

Donald Trump Tweet Author Analysis

This personal project looks at data from president Trump's Twitter. Donald Trump used to tweet from an iPhone while his campaign staff tweeted through his account on Android, so this data can be used to train a neural network to predict who wrote his newer tweets based on a set of features to include how often he used certain words and how long the tweets were.

Load Data

```
% Uncomment this on first run to load in data, takes a while
% opts = detectImportOptions('dataset.xlsx');
% opts.SelectedVariableNames = [8 12:334];
% X = readmatrix('dataset.xlsx',opts);
% opts.SelectedVariableNames = 2;
% Y = readmatrix('dataset.xlsx',opts);
load('X.mat');
load('Y.mat');
s = string(Y);
id = s == "Twitter for iPhone" | s == "Twitter for Android";
s = s(id);
X = X(id,:);
Y = s == "Twitter for iPhone";
```

Mean normalize X

```
non_present = max(X) == 0;
X = X(:,~non_present);
X = (X - repmat(mean(X), size(X,1), 1)) ./ repmat(2*std(X), size(X,1), 1);
X = [ones(size(X,1),1) X];
```

Divide data

```
m = size(X,1);
[trainInd,valInd,testInd] = dividerand(m,.6,.2,.2);
X_train = X(trainInd,:);
Y_train = Y(trainInd,:);
X_val = X(valInd,:);
Y_val = Y(valInd,:);
X_test = X(testInd,:);
Y_test = Y(testInd,:);
```

Logistic Regression - Train theta

Here I will use logistic regression, then I will try a neural network.

```
initial_theta = zeros(size(X_train,2),1);
options = optimoptions(@fminunc,'Algorithm','Quasi-Newton','GradObj','on','MaxIter',400);
[theta, cost] = fminunc(@(t)(costFunctionReg(t, X_train, Y_train,1)), initial_theta, options);
```

Local minimum possible.

fminunc stopped because it cannot decrease the objective function along the current search direction.

<stopping criteria details>

```
p = predict(theta, X_val);
success = p == Y_val;
success = success(success);
success_rate = length(success)/length(p);
fprintf('Logistic Regression Validation Set Accuracy = %f\n', success_rate);
```

Logistic Regression Validation Set Accuracy = 0.854509

Neural Network

The neural network is a simple 1-hidden layer neural network.

Initialize Thetas

```
% nn script
hidden_layer_size = 40;
input_layer_size = 308;

initial_Theta1 = randInitializeWeights(input_layer_size, hidden_layer_size);
initial_Theta2 = randInitializeWeights(hidden_layer_size, 1);

% Unroll parameters
initial_nn_params = [initial_Theta1(:) ; initial_Theta2(:)];
```

Train

```
options = optimset('MaxIter', 60);
lambda = 3;

% Create "short hand" for the cost function to be minimized
costFunction = @(p) nnCostFunction(p, input_layer_size, hidden_layer_size, X_train, Y_train, lambda);

% Now, costFunction is a function that takes in only one argument (the
% neural network parameters)
[nn_params, ~] = fmincg(costFunction, initial_nn_params, options);
```

Iteration	1		Cost: 7.003901e-01
Iteration	2		Cost: 6.277823e-01
Iteration	3		Cost: 4.399062e-01
Iteration	4		Cost: 3.655925e-01
Iteration	5		Cost: 3.611915e-01
Iteration	6		Cost: 3.267465e-01
Iteration	7		Cost: 3.231518e-01
Iteration	8		Cost: 3.181216e-01
Iteration	9		Cost: 3.179061e-01
Iteration	10		Cost: 3.155242e-01
Iteration	11		Cost: 3.140591e-01
Iteration	12		Cost: 3.106492e-01
Iteration	13		Cost: 3.095375e-01
Iteration	14		Cost: 3.061780e-01

Iteration	15	Cost: 3.053409e-01
Iteration	16	Cost: 3.028552e-01
Iteration	17	Cost: 3.015783e-01
Iteration	18	Cost: 2.981240e-01
Iteration	19	Cost: 2.962275e-01
Iteration	20	Cost: 2.939230e-01
Iteration	21	Cost: 2.888381e-01
Iteration	22	Cost: 2.884399e-01
Iteration	23	Cost: 2.865061e-01
Iteration	24	Cost: 2.851014e-01
Iteration	25	Cost: 2.833821e-01
Iteration	26	Cost: 2.819007e-01
Iteration	27	Cost: 2.784343e-01
Iteration	28	Cost: 2.775216e-01
Iteration	29	Cost: 2.711898e-01
Iteration	30	Cost: 2.668222e-01
Iteration	31	Cost: 2.659855e-01
Iteration	32	Cost: 2.586643e-01
Iteration	33	Cost: 2.580030e-01
Iteration	34	Cost: 2.547711e-01
Iteration	35	Cost: 2.542652e-01
Iteration	36	Cost: 2.527636e-01
Iteration	37	Cost: 2.522833e-01
Iteration	38	Cost: 2.501889e-01
Iteration	39	Cost: 2.500489e-01
Iteration	40	Cost: 2.483174e-01
Iteration	41	Cost: 2.470006e-01
Iteration	42	Cost: 2.451896e-01
Iteration	43	Cost: 2.440025e-01
Iteration	44	Cost: 2.402973e-01
Iteration	45	Cost: 2.358044e-01
Iteration	46	Cost: 2.351440e-01
Iteration	47	Cost: 2.344891e-01
Iteration	48	Cost: 2.336491e-01
Iteration	49	Cost: 2.327416e-01
Iteration	50	Cost: 2.320184e-01
Iteration	51	Cost: 2.293800e-01
Iteration	52	Cost: 2.281883e-01
Iteration	53	Cost: 2.233164e-01
Iteration	54	Cost: 2.215880e-01
Iteration	55	Cost: 2.191914e-01
Iteration	56	Cost: 2.181871e-01
Iteration	57	Cost: 2.148518e-01
Iteration	58	Cost: 2.132022e-01
Iteration	59	Cost: 2.131667e-01
Iteration	60	Cost: 2.119214e-01

```
Theta1 = reshape(nn_params(1:hidden_layer_size * (input_layer_size + 1)), hidden_layer_size, (input_layer_size + 1));
Theta2 = reshape(nn_params((1 + (hidden_layer_size * (input_layer_size + 1))):end), 1, (hidden_layer_size + 1));
```

```
pred_val = predict_nn(Theta1, Theta2, X_val);
val_accuracy = mean(double(pred_val == Y_val));
fprintf('Neural Network Validation Set Accuracy = %f\n', val_accuracy);
```

Neural Network Validation Set Accuracy = 0.863616