

Capstone Project Report for legal document text classification

Project Overview

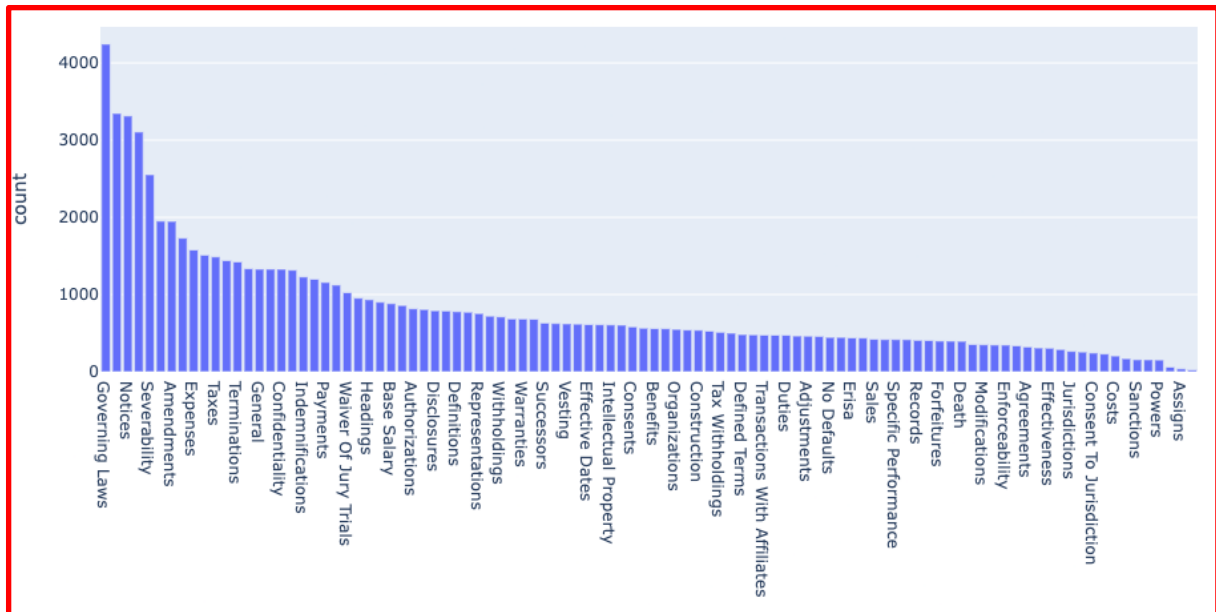
Unstructured free text is a prosperous source of information; however, extracting insights from it can be a complicated and highly time-consuming task. In recent years, thanks to many breakthroughs in the field of Natural Language Processing (NLP), which is a subfield of Machine Learning (ML), extracting actionable insight has significantly simplified. One of the essential parts of NLP is text classification which allows Data scientists to transform an unstructured corpus of text data into a structured format. During the past couple of years, novel algorithms such as Transformers[1], Universal Language Model Fine-tuning for Text Classification[2] and pre-trained language models such as Universal Sentence Encoder[3] and Pre-training of Deep Bidirectional Transformers for Language Understanding (BERT)[4] has improved the quality of text classification tasks while requiring much less tagged data to achieve state of the art results.

Problem statement

In this project, I created a SageMaker pipeline for a text classification application. The pipeline takes raw text data with its corresponding tag in a parquet format and uses varieties of NLP tools to train a text classification model. And at the end, the best-trained model that was selected based on the F1-score from the hyperparameter tuning step was selected and deployed to a SageMaker endpoint with autoscaling for real-time inference. Bespoke docker images for each node were developed and used to maximize flexibility.

Dataset

In this work I have used LEDGAR[1] dataset which is a multilabel corpus of legal provisions in contracts with label set of over 12'000 labels annotated in almost 100'000 provisions in over 60'000 contracts. The data was obtained from the US Securities and Exchange Commission (SEC) filings, that are publicly available. Each label represents the single main topic of the relevant contract.



