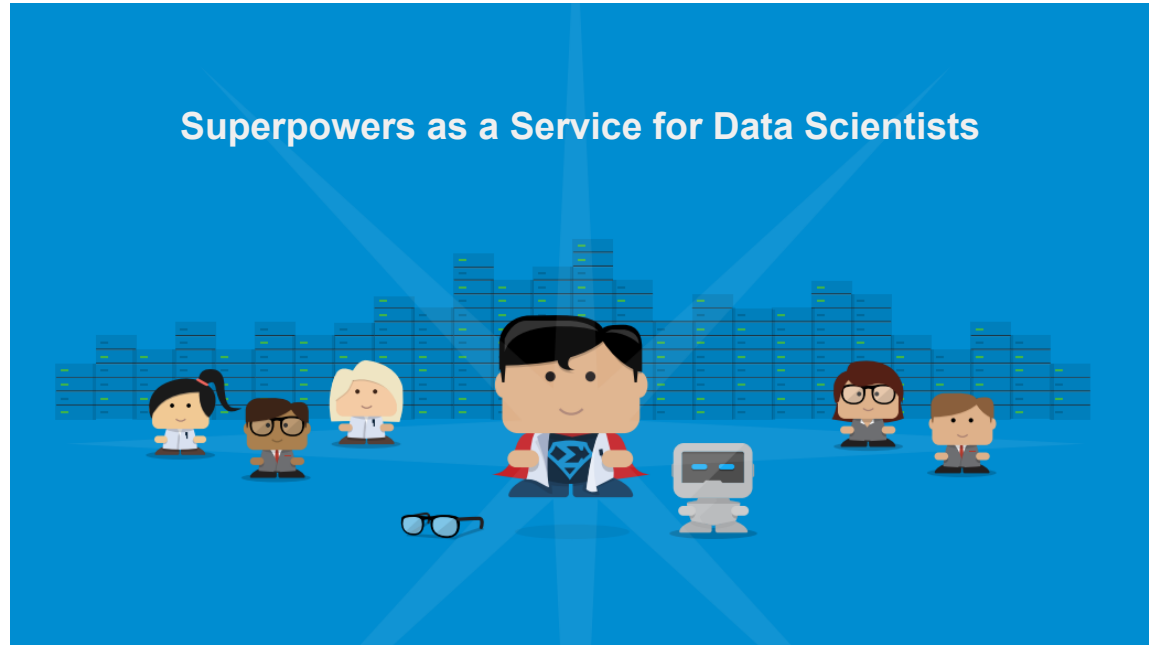# DataRobot

Winning the Liberty Mutual Kaggle Competition

Disclaimer:

- We do NOT state or imply Liberty Mutual endorses me, our team, our submission, Data Robot, Inc. or other companies I am involved with, or any of its or their products

- We thank Liberty Mutual for hosting the competition and identifying us as the winner of this competition.

# AGENDA

1. Team

2. Competition Summary

3. Challenges

4. Our Approach

5. Key Takeaways



**Superpowers as a Service for Data Scientists**

# DATAROBOT TEAM

### Xavier Conort

Chief Data Scientist

**Kaggle Rank**: #1    2012-2013

**Experience**:
Principal Research Engineer, I2R
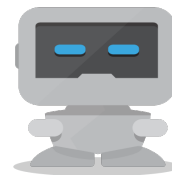Actuary, Risk Manager and CFO at CNP Assurances and AXA

### Owen Zhang

Chief Product Officer

**Kaggle Rank**: #1    2013-present

**Experience**:
VP,  Science, AIG
Sr. Director, Analytics at Travelers Insurance
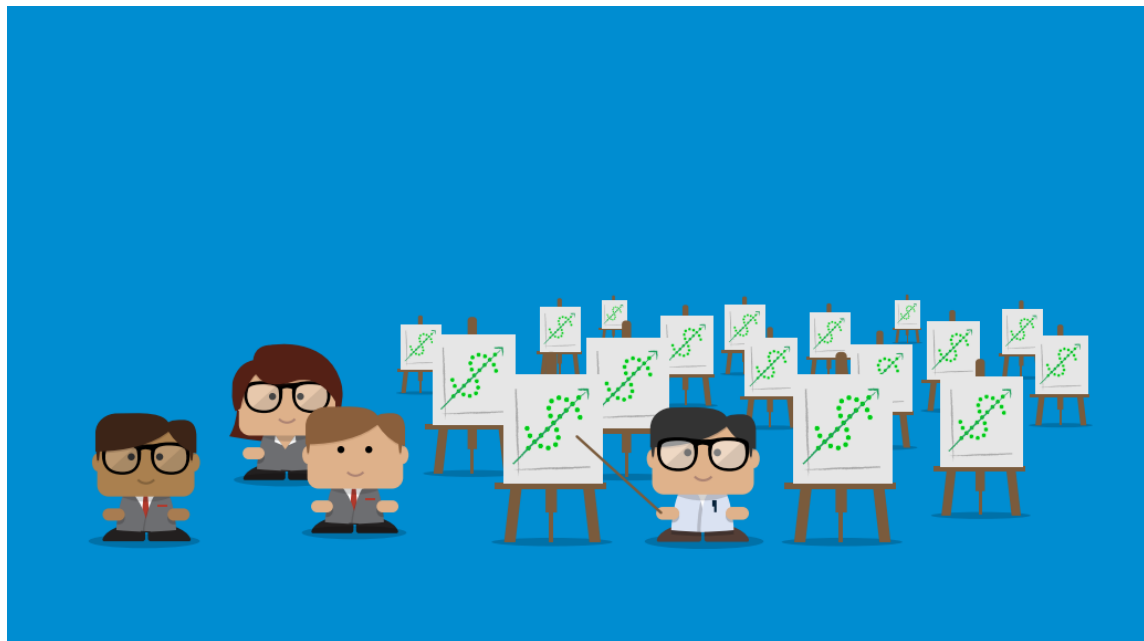
### DataRobot

Data Science Master

**Kaggle Rank**:   Master (top 1%)

**Experience**:
*10+ million* predictive models and counting ...

# AGENDA

1. Team

2. Competition Summary

3. Challenges

4. Our Approach

5. Key Takeaways

# COMPETITION SUMMARY

**Fire Peril Loss Cost**:

- Organized by Liberty Mutual

- Business problem: Predict expected fire losses for insurance policies
  - Significant portion of property losses
  - Volatile and hard to model correctly

- Started **July 8, 2014**

- Finished **Sept 2, 2014**



Completed • $25,000 • 634 teams

**Liberty Mutual Group - Fire Peril Loss Cost**

Tue 8 Jul 2014 – Tue 2 Sep 2014 (12 days ago)

Dashboard

Home
Data
Make a submission

Information

Description
Evaluation
Rules
Prizes
Timeline
Winners

Forum

Leaderboard
Public
Private

Competition Details » Get the Data » Make a submission

## Predict expected fire losses for insurance policies

A Fortune 100 company, Liberty Mutual Insurance has provided a wide range of insurance products and services designed to meet our customers' ever-changing needs for over 100 years.

Within the business insurance industry, fire losses account for a significant portion of total property losses. High severity and low frequency, fire losses are inherently volatile, which makes modeling them difficult. In this challenge, your task is to predict the target, a transformed ratio of loss to total insured value, using the provided information. This will enable more accurate identification of each policyholder's risk

Leaderboard

1. DataRobot
2. Ivanhoe
3. barisumog
4. datalab.se
5. paulperry

**Team DataRobot finished 1st out of 634 teams competing globally!**

# DATA OVERVIEW

- ~1 million insurance records

- 300 variables:

  **target** : The transformed ratio of loss to total insured value

  **id** : A unique identifier of the data set

  **dummy** : Nuisance variable used to control the model, but not a predictor

  **var1 – var17** : A set of normalized variables representing policy characteristics

  **crimeVar1 – crimeVar9** : Normalized Crime Rate variables

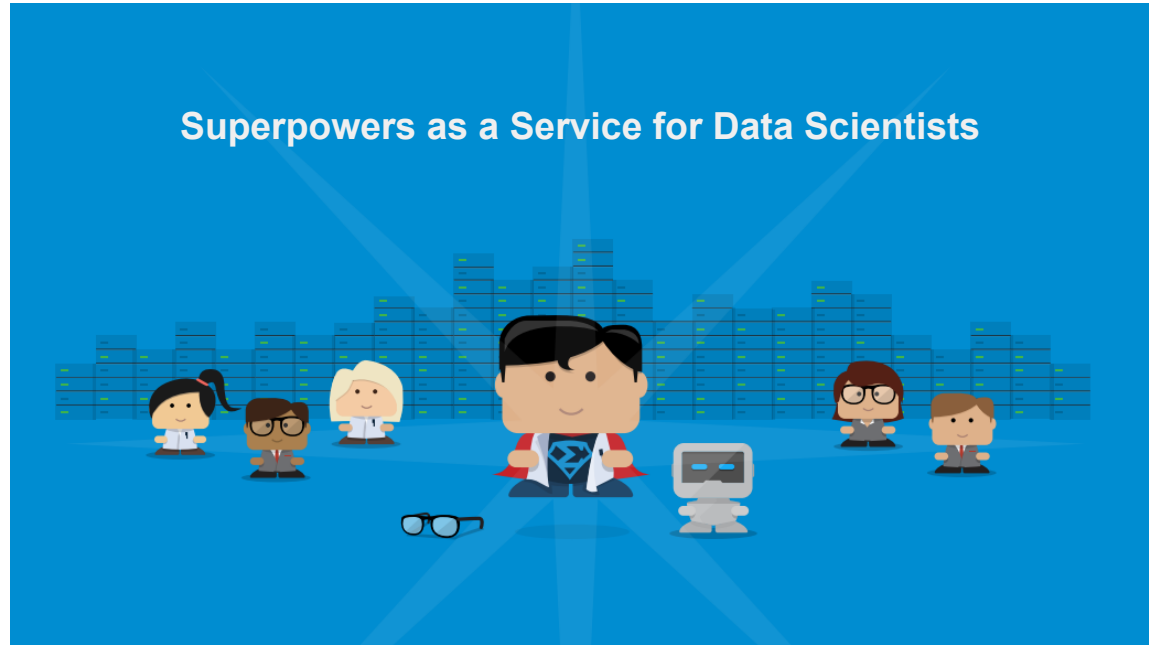  **geodemVar1 – geodemVar37** : Normalized geodemographic variables

  **weatherVar1 – weatherVar236** : Normalized weather station variables

| Numeric Variable Name | Variable Type |
| --- | --- |
| target | Continuous |
| id | Discrete |
| var10 | Continuous |
| var11 | Continuous |
| var12 | Continuous |
| var13 | Continuous |
| var14 | Continuous |
| var15 | Continuous |
| var16 | Continuous |
| var17 | Continuous |
| crimeVar1 – crimeVar9 | Continuous |
| geoDemVar1 – geoDemVar37 | Continuous |
| weatherVar1 – weath | |

| Categorical Variable Name | Variable Type | Possible Values |
| --- | --- | --- |
| var1 | Ordinal | 1, 2, 3, 4, 5, Z* |
| var2 | Nominal | A, B, C, Z* |
| var3 | Ordinal | 1, 2, 3, 4, 5, 6, Z* |
| var4[+] | Nominal | A1, B1, C1, D1, D2, D3, D4, E1, E2, E3, E4, E5, E6, F1, G1, G2, H1, H2, H3, I1, J1, J2, J3, J4, J5, J6, K1, L1, M1, N1, O1, O2, P1, R1, R2, R3, R4, R5, R6, R7, R8, S1, Z* |
| var5 | Nominal | A, B, C, D, E, F, Z* |
| var6 | Nominal | A, B, C, Z* |
| var7 | Ordinal | 1, 2, 3, 4, 5, 6, 7, 8, Z* |
| var8 | Ordinal | 1, 2, 3, 4, 5, 6, Z* |
| var9 | Nominal | A, B, Z* |
| dummy | Nominal | A, B |

# AGENDA

1.  Team

2.  Competition Summary

3.  Challenges

4.  Our Approach

5.  Key Takeaways



**Superpowers as a Service for Data Scientists**

# CHALLENGES

Low frequency
Very few homes get burned down!

High severity
A burnt down home may cost a lot of $$$

Many features

300 different variables!

Common Challenge for P&C Insurance dataset

# AGENDA

# OUR APPROACH

| High severity (presence of large claims) | Low frequency (few claims events) | Many Features |
|---|---|---|

| Censor large claims | Downsample non-claims | Reduce noise through feature selection |
|---|---|---|

| Explore best pre-processing and algorithms | Explore other models manually |
|---|---|

# CENSORING LARGE CLAIMS

Technique

- To censor large claims, we capped the target at the 80% quantile of non-zero losses

Impact

- More accurate and robust results than modeling with the raw target or its log transformation

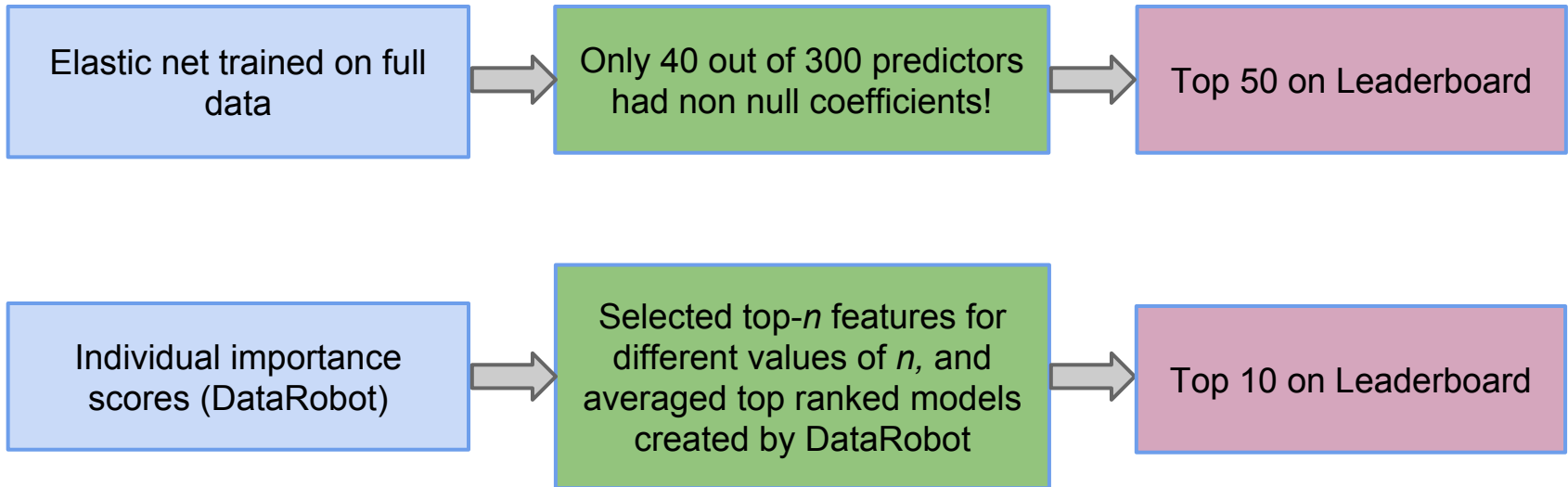# DOWNSAMPLING MAJORITY EVENTS

Technique

- Kept all non-0 records but only small % of 0 records from the training dataset

Impact

- Sped up model training significantly, and improved accuracy of certain ML algorithms (RandomForest, ExtraTrees, etc.)
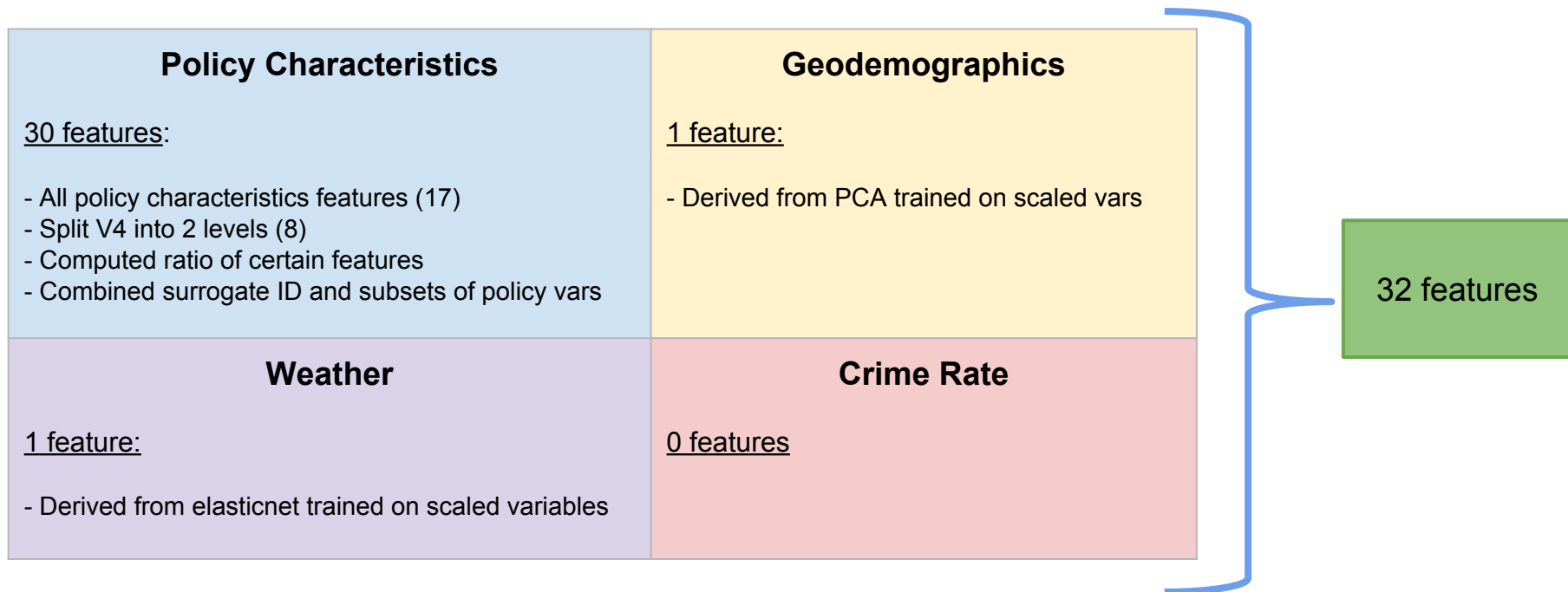
# REDUCING NOISE

- Ran 2 automated feature selections to detect large presence of noise
- Produced our own feature list with significantly low noise

| | | |
|---|---|---|
| Elastic net trained on full data | Only 40 out of 300 predictors had non null coefficients! | Top 50 on Leaderboard |

| | | |
|---|---|---|
| Individual importance scores (DataRobot) | Selected top-*n* features for different values of *n,* and averaged top ranked models created by DataRobot | Top 10 on Leaderboard |

# FEATURE ENGINEERING
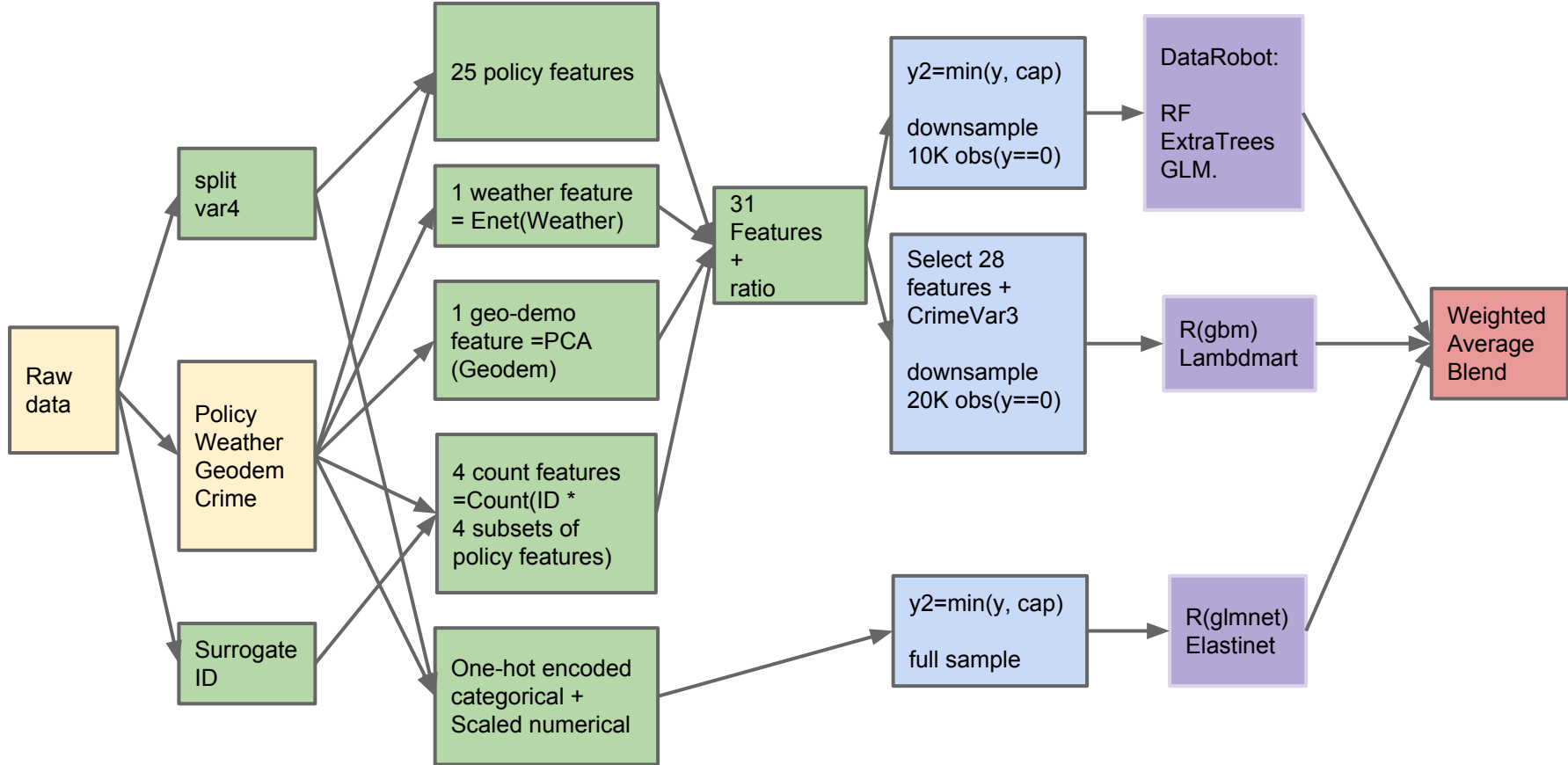
- Broke feature set into 4 components
- Created surrogate ID based on identical crime, geodemographics and weather variables

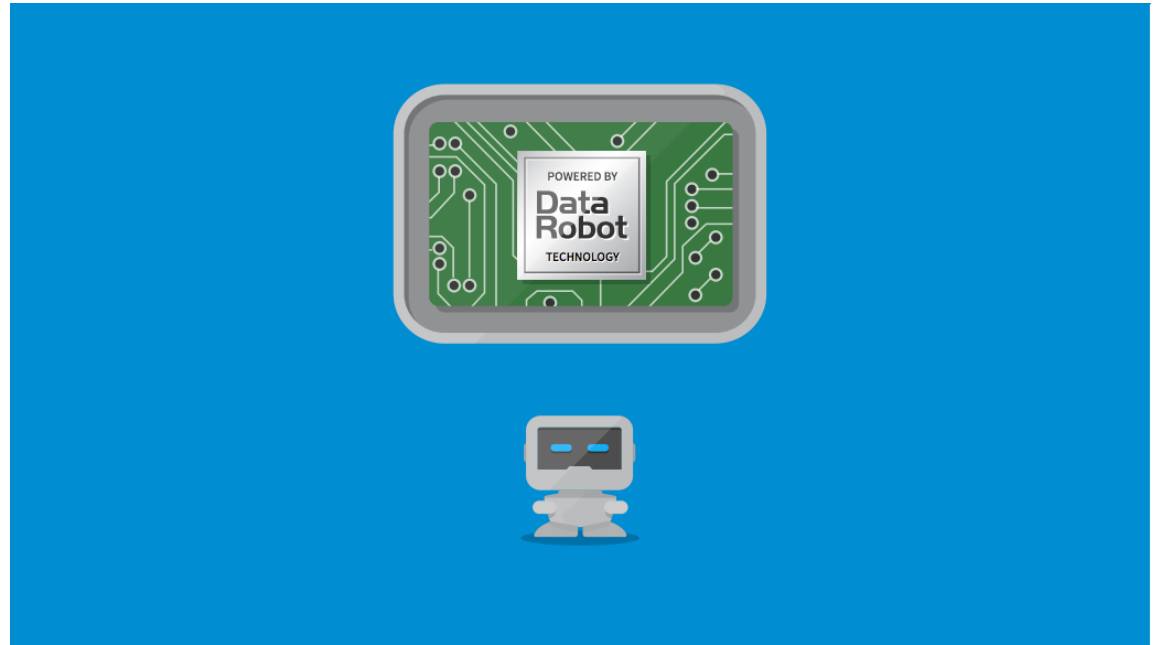| Policy Characteristics | Geodemographics |
|---|---|
| **30 features:**<br><br>- All policy characteristics features (17)<br>- Split V4 into 2 levels (8)<br>- Computed ratio of certain features<br>- Combined surrogate ID and subsets of policy vars | **1 feature:**<br><br>- Derived from PCA trained on scaled vars |
| **Weather** | **Crime Rate** |
| **1 feature:**<br><br>- Derived from elasticnet trained on scaled variables | **0 features** |

**32 features**

# MODELING WITH DATAROBOT

- Uploaded downsampled dataset with 32 features and capped target

- Ran multiple models in parallel with Normalized Gini metric to rank them

- Selected 3 promising models based on their 5-CV scores and trained them on 100% data

  - ExtraTrees Regressor

  - RandomForest Regressor

  - Ordinary Least Squares Regressor

# FINAL SOLUTION SUMMARY

# AGENDA

1. Team

2. Competition Summary

3. The Challenge

4. Our Approach

5. Key Takeaways

# KEY TAKEAWAYS

1.  Classic insurance problem set - keys to addressing this problem were:

    ○  Capping

    ○  Downsampling

    ○  Feature selection to reduce noise

    ○  Extracting value from blocks of features

2.  Capabilities of DataRobot

    ○  Rapidly explore diverse combinations of feature transformations and ML algorithms

    ○  Build highly accurate models from this search space automatically