

# CSE 6748: Applied Analytics Practicum

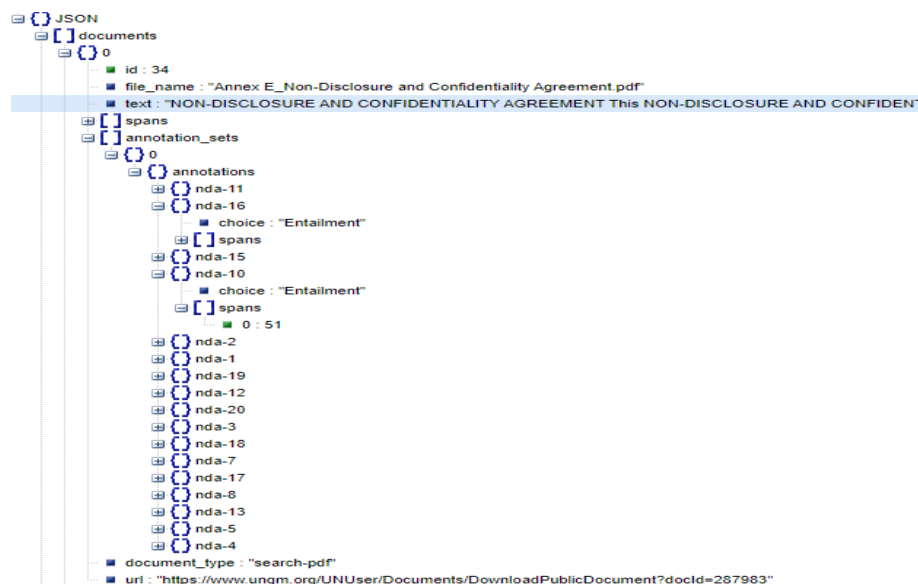
Siddharth Gudiduri  
GT ID: 903464761

## Overview

I was tasked with developing an end-to-end process to automate contract reviews for the CSE 6748 project. Our goal was to create a process that resembled that of an experienced attorney, but with increased speed and accuracy. This type of service will not only maximize document review but will also reduce the risk of human error by mitigating unacceptable or missing clauses. Demonstration of concept The source code for this project can be found at <https://github.com/sgudiduri/CSE-6748>. Because there are many details within the algorithm and check-in omitted, this code is only for your preview and not the company's active repo check-in. Contract Review Automation (CRA) is divided into three mandatory parts and three optional parts. CRA stages include data analysis, model development, productionizing code, and creating model API, Deploy to PaaS, Testing. Bonus tasks involve implementing Pocket base, Caching via Redis and Scaling concepts.

## 1. Data Analysis

Data is received json format. Training data contains 423 documents, and test 123 documents. structure shown below



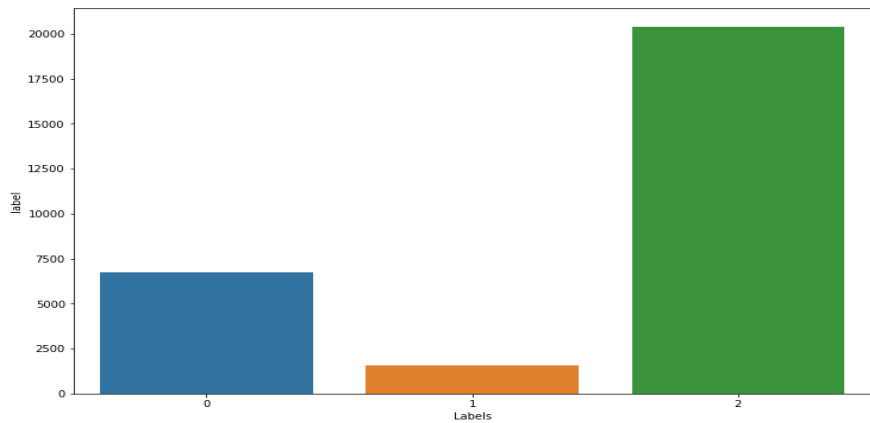
At the document level, we have the "text" key, which contains the entire document text, and the "spans" key, which divides the text into a list of premises. The key term "annotation sets" refers to a list of multiple annotations for a given document. At the annotation level, each key "nda-1," "nda-2," and so on represents a hypothesis that either implies, contradicts, or is neutral to the

given document. The "spans" key is indexed at the document level under each hypothesis. For example, "nda-1" includes the spans 1, 13, and 91. At the document level, span 1 corresponds to sentence text indexed between characters [25, 89]. Each hypothesis' text sequence is described by the "labels" key.

