

Data Cleaning

```
In [5]: def data_cleaning(df):  
        print("##### Cleaning Data #####")  
  
        # Removing empty cells  
        df.dropna(inplace=True)  
  
        # Pre-processing Text Reviews  
  
        # Lowercase Words  
        df['review_body'] = df['review_body'].apply(  
            lambda x: x.lower()  
        )  
  
        print("\n##### Lowercase Complete #####")
```

```

lambda x: x.replace("<chr /><br /></>", "")

print("\n##### Remove Stop Words Complete #####")

# Remove Punctuations
tokenizer = RegexTokenizer(r'\w+')
df['review_body'] = df['review_body'].apply(
    lambda x: ' '.join([word for word in tokenizer.tokenize(x)])

print("\n##### Remove Punctuation Complete #####")

# Lemmatization using .lemmas
nlp.spacy.load('en_core_web_sm', disable=['parser', 'ner'])
df['review_body'] = df['review_body'].apply(
    lambda x: ' '.join([token.lemma_ for token in nlp(x)])

print("\n##### Data Cleaning Complete #####")

return df

```

```

In [6]: apparel_cleaned = data_cleaning(apparel_small)
# apparel_cleaned

##### Cleaning Data #####

##### Lowercase Complete #####

##### Remove Stop Words Complete #####

##### Remove Punctuation Complete #####

##### Data Cleaning Complete #####

```

Add Bigrams

```

In [7]: # https://stackoverflow.com/questions/48331315/how-to-extract-all-the-ngrams-from-a-text-dataframe-could-
mn-in-different-order-in

from collections import Counter
from nltk import ngrams
from itertools import chain

def find_ngrams(input_list, n):
    return list(zip(*[input_list[i:] for i in range(n)]))

```

```

from nltk import ngrams
from itertools import chain

def find_ngrams(input_list, n):
    return list(zip(*[input_list[i:] for i in range(n)]))

apparel_cleaned['bigrams'] = apparel_cleaned['review_body'].map(lambda x: find_ngrams(x.split(" "), 2))
# apparel_cleaned

In [8]: apparel_un_verified = apparel_cleaned[apparel_cleaned['verified_purchase'] == 'N']
apparel_verified = apparel_cleaned[apparel_cleaned['verified_purchase'] == 'Y']

# apparel_un_verified

In [9]: apparel_verified_bigrams = apparel_verified['bigrams'].tolist()
verified_bigrams = list(chain(*verified_bigrams))

verified_bigram_counts = Counter(verified_bigrams)
verified_bigram_counts.most_common(20)

Out[9]: (('I', 'm'), 928),
 (('love', 'it'), 419),
 (('fit', 'perfectly'), 238),
 (('fit', 'great'), 224),
 (('look', 'like'), 217),
 ('m', '5'), 215),
 (('fit', 'well'), 209),
 (('good', 'quality'), 182),
 (('fit', 'perfect'), 172),
 (('can', 't'), 168),
 (('well', 'make'), 164),
 (('run', 'small'), 163),
 (('order', 'size'), 156),
 (('like', 'picture'), 148),
 (('look', 'great'), 142),
 (('I', 've'), 138),
 (('love', 'shirt'), 127),
 (('fit', 'like'), 125),
 (('order', 'large'), 123),
 (('love', 'dress'), 120)]

In [10]: un_verified_bigrams = apparel_un_verified['bigrams'].tolist()
un_verified_bigrams = list(chain(*un_verified_bigrams))

un_verified_bigram_counts = Counter(un_verified_bigrams)
un_verified_bigram_counts.most_common(20)

Out[10]: (('exchange', 'honest'), 446),
 (('I', 'm'), 442),
 (('receive', 'product'), 359),
 (('honest', 'review'), 340),
 (('unbiased', 'review'), 287),
 (('I', 'receive'), 282),
 (('honest', 'unbiased'), 264),
 (('product', 'discount'), 262),
 ('br', 'I'), 240),
 (('discount', 'exchange'), 219),
 (('love', 'it'), 180),
 (('well', 'make'), 173),
 (('discount', 'price'), 153),
 (('fit', 'perfectly'), 148),
 (('I', 've'), 135),
 (('fit', 'well'), 130),
 (('t', 'shirt'), 119),
 (('price', 'exchange'), 112),
 (('good', 'quality'), 112),
 (('look', 'great'), 101)]

