

CSC401/2511 Assignment 1 Tutorial 1

Winston (Yuntao) Wu

University of Toronto

Overview

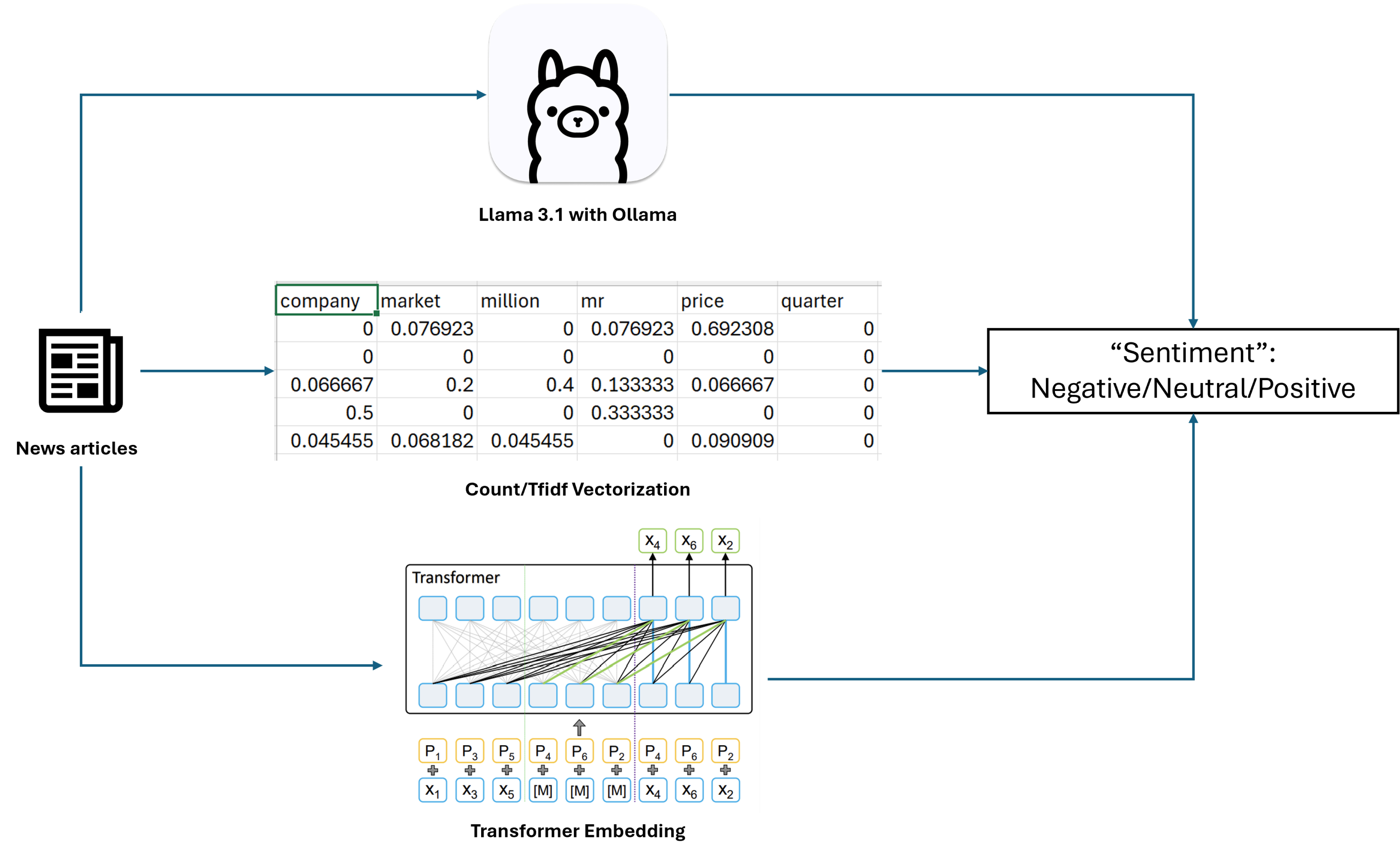


Image Source1 Image Source2

Dataset (FPB)

Financial Phrase Bank:

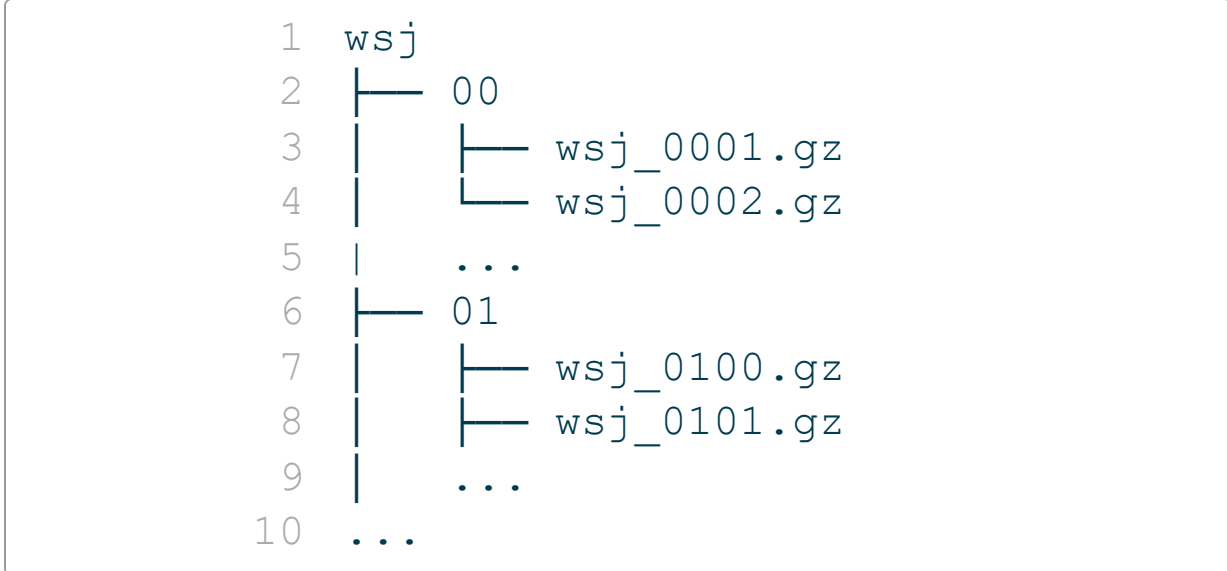
- Polar sentiment dataset of sentences from financial news
- 4k sentences from English language financial news rated by 5-8 annotators
- The dataset can be accessed from /u/cs401/A1/data/fpb_dataset.parquet

		text	label	label_numeric
0		According to Gran , the company has no plans t...	neutral	0
1		Technopolis plans to develop in stages an area...	neutral	0
2		The international electronic industry company ...	negative	-1
3		With the new production plant the company woul...	positive	1
4		According to the company 's updated strategy f...	positive	1

Dataset (WSJ89)

Wall Street Journal (1989):

- 2k articles from 1989 Wall Street Journal from the English Penn Treebank corpus. Dataset is at </u/cs401/A1/data/wsj.gz>
- We tag each article by the sign of 5-day log return of S&P500 on previous day ($\log \frac{p_{t-1}}{p_{t-6}}$). Tags are save in /u/cs401/A1/data/wsj89_labels.parquet
- You do not need to worry about WSJ89 in first part.



wsj.gz File Structure

```
.START

After a bad start, Treasury bonds were buoyed
by a late burst of buying to end modestly
higher.

"The market was pretty dull" for most of the
day, said Robert H. Chandross, vice president
at Lloyds Bank PLC.
He said some investors were reluctant to plunge
into the market ahead of several key economic
indicators due this week, especially Friday's
potentially market-moving employment report.
```

Sample article

	fn	label_numeric
0	wsj_0001.gz	-1.0
1	wsj_0002.gz	-1.0
2	wsj_0003.gz	-1.0
3	wsj_0004.gz	1.0
4	wsj_0005.gz	-1.0

wsj89_labels.parquet

Part 1 Llama 3.1 Analysis (due Sept 24)

