
DEMO BUILD

The Challenge

Predict Topics Deployed Stop Restart Deploy New Build

Overview Deployments Builds Monitoring Settings

Description p

Sample Code Shell Python R

```
curl -H "Content-Type: application/json" -X POST https://cdsw.cloudera.tellarius.eu/api/altus-ds-1/models/call-model -d '{"accessToken":"mec6eeusr8vk7xqrhlhxb7atoa9b7sf","request":{"phrase":"my credit report is completely inaccurate"}}'
```

Test Model

Input

```
{  "phrase": "my credit report is completely inaccurate"}
```

Test Reset

Result

Status ● success

Response

```
{  "topics": "[('1: Loans', 0.040193766), ('2: Debt', 0.040334601), ('3: Credit Report', 0.83919954), ('4: General', 0.040070638), ('5: Payments', 0.0402015)]"
```

Model Details

Model Id	6
Deployment	352
Build	67
Deployed By	jrothwell
Comment	
Kernel	python3
Engine Image	Base Image v5
File	predict_topics.py
Function	predict

Model Resources

Replicas	1
Total CPU	0.5 vCPUs
Total Memory	1.00 GiB

Project Flow

- Acquire Data
- Create new Python project and Upload data file
- Step 01 - populate source dataframe and series
- Step 02 - Pre-process the data: lowercase, tokenise, stop-words, stemming
- Step 03 - Build the LDA model
- Step 04 - Visualise the model
- Step 05 - Write the prediction function
- Step 06 - Deploy the model

Logistics

- Have you brought your laptop and want to do the hands on lab?

- Collect a UserID and Password from one of the team:

Trainingxx/ bwh3PcZX6ij6puz

- Data Science Workbench URL: <https://cdsw.cloudera.tellarius.eu>

- Github repo: <https://github.com/bograt71/Data-Innovation-Summit-2019>

<https://bit.ly/2UBfETt>

Acquire Data

- <https://catalog.data.gov/dataset/consumer-complaint-database>
- >600MB
- We will run with a subset for today : [consumer_complaints_16952.csv](#)
- Open the csv file and take a quick look
 - It is already labelled – we will be ignoring that for the demo
 - We are just going to use the Customer complaint narrative field
 - Notice there is a fair amount of superfluous text in there

Create a new Blank Project and open a Workbench

- Open a Python 3 Workbench
- Upload the source data file
- Create a Readme.md
- Create a setup.py to install the required packages
- setup.py is here: [setup.py](#)
- Create a “models” folder where we will store the trained model

Step 01 – populate source dataframe and series

- Create a file called **train_model.py**
 - this will hold all our code for data prep, model training and visualisation
 - We will be incrementally adding code fragments to this file
- Import packages
- Read complaints csv into a dataframe and extract the complaints narrative into a series
- Code is here: [train_model_01.py](#)
- Take a look at the dataframe and the series

```
df.head(5), doc_set, doc_set[0]
```


Step 02 - Pre-process data

- Build a list of tokenised (and “cleansed”) documents
 - Loop through each row in the `doc_set`
 - Lower case it
 - Tokenise it
 - Remove stop words (I, on, my, and ...)
 - Stem to the word root (change,changing,changed -> chang)
- Code is here: [train_model_02.py](#)
- Now compare your source complaint to the pre-processed one
`doc_set[0], texts[0]` (stopped words and stemmed)

Step 03 - Build the LDA model

- turn our tokenized documents into a id <-> term dictionary
- convert tokenized documents into a document-term matrix
- generate LDA model
- save the trained model
- Take a look at the term dictionary and document-term matrix
 - `dictionary[0], len(corpus), corpus[0]`
- Code is here: [train_model_03.py](#)

Step 04 - Visualise the model

- Look at the term composition of the main topics in the model
- Use pyLDAvis to visualise the model to aid labelling of the topics
- Code is here: [train_model_04.py](#)

Step 05 - Predict topics

- Create the `predict_topics.py` to receive a phrase and return the topic list
- Load saved model
- Pre-process the phrase
 - Lowercase, tokenise, stop_words, stemming
- Assemble a bag of words (bow) from the pre-processed phrase
- Score the bow with the LDA model
- Return list of topics
- Try the sample phrases

Step 06 - Deploy the predict_topics.py

- First create the `cdsw-build.sh` script
 - Include the contents of `setup.py` (i.e. import all necessary packages)
- Open a terminal to your K8S container and change permissions:

```
chmod u+x /home/cdsw/cdsw-build.sh
```
- Code fragment here: [cdsw-build.sh](#)

Step 06 - Deploy the predict_topics.py

- Name: Predict topics
- Description: Topic modelling with LDA
- File: predict_topics.py
- Function: predict
- Example Input: { "phrase": "my credit report is completely inaccurate"}
- Kernel: Python 3
- Engine Profile: 0.5 vCPU/1 GiB memory

WELL DONE