

Detecting JET Outliers in the Tokamak Database

Dan Segal
@djsegal

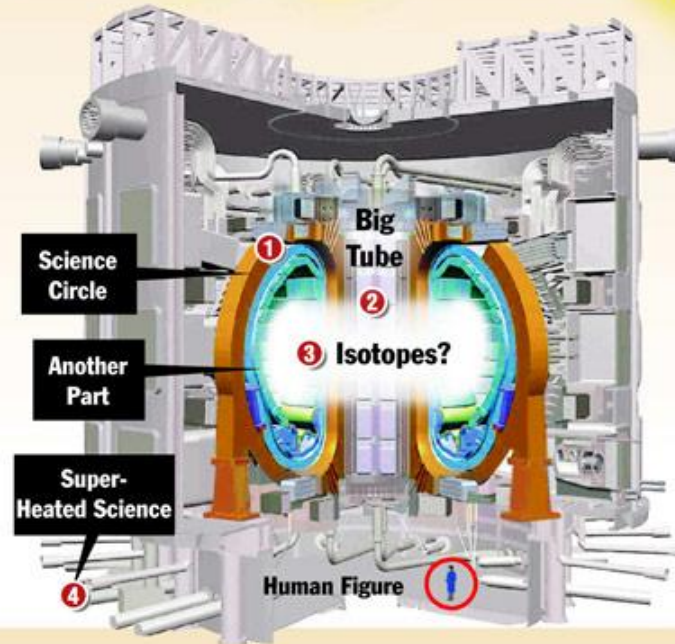
**Tokamaks will be fusion reactors
that compete with baseload
power sources – like coal**

ITER, at \$20B,
is the second
most costly
experiment in
the world

Onion Science Thursday

Giant Machine Creates Science

The Onion explains the inner workings of the complex, expensive science thing.



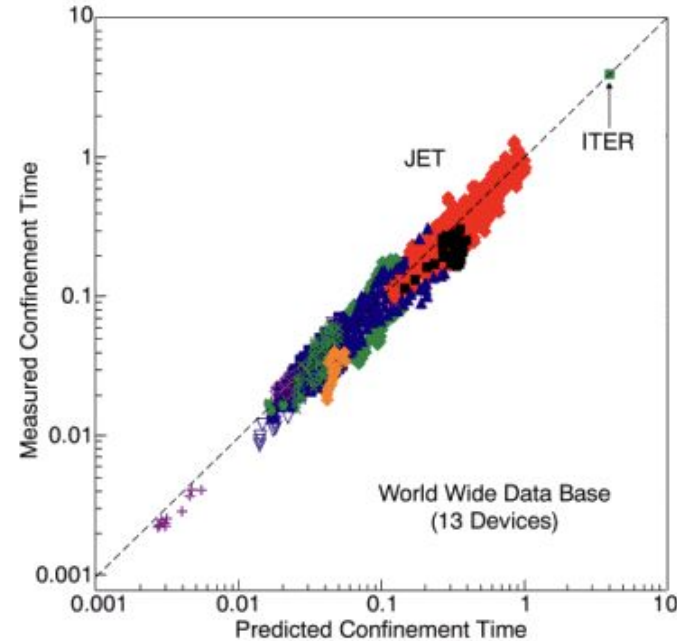
A Science Machine

The expensive device will test and execute more science than ever before

- 1 Scientists make sure machine's On/Off button is switched to On
- 2 Parts of the machine begin to move, at first slowly, and then rapidly
- 3 A lot of science begins to generate
- 4 Many things light up and sounds of thunder happen
- 5 Science ends

ITER was built
using a linear
regression on
confinement:

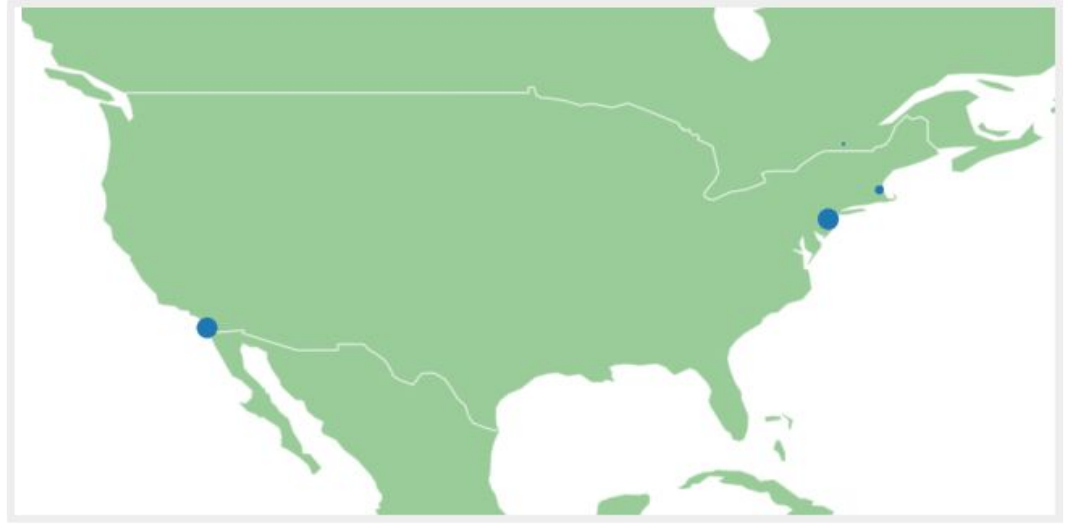
$$\tau_E^H = 0.145 H \frac{I_P^{0.93} R_0^{1.39} a^{0.58} \kappa^{0.78} \bar{n}^{0.41} B_0^{0.15} A^{0.19}}{P_{src}^{0.69}}$$