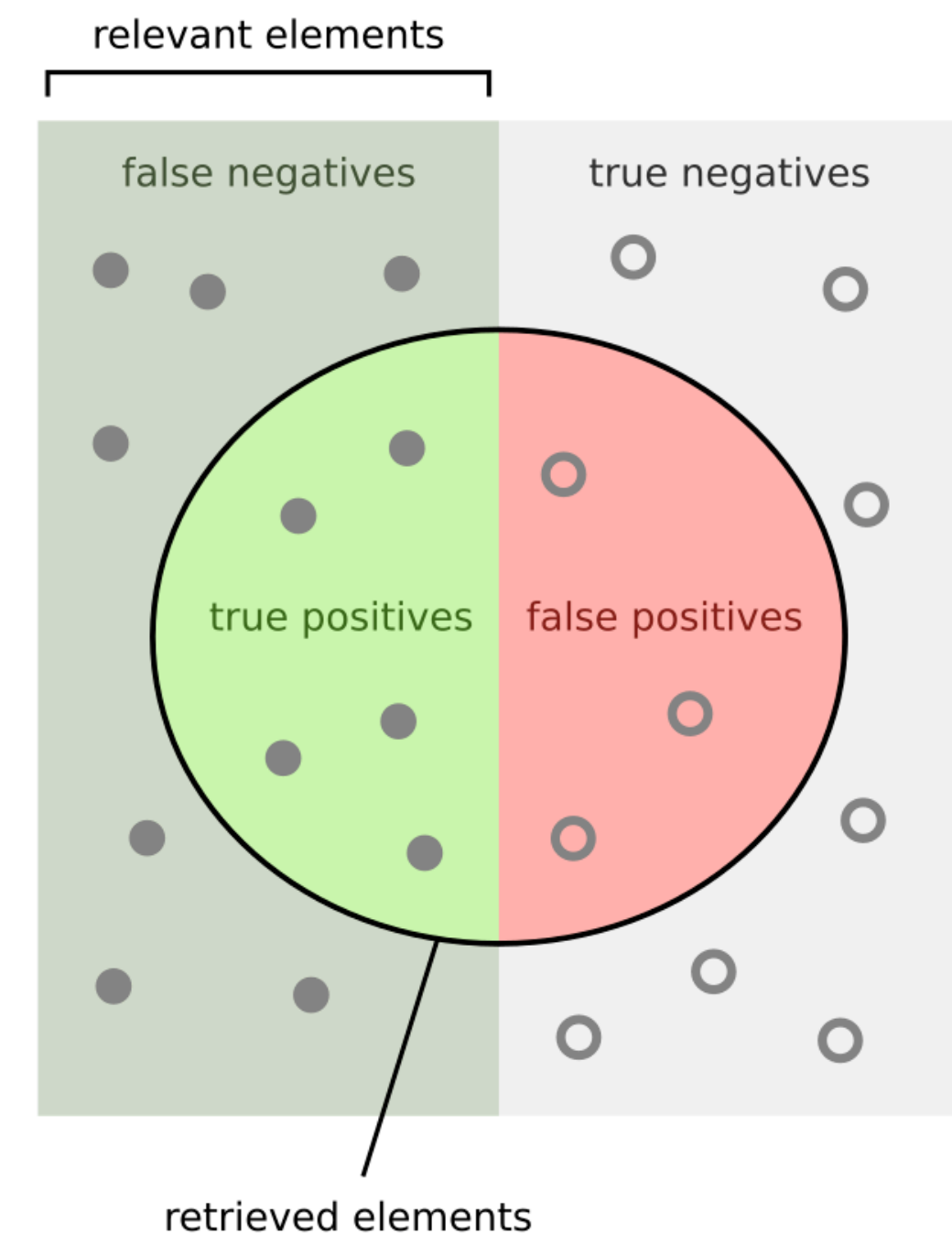
In this part, you are tasked with sentiment analysis using Llama 3.1. Submit on Markus (A1p1).

Deliverables:

- [a1_utils.py](#): Functions for computing metrics (accuracy, precision, recall)
- [a1_llama3.py](#): Select articles from FPB, query Llama 3.1, parse the response to get the sentiment scores, and compute metrics
- [a1_part1.txt](#) [Due Oct 8th with the full assignment]: Save the accuracy, precision, recall rate by the provided format string. Comment on the performance of Llama3.1, based on the accuracy, precision, recall rate, and samples of the generated explanation (at least one for correct and at least one for incorrect, with your comments)

Part 1: Metrics

- Accuracy: The total number of correctly classified instances over all classifications: $A = \frac{\sum_i c_{i,i}}{\sum_{i,j} c_{i,j}}$
- Recall: For each class κ , the fraction of cases that are truly class κ that were classified as κ , $R(\kappa) = \frac{c_{\kappa,\kappa}}{\sum_j c_{\kappa,j}}$.
- Precision: For each class κ , the fraction of cases classified as κ that truly are κ , $P(\kappa) = \frac{c_{\kappa,\kappa}}{\sum_i c_{i,\kappa}}$.
- You can use [sanity_check.py](#) to check your implementation.



How many retrieved items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are retrieved?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Part 1: Llama 3.1

We have deployed a [Llama 3.1 8B Instruct](#) model on teach.cs server. The model is loaded with 4-bit quantization using [Ollama](#).

