First attempt at a model - Group 4 (Mayuresh Mali, Divya G Tripathi, Devika Antarkar, Soma Ghosh, Kshitij Pathak) **Import Libraries** In [1]: import sklearn.feature\_extraction import pandas as pd import numpy as np import warnings warnings.filterwarnings("ignore", category=DeprecationWarning) import tldextract import matplotlib as plt from pylab import \* import re import math from itertools import groupby Read the data In [2]: train\_data = pd.read\_csv('train\_data.csv') train\_data\_copy = train\_data test\_data = pd.read\_csv('test\_data.csv') In [3]: import tldextract def domain\_extract(uri): ext = tldextract.extract(uri) if (not ext.suffix): return np.nan else: return ext.domain train\_data\_copy['domain'] = [ domain\_extract(DNS\_Address) for DNS\_Address in train\_data\_copy['DNS\_Address']] del train\_data\_copy['DNS\_Address'] train\_data\_copy.head() Out[3]: DNS\_Type domain dr-beckmann lifeadvicedaily benign benign aixenprovencetourism benign ayto-alcaladehenares In [4]: train\_data\_copy['domain'] = train\_data\_copy['domain'].astype(str) **Feature Engineering- Adding few more features** Let me create some features which can be useful such as length of domain name, number of vowels in it, consecutive consonants and digits in the domain name and also since XGBoost needs numeric features, we need to encode the data as numeric. In [5]: # Add a length field for the domain train\_data\_copy['length'] = [len(x) for x in train\_data\_copy['domain']] In [6]: train\_data\_copy.head() Out[6]: DNS\_Type domain length dr-beckmann 11 lifeadvicedaily 15 benign 20 aixenprovencetourism benign ayto-alcaladehenares In [7]: # Calculating character entropy import math from collections import Counter def entropy(s): p, lns = Counter(s), float(len(s)) return -sum( count/lns \* math.log(count/lns, 2) for count in p.values()) In [8]: def tri\_gram(domain): s = []count = 2while count < len(domain):</pre> s.append(domain[count - 2] + domain[count - 1] + domain[count]) count = count + 1**return** s  $dataset = \{\}$  $sum\_of\_frequency = 0.0$ # load trigrams\_google into dataset with open("trigram\_google.txt") as f: **for** line **in** f: (key, val) = line.split() dataset[key] = float(val) sum\_of\_frequency += dataset[key] # print sum def calc\_freq(trigrams): frequency = sum([dataset.get(trigram, 0) for trigram in trigrams]) / sum\_of\_frequency return frequency In [9]: # calculate vowels def calc\_vowels(y):  $num\_vowel = 0$ vowels = list('aeiou') **for** char **in** y: if char in vowels: num\_vowel += 1 return num\_vowel #Calculate number of digits def calc\_digits(z):  $num\_digit = 0$ digits = list('0123456789') for char in z: if char in digits: num\_digit += 1 return num\_digit In [10]: # Maximum length of Consecutive consonants def consecutive\_consonants(string): from itertools import groupby is\_vowel = lambda char: char in "aAeEiIoOuU" best = 0listnames = ["".join(g) for v, g in groupby(string, key=is\_vowel) if not v] for index in range(len(listnames)): if len(listnames[index]) > best: best = len(listnames[index]) return best In [11]: # Add a entropy field for the domain train\_data\_copy['entropy'] = [entropy(x) for x in train\_data\_copy['domain']] In [12]: # Add a trigram field for the domain train\_data\_copy['trigrams'] = [tri\_gram(x) for x in train\_data\_copy['domain']] In [13]: # Add a trigram frequency field for the domain train\_data\_copy['trigram\_freq'] = [calc\_freq(x) for x in train\_data\_copy['domain']] In [14]: # Add a vowels field for the domain train\_data\_copy['vowels'] = [calc\_vowels(x) for x in train\_data\_copy['domain']] In [15]: # Add a digits field for the domain train\_data\_copy['digits'] = [calc\_digits(x) for x in train\_data\_copy['domain']] In [16]: # Add a consec\_consonants field for the domain train\_data\_copy['consec\_consonants'] = [consecutive\_consonants(x) for x in train\_data\_copy['domain']] In [17]: train\_data\_copy['threat\_type'] = train\_data\_copy.DNS\_Type.apply(lambda x: 0 if x=='benign' else 1) In [18]: train\_data\_copy.head() Out[18]: DNS\_Type domain length entropy trigrams trigram\_freq vowels digits consec\_consonants threat\_type dr-beckmann 11 3.277613 [dr-, r-b, -be, bec, eck, ckm, kma, man, ann] benign 15 3.056565 benign lifeadvicedaily [lif, ife, fea, ead, adv, dvi, vic, ice, ced, ... smiles 6 2.251629 [smi, mil, ile, les] benign aixenprovencetourism 20 3.684184 [aix, ixe, xen, enp, npr, pro, rov, ove, ven, ... 