The next step in data analysis is to process the data and extract the features that I will need for building. This feature engineer step was recorded and will be used later in the machine learning pipeline to process incoming data. The extracted Tibble used for Model building is shown below.

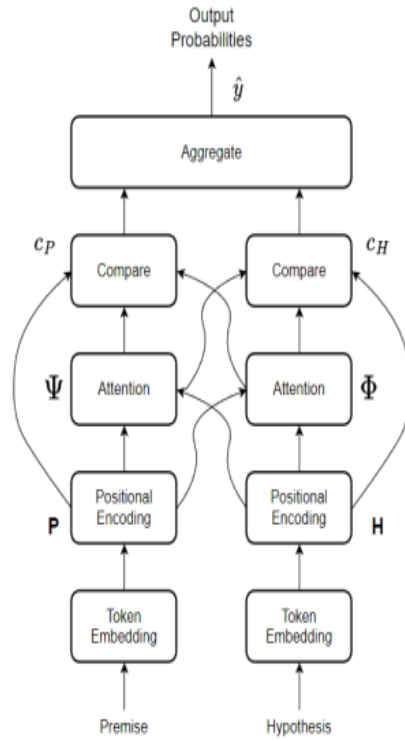
	hypothesis	premise	label
31	Receiving Party shall not reverse engineer any...	i obtained by the Recipient without restrictio...	2
56	Receiving Party shall not reverse engineer any...	For and on behalf of UNHCR For and on behalf o...	2
13	Receiving Party shall not reverse engineer any...	NOW THEREFORE the Parties agree as follows	2
39	Receiving Party shall destroy or return some C...	5 All Confidential Information in any form and...	0
40	Receiving Party shall destroy or return some C...	a if a business relationship is not entered in...	0
...	...	...	...
10	Receiving Party shall not use any Confidential...	a make use of any of the Provider s Confidenti...	0
16	Receiving Party shall not use any Confidential...	Neither the Recipient nor any of the Recipient...	0
89	Receiving Party shall not use any Confidential...	c irrevocably and unconditionally waives the r...	2
6	Receiving Party shall not use any Confidential...	This Agreement sets forth the Parties obligati...	2
48	Receiving Party shall not use any Confidential...	a of this Section 7	2

I discovered that the data was imbalanced, which will play a role in selecting hyperparameters for model building. The screenshot below depicts a higher frequency of neutral cases, followed by entitlement and contradiction. Here's a research code:



## 2. Model Building

I modified one of the existing architectures Decomposable Attention Model with some minor changes like adding additional layers, dropouts, and attention mechanism.



### Token Embedding:

continuing for word representation, I used Global Vectors (GloVe) embedding, an unsupervised learning algorithm for obtaining vector representations for words, with training performed on aggregated global word-word co-occurrence statistics from a corpus. GloVe 6B 100d is used to perform word embedding.

### Positional Encoding:

Below formula uniquely encodes information about the position of a token.

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

Here, d is the embedding dimension, pos is the position of the token in the sequence, and i maps to sin and cosine functions.