As we can see from the above image, the number of available text for each “clause_type” ranges from ~4000 to ~20. In this work I will make a text classification for the top 5 classes and will add an “other” class to capture anything else. The following is the list of “clause_type” that I’m going to train a text classification model for:

1. Governing Laws
2. Counterparts
3. Notices
4. Entire Agreements
5. Severability
6. Others

Evaluation metrics

The evaluation metric for this work is going to be a multi-class F1-Score. This metric is a harmonic mean of precision and recall, which is the preferred metric for unbalanced datasets like we have here.

$$\text{F-1 Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The above equation is used to calculate F1-Score.

Solution statement:

Figure 2 depicts that pipeline that was developed for this work. In this section I will go over what each section is designed to produce.

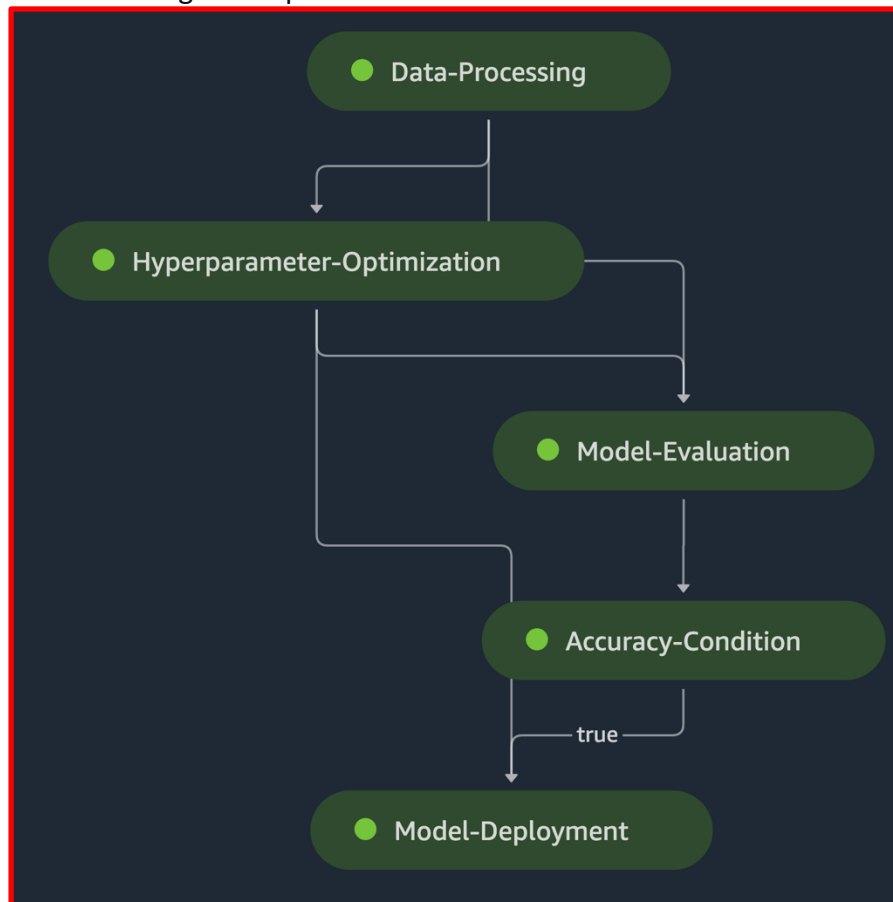


Figure2. SageMaker pipeline graph

Data-Processing:

Figure 3 shows an overview of the Data-Processing step which includes reading the data from an S3 bucket, process the data using the processing.py file, and uploading the processed data to a folder in S3 bucket that is then used in Hyperparameter-Optimization step and Model-Evaluation step.