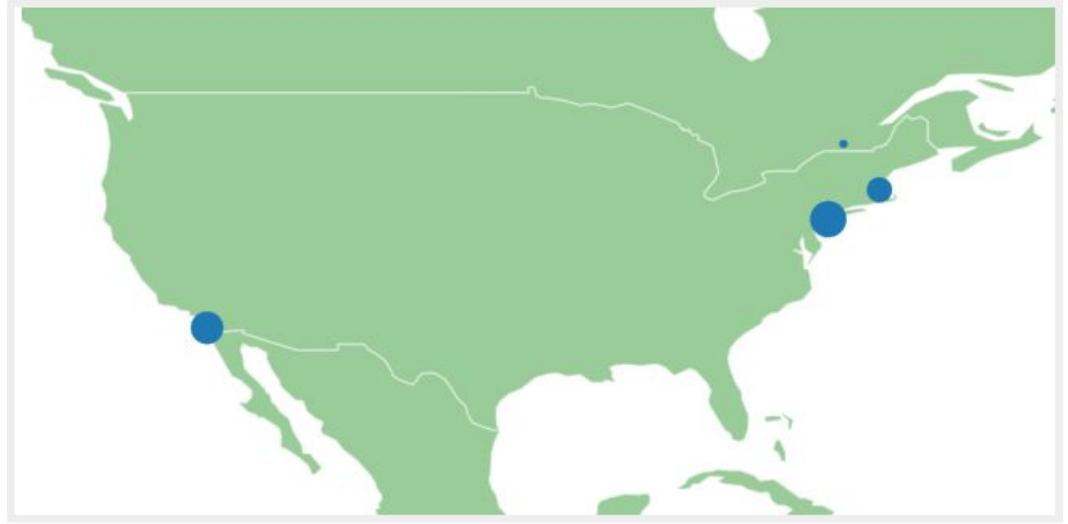
**Our goal is
to detect
JET outliers**

	Good	Bad
JET	1434	1710
Else	2097	913

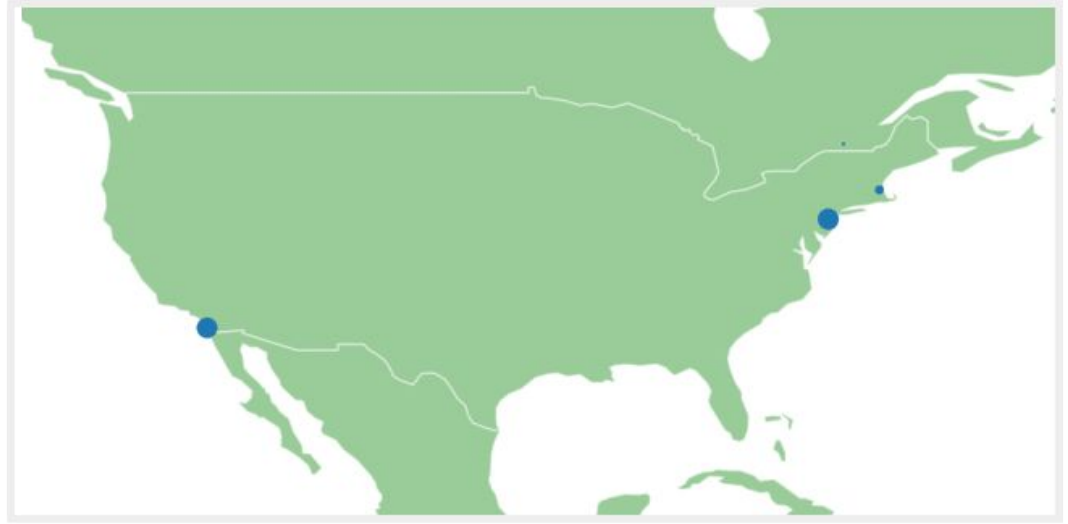
**The first step
is balancing
shots from all
the tokamaks**