```

```

authentic_reviews_df = df[df['verified_purchase'] == 'Y']
fake_reviews_df = df[df['verified_purchase'] == 'N']

authentic_reviews_us_df = authentic_reviews_df.sample(sample_size)
under_sampled_df = pd.concat([authentic_reviews_us_df, fake_reviews_df], axis=0)

print("Under-Sampled Verified", len(under_sampled_df[(under_sampled_df['verified_purchase'] == 'Y')]))
print("Under-Sampled Un-Verified", len(under_sampled_df[(under_sampled_df['verified_purchase'] == 'N')]))

# Graph of Data Distribution
fig, ax = plt.subplots(figsize=(6, 4))
sns.countplot(x='verified_purchase', data=under_sampled_df)
plt.title("Count of Reviews")
plt.show()
print("Under-Sampling Complete")
return under_sampled_df

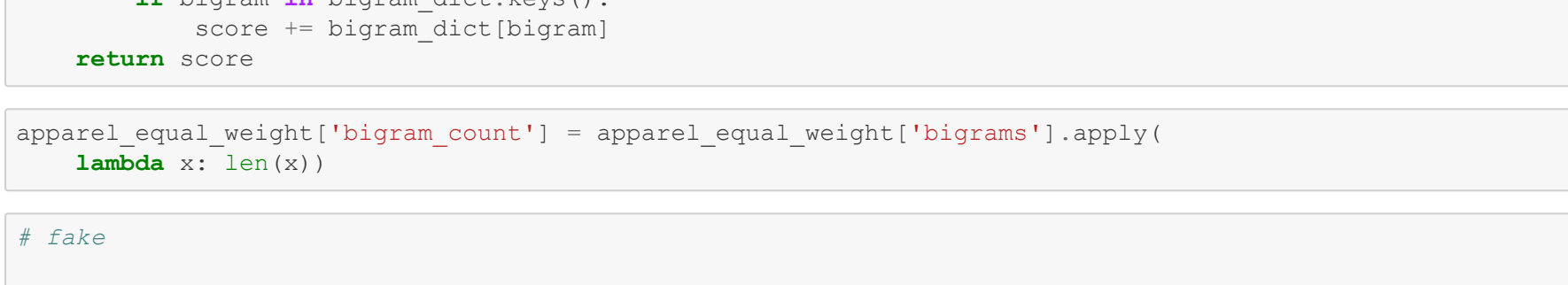
```

In [12]: `apparel_equal_weight = under_sampling(apparel_cleaned)`
`# apparel_equal_weight`

Under-Sampling Data
Verified: 7353
Un-Verified: 2646
Under-Sampled Verified 2646
Under-Sampled Un-Verified 2646

verified_purchase	count
Y	2646
N	2646

Under-Sampling Complete



```
apparel_equal_weight['fake_bigram_count']  
    lambda x: get_bigram_count(x, fake_bigram_count)
```

```
apparel_equal_weight['fake_bigram_percent'] = apparel_equal_weight['fake_bigram_percent'] + apparel_equal_weight['fake_bigram_count'] / apparel_equal_weight['fake_bigram_count']
```

```
lambda x: get_bigram_simple_score(x, fake_bigram_dict))

apparel_equal_weight['fake bigram normalized score'] = apparel_equal_weight['fake bigram si
```

```
/ apparel_equal_weight['bigram_count']
```

```
In [16]: # gold
gold_bigram_dict = dict(verified_bigram_counts) # gold_bigram_dict = dict(verified_bigram_counts.most_c
common(30))

apparel_equal_weight['gold_bigram_count'] = apparel_equal_weight['bigrams'].apply(
    lambda x: get_bigram_count(x, gold_bigram_dict))

apparel_equal_weight['gold_bigram_percent'] = apparel_equal_weight['gold_bigram_count'] / apparel_equal
_weight['bigram_count']

apparel_equal_weight['gold_bigram_simple_score'] = apparel_equal_weight['bigrams'].apply(
    lambda x: get_bigram_simple_score(x, gold_bigram_dict))

apparel_equal_weight['gold_bigram_normalized_score'] = apparel_equal_weight['gold_bigram_simple_score']
/ apparel_equal_weight['bigram_count']

In [22]: apparel_equal_weight = apparel_equal_weight.fillna(0)
apparel_equal_weight.head(1)
```

verified_purchase	review_body	bigrams	bigram_count	fake_bigram_count	fake_bigram_percent	fake_bigram_simple_score	fake
		[[would, ve), (ve, fit), (fit, accord), (accord, size), (size,					

