# Building a Custom Language Classification Model on Cloud AutoML NLP

This section will walk through creating a document dataset and building a custom language classification model on AutoML Vision.

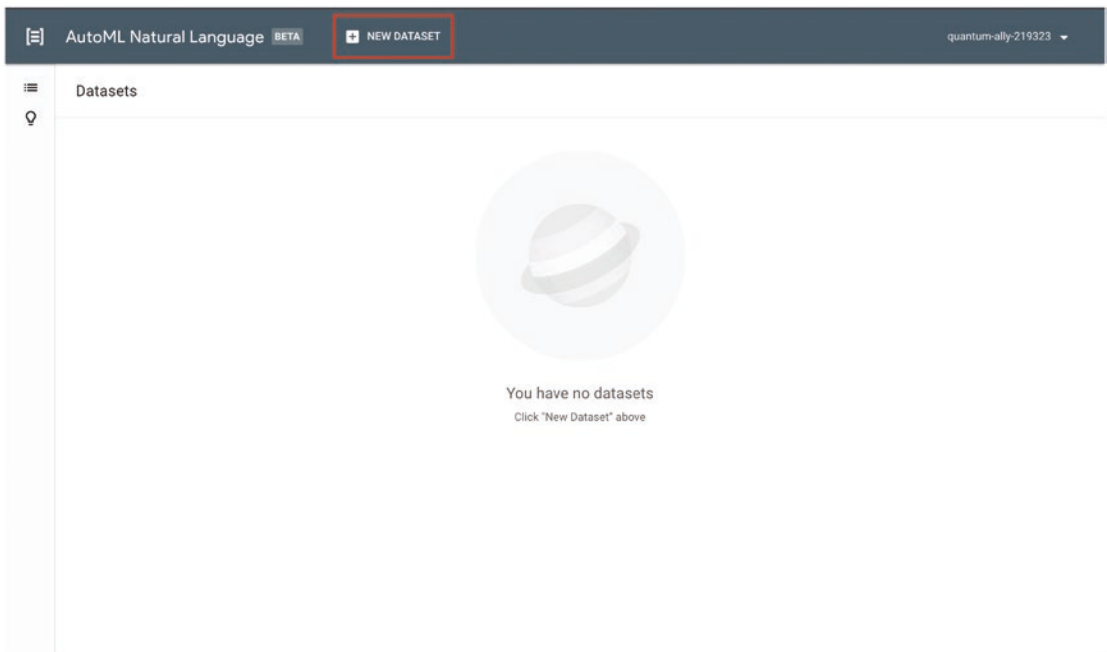1. From the Cloud AutoML NLP dashboard, click **NEW DATASET** as shown in Figure 43-5.



***Figure 43-5.*** *New Dataset on AutoML NLP*

2. To create a Dataset on Cloud AutoML NLP, set the following parameters as shown in Figure 43-6:

*Figure 43-6.*  *Create a Dataset on Cloud AutoML NLP*

a. Dataset name: toxicity_dataset.

b. Select a CSV file on Cloud Storage (this is the CSV file placed on the bucket created when Cloud AutoML was configured that contains the path to the text documents): gs://quantum-ally-219323-lcm/file/data.csv.

c. Click **CREATE DATASET** to begin importing images (see Figure 43-7).

*Figure 43-7.  Cloud AutoML NLP: Importing text items*



*Figure 43-8.  Cloud AutoML NLP: Imported text documents and their labels*

3.  After importing the Dataset, click **TRAIN** (see Figure 43-8) to initiate the process of building a custom language classification model.

4.  In this example, we have a good enough number of training examples as seen in Figure 43-9, so hopefully, it makes sense to expect a good language classification model. Click **START TRAINING** to begin the training job.

