| | |
|---|---|
| Data Exploration: Analyzing Product Sentiment | Read some product review data (product)<br>- product.head() |
| Creating Word Count Vector | Add a new word_count column that counts the number of words in the review column<br>- Output: count how many times each words appear in the review |
| Find Most Popular Product | 1. Create a histogram via product.show() to see the most frequent item<br>2. Create a dataset to contain the reviews of only the most popular product<br>3. Order the created dataset in order of ratings and use .show() to see which one has the highest frequency<br>4. Can determine if this popular product has good or bad reviews |
| Sentiment classification (positive or negative) | 1. Assign the ranks<br>    a. Unclear Ratings = 3 we ignore for now<br>    b. Positive Ratings = 4 or 5, value is 1 for true<br>    c. Negative Ratings= 2 or 1 value is 0 for false<br>2. Put the ratings of 1 or 0 in a new column called sentiment |
| Train the sentiment classifier | 1. Assign train data and test data using *random_split* with 80% for training and 20% for testing<br>2. Use *logistic_classifier.create()(data _set, target, features, validation_set)*<br>    a. Data set = train data<br>    b. Target = name of column containing target variable (sentiment)<br>    c. Features = name of column containing features (word count)<br>    d. Validation_set = test_data<br>3. Output should show the number of iterations and the accuracy increasing with each iteration |
| Evaluate the Classifier with ROC curve. | 1. ROC_Curve is a way to explore false positive/negatives<br>2. Use .evaluate(dataset, metric)<br>    a. Dataset = test_data<br>    b. metric=roc_curve<br>3. Use .show(view='Evaluation") to visualize the .evaluate outputs<br>    a. Output should show a curved line that evaluates false/true positives/negatives<br><br>ROC Curves<br>- Shows the probability of a true positive or false positive based on where the data hits on the curve. |

| | |
|---|---|
| | - Ex: value where FP=0.05 shows 0.2 chance of being a true positive<br><br>True Positive Rate<br><br>(0.95 , 0.98 )<br><br>False Positive Rate<br><br>True Positive 26595   False Negative 1459   Accuracy 0.914   Precision 0.949   Threshold 0.46<br>False Positive 1426   True Negative 3896   Recall 0.948   F1 Score 0.949 |
| Applying the Model | Applying the learned model to understand sentiment for the most popular product (Giraffe)<br>1. Extract the predicted sentiment for giraffe reviews from the sentiment_model into a new dataset called giraffe_reviews<br>2. Sort the reviews based on the predicted sentiment and explore using giraffe_reviews.sort('sorted item',ascending=false) |
| Exploring the most positive/negative aspect of a product | Dataset[0][column]  => first time on the table<br>Dataset[1]][column]   => first time on the table<br>.<br>.<br>.<br>Dataset[-1]][column]  => last item on the table<br>Dataset[-2]][column]  => second to last item on the table |