



```
In [ ]: print(df_train.info())
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 26707 entries, 0 to 26706
Data columns (total 37 columns):
 # Column          Non-Null Count Dtype  
 --- -----
 0 h1n1_concern      26707 non-null float64
 1 h1n1_knowledge     26707 non-null float64
 2 behavioral_antiviral_meds 26707 non-null float64
 3 behavioral_avoidance 26707 non-null float64
 4 behavioral_face_mask 26707 non-null float64
 5 behavioral_wash_hands 26707 non-null float64
 6 behavioral_large_gatherings 26707 non-null float64
 7 behavioral_outside_home 26707 non-null float64
 8 doctor_recc_h1n1    26707 non-null float64
 9 chronic_med_condition 26707 non-null float64
 10 child_under_6_months 26707 non-null float64
 11 health_worker      26707 non-null float64
 12 health_insurance   26707 non-null float64
 13 opinion_h1n1_vac_effective 26707 non-null float64
 14 opinion_h1n1_risk   26707 non-null float64
 15 opinion_h1n1_sick_from_vacc 26707 non-null float64
 16 opinion_seas_vacc_effective 26707 non-null float64
 17 opinion_seas_risk    26707 non-null float64
 18 opinion_seas_sick_from_vacc 26707 non-null float64
 19 age_group          26707 non-null float64
 20 education          26707 non-null float64
 21 race               26707 non-null float64
 22 income_poverty     26707 non-null float64
 23 marital_status     26707 non-null float64
 24 rent_or_own        26707 non-null float64
 25 hhs_geo_region     26707 non-null float64
 26 census_msa         26707 non-null float64
 27 household_adults   26707 non-null float64
 28 employment_industry 26707 non-null float64
 29 h1n1_vaccine       26707 non-null float64
 30 seasonal_vaccine   26707 non-null float64
 31 household_children  26707 non-null float64
 32 employment_occupton 26707 non-null float64
 33 age_group          26707 non-null float64
 34 education          26707 non-null float64
 35 race               26707 non-null float64
 36 sex                26707 non-null float64
dtypes: float64(37)
memory usage: 7.57 MB
None
```

```
Correlation Heatmap
```

```
In [ ]: h_labels = [x.replace('_train', '_').replace('.d', '.t') for x in list(df_train.select_dtypes(include=['number', 'bool']).columns.values)]
fig, ax = plt.subplots(figsize=(10,6))
sns.heatmap(df_train.corr(), annot=True, xticklabels=h_labels, yticklabels=h_labels)
```

```
Correlation Table
```

```
In [ ]: df_train.corr()
```

```
Out[ ]:
```

	h1n1_concern	h1n1_knowledge	behavioral_antiviral_meds	behavioral_avoidance
h1n1_concern	1.000000	0.076568	0.085562	
h1n1_knowledge	0.076568	1.000000	-0.001816	
behavioral_antiviral_meds	0.085562	-0.001816	1.000000	
behavioral_avoidance	0.050465	0.095192	0.051383	
behavioral_face_mask	0.51493	0.29994	0.139040	
behavioral_wash_hands	0.289809	0.094993	0.065579	
behavioral_large_gatherings	0.250078	-0.039979	0.102030	
behavioral_outside_home	0.240748	-0.057778	0.111666	
doctor_recc_h1n1	0.246532	0.094727	0.069112	
chronic_med_condition	0.267093	0.092033	0.047340	
child_under_6_months	0.032723	0.070276	0.004767	
health_worker	0.187183	-0.10431	-0.07178	
opinion_h1n1_vac_effective	0.213980	0.135095	0.016339	
opinion_h1n1_risk	0.351031	0.088906	0.085836	
opinion_h1n1_sick_from_vacc	0.353772	0.002560	0.064301	
opinion_seas_vacc_effective	0.202705	0.114322	-0.01074	
opinion_seas_risk	0.307118	0.098998	0.063913	
opinion_seas_sick_from_vacc	0.208736	-0.029376	0.066733	
age_group	0.051225	-0.051356	-0.096319	
education	0.014636	-0.116299	-0.110590	
race	0.242125	-0.052220	0.022828	
income_poverty	0.123978	-0.035413	-0.054453	
marital_status	0.353413	-0.138041	-0.054453	
rent_or_own	0.053092	0.060474	0.081529	
hhs_geo_region	0.008409	-0.015268	0.007902	
census_msa	0.030374	-0.059151	0.010412	
household_adults	-0.011020	0.033406	0.040948	
household_children	0.053092	0.060474	0.081529	
employment_industry	0.2569823	0.098985	0.019574	
employment_occupton	0.167323	0.248827	0.013989	
h1n1_vaccine	0.120970	0.115295	0.033279	
seasonal_vaccine	0.152583	0.116014	0.000779	