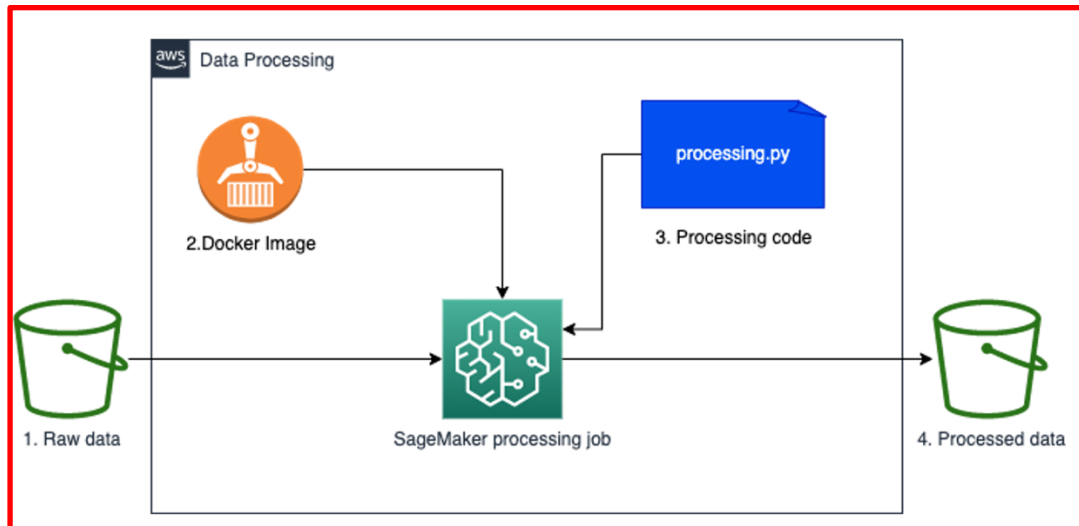


Figure3. BYOC Processing Job

The processing.py file is performing the following:

1. Load the parquet data file to a Panda DataFrame
2. Split the data into train/dev/test
3. Remove some stop words and clean the text
4. Reduce the number of target variables from 100 to six
5. Transform the text to vector embedding, using one of the Spacy language models
6. Save the DataFrame to a parquet file

Hyperparameter-Optimization

In this work the best model performance was selected using a HyperparameterTuner method from SageMaker pipeline. Figure 4 shows the overview of this step.

1. Processed data:
 - S3 bucket or feature store
2. Docker image on ECR
 - Python image contains all the libraries for model training
 - train.py: training workflow
3. Export model artifact to a dedicated S3 bucket

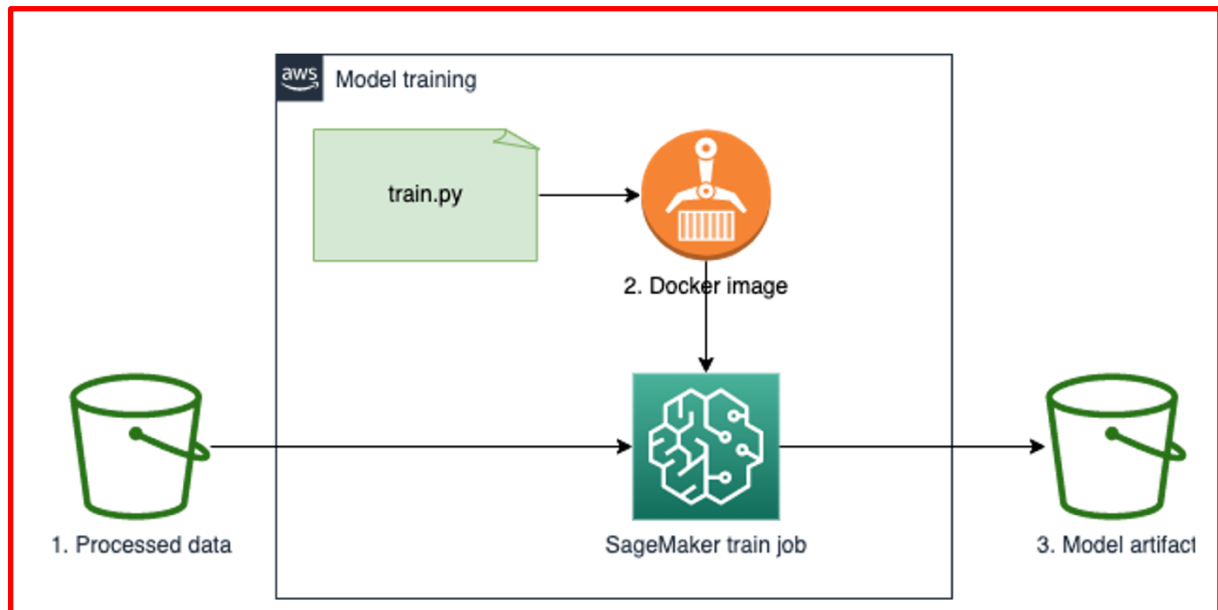


Figure4, model training and selection

To select the best model, I used the bayesian optimizer module for the hyperparameter tuning job in SageMaker and searched two parameters, `n_estimators` and `min_sample_split`. Figure5 shows that we have trained four models, and the one with the golden star is the best performing model, with an oval F1 score of 96%.