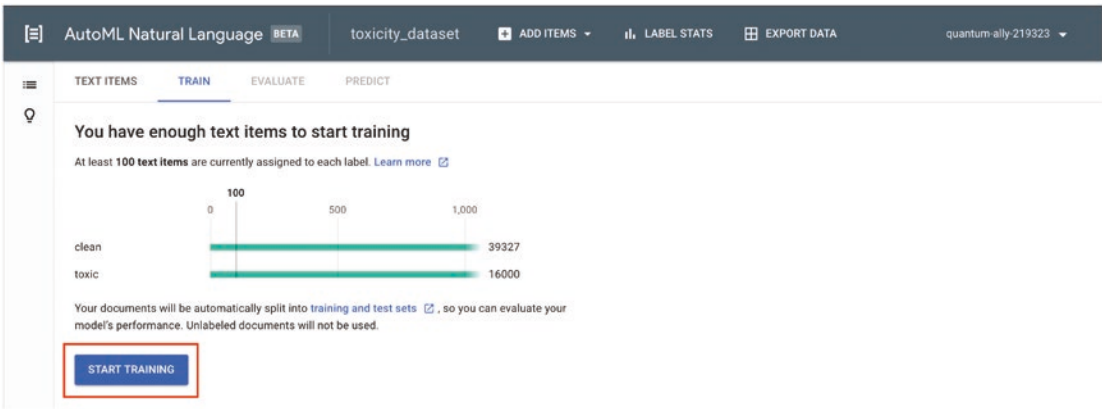*Figure 43-9.*  *Cloud AutoML NLP checking the adequacy of training examples*

5. Accept the default model name, and click **START TRAINING** (see
   Figure 43-10) to begin building the model as seen in Figure 43-11.
   Note that this training might take about an hour to complete. When
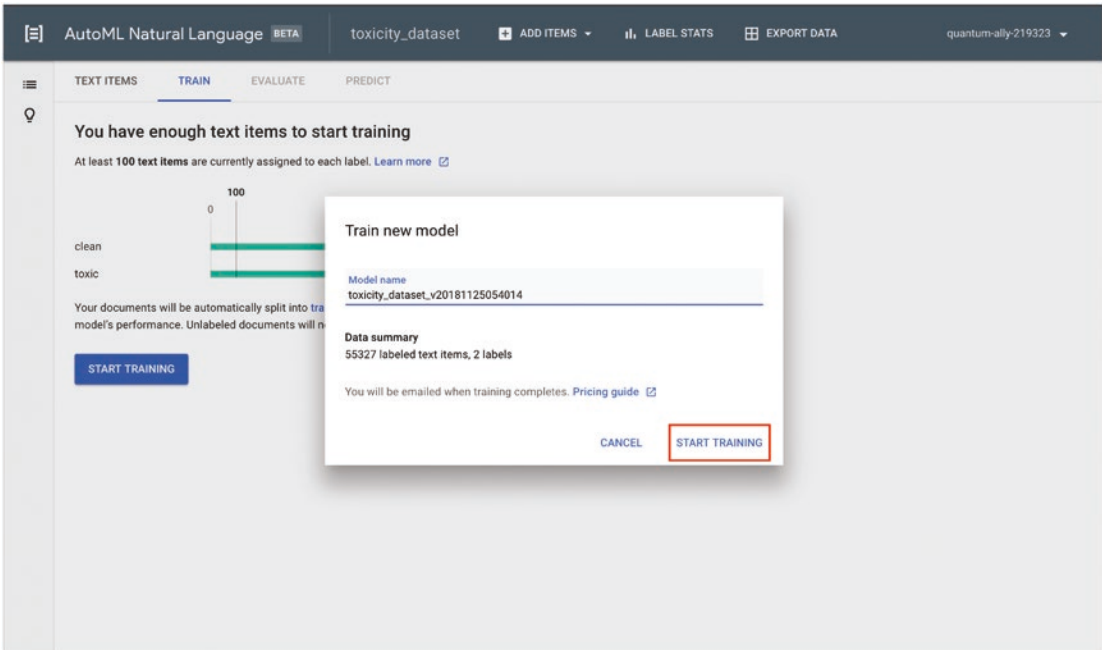   done, the user will get an email of completion.



