transition to skills
5     skills = "Delving into technicalities, " + lexical_dev["skills"]
6
7     # Seamless connection between past experience and accolades
8     experience = lexical_dev["experience"].replace(", they have", ". Additionally
, they have")
9
10    # Highlight the present role

```

```

11     current_work = "On the present front, " + lexical_dev["current_work"]
12
13     # Conclude with personal interests
14     personal = "When away from the code, " + lexical_dev["personal"]
15
16     return f"{intro} {skills} {experience} {current_work} {personal}"

```

We then need to call all the functions in the correct order, resulting in the following code.

```

1 for dev in developers:
2     content = content_determination(dev)
3     structured = document_structuring(content)
4     aggregated = aggregation(structured)
5     lexically_enriched = lexical_choice(aggregated)
6     referred = referring_expression(lexically_enriched)
7     final_output = realisation(referred)
8     print(final_output, "\n")

```

4 Exercises

Exercises are basic tasks that can help you better understand the content. You are not expected to finish everything in the lab.

1. Print some statistics about your document: number of sentences, number of tokens, average number of tokens per sentence.
2. Print some statistics about your language models: 10 most common unigrams/bigrams/trigrams, 10 least common unigrams/bigrams/trigrams, and their probability of appearing in the text.
3. Using data from the Gutenberg library, build a language model that can write in the style of your favourite author.
4. Using data from the IMDB movie review dataset, build two language models: one to generate positive reviews, and one to generate negative reviews.

5 Tasks

Tasks are much more time-consuming than exercises and allow you to go much deeper into content mastery. We do not expect you to finish all of them in the lab, but they can be good toy problems to practice with.

1. [★] Build yourself a training and a test set, and use perplexity to compare the performance of a basic, non-smoothed language model with a language model using the smoothing function of your choice.
2. [★ ★] The LM module allows you to calculate the perplexity of your language model. One common use of language model perplexity is to find out the author of a text. Implement a plagiarism detector using language models (without googling “perplexity author detection”).
3. [★ ★ ★] Can you make two language models talk to each other? Pair up with a friend (or with yourself) and train two language models on different datasets (there are several open-source books on the Gutenberg library, but you can find resources elsewhere if you want). Then, work together to make your language models speak to each other, in turns.
4. [★ ★ ★] Bootstrap your own dataset by building an object similar to the one described in 3.3.1, containing the details of your CV, and develop a basic NLG system which outputs your bio and adapts the language, structure, and amount of details based on the following criteria:
 - technical vs non-technical audience (vocabulary used);
 - complete bio vs shortened bio (amount of details);
 - purpose of the bio (job application, social event, dating profile).
5. [★ ★] Building on the previous task, embed your NLG system in an I/O loop as seen in the previous lab, and make it progressively build a profile of the user that can be used to generate the bio.

References

- [1] DALE, R., AND REITER, E. Computational interpretations of the gricean maxims in the generation of referring expressions. *Cognitive science* 19, 2 (1995), 233–263.
- [2] FELLBAUM, C. *WordNet: An electronic lexical database*. MIT press, 1998.