| Hyperparameter-Optimization | | X | |
|-------------------------------|--------|----------------------|-------------|
| Input | Output | Logs | Information |
| Top Training Jobs ⓘ | | Tuning Job Objective | |
| ★ 52et02alu9lu-Hyperpa-H75... | | 0.960152804851532 | |
| 52et02alu9lu-Hyperpa-H75oM... | | 0.9587489366531372 | |
| 52et02alu9lu-Hyperpa-H75oM... | | 0.9582776427268982 | |
| 52et02alu9lu-Hyperpa-H75oM... | | 0.9559884667396545 | |
| Overall Best Training Job | | Tuning Objective | |
| 52et02alu9lu-Hyperpa-H75oM... | | 0.960152804851532 | |
| Tuned Hyperparameters ⓘ | | Value | |
| min_samples_split | | 4 | |
| n_estimators | | 100 | |
| Training Job Status | | Number of Job | |
| Completed | | 4 | |
| InProgress | | 0 | |
| RetryableError | | 0 | |
| NonRetryableError | | 0 | |
| Stopped | | 0 | |

Figure5, Hyperparameter Optimization

| Training job status counter | | | | | | |
|---|-----------|------------------------|------------------------|-------------------|--|--|
| Completed 1 In Progress 0 Stopped 0 Failed 0 (Retriable: 0, Non-retriable: 0) | | | | | | |
| Training jobs | | | | | | |
| Sorting by objective metric value will display only jobs that have metric values. | | | | | | |
| <input type="text" value="Search training jobs"/> View logs View instance metrics Stop Create model | | | | | | |
| Name | Status | Objective metric value | Creation time | Training Duration | | |
| S2et02aUuH-Hyperpa-H75oMR6de-004-5ec66f1 | Completed | 0.958748936651372 | Aug 16, 2022 10:54 UTC | 3 minute(s) | | |
| S2et02aUuH-Hyperpa-H75oMR6de-003-c11a3cd | Completed | 0.9559884667396545 | Aug 16, 2022 10:54 UTC | 2 minute(s) | | |
| S2et02aUuH-Hyperpa-H75oMR6de-002-5c9a2640 | Completed | 0.9582776427268982 | Aug 16, 2022 10:54 UTC | 3 minute(s) | | |
| S2et02aUuH-Hyperpa-H75oMR6de-001-b89aa556 | Completed | 0.960152804851532 | Aug 16, 2022 10:54 UTC | 2 minute(s) | | |

Figure 6. shows the training jobs

Model evaluation & Accuracy check

After identifying the best model performance, we need to ensure that the model performance is good enough for the model to be deployed or update the currently in-production SageMaker endpoint.