**The first step
is balancing
shots from all
the tokamaks**

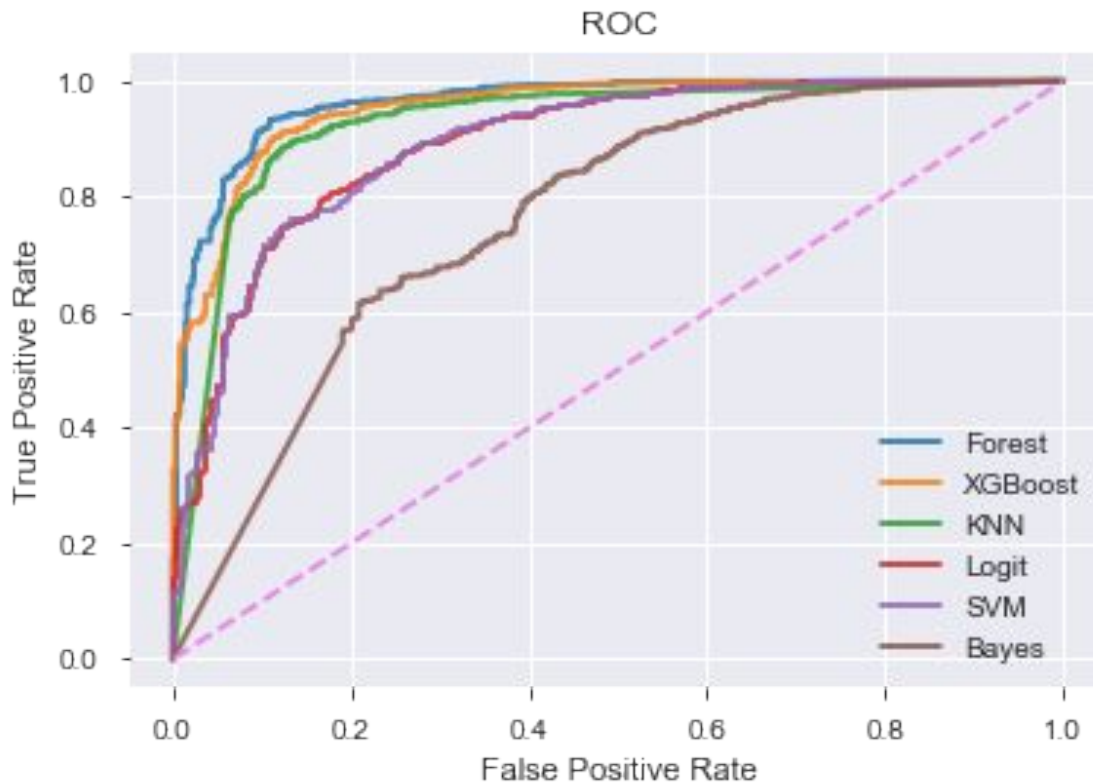


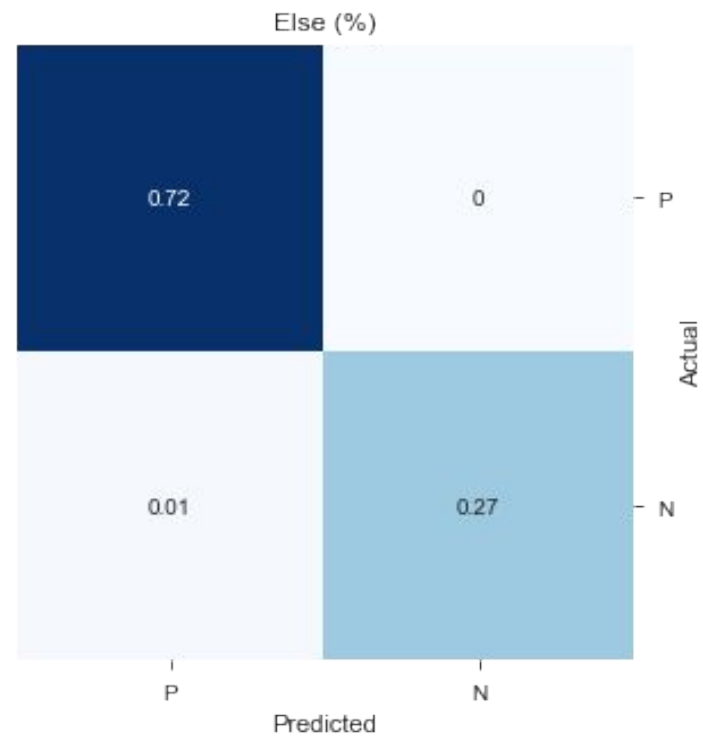
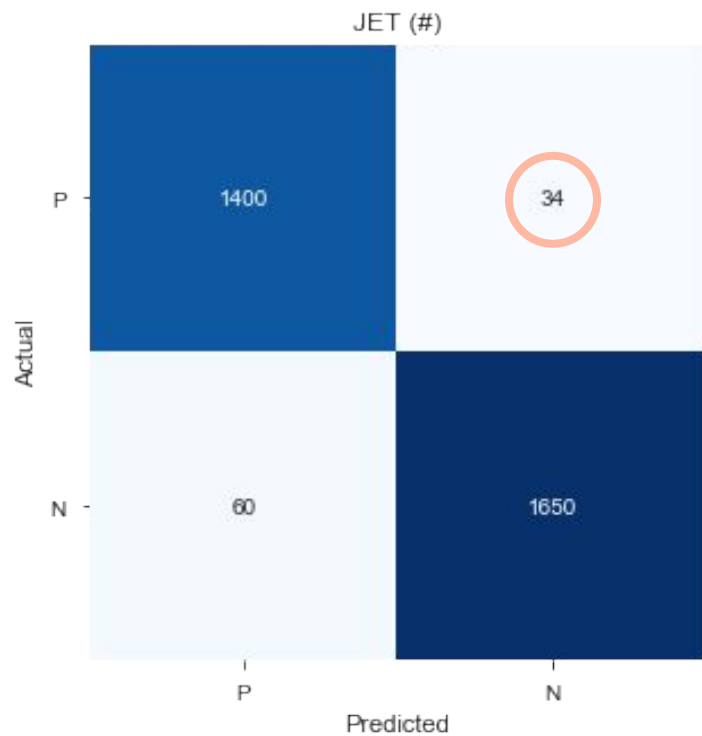
**The first step
is balancing
shots from all
the tokamaks**



