# Data Science Career Track
## Model Metrics Exercise

1. Look at the table below. If the goal is to optimize the True Positives, which model would you choose and why?  I chose the Logistic model.   This model strikes a balance between a good score for ability to capture all of the True Positives (Recall 0.746) with the highest and good score for correctly identifying the True Positives (Precision 0.775).  As a result the F1 score is the highest as it derives its value from the mean of Recall and Precision.  The .999 Accuracy may indicate overfitting, if so the Logistic with auto threshold is a second choice.

| Model | Recall | Precision | Accuracy | F1 |
|---|---|---|---|---|
| Logistic | 0.746 | 0.775 | 0.999 | 0.761 |
| Logistic with auto threshold | 0.891 | 0.061 | 0.976 | 0.114 |
| Logistic with class weights | 0.878 | 0.110 | 0.988 | 0.195 |
| Hinge with auto threshold | 0.905 | 0.014 | 0.890 | 0.028 |
| Hinge with class weights | 0.878 | 0.103 | 0.987 | 0.185 |

2. Calculate the F-1 scores for each model and identify the best model based on the F1 score.  Deep NN has the highest F1 score (.805).

| Model | Recall | Precision | F1 | Auc/Roc |
|---|---|---|---|---|
| Deep NN | 0.79 | 0.82 | 2*0.82*0.79/(0.82+0.79)<br>0.805 | 0.92 |
| Logistic Regression | 0.75 | 0.79 | 2*0.79*0.75/(0.79+0.75)<br>0.769 | 0.90 |
| Random Forest | 0.80 | 0.66 | 2*0.66*0.80/(0.66+0.80)<br>0.723 | 0.90 |
| LinearSVC | 0.74 | 0.75 | 2*0.75*0.74/(0.75+0.74)<br>0.745 | 0.82 |

3. Identify the best parameter values for 'alpha' and 'L1-ratio' based on the above comparison. The values of 0.5 and 0.2 in Linear Regression (row 1) have a good MAE score (84.27) coupled with the highest R-squared value. The R-squared value is not great, but based upon it and the RMSE scores it appears the data has some outliers.

| Model | Parameter | Parameter | Metric | Metric | Metric |
|---|---|---|---|---|---|
| | Alpha | L1-ratio | MAE | R-squared | RMSE |
| Linear Regression | 0.5 | 0.2 | 84.27 | 0.277 | 158.1 |
| Linear Regression | 0.2 | 0.5 | 84.08 | 0.264 | 159.6 |
| Linear Regression | 0.5 | 0.5 | 84.12 | 0.272 | 158.6 |
| Linear Regression | 0 | 0 | 84.49 | 0.249 | 161.2 |