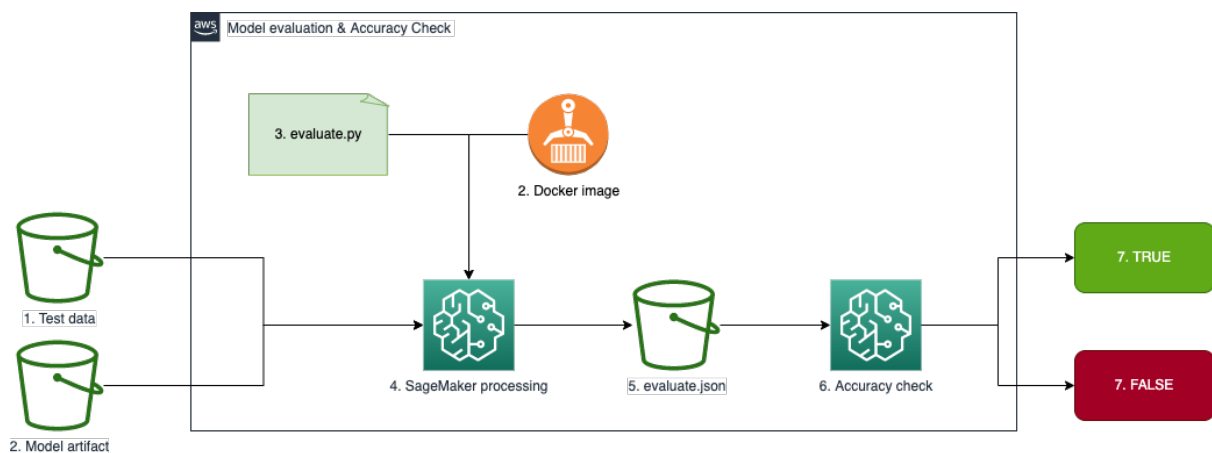


Figure7. model evaluation and check

Figure6, shows the steps that are part of the model evaluation workflow:

1. Load test data
2. Load the trained model
3. evaluate.py: contains the test logic: F1, edge cases, etc.
4. Save the output test result to evaluate.json
5. SageMaker SDK to check the evaluate.json
6. Provide a binary logic to move forward

The evaluation.py contains the code for:

1. Load test data to a DataFrame
2. Load pre-trained model from the training step
3. Make inference on the test data
4. Save the metrics such as Accuracy, F1, recall, etc. to the evaluate.json file

Endpoint Deployment

Figure 7 shows the diagram for the creating a SageMaker endpoint.

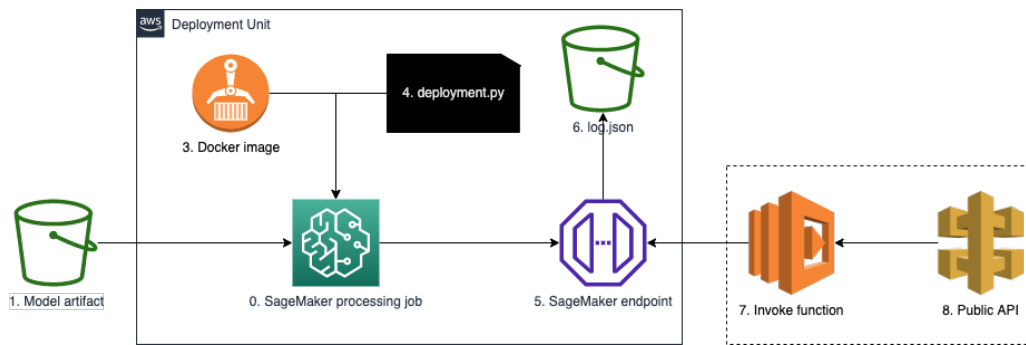


Figure 7. model deployment

1. Copy model artifacts to the container
2. Docker image with python endpoint
3. deployment.py contains all the deployment steps
4. SageMaker Endpoint
5. API logs saved to the S3 bucket for model monitoring
6. Lambda function that invokes the endpoint
7. API gateway to monitor and manage the API calls

Benchmark:

To save cost and training time, I have decided to reduce the number of classes from 100 to 6 in the current work. As a result of that, it is not relevant to compare the F1-Score of the current project with the literature. For instance, reference 6 reports an F1-Score between 87.7 to 88, which is much lower than the 96% score we achieved in this work. I did test some extracting embedding using a smaller word embedding from Spacy, and the result model performance was not as good as the model that was implemented in the pipeline. I recorded an F1-Score of ~92%. However, since we have six classes to identify in this dataset, if we assume that we have a random classifier to pick each of these classes at random with a normal distribution, we get to the accuracy of ~ 16%. The current f1-score after the current workflow is 96%.

References

1. Attention Is All You Need (31st Conference on Neural Information Processing Systems (NIPS 2017))(
<https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>)

2. Universal Language Model Fine-tuning for Text Classification
[<https://arxiv.org/abs/1801.06146>]
3. Universal Sentence Encoder[<https://arxiv.org/abs/1803.11175>]
4. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
[<https://arxiv.org/abs/1810.04805>]
5. LEDGAR: A Large-Scale Multi-label Corpus for Text Classification of Legal Provisions in Contracts [Don Tuggener, Pius von Däniken, Thomas Peetz, Mark Cieliebak]
(<https://aclanthology.org/2020.lrec-1.155/>)
6. <https://eprints.whiterose.ac.uk/178919/1/2110.00976v1.pdf>