3 20 3.403702 [ayt, yto, to-, o-a, -al, alc, lca, cal, ala, ... benign ayto-alcaladehenares In [19]: train\_data\_copy.tail() Out[19]: domain length entropy DNS\_Type trigrams trigram\_freq vowels digits consec\_consonants threat\_type [und, nde, der, erc, rci, cir, irc, rcu, 25 3.463465 999995 dga undercircumstancestyranny [suc, uch, cho, hof, ofc, fco, con, 999996 suchofconsentourabuses 22 3.390805 [eme, mek, ekc, kcw, cwk, wkw, 999997 16 2.602217 emekcwkwkmwkykca [cdk, dko, kok, oko, koe, oed, edn, 999998 cdkokoednvlk 12 2.855389 dga [pbo, boo, oow, own, wnx, nxe, xee, 999999 15 3.506891 dga pboownxeeayddcg In [20]: #Rearranging the order of columns train\_data\_copy = train\_data\_copy[['domain','DNS\_Type','threat\_type','length','entropy','trigrams','trigram\_freq','v owels','digits','consec\_consonants']] In [21]: | train\_data\_copy.head() Out[21]: DNS\_Type threat\_type length entropy trigrams trigram\_freq vowels digits consec\_consonants dr-beckmann 0 11 3.277613 [dr-, r-b, -be, bec, eck, ckm, kma, man, ann] 6 2.251629 [smi, mil, ile, les] smiles 3 aixenprovencetourism 20 3.684184 [aix, ixe, xen, enp, npr, pro, rov, ove, ven, ... 20 3.403702 4 ayto-alcaladehenares [ayt, yto, to-, o-a, -al, alc, lca, cal, ala, ... In [22]: # # selecting rows based on condition # rslt\_df = train\_data\_copy[train\_data\_copy['trigram\_freq'] > 0] # print('\nResult dataframe :\n', rslt\_df) **Model attempt 1** In [25]: **from xgboost import** XGBClassifier from sklearn.model\_selection import train\_test\_split from sklearn.metrics import accuracy\_score, f1\_score X = train\_data\_copy.as\_matrix(['length', 'entropy', 'vowels', 'digits','consec\_consonants']) y = np.array(train\_data\_copy['threat\_type'].tolist()) # Train on a 80/20 split from sklearn.model\_selection import train\_test\_split X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2,random\_state=0) model = XGBClassifier(silent=False, scale\_pos\_weight=1, learning\_rate=0.1, colsample\_bytree = 0.4, subsample = 0.8, objective='binary:logistic', n\_estimators=100,  $reg_alpha = 0.3,$ max\_depth=4, gamma=10) model.fit(X\_train, y\_train) C:\Users\mayur\AppData\Local\Continuum\anaconda3\lib\site-packages\ipykernel\_launcher.py:5: FutureWarning: Method .as \_matrix will be removed in a future version. Use .values instead. Out[25]: XGBClassifier(base\_score=0.5, booster='gbtree', colsample\_bylevel=1, colsample\_bynode=1, colsample\_bytree=0.4, gamma=10, learning\_rate=0.1, max\_delta\_step=0, max\_depth=4, min\_child\_weight=1, missing=None, n\_estimators=100, n\_jobs=1, nthread=None, objective='binary:logistic', random\_state=0, reg\_alpha=0.3, reg\_lambda=1, scale\_pos\_weight=1, seed=None, silent=False, subsample=0.8, verbosity=1) In [26]: ##Testing the model y\_pred = model.predict(X\_test) predictions = y\_pred In [27]: #Evaluating predictions accuracy = accuracy\_score(y\_test, predictions) f\_score = f1\_score(y\_test, predictions) print("Accuracy: %.3f%%" % (accuracy \* 100.0)) print("F1\_score:", f\_score) Accuracy: 88.599% F1\_score: 0.8885478273620412 **Model 1: Predictions on test data** In [28]: def test\_it(domain): In [29]: # Extracting just the second-level domain name from test data set rows. test\_data['newdomain'] = [ domain\_extract(x) for x in test\_data['domain']] #del test\_data['domain'] test\_data.head() # Adding the features to the test data set testdomain = test\_data.newdomain.apply(lambda x : test\_it(x)) In [30]: X2 = np.array(testdomain.tolist()) In [31]: ##Testing the model y\_pred2 = model.predict(X2) predictions = [round(value) for value in y\_pred2] y\_pred2 Out[31]: array([0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0]) In [32]: test\_data['threat']= y\_pred2 In [33]: test\_data.threat.value\_counts() Out[33]: 0 303 1 197 Name: threat, dtype: int64 Saving the data with predictions to a file In [34]: #Re-encode dga and benign labels test\_data['threat'] = test\_data.threat.apply(lambda x: 'benign' if x==0 else 'dga') del test\_data['newdomain'] In [36]: #Save output to file #os.remove("submission\_output3.csv") test\_data.to\_csv('submission\_output3.csv', index=False) We also tried another model with character encoding for letters a-z, numerals 0-9 and special characters to make 40 numbers denoting one for every character of the domain name. But we prefer this model as the accuracy and F1\_score are slightly better in this model. We will try different ensemble methods

and feature engineering techniques (such as n-gram probabilities and dictionary check or letter label encodings etc.) in future to improve accuracy of the

model.