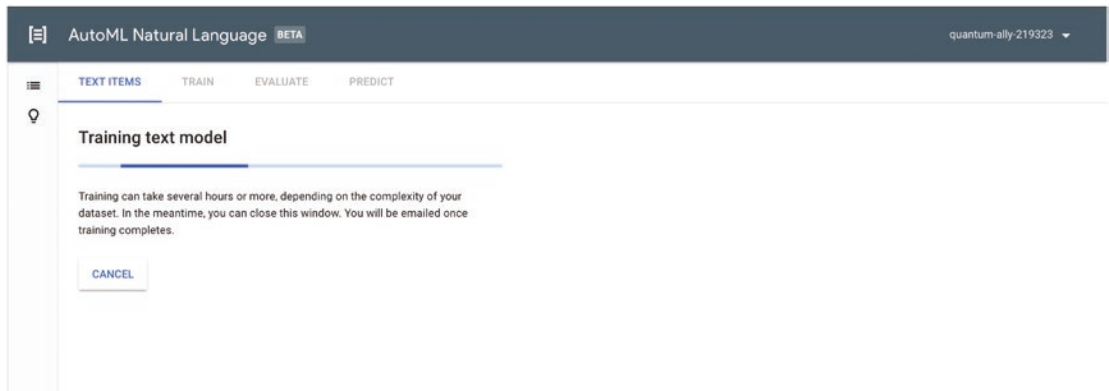*Figure 43-10.*  *Accept the Model name and click on "Start Training"*

***Figure 43-11.*** *Training the text classification model on Cloud AutoML NLP*

6. The training summary is shown in Figure 43-12. The training phase lasted for approximately 4 hours and 45 minutes.
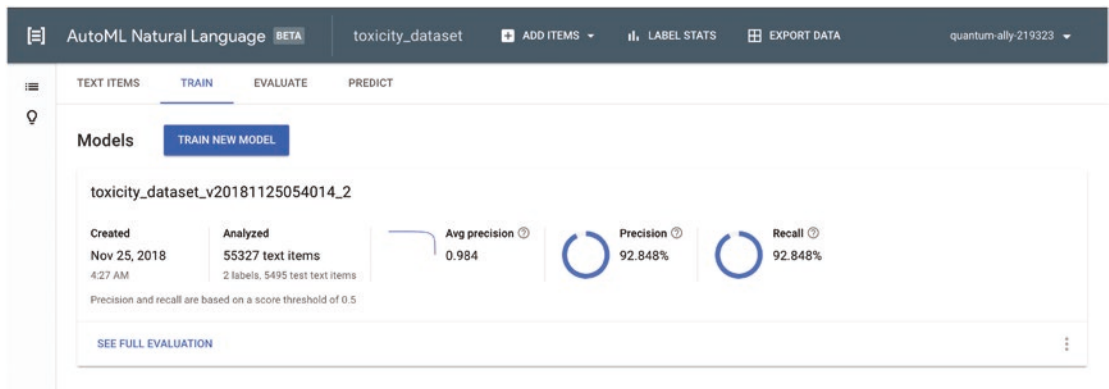