37 rows × 37 columns

Variance

```
In [ ]: df_train.var()
```

```
Out[ ]:
```

	h1n1_concern	h1n1_knowledge	behavioral_antiviral_meds	behavioral_avoidance	behavioral_f
0	0.849349	0.402587	0.402587	0.076568	0.085562
1	0.402587	0.402587	0.402587	0.076568	-0.001816
2	0.402587	0.402587	0.402587	0.076568	1.000000
3	0.402587	0.402587	0.402587	0.076568	0.051383
4	0.402587	0.402587	0.402587	0.076568	0.139040
...
26702	0.228338	0.230597	0.233459	0.032723	0.102030
26703	0.102030	0.102030	0.102030	0.032723	0.111666
26704	0.102030	0.102030	0.102030	0.032723	0.069112
26705	0.0	0.0	0.0	0.0	0.0
26706	0.0	0.0	0.0	0.0	0.0

26707 rows × 37 columns

<

```
In [ ]: df_train.dtypes
```

```
Out[ ]:
```

	h1n1_concern	h1n1_knowledge	behavioral_antiviral_meds	behavioral_avoidance	behavioral_f
0	float64	float64	float64	float64	float64
1	float64	float64	float64	float64	float64
2	float64	float64	float64	float64	float64
3	float64	float64	float64	float64	float64
4	float64	float64	float64	float64	float64
...
26702	float64	float64	float64	float64	float64
26703	float64	float64	float64	float64	float64
26704	float64	float64	float64	float64	float64
26705	float64	float64	float64	float64	float64
26706	float64	float64	float64	float64	float64

Balancing/Oversampling minority classes using SMOTE.

```
In [ ]: df_train['h1n1_vaccine'] = df_train['h1n1_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
df_train['seasonal_vaccine'] = df_train['seasonal_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
```

```
Out[ ]:
```

	h1n1_vaccine	seasonal_vaccine
0	0.849349	0.402587
1	0.402587	0.402587
2	0.402587	0.402587
3	0.402587	0.402587
4	0.402587	0.402587
...
26702	0.228338	0.230597
26703	0.102030	0.102030
26704	0.102030	0.102030
26705	0.0	0.0
26706	0.0	0.0

26707 rows × 37 columns

```
< 
```

```
In [ ]: df_train.dtypes
```

```
Out[ ]:
```

	h1n1_vaccine	seasonal_vaccine
0	float64	float64
1	float64	float64
2	float64	float64
3	float64	float64
4	float64	float64
...
26702	float64	float64
26703	float64	float64
26704	float64	float64
26705	float64	float64
26706	float64	float64

26707 rows × 37 columns

```
< 
```

```
In [ ]: df_train['h1n1_vaccine'] = df_train['h1n1_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
df_train['seasonal_vaccine'] = df_train['seasonal_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
```

```
Out[ ]:
```

	h1n1_vaccine	seasonal_vaccine
0	0.849349	0.402587
1	0.402587	0.402587
2	0.402587	0.402587
3	0.402587	0.402587
4	0.402587	0.402587
...
26702	0.228338	0.230597
26703	0.102030	0.102030
26704	0.102030	0.102030
26705	0.0	0.0
26706	0.0	0.0

26707 rows × 37 columns

```
< 
```

```
In [ ]: df_train['h1n1_vaccine'] = df_train['h1n1_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
df_train['seasonal_vaccine'] = df_train['seasonal_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
```

```
Out[ ]:
```

	h1n1_vaccine	seasonal_vaccine
0	0.849349	0.402587
1	0.402587	0.402587
2	0.402587	0.402587
3	0.402587	0.402587
4	0.402587	0.402587
...
26702	0.228338	0.230597
26703	0.102030	0.102030
26704	0.102030	0.102030
26705	0.0	0.0
26706	0.0	0.0

26707 rows × 37 columns

```
< 
```

```
In [ ]: df_train['h1n1_vaccine'] = df_train['h1n1_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
df_train['seasonal_vaccine'] = df_train['seasonal_vaccine'].apply(lambda x: 1 if x == 'Yes' else 0)
```

```
Out[ ]:
```