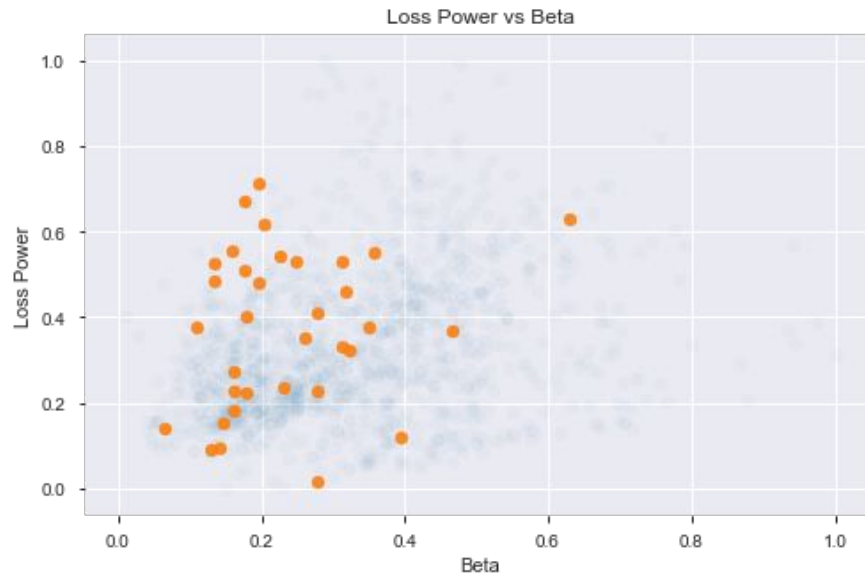
**Next we will plug 10k rows
into 6 classifiers using
the 50 strongest features**