### Attending:

We perform soft alignment of the premise and hypothesis essentially achieved by passing the input premise and hypothesis through a multi-layer perceptron and then computing soft attention weights

$$w_{ij} = F(p_i)^T F(h_j)$$

where F is the multi-layer perceptron with ReLU nonlinear activation that maps  $p_i, h_j$  to a hidden dimension space. This allows us to calculate the projection of the premise over the hypothesis. The intuition behind the alignment model is based on a bidirectional RNN used as an encoder and decoder

$$\Phi_i = \sum_{j=1}^b \frac{\exp(w_{ij})}{\sum_{k=1}^b \exp(w_{ik})} h_j$$

$$\Psi_j = \sum_{i=1}^a \frac{\exp(w_{ij})}{\sum_{k=1}^a \exp(w_{kj})} p_i$$

### Comparing:

In the compare section, all the tokens from one sequence, with their corresponding weights are compared with a token in the other sequence.

$$c_{P,i} = G([p_i, \Phi_i]), i = 1 \dots a$$

$$c_{H,j} = G([h_j, \Psi_j]), j = 1 \dots b$$

The representation is the concatenation of premise token  $p_i$  and the softly aligned weight representation for that token  $\Phi_i$ . A similar operation is performed for the hypothesis as well. As the concatenation operation is performed along the embedding dimension, the multi-layer perceptron G maps input dimension equal to twice the embedding dimension, to the number of hidden units.

### Aggregating:

The final step performed by the decomposable attention model is aggregating the information obtained from the comparison step. The information in the comparison vectors is aggregated through a summation operation. The summed-up results are now fed into a multi-layer perceptron H and are mapped to the number of outputs - Entailment, Contradiction and Neutral. Below are learnable parameters:

$$c_P = \sum_{i=1}^a c_{P,i}$$

$$c_H = \sum_{j=1}^b c_{H,j}$$

$$\hat{y} = H([c_P, c_H])$$

## Focal Loss:

While training the decomposable attention model, we use focal loss as there exists class imbalance among the 3 classes - Entailment, Contradiction and Neutral.

$$CB_{focal}(z, y) = -\frac{1-\beta}{1-\beta^{n_y}} \sum_{i=1}^C (1-p_i^t)^\gamma \log(p_i^t)$$

The  $\beta$  hyper-parameter can be tuned to perform reweighting. When  $pt$  is small and consequently,  $(1-pt)^\gamma$  is close to 1, then Focal loss becomes classic cross entropy, and would result in incorrect classification by the model. As the model adjusts its weights, Focal Loss scales down the contribution of easy examples during training and instead focuses on the harder examples, resulting in an improvement in prediction accuracy for the minor classes.

## Results:

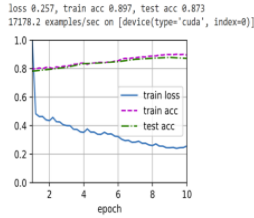


Figure 3. Loss (CE) and Accuracy Curves (Standard Dataset)

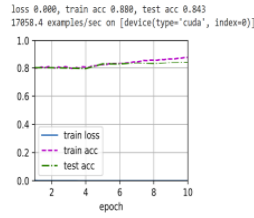


Figure 4. Loss (Focal) and Accuracy Curves (Standard Dataset)

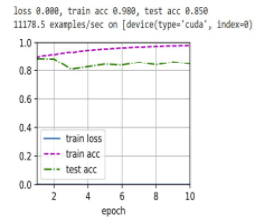


Figure 5. Loss (Focal) and Accuracy Curves (Standard Dataset with Faiss)

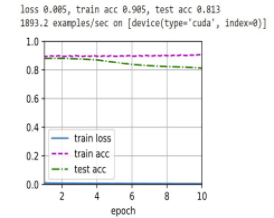


Figure 6. Loss (Focal) and Accuracy Curves (Standard Dataset with Faiss)

Optimizer [SGD, ADAM]	Batch Size [64, 128, 256]	Learning Rate [1e-2, 1e-3, 1e-4]	Regularization [1e-3, 1e-4, 1e-5]	Entailment F1 Score	Contradiction F1 Score	Neutral F1 Score	Accuracy
ADAM	256	1e-2	1e-3	0.08	0.28	0.54	0.88

Table 1. Optimal Hyper-parameters and Validation Results for Standard Contract NLI Dataset using Cross Entropy Loss