***Figure 43-12.*** *Cloud AutoML NLP: Training summary*

7. AutoML NLP sets aside a portion of the documents as a test set in order to evaluate the quality of the model after training (see Figure 43-13). The F1 plot shows the trade-off between precision and recall. Also, a confusion matrix provides further insight into the model quality (see Figure 43-14).
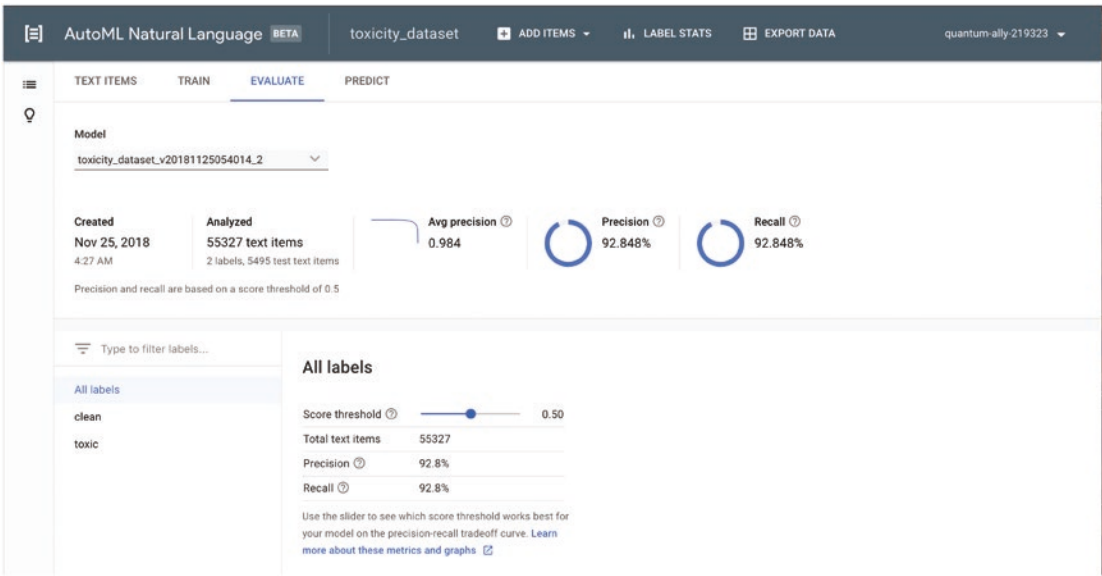
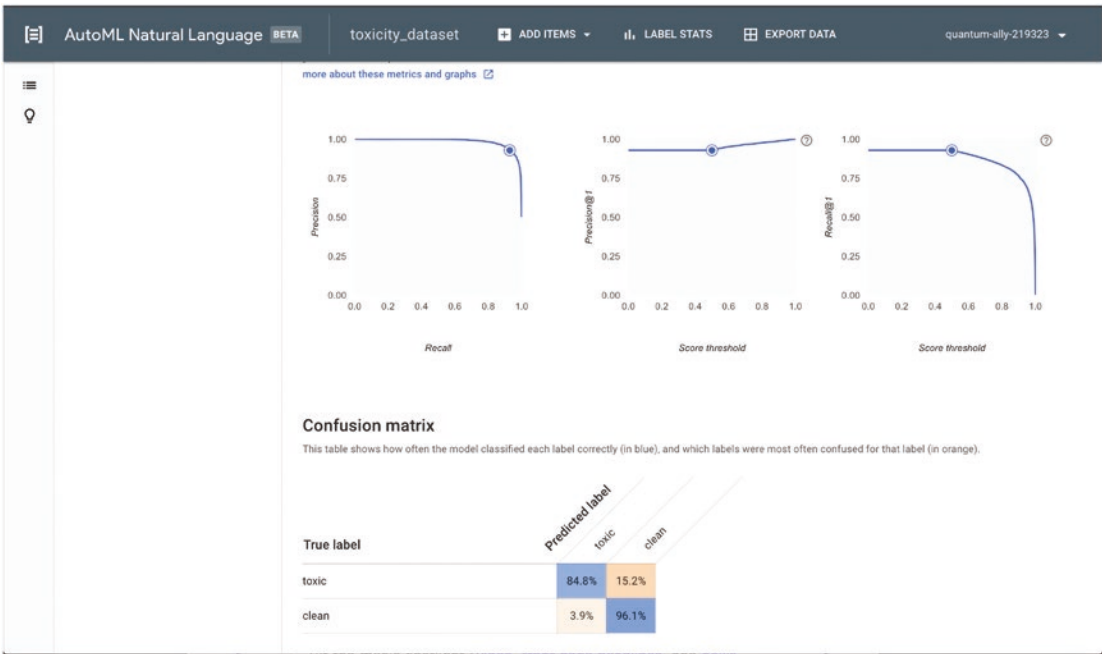***Figure 43-13.*** *Cloud AutoML NLP: Model evaluation*



