

Fifty shades of....

*Automated
Machine Learning*

Rui Quintino, Data Research @ DevScope
Dewan Fayzur, Data Scientist @ DevScope



[DSPT #27, 20180417](#)

devscope

Rui Quintino

Data R&D @ DevScope

#PowerBI #SQLServer #Web
#Analytics #Azure #Microsoft

#MachineLearning #AutoML
#R #Linux #Dataiku #Docker
#Python #Coaching #Learning

also very often a #DataSkeptic ☺

twitter.com/rquintino

rquintino.wordpress.com

rquintino@gmail.com



“jack of all trades (and master of none)”

1. a person who can do many different types of work but who is not (necessarily...) very competent at any of them...

SmartDocumentor ReviewStation - version 3.1.315.0

FISCALNOVA (0)

No templates used

Review

Document Properties

Documentos contabilísticos

Validation OK

Classe Documental: Selecione uma classe documental

Fornecedor: Num. Contribuinte 505207583

Fatura

- Nº Documento: 14192
- Data Documento: 2014-09-04
- Data Vencimento: 2014-09-04
- Prazo Pagamento: Pronto Pagamen
- Base Incidência IVA: 227.64

Taxas IVA

Taxa	Base	Valor
Taxa 1	0.23	227.64
Taxa 2		52.36

QUINTA DE S. JOSÉ

João Brito e Cunha, Lda
Rua Augusto César, 99
5000-591 Vila Real
Telephone: 259 325 147
Fax: 259 325 147
Email: joao.britoecunha@quintasjose.com
URL: www.quintasjose.com

Capital Social: 10.000,00 €
C. R. C. Vila Real n.º: 505207583
NIF/VAT: 505207583

RECOBRO DE CLIENTE

Nº Documento: 14192
Data: 2014-09-04
ORIGINAL

DevScope, SA
Rua Passos Manuel 223, 4º
4000-385 PORTO

DIGITS smartdoc-sgennl-model-digits5 Test One

Infer One Image

Job Status Done

- Initialized at 12:55:58 AM (1 second)
- Running at 12:55:59 AM (3 seconds)
- Done at 12:56:02 AM (Total - 4 seconds)

Infer Model Done

Source image **Inference visualization**

Generic Image Model

devscope

Getting Value from Machine Learning Isn't About Fancier Algorithms — It's About Making It Easier to Use

by [Ben Schreck](#), [Max Kanter](#), [Kalyan Veeramachaneni](#), [Sanjeev Vohra](#), and [Rajendra Prasad](#)