Optimizer [SGD, ADAM]	Batch Size N	Learning Rate $\alpha$	Regularization $\lambda$	Class Re-balance Factor $\beta$	Modulating Factor $\gamma$	Entailment F1 Score	Contradiction F1 Score	Neutral F1 Score	Accuracy
ADAM	256	1e-2	1e-3	0.9	2	0.46	0.22	0.86	0.84

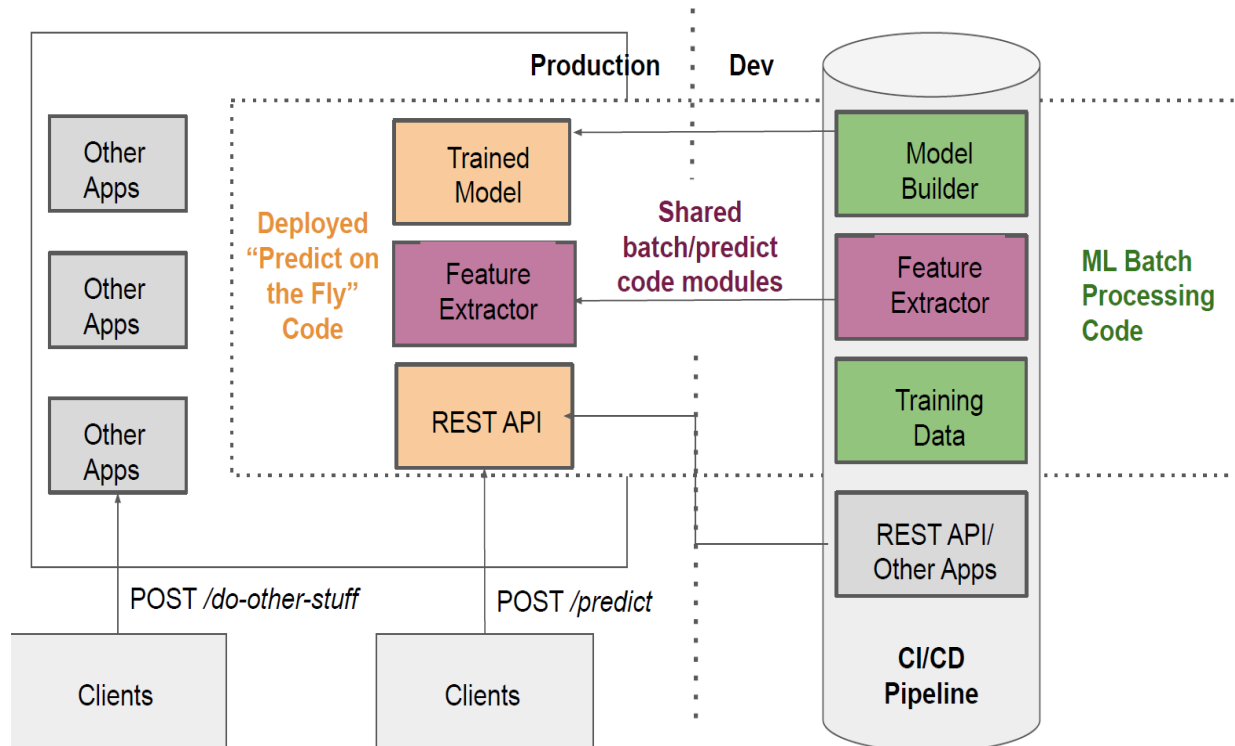
Table 2. Optimal Hyper-parameters and Validation Results for Standard Contract NLI Dataset using Focal Loss

using the above hyperparameter we got the best precision, recall, f1-score along with accuracy and loss. After selecting model and hyperparameter our next task was to productionizing code.