[illegible]

(way, big),
(big
reorder),
(reorder,
time),
(time,
probably),
(probably,
go), (go,
two), (two,
size),
(size,

```

small])

In [23]: apparel_equal_weight[apparel_equal_weight['verified_purchase'] == 'N'].head(1)

Out[23]:
```

verified_purchase	review_body	bigrams	bigram_count	fake_bigram_count	fake_bigram_percent	fake_bigram_simple_score	fake_
		[(second, tank), (tank, order)]					

28	N	second tank order wear often wash well receive many compliment they	(order, wear), (wear, often), (often, wash), (wash, well), (well, receive), (receive, many).	10	10	1.0	74
----	---	---	--	----	----	-----	----

```
dt = dt.copy()

df_unlabeled = df[['fake_bigram_count', 'fake_bigram_percent', 'fake_bigram_simple_score', 'fake_bigram_normalized_score',
                  'gold_bigram_count', 'gold_bigram_percent', 'gold_bigram_simple_score', 'gold_bigram_normalized_score']]

df['verified_purchase'] = df['verified_purchase'].apply(lambda x: 1 if x == 'Y' else 0)
print("Training a algorithm with Model1")
labels = df['verified_purchase']

train_data, test_data, train_label, test_label = train_test_split(df_unlabeled, labels, test_size=0.25, random_state=42)

test_data_copy = test_data.copy()
test_label_copy = test_label.copy()

all_labeled = False

current_iteration = 0

pbar = tqdm(total=iterations)

while not all_labeled and (current_iteration < iterations):
```

```
current_iteration += 1
model.fit(train_data, train_label)

probabilities = model.predict_proba(test_data)
pseudo_labels = model.predict(test_data)

indices = np.argmax(probabilities > threshold)

for item in indices:
```

```

        train_data.loc[test_data.index[item[0]]] = test_data.loc[item[0]]
        train_label.loc[test_data.index[item[0]]] = pseudo_labels[item[0]]
    test_data.drop(test_data.index[indices[:, 0]], inplace=True)
    test_label.drop(test_label.index[indices[:, 0]], inplace=True)

    print("--" * 20)

    if len(test_data) == 0:
        print("Exiting...True")
    all_labeled = True

```

```

        pbar.update(1)

pbar.close()
predicted_labels = model.predict(test_data_copy)

print(algorithm + ' Model Results')
print('--> 20)')
print('Accuracy Score : ' + str(accuracy_score(test_label_copy, predicted_labels)))
print('Precision Score : ' + str(precision_score(test_label_copy, predicted_labels, pos_label=1)))

```

```
print('Recall Score : ' + str(recall_score(test_label_copy, predicted_labels, pos_label=1)))
print('F1 Score : ' + str(f1_score(test_label_copy, predicted_labels, pos_label=1)))
print('Confusion Matrix : \n', str(confusion_matrix(test_label_copy, predicted_labels)))
plot_confusion_matrix(test_label_copy, predicted_labels, classes=[1, 0],
                       title=algorithm + ' Confusion Matrix').show()

def plot_confusion_matrix(y_true, y_pred, classes, title=None, cmap=plt.cm.Blues):
    # Compute confusion matrix
    cm = confusion_matrix(y_true, y_pred)
    plot_confusion_matrix(cm, classes,
                          title=title, cmap=cmap)
```

```
# Only use the labels that appear in the data

fig, ax = plt.subplots()
im = ax.imshow(cm, interpolation='nearest', cmap=cmap)
ax.figure.colorbar(im, ax=ax)
# We want to show all ticks
ax.set(xticks=np.arange(cm.shape[1]),
        yticks=np.arange(cm.shape[0]),
        xticklabels=classes,
```

```

        yticklabels=classes,
        title=title,
        ylabel='True label',
        xlabel='Predicted label')

# Rotate the tick labels and set their alignment.
plt.setp(ax.get_xticklabels(), rotation=45, ha="right",
         rotation_mode="anchor")

# Loop over data dimensions and create text annotations