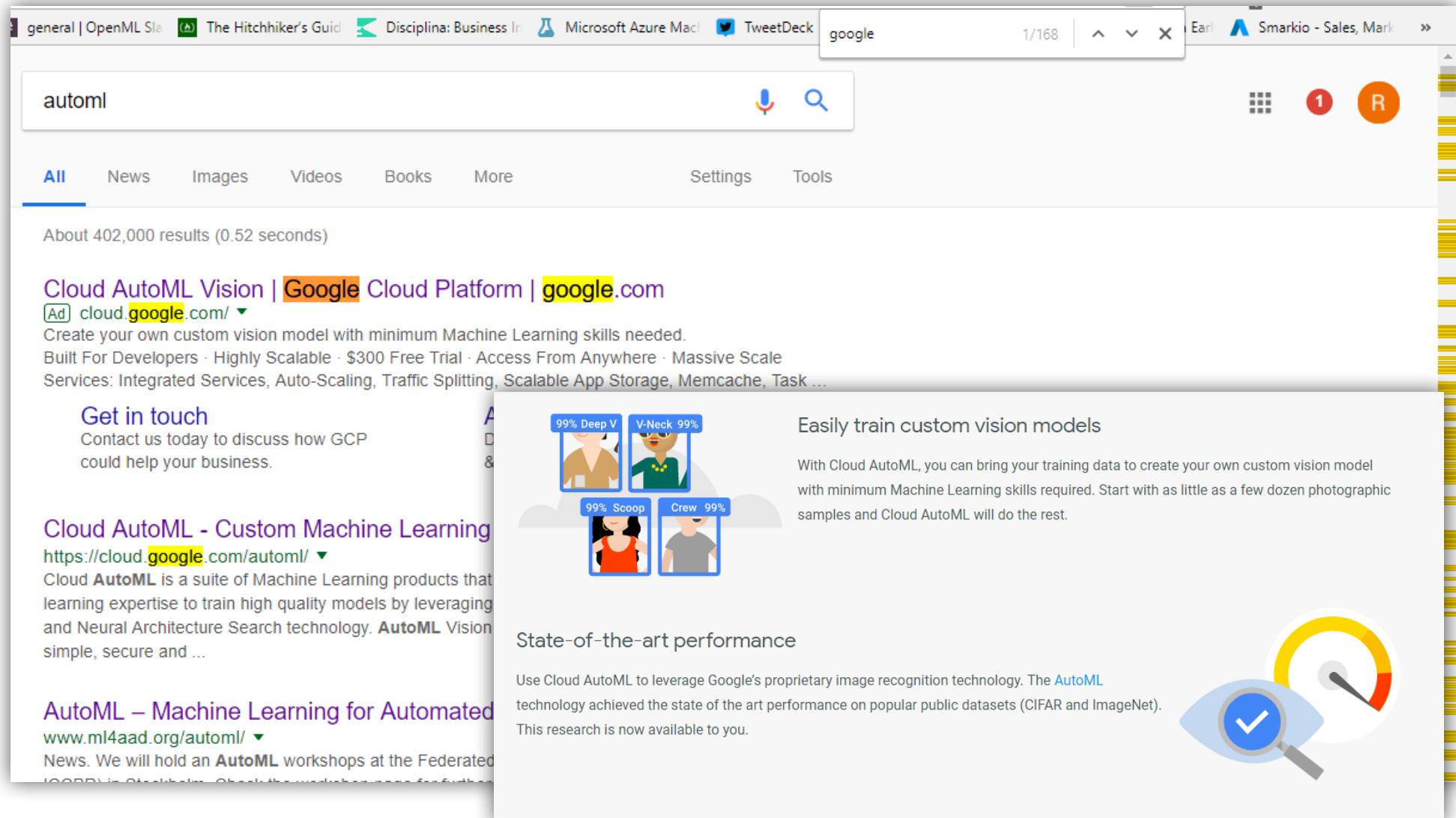
MARCH 06, 2018

Increasing ROI with AutoML

- Productivity (less time/less errors/more learning)
- Discovery
- Democratization/Massification
- Better Quantify Uncertainty
- Optimization (better results)
- -> Free more time to higher level tasks!

Tiny AutoML R&D Tip...



general | OpenML Slides | The Hitchhiker's Guide | Disciplina: Business Law | Microsoft Azure Marketplace | TweetDeck | google | 1/168 | 

automl

All News Images Videos Books More Settings Tools

About 402,000 results (0.52 seconds)

Cloud AutoML Vision | Google Cloud Platform | google.com

Ad [cloud.google.com/](https://cloud.google.com/automl/) ▾

Create your own custom vision model with minimum Machine Learning skills needed.

Built For Developers · Highly Scalable · \$300 Free Trial · Access From Anywhere · Massive Scale

Services: Integrated Services, Auto-Scaling, Traffic Splitting, Scalable App Storage, Memcache, Task ...

Get in touch
Contact us today to discuss how GCP could help your business.

Cloud AutoML - Custom Machine Learning
<https://cloud.google.com/automl/> ▾

Cloud AutoML is a suite of Machine Learning products that bring learning expertise to train high quality models by leveraging and Neural Architecture Search technology. AutoML Vision is simple, secure and ...

AutoML – Machine Learning for Automated
www.ml4aad.org/automl/ ▾

News. We will hold an AutoML workshops at the Federated ICML 2019. Check them out for further information.

Easily train custom vision models

With Cloud AutoML, you can bring your training data to create your own custom vision model with minimum Machine Learning skills required. Start with as little as a few dozen photographic samples and Cloud AutoML will do the rest.

State-of-the-art performance

Use Cloud AutoML to leverage Google's proprietary image recognition technology. The AutoML technology achieved the state of the art performance on popular public datasets (CIFAR and ImageNet). This research is now available to you.



automl -google

All News Images Videos Books More Settings Tools

About 226,000 results (0.58 seconds)

Cloud AutoML Vision | Google Cloud
Ad cloud.google.com/ ▾
Create your own custom vision model with minimum effort.
Built For Developers · Deploy At Google Scale · Machine Learning Services: Integrated Services, Auto-Scaling, Traffic

All Products
Discover Compute, Storage, Big Data & More Products On Google's Cloud

Information about Automated Machine Learning
[https://automl.info/ ▾](https://automl.info/)
Information about automated machine learning (AutoML)

AutoML: Automatic Machine Learning
<https://docs.h2o.ai/h2o/latest-stable/h2o-docs/automl/>
H2O's **AutoML** can be used for automating the machine learning process by automatically training and tuning of many models within a user-specified time limit. It can automatically train on collections of individual machine learning models which, in most cases, ...