## 3. Architecture Component Breakdown

Next our goal is building the below architecture by creating model package, web api and CI/CD pipelines for package and api. Below is architecture breakdown

**Diagram: Train by batch, predict on the fly, serve via REST API**



## 4. Productionizing Model package

Continuing with the project, the next step was to write production code that would be deployed to end users. I kept my attention focused on testability, maintainability, scalability, performance, and reproducibility. The package structure created below breaks down research code into separation of concern components, implying that each module has a single responsibility in performing its function. As a result, the config package will only contain configuration modules. Only modules designed for testing will be included in the testing package. NOTE: The company code contains a lot more unit tests and an ensemble of models with different hyperparameters.

```

└─ Production/
    └─ contract_nli/
        └─ config/
            ├── __init__.py
            └─ core.py
        └─ data/
            ├── test.csv
            └─ train.csv
        └─ model/
            ├── __init__.py
            ├── aggregate.py
            ├── attend.py
            ├── compare.py
            ├── decomposable_attention.py
            └─ mlp.py
        └─ preprocessing/
            ├── __init__.py
            ├── data_management.py
            └─ snli_dataset.py
        └─ tests/
            ├── __init__.py
            ├── conftest.py
            ├── test_config.py
            ├── test_predict.py
            └─ test_validation.py
        └─ trained_model/
            ├── __init__.py
            ├── model.pth
            └─ model.pth
    └─ requirements/
        └─ requirements.txt
    └─ MANIFEST.in
    └─ mypy.ini
    └─ pyproject.toml
    └─ requirements.txt
    └─ setup.py

```

As you can see from the folder structure above, I created various model modules to build a deep learning package. The dependency injection pattern, which involves passing objects that objects require rather than creating them, aided in the creation of scalable and testable code. Below is an implementation of pydantic, a Python library that simplifies configuration code and compiles it into an object that can be passed into various parts of the application.

```

1  from pathlib import Path
2  import typing as t
3  from pydantic import BaseModel, validator
4  from strictyaml import load, YAML
5  import os
6
7
8  # Project Directories
9  PWD = os.path.dirname(os.path.abspath(__file__))
10 PACKAGE_ROOT = Path(os.path.abspath(os.path.join(PWD, '..')))
11 ROOT = PACKAGE_ROOT.parent
12 CONFIG_FILE_PATH = PACKAGE_ROOT / "config.yaml"
13 TRAINED_MODEL_DIR = PACKAGE_ROOT / "trained_model"
14 DATA_DIR = PACKAGE_ROOT / "data"
15
16 class AppConfig(BaseModel):
17     """
18     Application-level config.
19     """
20     package_name: str
21     train_path: str
22     test_path: str
23     vocab_path: str
24     model_path: str
25
26
27 class ModelConfig(BaseModel):
28     """
29     All configuration relevant to model
30     training and feature engineering.
31     """
32     num_step: int
33     batch_size: int
34     learning_rate: float
35     embed_size: int
36     num_hiddens: int
37     epochs: int
38     save_best: bool
39     trainer: str
40     loss: str
41

```

Below are some screenshots for Data Service class

29 lines (26 sloc)
1.17 KB

Raw
Blame

```

1  import pandas as pd
2  import torch
3  import d2l
4  from .snli_dataset import SNLIDataset
5
6  class DataService():
7      def __init__(self):
8          pass
9
10     #loads csv into pandas dataframe.
11     def load_data(self, train_path, test_path):
12         train = pd.read_csv(train_path)
13         test = pd.read_csv(test_path)
14         train["label"] = train["label"].astype(int)
15         test["label"] = test["label"].astype(int)
16         return train, test
17
18     #load pandas dataframe to dataset.
19     def create_snli_dataset(self, train, test, num_steps = 50, batch_size = 256, num_workers = 4):
20         train_set = SNLIDataset(train, num_steps)
21         test_set = SNLIDataset(test, num_steps, train_set.vocab)
22         vocab = train_set.vocab
23         train_iter = torch.utils.data.DataLoader(train_set, batch_size,
24                                                 shuffle=False,
25                                                 num_workers=num_workers)
26         test_iter = torch.utils.data.DataLoader(test_set, batch_size,
27                                                 shuffle=False,
28                                                 num_workers=num_workers)
29         return train_iter, test_iter, vocab

```



Below are some screenshots for testing code modules

 main ▾ CSE-6748 / production / contract\_nli / tests / test\_predict.py / <> Jump to ▾

 sguididuri Create/UpdateUnit tests ...