Model	AUC (%)
Forest	96
XGBoost	95
KNN	93
Logit	89
SVM	89
<i>Bayes</i>	<i>77</i>



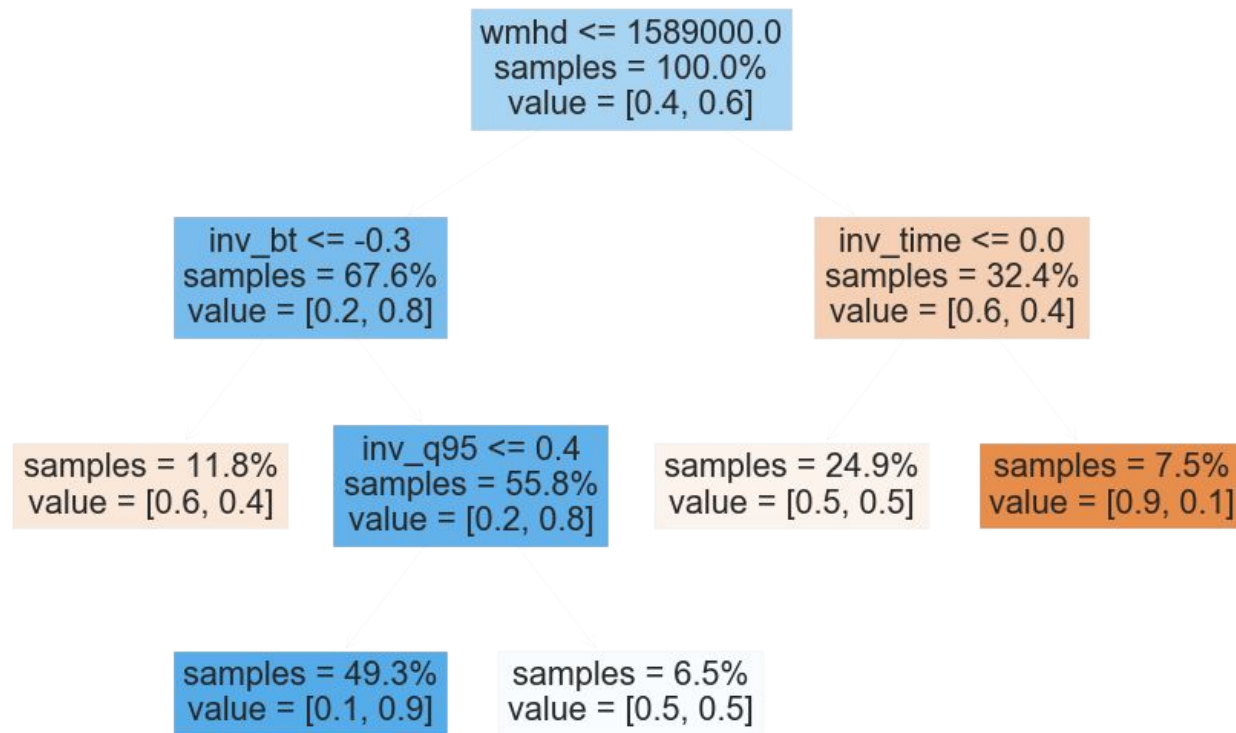


**The 34 outliers
tend to have low
betas and high
loss powers**



Features

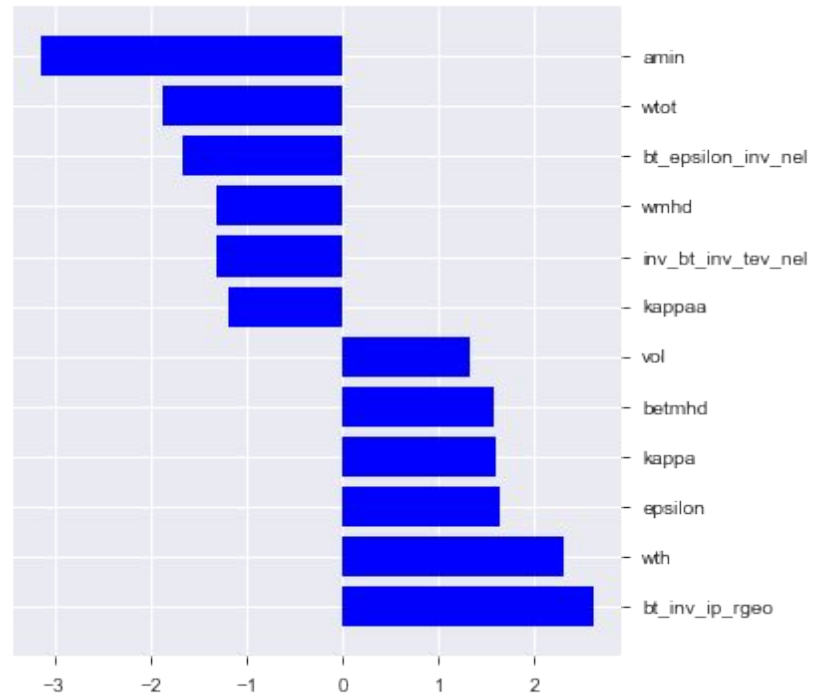
- Started with ~100 numeric parameters
- Imputed NaNs with median or predicted
- Expanded to ~1000 parameters w/ polynomials
- Removed ~500 highly correlated variables
- Eliminated ~200 variables with random forests
- Dwindled to ~50 variables with L1 classifiers



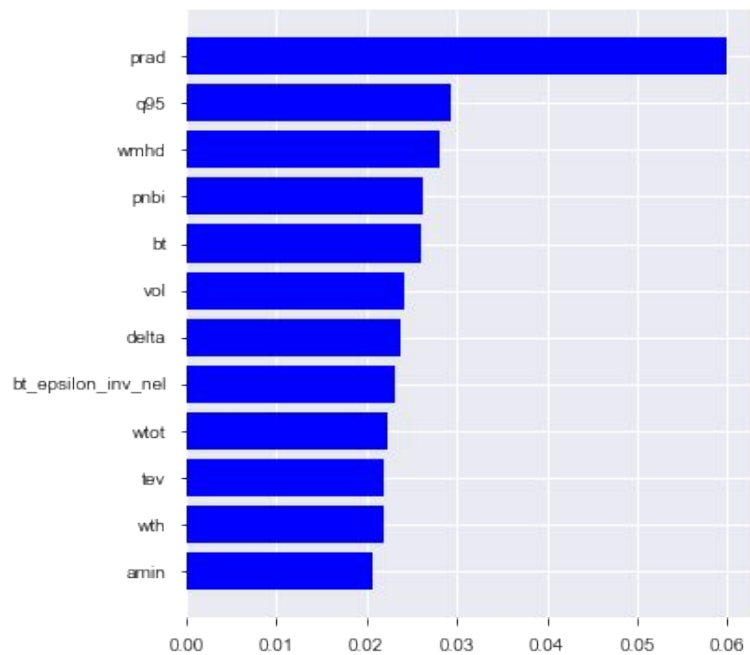
Logistic Features



SVM Features



Random Forests



XGBoost