Auto M L - Reparações Gerais de Automóveis, Lda - Portugalio
[https://www.portugalio.com/auto-m-l-reparacoes-gerais-de-automo... ▾](https://www.portugalio.com/auto-m-l-reparacoes-gerais-de-automo...) Translate this page
Auto M L - Telefone (2558153...), Fax, empresa situada em Lustosa, Lousada, Porto. Manutenção e reparação dos veículos automóveis. Automóveis.

Запознаване с AutoML - Софтуерен университет - SoftUni
[https://softuni.bg/forum/16064/zapoznavane-s-automl ▾](https://softuni.bg/forum/16064/zapoznavane-s-automl) Translate this page
Jun 5, 2017 - 1 post
Запознаване с AutoML. Очевидно DL се развива с бесни темпове и вече има AutoML -
<https://themerkle.com/googles-ai-is-creating-ai-and-its-better-than-company-engineers-at-it/> Въпросът ми е по - скоро молба. Ще можем ли да вместим в програмата на курса и кратко запознаване с AutoML?

AutoML 2015 workshop @ ICML 2015 (11 July 2015) · LAL Events ...
<https://indico.lal.in2p3.fr/event/2914/sessions/1168/>
18:00. Joaquin Vanschoren. Invited Talk: OpenML: A Foundation for Networked & Automatic Machine Learning. 16:30 - 17:10. Lille Grand Palais. Marc Boule. **AutoML** Challenge. 17:10 - 17:30. Lille Grand Palais. Panel Discussion: Next steps for **AutoML**. 17:30 - 18:00. Lille Grand Palais. Updating the timetable... Indico.

Automated Machine Learning (in this session)

- (for Data Scientists/ML Experts)
- Common tasks, Repetitive tasks (free us)
- Search tasks/ Even “Reasoning” (augment us)
- *Ps-Not Automated Data Science*

A few shades of AutoML...



Machine Learning GUIs - Productivity



```
>docker run -it -p 10000:10000 dataiku/dss
```

The screenshot shows the Dataiku DSS interface with the following details:

- Project:** titanic
- Analysis:** Insights into Survived for titanic_train
- Script Tab:** Active (highlighted in blue)
- Summary, Charts, Models:** Other tabs available
- Actions:** Deploy Script button and Actions dropdown
- Display:** Options for grid, list, and chart
- Search:** Search bar at the top left
- Rows:** 891 matching rows
- Table View:** Shows the structure of the titanic dataset with 13 columns:
 - PassengerId (Integer)
 - Pclass (Integer)
 - Name (Natural lang.)
 - Age (Decimal)
 - SibSp (Integer)
 - Parch (Integer)
 - Fare (Decimal)
 - Embarked (Text)
 - Survived (Integer)
- Data Preview:** A small preview of the first few rows is visible.

titanic ➔ 📈 ⚙️ 🔍 ANALYSES

INSIGHTS Insights into Survived for titanic_train 🔄

Predict Survived (Binary classification) 🛡️ DESIGN RESULT SAVED TRAIN

Target

Train / Test Set

Python environment

FEATURES

Features handling

Feature generation

Feature reduction

MODELING

Algorithms

Hyperparameters

EVALUATION

Metric

Algorithms

Algorithms

- Random Forest
- Gradient tree boosting
- Logistic Regression
- XGBoost
- Decision Tree
- Support Vector Machine
- Stochastic Gradient Descent
- KNN
- Extra Random Trees
- Neural Network
- Lasso Path

Random Forest

A Random Forest is made of many decision trees. Each tree in the forest predicts a record, and each tree "votes" for the final answer of the forest.

Show more...

Numbers of trees Number of trees in the forest.

Feature sampling strategy Adjusts the number of features to sample at each split.

Maximum depth of tree Maximum depth of each tree in the forest. Higher values generally increase the quality of the prediction, but can lead to overfitting. High values also increase the training and prediction time. Use 0 for unlimited depth (ie, keep splitting the tree until each node contains a single target value).

Minimum samples per leaf Minimum number of samples required in a single tree node to split this node. Lower values increase the quality of the prediction (by splitting the tree more), but can lead to overfitting and increased training and prediction time.

Parallelism

titanic ANALYSES

Insights into Survived for titanic_train

Predict Survived (Binary classification)

DESIGN **RESULT**

SAVED **TRAIN**

Target

Train / Test Set

Python environment

FEATURES

- Features handling** (selected)
- Feature generation
- Feature reduction

MODELING

- Algorithms
- Hyperparameters

EVALUATION

- Metric

Features Handling

Dataset: Filter

Feature	Description	Status
A Gender	Dummy-encode	ON
A Family	Dummy-encode	ON
# PassengerId	Reject	OFF
A Pclass	Dummy-encode , impute missing	ON
I Name	Reject	OFF
# Age	Avg-std rescaling	ON
# SibSp	Avg-std rescaling	ON
# Parch	Avg-std rescaling	ON
# Fare		ON

Handling of "Gender"

Role: Input

Variable type: A Categorical

Category handling: Dummy-encoding (vectorization)

Drop dummy: Let DSS decide

Clipping: Max nb. categories

Max. Nb. Categories: 100

Missing values: Treat as a regular value

2 distinct values, with 0.0% empty cells

Gender	Percentage
Men	64.9%
Women	35.1%

titanic ANALYSES

INSIGHTS

Insights into Survived for titanic_train

Predict Survived (Binary classification)

DESIGN RESULT ACTIONS SAVED TRAIN

Target

Train / Test Set

Python environment

FEATURES

Features handling

Feature generation

Feature reduction

MODELING

Algorithms

Hyperparameters

EVALUATION

Metric

Feature generation

FEATURE INTERACTIONS

This will generate interactions between features:

- Numerical features will be multiplied
- Numerical and categorical features will produce a dummies multiplied by the numerical feature.
- Two categorical features will produce dummies in the cross-product of the two features

+ ADD INTERACTION

AUTOMATIC GENERATION OF NUMERICAL FEATURES

Pairwise linear combinations Generates A+B and A-B for pairs of numerical features

Polynomial combinations Generates A*B for pairs of numerical features

The screenshot shows a software interface for analyzing the 'titanic' dataset. The main window has tabs for DESIGN and RESULT, with DESIGN currently selected. On the left, there's a sidebar with categories like FEATURES, MODELING, and EVALUATION, each with sub-options. The FEATURES category has 'Feature generation' selected. The central area displays the 'Feature generation' configuration, which includes a 'FEATURE INTERACTIONS' section describing how it generates interactions between features (multiplication of numerical and categorical features) and a 'AUTOMATIC GENERATION OF NUMERICAL FEATURES' section with options for pairwise linear combinations and polynomial combinations. A large blue button labeled '+ ADD INTERACTION' is visible. The top navigation bar includes icons for file operations and a search function, along with tabs for Summary, Script, Charts, and Models.



▼ Predict Survived (Binary classification)

DESIGN

RESULT

SAVED

TRAIN



Search...

Filter

Metric: ROC AUC



SESSIONS

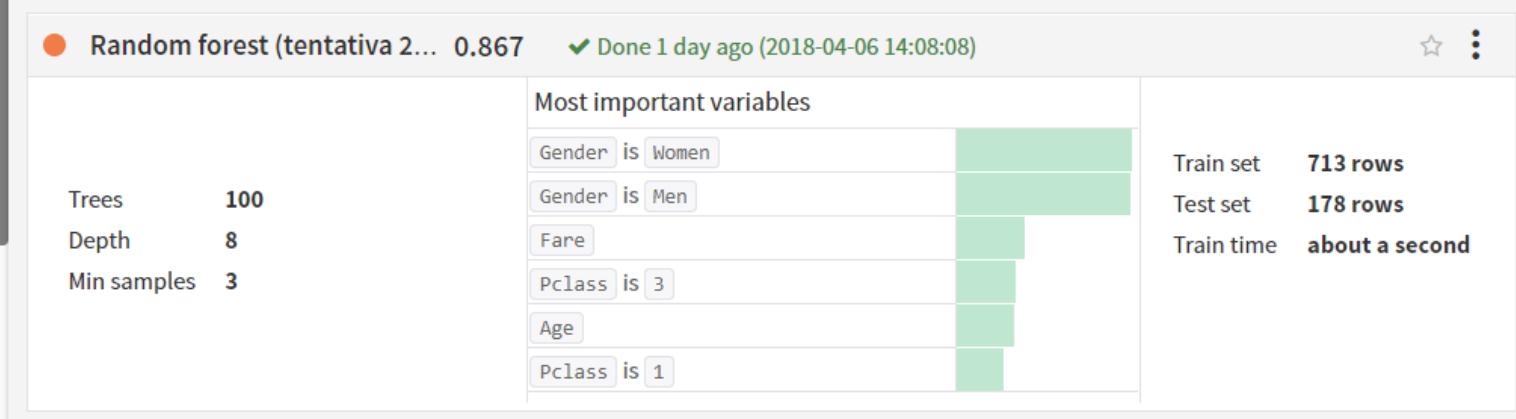
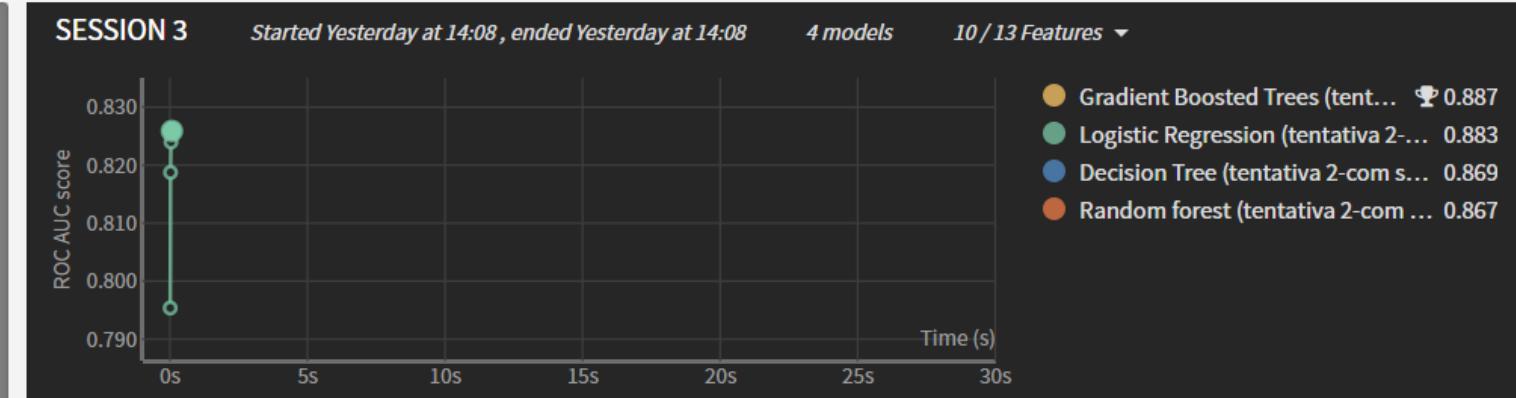
MODELS

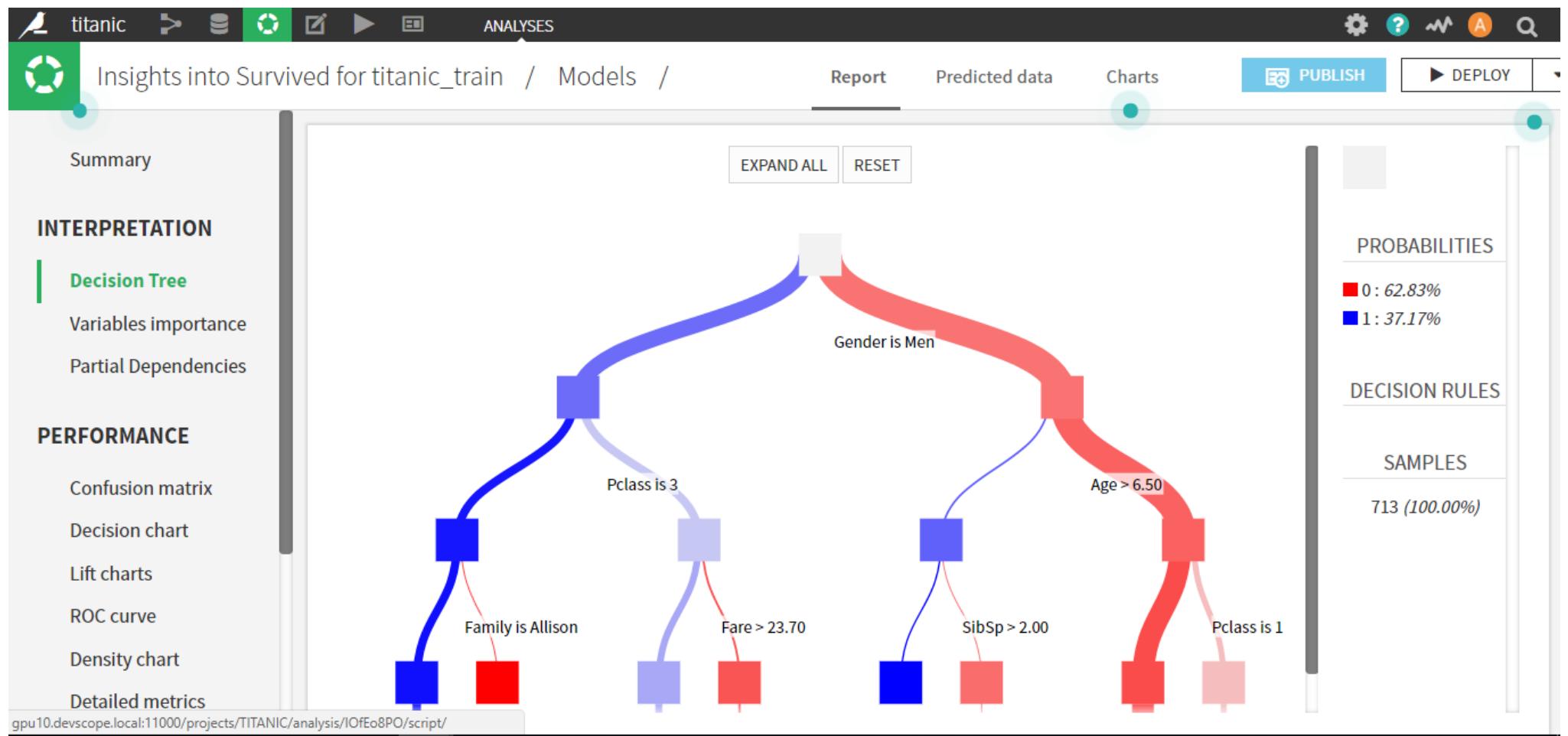
TABLE

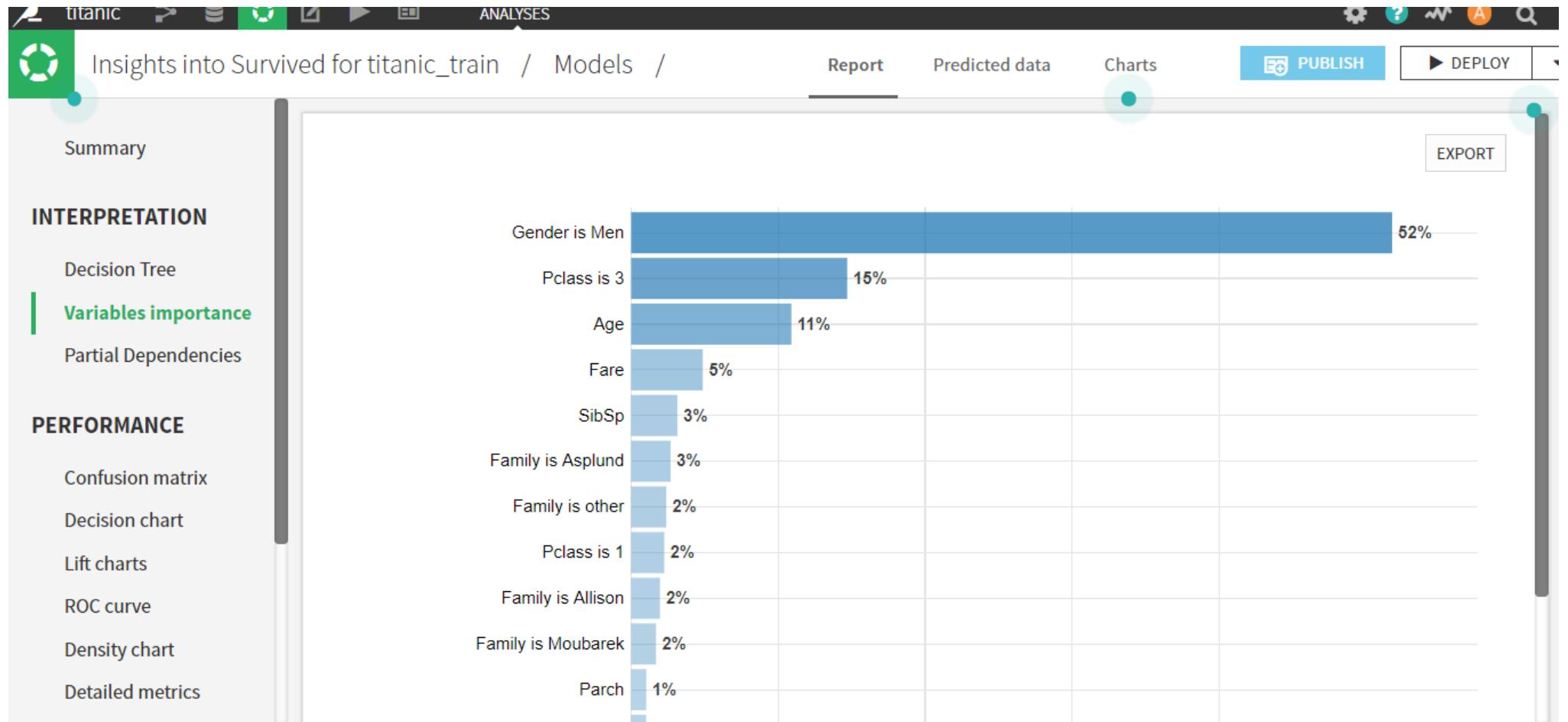
SESSION 3		
	Random forest ...	0.867
	Gradient Bo...	0.887
	Logistic Regres...	0.883
	Decision Tree (t...	0.869

SESSION 2		
	Random for...	0.791
	Gradient Boost...	0.774
	Logistic Regres...	0.760
	Decision Tree	0.777

SESSION 1		
	Random for...	0.805







Insights into Survived for titanic_train / Models /

Report Predicted data Charts PUBLISH DEPLOY

Summary

INTERPRETATION

- Decision Trees
- Variables importance
- Partial Dependencies**

PERFORMANCE

- Confusion matrix
- Decision chart
- Lift charts
- ROC curve
- Density chart
- Detailed metrics

MODEL INFORMATION

- Data preparation
- Features

Select a feature : Age

Partial Dependency

A partial dependency plot showing the relationship between Age (X-axis, ranging from 0.42 to 80) and Partial dependency (Y-axis, ranging from -1.27 to 1.99). The plot shows a piecewise constant function that starts at approximately 1.99 for very young ages, drops sharply to about 0.5 around age 10, then to near zero at age 18, and remains relatively flat until age 35. It then fluctuates between -0.2 and -0.6 until age 45, where it drops sharply again to approximately -1.27. This low value is maintained until age 62, where it rises sharply to about -0.5, and finally levels off at approximately -0.3 for older ages.

Reading tips

A partial dependency plot shows the dependence of the predicted response on a single feature. The x axis displays the value of the selected feature, while the y axis displays the partial dependence.

The value of the partial dependence is by how much the log-odds are higher or lower than those of the average probability.

Note : the log-odds for a probability p are defined as $\log(p / (1 - p))$. They are strictly increasing, ie. higher log odds mean higher probability.

titanic ► ▶ NOTEBOOKS WEB APPS LIBRARIES RMARKDOWN REPORTS

Predict Survived in titanic_train

jupyter Predict Survived in titanic_train (unsaved changes)

File Edit View Insert Cell Kernel Help Python 2

building tree 85 of 100
building tree 86 of 100
building tree 87 of 100
building tree 88 of 100
building tree 89 of 100
building tree 90 of 100
building tree 91 of 100
building tree 92 of 100
building tree 93 of 100
building tree 94 of 100building tree 95 of 100

building tree 96 of 100
building tree 97 of 100
building tree 98 of 100
building tree 99 of 100
building tree 100 of 100
CPU times: user 252 ms, sys: 28 ms, total: 280 ms
Wall time: 275 ms

[Parallel(n_jobs=4)]: Done 33 tasks | elapsed: 0.0s
[Parallel(n_jobs=4)]: Done 100 out of 100 | elapsed: 0.1s finished