Latest commit afc0a00 2 weeks ago  History

 1 contributor

19 lines (16 sloc) | 1.02 KB

Raw Blame    

```
1  '''
2  Here is where we want to test inputs, outputs and model quality
3  below are sample tests
4  '''
5
6  #Testing for entailment condition
7  def test_prediction_quality_against_benchmark_entailment(load_predict_class, sample_input_data_1):
8      res = load_predict_class.make_single_prediction(sample_input_data_1["premise"].split(), sample_input_data_1["hypothesis"].split())
9      assert res == sample_input_data_1["result"]
10
11 #Testing for Contradiction condition
12 def test_prediction_quality_against_benchmark_entailment(load_predict_class, sample_input_data_2):
13     res = load_predict_class.make_single_prediction(sample_input_data_2["premise"].split(), sample_input_data_2["hypothesis"].split())
14     assert res == sample_input_data_2["result"]
15
16 #Testing for neutral condition
17 def test_prediction_quality_against_benchmark_entailment(load_predict_class, sample_input_data_3):
18     res = load_predict_class.make_single_prediction(sample_input_data_3["premise"].split(), sample_input_data_3["hypothesis"].split())
19     assert res == sample_input_data_3["result"]
```

```
63 lines (53 sloc) | 1.96 KB
1  #!/usr/bin/python
2
3  '''
4  https://stackoverflow.com/questions/34466027/in-pytest-what-is-the-use-of-conf-test-py-files
5  conf-test.py is used to define
6  fixture used to define static data used by tests
7  External plugin
8  Hooks
9  Test root path
10 '''
11
12 import pytest
13 from contract_nli.config.core import config, TRAINED_MODEL_DIR
14 from contract_nli.predict import Predict
15
16 trained_model_dir_path = TRAINED_MODEL_DIR.as_posix()
17 model_config = config.model_config
18 embed_size=model_config.embed_size
19 num_hidden=model_config.num_hidden
20
21 @pytest.fixture()
22 def raw_app_config():
23     #For larger datasets, here we would use a testing sub-sample.
24     return config.app_config
25
26 @pytest.fixture()
27 def raw_model_config():
28     return model_config
29
30 @pytest.fixture()
31 def load_predict_class():
32     model_path = f'{trained_model_dir_path}/{config.app_config.model_path}'
33     vocab_path = f'{trained_model_dir_path}/{config.app_config.vocab_path}'
34     pr = Predict(model_config.embed_size,model_config.num_hidden, model_path, vocab_path)
35     return pr
36
37 @pytest.fixture()
38 def sample_input_data_1():
39     row_1 = {
40         "hypothesis": "Receiving Party shall destroy or return some Confidential Information upon the termination of Agreement",
41         "premise": "I the completion or termination of the dealings between the parties contemplated hereunder or",
42         "result": "Entailment"
43     }
44     return row_1
45
46 @pytest.fixture()
47 def sample_input_data_2():
48     row_2 = {
49         "hypothesis": "All Confidential Information shall be expressly identified by the Disclosing Party",
50         "premise": "I marked confidential or proprietary or",
51         "result": "Contradiction"
52     }
53     return row_2
54
55 @pytest.fixture()
56 def sample_input_data_3():
57     row_3 = {
58         "hypothesis": "Receiving Party shall not reverse engineer any objects which embody Disclosing Party's Confidential Information",
59         "premise": "6 Compelled Disclosure of Confidential Information",
60         "result": "Neutral"
61     }
62     return row_3
```

[Give feedback](#)

After finishing the packaging module, I began integrating Azure pipelines for CI/CD, which stands for continuous integration, continuous delivery, and continuous deployment. This means that when a developer, such as myself, submits code for review and checks in after approval, the code goes through a process of building, testing, and publishing files to a private server. Instead of creating a monolithic application, this is done so that the machine learning model can be integrated with a website or a web API. Here's an example from my company's pipeline after a feature for this project was checked in.

← Jobs in run #20230221.3.0

Build CPS

✓	Build application	58s
✓	Initialize job	5s
✓	Checkout r...	3s
✓	Use Visual Studio 5.5.1	1s
✓	Restore packages	7s
✓	Build solution	16s
✓	Test Assemblies	14s
✓	Publish symbols path	3s
⌚	Copy Files to: D:\a\1\...	<1s
✓	Publish Artifact	<1s
✓	Post-job: Checkout .NET Core...	<1s
✓	Finalize Job	<1s
✓	Report build status	<1s

✓ Build application

1 Pool: [Azure Pipelines](#)

2 Image: windows-2019

3 Agent: Hosted Agent

4 Started: Today at 5:06 PM

5 Duration: 58s

6

7 ▶ Job preparation parameters

8 ▶ /x 7 queue time variables used

9 📊 100% tests passed

Note, I created a similar example for this class as a POC before integrating with company code. This deployed on pypi is an experimental version and not the model package used at my company. Link can be found [here](#)

## 5. Octopus Private model server

Octopus is a repository of software packages for the various programming language. It is a central location where developers can upload and share their software packages, making it easy for other developers to install and use them in their own projects. Octopus is not specifically a server for deploying machine learning models, but it is commonly used for hosting and distributing Python packages that contain machine learning models, as well as other types of software tools and libraries.

To deploy a machine learning model using Octopus, you would first need to create a Python package that includes the model and any necessary dependencies. You can then upload the package to Octopus using tools like VSTS build or the Octopus web interface. Once the package is uploaded, other developers can install it using package management tools.

Screenshot of library below

Library

Certificates

External Feeds

Git Credentials

Lifecycles

Packages

Build Information

Script Modules

Step Templates

Tenant Tag Sets

Variable Sets

Built-in Package Repository ⓘ

UPLOAD PACKAGE

Octopus Built-in Package Repository

Octopus does not allow NuGet clients to connect to this feed to retrieve packages; the feed only supports deployable packages, and doesn't allow NuGet packages to be consumed from Visual Studio and other tools.

Search...

ID	Highest version	Description
ArchReport	20230321.1	
AdminTools_DBMigrations	20230328.2-master	The AdminTools_DBMigrations deployment package, built on 3/28/2023
AnomalyDetection	20200429.1	
AnomalyDetection.AlertGenerator	20190828.1	
AnomalyDetection.Distribution	20200220.4	A deployment package created from files on disk.
AnomalyDetection.Predictor	20190419.2	
AnomalyDetection.WebApi	20200430.1	
ArchClient_DBMigrations	20230328.2-master	The ArchClient_DBMigrations deployment package, built on 3/28/2023
ARCHWebAdmin	20230328.1-master	The ARCHWebAdmin deployment package, built on 3/28/2023
ARCommonV2	20230328.1-master	The ARCommonV2 deployment package, built on 3/28/2023

Repository Retention

Permission required

The ConfigureServer permission is required to view and change the repository retention settings

Package Indexing

Number of packages: **10668**

Last re-index: **Unknown**

Unknown

The built-in package repository **will not** be re-indexed at startup.

Permission required

The ConfigureServer permission is required to change the repository indexing status

Octopus can automatically re-index the built-in repository at startup to ensure that it is in sync. [Learn more](#)

## 6. API to interact with Model

To create Web API to we imported libraries from the private server. Then we defined the input and the output formats for our model. This is a JSON request with specific keys. Next we defined end points using Fast API decorator. Below are the screenshot of the a) file structure b) main.py c) on the web.