The following is a replication code for running Llama 3.1 using pure huggingface transformer library:

```

1 # Load model
2 model_id = "./models/Meta-Llama-3.1-8B-Instruct"
3 tokenizer = AutoTokenizer.from_pretrained(model_id, trust_remote_code=True)
4 model = AutoModelForCausalLM.from_pretrained(
5     model_id,
6     quantization_config=BitsAndBytesConfig(load_in_4bit=True),
7     device_map="auto",
8     trust_remote_code=True
9 ).eval()
10
11 # Tokenize chat input
12 messages = [
13     {"role": "system", "content": "You will compute the sentiment score of articles. Start with the classification, capitalized. Then"},
14     {"role": "user", "content": input_text},
15 ]
16 input_ids = tokenizer.apply_chat_template(
17     messages,
18     add_generation_prompt=True,
19     return_tensors="pt"

```

You don't have to worry too much about the parameter settings. You will learn more about [top_k](#) and [temperature](#) in this class.

If you want to change any parameter or the system prompt, you need to specify what you changed, and the results before/after the change in your written analysis ([a1_part1.txt](#) and [a1_compare.txt](#))

Part 1: Server Setup

The Llama 3.1 is setup as a 4-bit quantized model on teach.cs server through Ollama.

It is wrapped by a queue-based (FIFO) request processing system that checks your user id. The logic of `send_request(user_id, input_text)` is as follows:

- Submit request to `http://csc401:8000/submit_request_llama3`. This will give you a `task_id` that's put onto the queue.

```
1 messages = [  
2     {"role": "system", "content": "You will compute the sentiment score of articles. Start with the classification, capitalized. Ther  
3     {"role": "user", "content": input_text},  
4 ]  
5 req_body = {  
6     "user_id": user_id,  
7     "messages": messages,  
8     "max_new_tokens": 256,  
9     "temperature": 0.5,  
10    "top_p": 0.95,  
11    "top_k": 10,  
12 }  
13 res = requests.post("http://csc401:8000/submit_request_llama3", json=req_body, headers={"Content-Type": "application/json"})
```

- Then use this `task_id` to check your request status through:

```
1 res = requests.get(f"http://csc401:8000/check_request_result?user_id={user_id}&task_id={task_id}", headers={"Content-Type": "applicat
```

Part 1: Requesting the Server

Once the request is processed, the expected output by calling `send_request` in `a1_llama3.py` is:

```

1 {
2     "model": "registry.ollama.ai/library/llama3:latest",
3     "created_at": "2023-12-12T14:13:43.416799Z",
4     "message": {
5         "role": "assistant",
6         "content": response
7     },
8     "done": true,
9     "total_duration": time spent in nanoseconds generating the response,
10    "load_duration": time spent in nanoseconds loading the model,
11    "prompt_eval_count": number of tokens in the prompt,
12    "prompt_eval_duration": time spent in nanoseconds evaluating the prompt,
13    "eval_count": number of tokens in the response,
14    "eval_duration": time in nanoseconds spent generating the response,
15    "status": "done"
16 }
```

You need to parse the response by `parse_response` in `a1_llama3.py` to the following form:

```

1 {
2     "label": extracted text label (POSTIVE, NEUTRAL, NEGATIVE) from the response,
3     "label_numeric": 1/0/-1 for positive, neutral, negative,
4     "raw_result": ["message"]["content"],
5     "compute_time": total_duration in second
6 }
```

Part 1: Process Dataframe

You need to write the code for `process_df` in `a1_llama3.py`:

- Use `np.random.choice(len(df), sample_num, replace=False)`, with the specific seed we set using your UTORid
- Get the text, and label at the specific indices, submit the request using `send_request`. The request may not always be successful. Error handling is expected (e.g. resend the request, or save a failed value and try later).
- Parse the response using `parse_response`. You can assume that at least one of `POSTIVE`, `NEUTRAL`, `NEGATIVE` exist in Llama3.1 response. However, you can also manually check the response and edit in the result dataframe before computing the metrics.
- Any accuracy rate is expected, based on your student number. We are not looking for the best accuracy.
- You are expected to provide some examples where Llama predicts correctly and where Llama predicts incorrectly, with your comment (at least one for correct and at least one for incorrect)

Part 1: Submission

You need to submit the following files by Sept 24, which is 7/10 for part 1. If you do not submit them on time, they will not be graded in the final submission.

- [a1_utils.py](#)
- [a1_llama3.py](#)

You should resubmit the above two files and the following file for part 1 when the assignment is due on Oct 8.

- [a1_part1.txt](#): accuracy, precision, recall rate by the provided format string. Comment on the performance of Llama3.1, based on the accuracy, precision, recall rate, and samples of the generated explanation (at least one for correct and at least one for incorrect, with your comments)

Part 1: Notes

- Computing sentiment scores using Llama 3.1 takes around 4 seconds per article. The server operates on a queuing system, so you may experience delays. Cache the results in a dataframe and save them locally to avoid repeated endpoint calls.
- Only one Llama 3.1 server is available for CSC401/2511 and CSC485. Start early to avoid long wait times in the request queue.
- Use the provided server exclusively for assignments, not personal projects. Each UTORid will be rate-limited.
- Use your own UTORid in your requests. Only requests from registered UTORids will be processed. We will track the total number of requests per UTORid as partial validation of your work.
- Llama 3.1 request may fail the first time, but should work properly afterwards. If Llama 3.1 request constantly fails, firstly try to increase the initial wait time `time.sleep(2)` to a longer time. If it does not help, contact the instructors on Piazza.
- Use `python3.10` instead of `python3` for consistent versioning with later parts.

VS Code Debugging

Install the Python extension for VSCode if you have not yet.

In Debug panel, choose `create a launch.json file` with Debug configuration set to Python File.

The screenshot shows the VS Code interface with a debug configuration file named `launch.json` open. The file is located in the `Launch Targets` folder. The configuration is as follows:

```

2 // Use IntelliSense to learn about possible attributes.
3 // Hover to view descriptions of existing attributes.
4 // For more information, visit: https://go.microsoft.com/fwlink/?linkid=830387
5 "version": "0.2.0",
6 "configurations": [
7   {
8     "name": "a1_preprocess",
9     "type": "debugpy",
10    "python": "python3.10",
11    "request": "launch",
12    "cwd": "${workspaceFolder}/demo",
13    "program": "a1_preprocess.py",
14    "args": [
15      "--a1_dir", "data",
16      "--filename_prefix", "wsj89",
17    ],
18    "console": "integratedTerminal",
19    "justMyCode": true
20  },
21
22  "name": "Python: Current File",
23  "type": "debugpy",
24  "request": "launch",
25  "program": "${file}",
26  "console": "integratedTerminal",
27  "justMyCode": true
28 ],
29 ]
30 }

```

The left sidebar shows the Debug Console with the following sections:

- VARIABLES**: Empty.
- WATCH**: Empty.
- CALL STACK**: Empty.
- BREAKPOINTS**:
 - Raised Exceptions
 - Uncaught Exceptions
 - User Uncaught Exceptions
 - a1_extract_features.py ~/a1_test (92)
 - a1_preprocess.py ~/student_submissi... (135)
 - a1_preprocess.py ~/demo (112)
 - test_34_accuracies.py ~/submissions2/... (8)
 - test_34_accuracies.py (8)
 - test_malloc.c (5)

The status bar at the bottom indicates the current file is `a1_preprocess (t4wuyunt)` and the cursor is at line 21, column 10. The status bar also shows the file encoding as UTF-8 and line endings as LF.

Debug and Watch Function

You can add breakpoints and start the debug function as highlighted in the image. Be sure to select a proper python interpreter (3.10.13 for A1) before you start.

Watch function helps you evaluate any expressions as you debug

The screenshot shows the VS Code interface with the following components:

- Run and Debug Sidebar:**
 - VARIABLES:** Shows local variables like `special variables`, `function variables`, and `args`.
 - WATCH:** Shows expressions being evaluated, such as `arr[0,1]` and `np.all(arr==0)`. A note says: "Watch panel, click the plus symbol to add any values you want to evaluate".
 - CALL STACK:** Shows the current step: `<module> a1_extract_features.py 91:1`.
 - BREAKPOINTS:** Shows a breakpoint set at `a1_extract_features.py ~/a1_test 90`.
- Main Editor:** Shows a Python script `a1_extract_features.py` with a breakpoint at line 91. A note says: "Region to add breakpoints".
- Terminal:** Shows the command: `bash a1_test Python Deb...`. A note says: "Here, all the parameters are properly set by the previous step configuration." and "Make sure you have the correct python interpreter".
- Status Bar:** Shows the current Python interpreter: `Python 3.11.6 64-bit`.