	h1n1_vaccine	seasonal_vaccine
0	0.849349	0.402587
1	0.402587	0.402587
2	0.402587	0.402587
3	0.402587	0.402587
4	0.402587	0.402587
...
26702	0.228338	0.230597
26703	0.102030	0.102030
26704	0.102030	

```
#Gradient Boosting Classifier
gb.classif = GradientBoostingClassifier()
gb.classif.fit(X_train, y_train_swine)
y_pred = gb.classif.predict(X_test)
f1 = round(f1_score(y_test_swine, y_pred, average='macro'), 3)
print("\nGradientBoosting (using scalar inputs) f1 score: ", f1)
classificationSummary(y_test_swine, y_pred)

#Gaussian Naive Bayes model
gnb = GaussianNB()
gnb.fit(X_train, y_train_swine)
y_pred = gnb.predict(X_test)
f1 = round(f1_score(y_test_swine, y_pred, average='macro'), 3)
print("\nNaive Bayes (using scalar inputs) f1 score: ", f1)
classificationSummary(y_test_swine, y_pred)

#K-Nearest Neighbors
knn = KNeighborsClassifier(n_neighbors=3, weights = 'distance')
knn.fit(X_train, y_train_swine)
y_pred = knn.predict(X_test)
f1 = round(f1_score(y_test_swine, y_pred, average='macro'), 3)
print("\nKNN f1 score: ", f1)
classificationSummary(y_test_swine, y_pred)

#Random Forest model
rf = RandomForestClassifier(max_depth=2, random_state = 12345)
rf.fit(X_train, y_train_swine)
y_pred = rf.predict(X_test)
f1 = round(f1_score(y_test_swine, y_pred, average='macro'), 3)
print("\nRandom Forest f1 score: ", f1)
classificationSummary(y_test_swine, y_pred)

Neural Network (using scalar inputs) f1 score: 0.331
Confusion Matrix (Accuracy 0.4945)

Prediction
Actual 0 1
0 6380
1 6240

Logistic Regression (using scalar inputs) f1 score: 0.785
Confusion Matrix (Accuracy 0.7848)

Prediction
Actual 0 1
0 1334
1 458
C:\users\eliefek\appdata\roaming\python\python39\site-packages\sklearn\svm\_base.py:122
5: ConvergenceWarning: Liblinear failed to converge, increase the number of iteration
S:
warnings.warn(
SVM (using scalar inputs) f1 score: 0.799
Confusion Matrix (Accuracy 0.7899)

Prediction
Actual 0 1
0 5137 1243
1 4584 4831

GradientBoosting (using scalar inputs) f1 score: 0.907
Confusion Matrix (Accuracy 0.9076)

Prediction
Actual 0 1
0 6033 347
1 819 5421

Naive Bayes (using scalar inputs) f1 score: 0.752
Confusion Matrix (Accuracy 0.7517)

Prediction
Actual 0 1
0 4584 1796
1 1338 4902

KNN f1 score: 0.805
Confusion Matrix (Accuracy 0.8101)

Prediction
Actual 0 1
0 4115 2265
1 131 6109

Random Forest f1 score: 0.827
Confusion Matrix (Accuracy 0.8269)

Prediction
Actual 0 1
0 4584 1796
1 1243 4997

In [ ]: #Find best parameters of Adaboost TODO using random search instead of grid search
boost_grid = {
    'n_estimators': [50, 100, 200, 400],
    'learning_rate': [0.01, 0.1, 0.5, 1.0, 1.5, 2.0]
}

boost_gridSearch = GridSearchCV(AdaBoostClassifier(), boost_grid, cv = 3)
boost_gridSearch.fit(X_train, y_train_swine)

print('Initial Adaboost Parameters:', boost_gridSearch.best_params_)

Initial Adaboost Parameters: {'learning_rate': 1.0, 'n_estimators': 400}
Confusion Matrix of Best Performing Model (Adaboost or Gradientboost I think but we need to
tinker with the hyperparameters).

In [ ]: #Boosted Classifier
adaboost = AdaBoostClassifier(DecisionTreeClassifier(max_depth = 2), learning_rate = 0.01)
adaboost.fit(X_train, y_train_swine)
y_pred = adaboost.predict(X_test)
f1 = round(f1_score(y_test_swine, y_pred, average='macro'), 3)
print("Adaboost f1 score: ", f1)
classificationSummary(y_test_swine, y_pred)

Adaboost f1 score: 0.888

In [ ]: #Create plot of the confusion matrix, showing performance of the model
cm = confusion_matrix(y_test_swine, y_pred)
group_names = ['True Neg', 'False Pos', 'False Neg', 'True Pos']
group_counts = ["{:0.2f}%".format(value) for value in
    cm.flatten()]
group_percentages = ["{:0.2%}".format(value) for value in
    cm.flatten()/np.sum(cm)]
labels = [(f'{v1}\n{v2}', f'{v3}\n{v4}') for v1, v2, v3, v4 in
    zip(group_names, group_counts, group_percentages)]
labels = np.asarray(labels).reshape(2,2)
sns.heatmap(cm, annot=labels, fmt='', cmap='Blues')

Out[ ]: <AxesSubplot:>



```