

Data Cleaning & Text Processing

Exercise 1

Pranav, Carolin and Karina

pranav.agrawal@uni-hamburg.de, karina.vida@uni-hamburg.de

1 Data Cleaning

One of the most pain-staking tasks in NLP is cleaning the data. Have a look at this table. One of our companies set up a form to collect the details of the donations of their fundraiser and the collected data is really messy. Here are some of the rows:

Email	Phone Number	Country	Donation
abc@gmailcom	04288221533	USA	20\$
xyz[at]gmail.com	+497161768885	Germany	50€
hello@uni-hamburg.de	0232 4495146	U.S.A	25.0
hamburg @uni-hamburg.de	02841371135a	Deutschland	7
g123@gmail.com	02773 919199	HK	500.0 HKD
dayta@co.de	06181965021510293	Hong Kong	\$ 70

1. Point out where the data is messing up and suggest the ways to clean it up.
2. Write a one-shot prompt for a LLM (Large Language Model) to clean and normalize the country column.
3. Write a one-shot prompt for a LLM (Large Language Model) to clean and normalize the donation column.
4. The best way to avoid irregularities in the phone number column is to provide a regular expression-based validation. Write a regular expression for a German phone number.
5. Similarly, the best way to avoid irregularities in the email column is to provide a regular expression-based validation. Write a basic regular expression for an email.
6. Once the data is cleaned, you need to feed it to an API which accepts the data in JSON form. Write a prompt with a template to convert this dataset into JSON form.

2 Regular Expressions and Text Analysis

Your TA, Pranav, has recently started learning German. As a text analysis nerd, he plans to use his NLP skills to improve his German vocabulary. Check out this beginner-level paragraph:

Der Hund bellt laut. Der große Hund spielt im Garten. Die Katze schläft auf dem Sofa. Der braune Hund jagt die Katze. Das kleine Kind streichelt den Hund. Die Katze trinkt Milch. Der Hund frisst sein Futter. Die Sonne scheint hell. Das Kind geht in die Schule. Die Schule ist groß. Der Lehrer unterrichtet in der Schule. Das Kind lernt viel. Der Hund wartet zu Hause. Die Familie kommt nach Hause. Der Hund freut sich. Die Katze bleibt auf dem Sofa. Am Abend essen alle zusammen. Der Hund liegt unter dem Tisch.

1. One of the biggest challenges for people who are learning German is to learn the grammatical gendered article of the words (like der, das, die). Write a regular expression to extract the German article followed by the word (for example: Der Hund, die Katze, das kind).
2. It is important for Pranav to learn the most frequent collocations to improve his vocabulary. Write a one-shot prompt and expected output for the LLM to find out at least 4 collocations here.
3. Pranav will build a German-English vocab list. For that, he wants to extract nouns, verbs and adjectives from the German text and translate them into the English. He also wants to include variants of the verbs and plurals of the nouns in his vocab list. See the below picture of how it should look like. Write a few-shot prompt to output a similar list.

der Anzug , "-e	suit
aus gehen , er geht aus , ist ausgegangen (<i>Abends beim Ausgehen haben wir viel Spaß.</i>)	to go out (<i>We have a lot of fun when we go out in the evening.</i>)
das Kleid , "-er	dress
die Kleidung (Sg.)	clothes
die Krawatte, "-n	tie
der Pullover , "-	jumper
tragen , er trägt, hat getragen	to wear
die Bluse , "-n	blouse
die Jeans , "-	jeans
der Mantel , "-	coat
die Mütze , "-n	hat

3 More Data Cleaning

1. Check out these tweets and texts. What is the best way to clean this data? What tricks would you use to remove repeated letters, emojis and confidential information?

- SLAYYYYYY!!!! ST. PAULI just SLAYED THAT GAME YAAASSS.



- - No worries, you can use my password which is katy-perry-flop.
 - John Doe, 41, from Hamburg, whose salary is 5000 euros per month, recently found out that he has diabetes.
 - I recently found out Stephanie and Christina in our school are gay but they haven't told anyone but me.
2. Many times, the texts might contain swear words. Would you consider removing those words? When would you do that, explain with examples.
 3. In order, to combat swear words, the earlier versions of OpenAI models used a simple blocklist to clean the data.¹ This blocklist contained many swear words but also contained words like *transgender*, *twink*, *homosexual* and so on. How this is problematic? Would you still recommend going for blocklists in cleaning the data?

¹<https://arxiv.org/pdf/2104.08758>