

# Predict MAX GB for AIS: RMSE Regression, Aug. 2016

## loadData

In [1]:

```
import graphlab
```

A newer version of GraphLab Create (v2.1) is available! Your current version is v2.0.1.  
You can use pip to upgrade the graphlab-create package. For more information see <https://turi.com/products/create/upgrade>.

In [8]:

```
df = graphlab.SFrame('data.csv')
df_no_missing = df.dropna()
df_no_missing
```

Read 52 lines. Lines per second: 3373.99

Finished parsing file /home/ubuntu/data.csv

Parsing completed. Parsed 52 lines in 0.017959 secs.

Finished parsing file /home/ubuntu/data.csv

Parsing completed. Parsed 52 lines in 0.016479 secs.

-----  
Inferred types from first 100 line(s) of file as  
column\_type\_hints=[str,str,str,str,int]  
If parsing fails due to incorrect types, you can correct  
the inferred type list above and pass it to read\_csv in  
the column\_type\_hints argument  
-----

Out[8]:

TABLE_NAME	ALL	ACTIVE	Var	MAX GB, Including Index
ALA_BASE	2,635,755	2,635,755	-	1
CEP_BASE	53,958	53,993	-35	1
CMP_BASE	372	372	-	1
COL_BASE	1,211,233	1,211,515	-282	1
CUC_BASE	100,794	81,480	19,314	1
CUS_BASE	6,385,988	3,247,298	3,138,690	2
CWS_BASE	55,978,362	55,978,362	-	0
DSC_BASE	1,502,131	1,502,131	-	1
EQC_BASE	2,846,987	2,086,960	760,027	1
EQE_BASE	43,356	86	43,270	1

[? rows x 5 columns]

Note: Only the head of the SFrame is printed. This SFrame is lazily evaluated.

You can use sf.materialize() to force materialization.

In [9]:

```
df.head(1)
```

Out[9]:

TABLE_NAME	ALL	ACTIVE	Var	MAX GB, Including Index
ALA_BASE	2,635,755	2,635,755	-	1

[1 rows x 5 columns]

## buildModel

In [12]:

```
train_data, test_data = df_no_missing.random_split(.8, seed=0)
```

In [13]:

```
reg_model = graphlab.linear_regression.create(train_data, target='MAX GB, Including Index', features=['ALL', 'ACTIVE', 'Var'])
```

WARNING: The number of feature dimensions in this problem is very large in comparison with the number of examples. Unless an appropriate regularization value is set, this model may not provide accurate predictions for a validation/test set.

Linear regression:

Number of examples : 37

Number of features : 3

Number of unpacked features : 3

Number of coefficients : 104

Starting Newton Method

Iteration	Passes	Elapsed Time	Training-max_error	Training-rmse
1	2	1.000626	0.171885	0.032390

SUCCESS: Optimal solution found.

In [14]:

```
reg_model.get('coefficients').print_rows(num_rows=18, num_columns=3)
```

name	index	value	...
(intercept)	None	1.17188459835	...
ALL	53,958	-0.0572895654035	...
ALL	372	-0.0859303731276	...
ALL	1,211,233	-0.0572895654156	...
ALL	100,794	-0.0572895654156	...
ALL	6,385,988	0.276012929135	...
ALL	55,978,362	-0.58586098905	...
ALL	1,502,131	-0.0859303731424	...
ALL	2,846,987	-0.0572895654194	...
ALL	43,356	-0.0572895654024	...
ALL	74,493,642	12.9415077214	...
ALL	173,199,033	9.60848277612	...
ALL	8,924,039	1.27592041273	...
ALL	5,891,071	-0.0572895654156	...
ALL	279,881,504	53.6044120547	...
ALL	40,722,585	96.6004338496	...
ALL	31,472,786	2.94243288543	...
ALL	204,729,741	25.9403050083	...

[104 rows x 4 columns]

In [15]:

```
print reg_model.evaluate(test_data)
```

```
{ 'max_error': 13.828115401649717, 'rmse': 5.6610034818192325 }
```

**applyModel**

In [16]:

```
def maxGB(x):  
    maxGB = df[df['TABLE_NAME']==x]  
    return reg_model.predict(maxGB)
```

In [17]:

```
maxGB('EQP_BASE')
```

Out[17]:

```
dtype: float  
Rows: 1  
[39.99640776272886]
```

## writeUp

This is a Root Mean Square Error (RMSE) regression model that results in predicted values close to the observed data. The RMSE value in the model results is an absolute measure of fit. One pitfall of RMSE is that it can incorporate too many variables, however this regression model has a Validation-rmse of 2.818729 and lower values indicate a better fit.