```

```
fnt = 'd'
thresh = cm.max() / 2.
for i in range(cm.shape[0]):
    for j in range(cm.shape[1]):
        ax.text(i, j, format(cm[i, j], fnt),
                ha="center", va="center",
                color="white" if cm[i, j] > thresh else "black")
fig.tight_layout()
```

```

        return pit

In [20]: start_time = time()
rf = RandomForestClassifier(random_state=42, criterion='entropy', max_depth=14, max_features='auto', n_estimators=500)
semi_supervised_learning(apparel_equal_weight, model=rf, threshold=0.7, iterations=15, algorithm='Random Forest')
end_time = time()

```

```
print("Time taken : ", end_time - start_time)
```

0% | 0/15 [00:00<, 2it/s]

Training Random Forest Model

7% | 1/15 [00:03<00:52, 3.73s/it]

13% | 2/15 [00:05<00:40, 3.08s/it]

15%	█	2/15 [00:00<00:40, 0.03s/it]
20%	█	3/15 [00:06<00:31, 2.63s/it]
27%	█	4/15 [00:08<00:25, 2.32s/it]

33%	██████████	5/15 [00:10<00:20, 2.09s/it]

40%	██████████	6/15 [00:11<00:17, 1.95s/it]

47%	██████████	7/15 [00:13<00:14, 1.84s/it]

53%	██████████	8/15 [00:14<00:12, 1.76s/it]
60%	██████████	9/15 [00:16<00:10, 1.71s/it]
67%	██████████	10/15 [00:17<00:08, 1.67s/it]

```

73% ██████████ | 10/15 [00:17<00:00, 1.18s/it]
-----
73% ██████████ | 11/15 [00:19<00:06, 1.64s/it]
-----
80% ██████████ | 12/15 [00:21<00:04, 1.62s/it]
-----

```

87%	██████████	13/15	[00:22<00:03, 1.60s/it]

93%	██████████	14/15	[00:24<00:01, 1.59s/it]

100%	██████████	15/15	[00:25<00:00, 1.72s/it]

```
-----
Random Forest Model Results
-----
Accuracy Score : 0.9606953892668179
Precision Score : 0.94044432132963989
Recall Score : 0.9869186046511628
F1 Score : 0.9631205673758866
Confusion Matrix :
[[592  43]
 [ 23  53]]
```

Random Forest Confusion Matrix

	Actual 0	Actual 1
Predicted 0	592	43
Predicted 1	10	457

The heatmap shows a strong bias towards predicting class 0. The diagonal elements (592 and 457) represent correct classifications, while the off-diagonal elements (43 and 10) represent misclassifications.

	Predicted label = 1	Predicted label = 0
True = 1	9	679
True = 0	0	0

```
Time taken : 26.026597499847412
```

```
In [21]: start_time = time()
nb = GaussianNB()
semi_supervised_learning(apparel_equal_weight, model=nb, threshold=0.7, iterations=15, algorithm='Naive Bayes')
end_time = time()

print("Time taken : ", end_time - start_time)
```

```
0%|          | 0/15 [00:00<?, ?it/s]

Training Naive Bayes Model

100%|██████████| 15/15 [00:02<00:00, 5.66it/s]

-----
-----
-----
```

```
-----
-----
-----
-----
-----
-----
-----
-----
```

```
-----
Naive Bayes Model Results
-----
Accuracy Score : 0.9138321995464853
Precision Score : 0.867948717948718
Recall Score : 0.9840116279069767
F1 Score : 0.9223433242506811
Confusion Matrix :
```

[[532 103]
[11 677]]

Naive Bayes Confusion Matrix

A heatmap visualization of the Naive Bayes Confusion Matrix. The x-axis is labeled '1' and the y-axis is labeled '0'. The matrix shows 532 true positives (dark blue) and 103 false positives (light blue) for class 1. The color scale on the right ranges from 500 to 600.

	0	1
0	677	11
1	103	532

	0	1
0	11	0
1	0	677

Predicted label:

Time taken : 2.770996570587158

In []: