# Datarobot, Automated ML

**Workflow and Benefits therein**

홍운표 🤖 **Data**Robot

# DataRobot

The world's most advanced Enterprise Machine Learning Automation platform
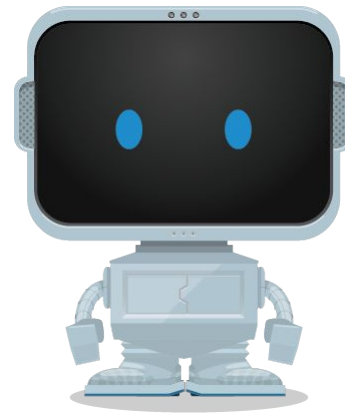
**2012**
Founded, HQ in Boston, MA

**$224M**
In funding

**1,000,000,000+**
Models built on DataRobot Cloud

**250+**
Data Scientists & Engineers (of 600+)

**4**
#1 ranked Data Scientists  kaggle

**50+**
Top 3 finishes  kaggle

INSURANCE     FINTECH     HEALTHCARE     MARKETING     BANKING     MANY MORE

DataRobot

# Best Practices and Technology

The top ranked Data Scientists in the world

**Owen Zhang**
Product Advisor
Highest: 1st
MASTER

**Xavier Conort**
Chief Data Scientist
Highest: 1st
MASTER

**Sergey Yurgenson**
Data Scientist
Highest: 1st
MASTER

**Amanda Schierz**
Data Scientist
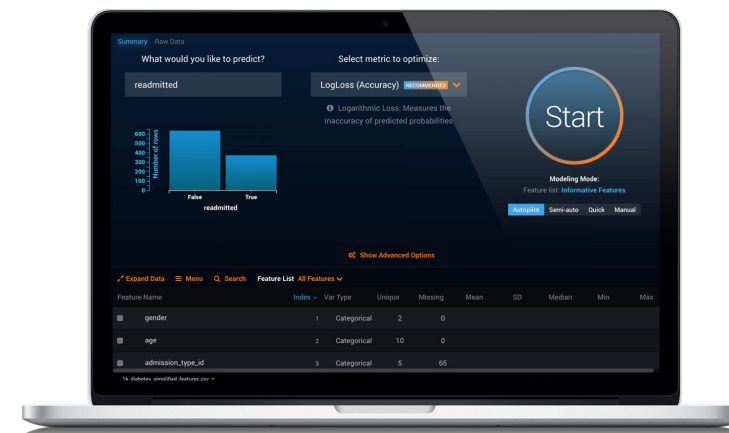Current: 1st Female, 1st in UK
MASTER

**Jeremy Achin**
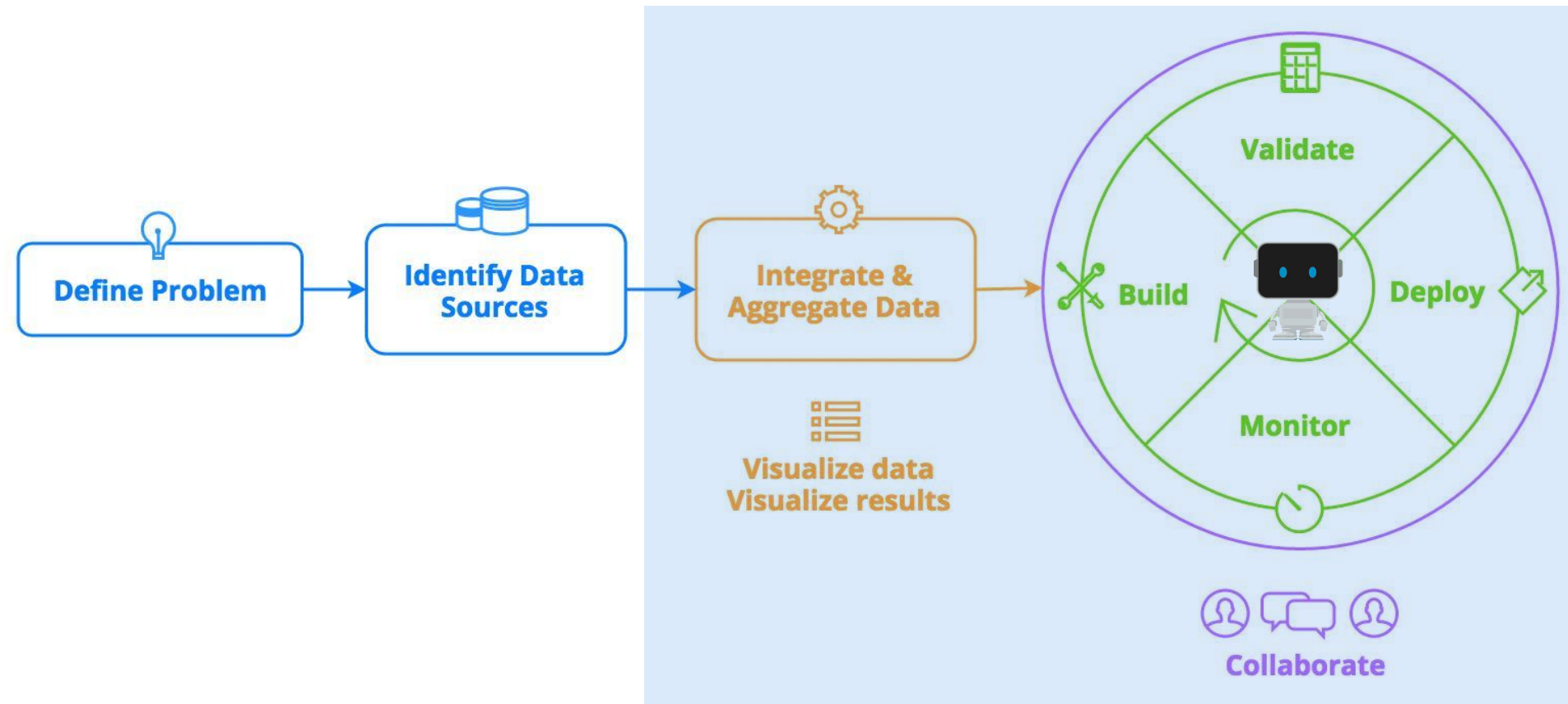CEO & Co-Founder
Highest: 20th
MASTER

**Tom de Godoy**
CTO & Co-Founder
Highest: 20th
MASTER

The best technologies in the world

# Data Science Methodology

Due to limited resource, call for amelioration



**No target goal**

**A few algorithms
& prone to overfit**

**Aging of model**

**Not sufficient Explanations**

DataRobot

# Motivations for AutoML

## Value of diverse set of algorithms

**Methodology driven**

⬇

**Problem driven**

TABLE 10.1. *Some characteristics of different learning methods. Key:* ▲ = *good,* ◆ = *fair, and* ▼ = *poor.*

| Characteristic | Neural Nets | SVM | Trees | MARS | k-NN, Kernels |
|---|---|---|---|---|---|
| Natural handling of data of "mixed" type | ▼ | ▼ | ▲ | ▲ | ▼ |
| Handling of missing values | ▼ | ▼ | ▲ | ▲ | ▲ |
| Robustness to outliers in input space | ▼ | ▼ | ▲ | ▼ | ▲ |
| Insensitive to monotone transformations of inputs | ▼ | ▼ | ▲ | ▼ | ▼ |
| Computational scalability (large $N$) | ▼ | ▼ | ▲ | ▲ | ▼ |
| Ability to deal with irrelevant inputs | ▼ | ▼ | ▲ | ▲ | ▼ |
| Ability to extract linear combinations of features | ▲ | ▲ | ▼ | ▼ | ◆ |
| Interpretability | ▼ | ▼ | ◆ | ▲ | ▼ |
| Predictive power | ▲ | ▲ | ▼ | ◆ | ▲ |

Source: http://statweb.stanford.edu/~tibs/ElemStatLearn/

**DataRobot**

# What is Automated Machine Learning

- 10 steps to building models
- An expert system that knows how to do each of these 10 steps, without human instructions
- Human friendly – not a black box
- Fast and accurate
- Replicable data science
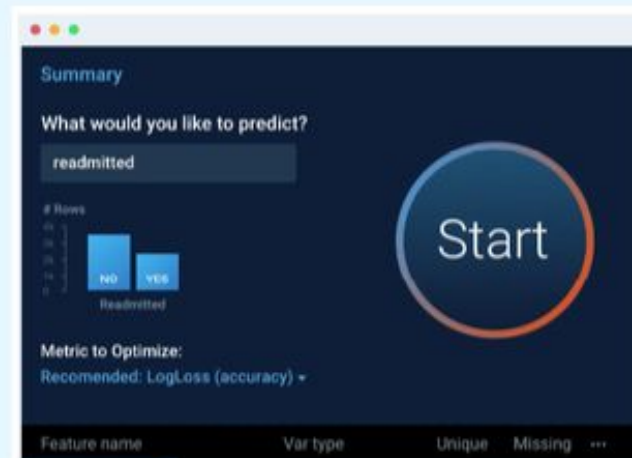


DataRobot

# What about DataRobot?



## Key Points

- End to end automated machine learning – all 10 steps are automated
- Hundreds of algorithms in the repository with new algorithms being added regularly
- Chooses the best algorithms for your data
- Best-in-class human-friendly insights
- Widest range of deployment options
- Enterprise ready
- Automatic model reports
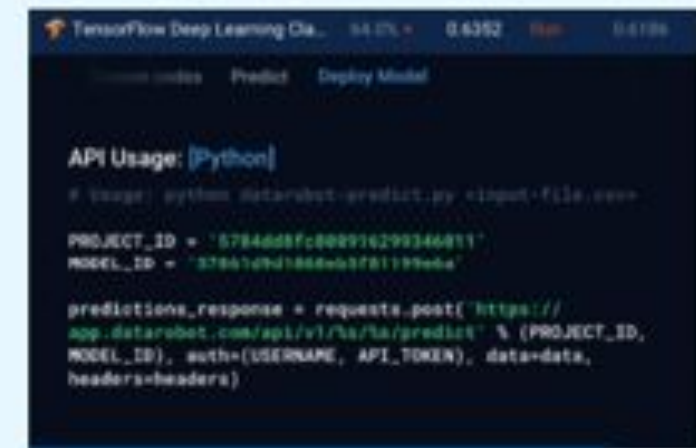- Large support team around the world

**DataRobot**

# DataRobot Workflow

# Different but powerful way of analysis

## A few perspectives (many more)

**Single model**

No need of **Hold-out** partition : just train/test or k-fold CV

Only **interpretable** algorithm is chosen ☐ Linear model is preferable

**Blending** starts from existing model

**Interaction** should be considered for model performance (linear model)

**parameter tuning** is limited for a model and time-consuming

**Multiple models**

**Hold-out** partition for evaluation of several models

**Interpretability** is model-agnostic

**Blending** is fair-basis reflecting multiple models performance, with speed vs accuracy data
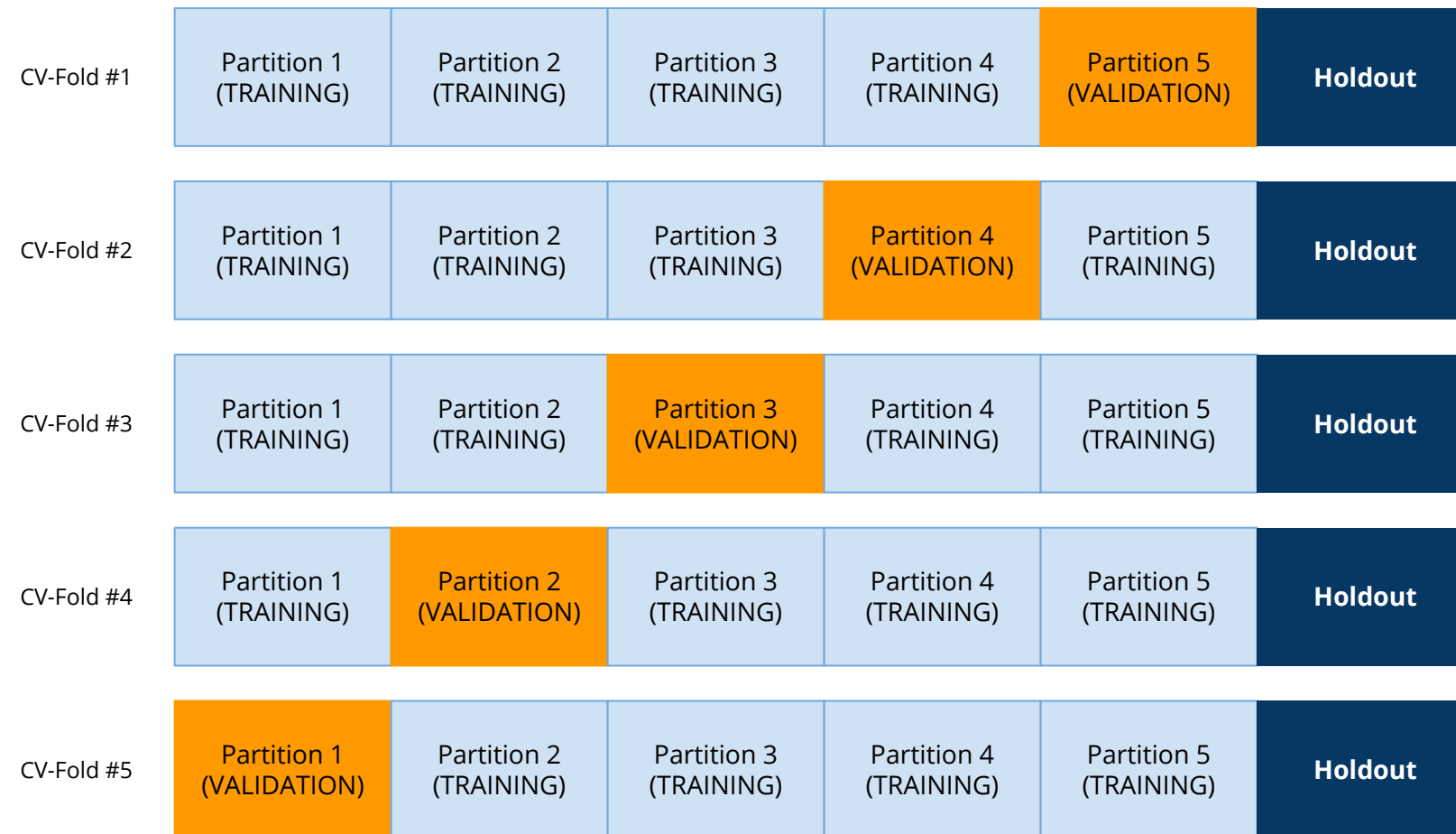
**Interaction** automatically reflected in tree-based algorithms. If interaction should be of importance, DR has GA2M model and R/Python api support for that

**Parameter tuning** is exhaustive for all candidate models.

One can easily confine the search space and quickly get the results

**DataRobot**

# Benefits : safer model

## Robust model free from the risk of overfitting

| | | | | | |
|---|---|---|---|---|---|
| CV-Fold #1 | Partition 1 (TRAINING) | Partition 2 (TRAINING) | Partition 3 (TRAINING) | Partition 4 (TRAINING) | Partition 5 (VALIDATION) | Holdout |

| | | | | | |
|---|---|---|---|---|---|
| CV-Fold #2 | Partition 1 (TRAINING) | Partition 2 (TRAINING) | Partition 3 (TRAINING) | Partition 4 (VALIDATION) | Partition 5 (TRAINING) | Holdout |

| | | | | | |
|---|---|---|---|---|---|
| CV-Fold #3 | Partition 1 (TRAINING) | Partition 2 (TRAINING) | Partition 3 (VALIDATION) | Partition 4 (TRAINING) | Partition 5 (TRAINING) | Holdout |

| | | | | | |
|---|---|---|---|---|---|
| CV-Fold #4 | Partition 1 (TRAINING) | Partition 2 (VALIDATION) | Partition 3 (TRAINING) | Partition 4 (TRAINING) | Partition 5 (TRAINING) | Holdout |

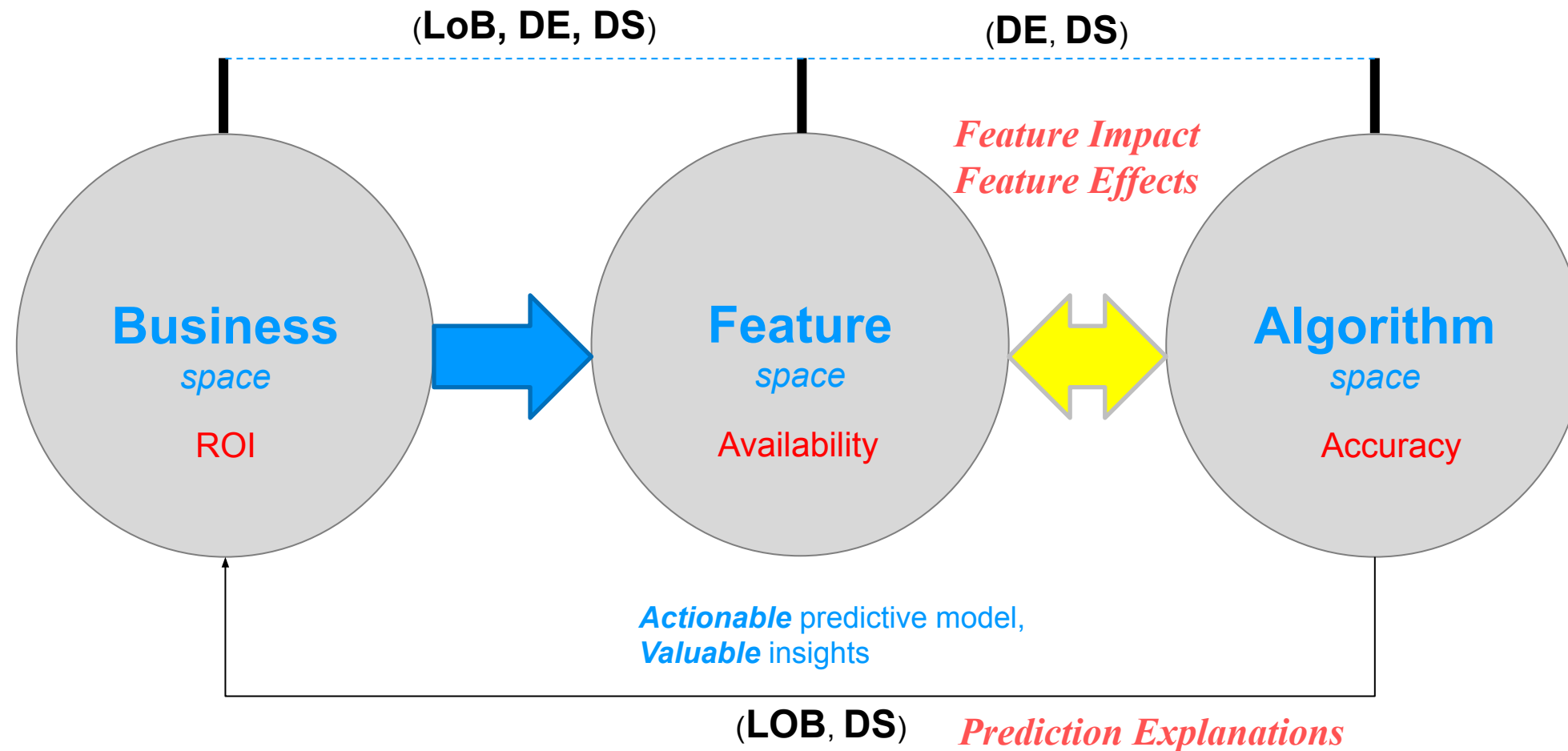| | | | | | |
|---|---|---|---|---|---|
| CV-Fold #5 | Partition 1 (VALIDATION) | Partition 2 (TRAINING) | Partition 3 (TRAINING) | Partition 4 (TRAINING) | Partition 5 (TRAINING) | Holdout |

The holdout is completely hidden from the models during the training process. After you have selected your optimal model, you can score your model on this to get your holdout score.

Average of these 5 validation scores is the cross validation score

**Data**Robot

# Benefits : more effort on feature space

"Feature engineering is the art of data science" (Sergey Yurgenson)



(**LoB, DE, DS**)

(**DE**, **DS**)

*Feature Impact*
*Feature Effects*

**Business** *space*

ROI

**Feature** *space*

Availability

**Algorithm** *space*

Accuracy

*Actionable* predictive model,
*Valuable* insights

(**LOB**, **DS**)   *Prediction Explanations*

**DataRobot**

# Benefits : Explainability
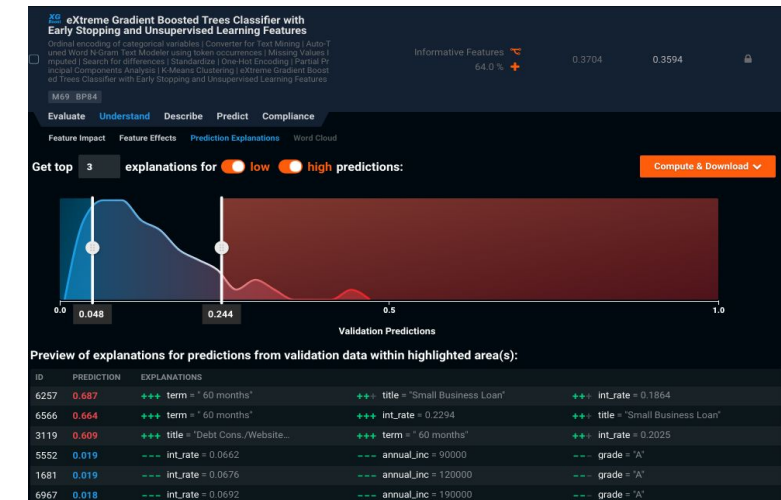
## Model-agnostic explanation



**[Feature Impact]**

- The importance of each feature
- Coincides with domain knowledge?
- Any new insights?

**[Feature Effect]**

- Relationship among target and a feature
- Relationship reflects domain knowledge?
- Any new insights or feature transform?

**[Prediction Explanation]**

- What is the basis of prediction?
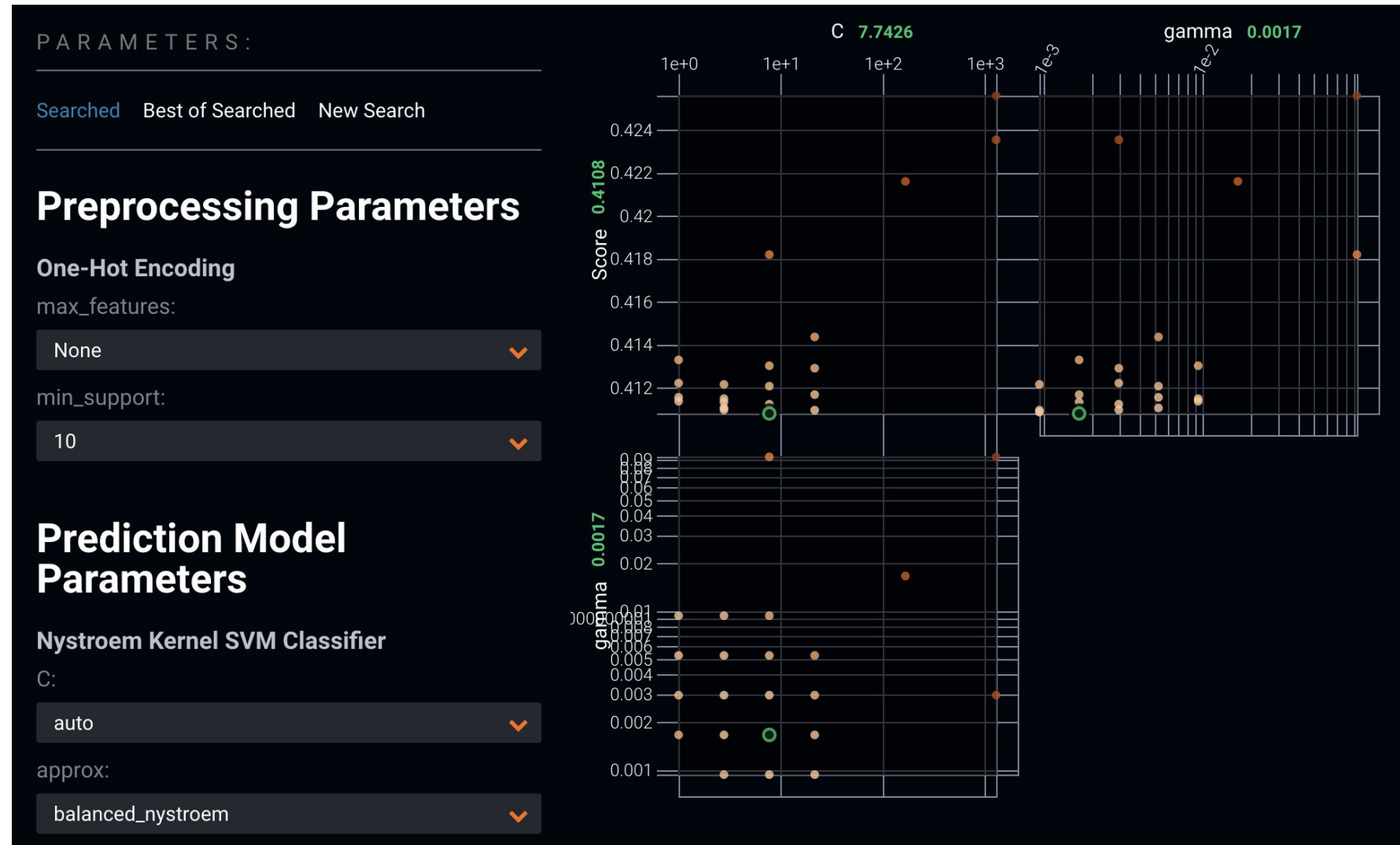- The predictions are reliablable to business people?

DataRobot

# Benefits : effective blending

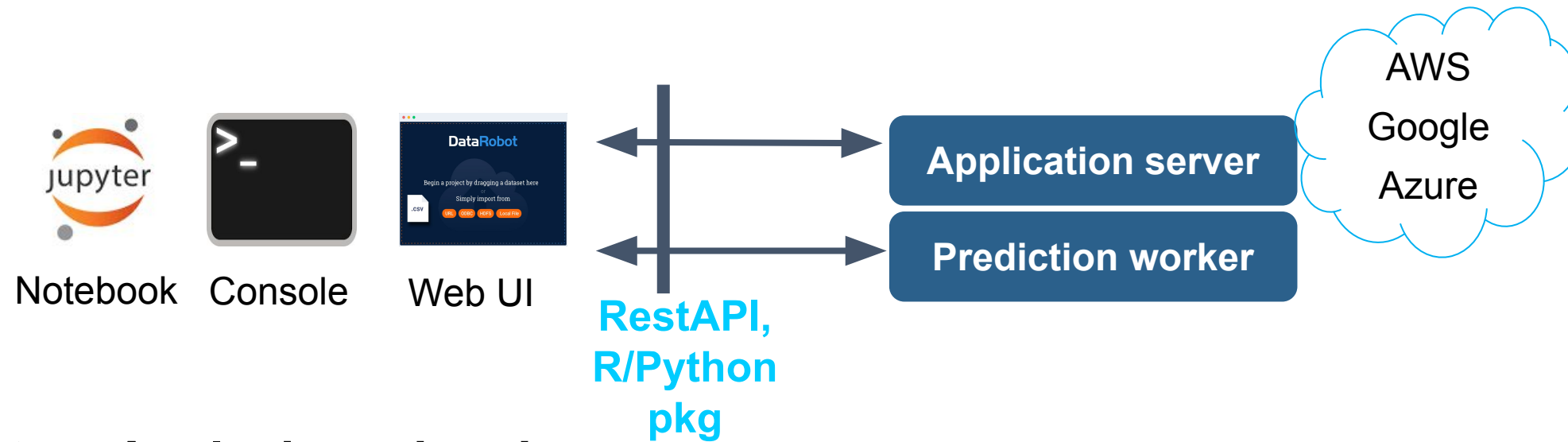Search over candidates which promises tangible improvement

# Benefits : Hyper-param Tuning

## Gradient-free and effective pattern search

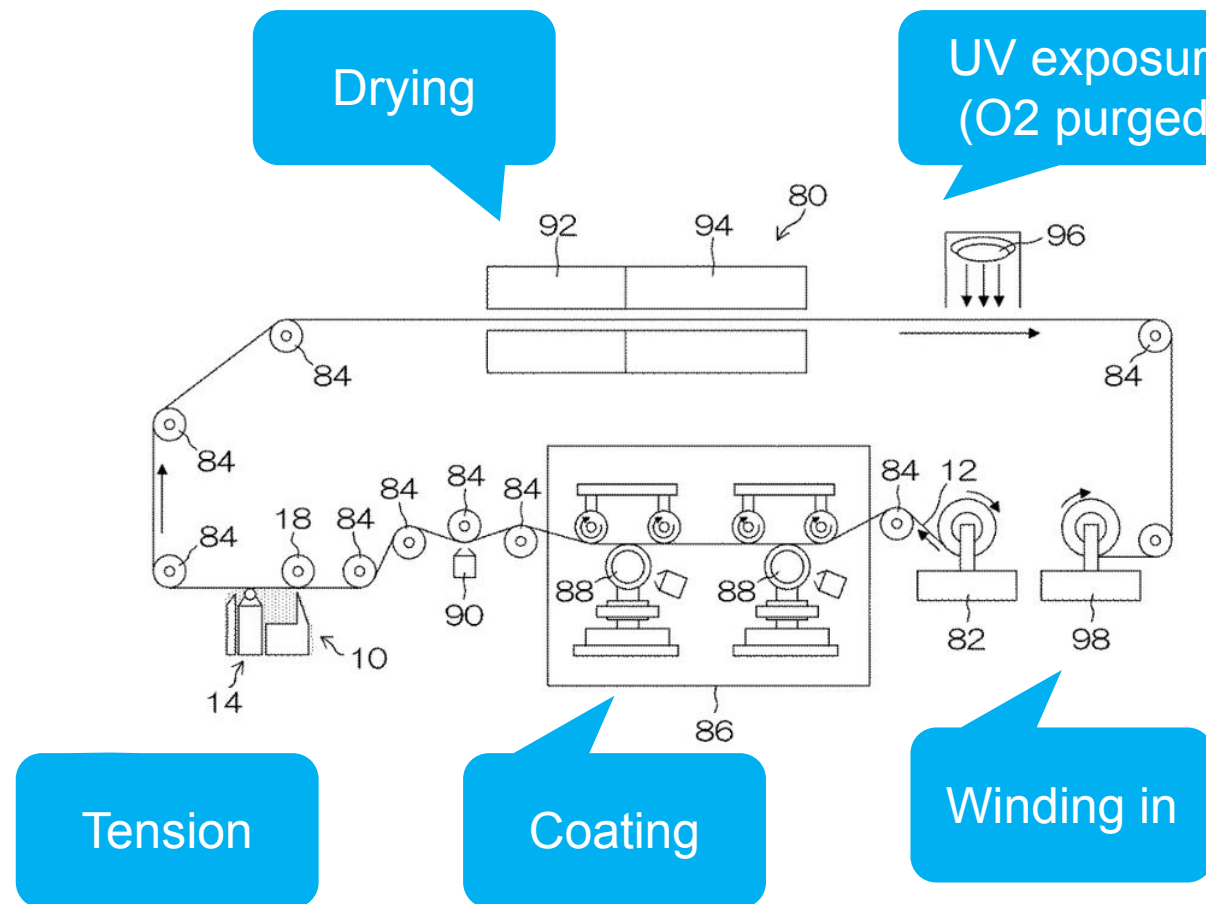# Benefits : API integration

data scientists and developers can use API



Notebook    Console    Web UI

**RestAPI, R/Python pkg**

**Application server**

**Prediction worker**

AWS
Google
Azure

<u>**Custom Analysis and various Analysis**</u>

| **Model Factory** | **Automatic Model Refresh** | **Model Diags & Viz** | **Feature Engineering** | **App. Integration** |

**DataRobot**

# Demo : Bleedout prediction

## Binary classification for QA



**Process:** Coating of thin film by covering the surface with coating solution and drying, followed by polymerizing with UV-light.

**Problem:** Unintended precipitation of powder such as unpolymerized monomer, antioxidant occurs causing "bleedout". It spoils the product and contaminates the production line.

**Data:**
- Material: length of film roll
- Project type: production vs experiment
- Control: winding tension, UV-exposure duration, O2 concentration etc

**DataRobot**

# Demo : Bleedout MFG process

1) Unwinding



2) Coating



3) Drying



4) UV exposure



5) Tension Control



6) Winding



**DataRobot**