***Figure 43-14.*** *F1 evaluation plot and confusion matrix on Cloud AutoML NLP*

8.  The custom text classification model is exposed as a REST or Python API for integration into software applications as a prediction service (see Figure 43-15). We can test our model by uploading a sample image for classification. Figure 43-16 passes a clean text example to the model and it predicts correctly with a probability of 98%, while Figure 43-17 passes a toxic text example to the model. This example is also correctly classified with a probability score of 99%.
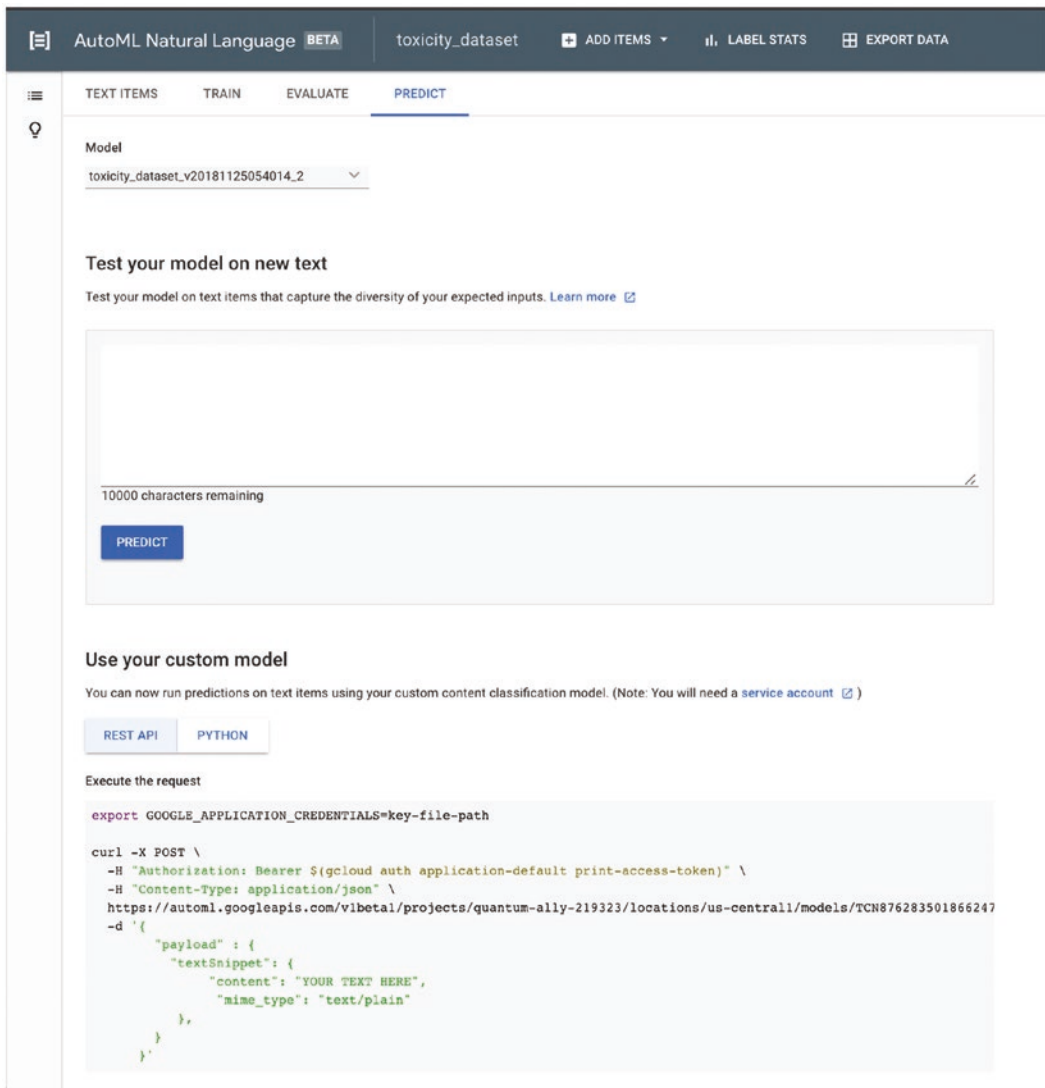


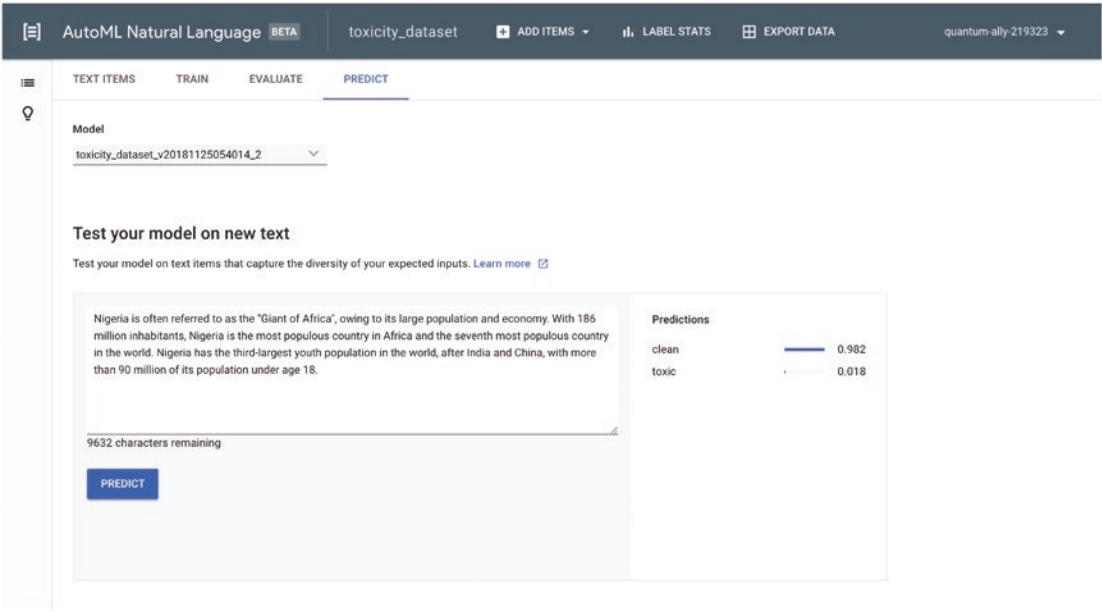*Figure 43-15.   Cloud AutoML NLP model as a prediction service*