a)

```
.
├── ContractNLI_API/
│   ├── app/
│   │   ├── __init__.py
│   │   ├── main.py
│   │   └── contract_nli/
│   │       └── app.py
│   ├── tests/
│   │   └── contract_nli/
│   │       └── app.py
│   ├── .dockerignore
│   ├── Dockerfile
│   ├── Procfile
│   ├── requirements.txt
│   ├── runtime.txt
│   ├── test_requirements.txt
│   └── tox.ini
```

b) partial implementation, real implementation is redacted.

```

1  from typing import Any
2
3  from fastapi import APIRouter, FastAPI, Request, File, UploadFile
4  from fastapi.middleware.cors import CORSMiddleware
5  from fastapi.responses import RedirectResponse, FileResponse
6  from loguru import logger
7  import numpy as np
8  from contract_nli.predict import Predict
9  from contract_nli.config.core import config, TRAINED_MODEL_DIR
10
11  favicon_path = 'favicon.ico'
12  api_router = APIRouter()
13
14  app = FastAPI(
15      title="Contract NLI V1", openapi_url=f"/api/v1/openapi.json"
16  )
17
18  def allowed_file(filename):
19      return '.' in filename and \
20          filename.rsplit('.', 1)[1].lower() in set(['png', 'jpg', 'jpeg'])
21
22
23  root_router = APIRouter()
24
25
26  @app.get('/favicon.ico', include_in_schema=False)
27  async def favicon():
28      return FileResponse(favicon_path)
29
30  @app.get("/")
31  async def docs_redirect():
32      return RedirectResponse(url='/docs')
33
34  @root_router.get("/health", response_model=dict, status_code=200)
35  async def health() -> dict:
36      """
37      API Endpoint to get health check
38      """
39      health = dict(
40          name="Contract NLI Language Inference API", api_version="1.0.0", model_version="1.0.7"
41      )
42
43      return health
44

```

```

45 @root_router.get("/PredictNLI/{premise}/{hypothesis}", response_model=dict, status_code=200)
46 async def PredictNLI(premise: str, hypothesis: str ) -> dict:
47     """
48     Given a premise and hypothesis predict Entailment, Contradiction or neutral
49     """
50     model_config = config.model_config
51     trained_model_dir_path = TRAINED_MODEL_DIR.as_posix()
52     vocab_path = f"{trained_model_dir_path}/{config.app_config.vocab_path}"
53     model_path = f"{trained_model_dir_path}/{config.app_config.model_path}"
54     p = Predict(model_config.embed_size, model_config.num_hiddens, model_path, vocab_path)
55     res = p.make_single_prediction(premise.split(), hypothesis.split())
56     return dict(Prediction=res)
57
58
59 app.include_router(root_router)
60
61 # #Set all CORS enabled origins
62 if __name__ == "__main__":
63     # Use this for debugging purposes only
64     logger.warning("Running in development mode. Do not run like this in production.")
65     import uvicorn
66

```

c) deployed version of contract NLI submissions and endpoint that returns prediction.



1.6.0-rc9

OAS3

[/openapi.json](#)

Authorize



default



POST

`/submissions` Predict

GET

`/submissions/{sub_id}` Get Result

DELETE

`/submissions/{sub_id}` Delete Result

GET

`/status/{sub_id}` Get Status

GET

`/artifacts/{sub_id}/{artifact_name}` Get Persistent Artifacts

POST

`/config/{document_type}/{config_name}` Update Config Map

GET

`/config` Get Config Map

GET

`/config/{sub_id}` Get Config Map Update Results



GET

/submissions/{sub\_id} Get Result



## Parameters

Try it out

Name	Description
------	-------------

**sub\_id** \* required

string(\$uuid)

(path)

## Responses

Code	Description	Links
------	-------------	-------

200

Successful Response

No links

Media type

Controls Accept header.

Example Value | Schema

```
{
  "sub_id": "string",
  "processingTime": 0,
  "usedOCR": true,
  "confidence": 0,
  "n_words": 0,
  "readingMethod": "string",
  "data": [
    {
      "actionName": "string",
      "actionType": "string",
      "columnIndex": 0,
      "output": [
        {
          "columnIndex": 0,
          "text": "string",
          "type": "string",
          "value": "string"
        }
      ]
    }
  ]
}
```

## **Next steps:**

Now that Machine learning model Web API is in UAT we are going to work on several other tasks. These tasks include monitoring and evaluation of machine learning model. Maintenance, Security and Scaling and then I will be finish with documenting model and its API so others may understand how to work with it.