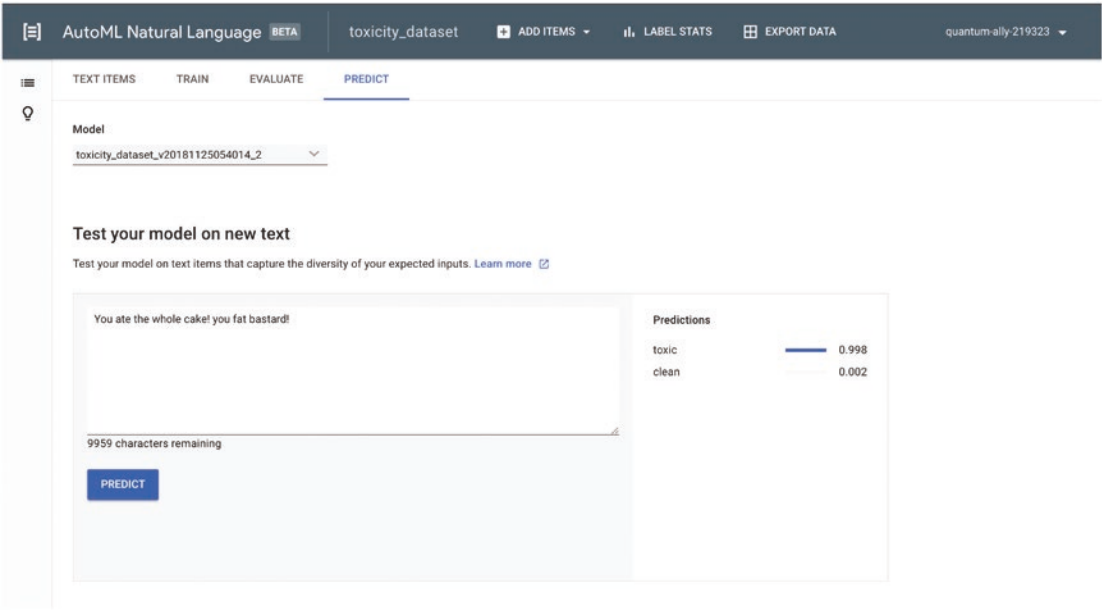*Figure 43-16.*   *Clean words example: AutoML NLP*



*Figure 43-17.*   *Toxic words example: AutoML NLP*

This chapter covered building and deploying custom text classification models using Google AutoML Cloud Vision. In the next chapter, we will build an end-to-end data science product on GCP.

# CHAPTER 44

# Model to Predict the Critical Temperature of Superconductors

This chapter builds a regression machine learning model to predict the critical temperature of superconductors. The features for this dataset were derived based on the following superconductor properties:

- Atomic mass

- First ionization energy

- Atomic radius

- Density

- Electron affinity

- Fusion heat

- Thermal conductivity

- Valence

And for each property, the mean, weighted mean, geometric mean, weighted geometric mean, entropy, weighted entropy, range, weighted range, standard deviation, and weighted standard deviation are extracted. Thus, this results in a total number of 8 x 10 = 80 features. In addition to this, a feature that contains the number of elements in the superconductor is added to the design matrix. The predictor variable is the critical temperature of the superconductor. Hence, the dataset has a total of 81 features and 21,263 rows.