Out[13]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
max_depth=8, max_features='auto', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=3, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=4,
oob score=False, random state=1337, verbose=2,

Model/Pipeline Selection, Hyper-Parameter Tuning



H2O AutoML (R/Python)

r python

```
library(h2o)

h2o.init()

# Import a sample binary outcome train/test set into H2O
train <- h2o.importFile("https://s3.amazonaws.com/erin-data/higgs/
test <- h2o.importFile("https://s3.amazonaws.com/erin-data/higgs/

# Identify predictors and response
y <- "response"
x <- setdiff(names(train), y)

# For binary classification, response should be a factor
train[,y] <- as.factor(train[,y])
test[,y] <- as.factor(test[,y])

aml <- h2o.automl(x = x, y = y,
                  training_frame = train,
                  max_runtime_secs = 30)
```

r python

```
import h2o
from h2o.automl import H2OAutoML

h2o.init()

# Import a sample binary outcome train/test set into H2O
train = h2o.import_file("https://s3.amazonaws.com/erin-data/higgs/higgs_train_10k.csv")
test = h2o.import_file("https://s3.amazonaws.com/erin-data/higgs/higgs_test_5k.csv")

# Identify predictors and response
x = train.columns
y = "response"
x.remove(y)

# For binary classification, response should be a factor
train[y] = train[y].asfactor()
test[y] = test[y].asfactor()

# Run AutoML for 30 seconds
aml = H2OAutoML(max_runtime_secs = 30)
aml.train(x = x, y = y,
```

H2O AutoML (R/Python)

```
# View the AutoML Leaderboard
lb <- aml@leaderboard
lb
```

#	model_id	auc	logloss
	StackedEnsemble_AllModels_0_AutoML_20171121_012135	0.788321	0.554019
	StackedEnsemble_BestOfFamily_0_AutoML_20171121_012135	0.783099	0.559286
	GBM_grid_0_AutoML_20171121_012135_model_1	0.780554	0.560248
	GBM_grid_0_AutoML_20171121_012135_model_0	0.779713	0.562142
	GBM_grid_0_AutoML_20171121_012135_model_2	0.776206	0.564970
	GBM_grid_0_AutoML_20171121_012135_model_3	0.771026	0.570270
	DRF_0_AutoML_20171121_012135	0.734653	0.601520
	XRT_0_AutoML_20171121_012135	0.730457	0.611706
	GBM_grid_0_AutoML_20171121_012135_model_4	0.727098	0.666513
	GLM_grid_0_AutoML_20171121_012135_model_0	0.685211	0.635138

H2O Flow UI / Leaderboard

The screenshot displays two panels from the H2O Flow UI.

Left Panel (Job):

- CS runAutoML {"training_frame": "higgs_train_10k.hex", "response_column": "response", "seed": -1, "max_models": 0, "max_runtime_secs": 30, "stopping_metric": "AUTO", "stopping_rounds": 3, "stopping_tolerance": 0.001}
- Job**
 - Run Time: 00:00:07.110
 - Remaining Time: 00:00:58.723
 - Type: Auto Model
 - Key: Q_AutoML_20170802_084713
 - Description: AutoML build
 - Status: RUNNING
 - Progress: 11%
 - Default Extremely Random T
 - Actions: [View](#) [Cancel Job](#)

Right Panel (Leaderboard):

- CS getLeaderboard "AutoML_20170608_075028"
- 60ms
- Leaderboard**
 - [Monitor Live](#)
 - MODELS**

models sorted in order of mean_residual_deviance, best first

model_id	mean_residual_deviance	rmse	mae	rmsle
0 StackedEnsemble_0_AutoML_20170608_075028	0.190473	0.436432	0.383892	0.307691
1 XRT_0_AutoML_20170608_075028	0.196350	0.443114	0.401439	0.313125
2 DRF_0_AutoML_20170608_075028	0.199015	0.446111	0.404586	0.314488
3 GBM_grid_0_AutoML_20170608_075028_model_0	0.209251	0.457440	0.390030	0.321598
4 GBM_grid_0_AutoML_20170608_075028_model_1	0.220931	0.470033	0.463941	0.332583
5 GLM_grid_0_AutoML_20170608_075028_model_0	0.223552	0.472812	0.450720	0.334355
6 GLM_grid_0_AutoML_20170608_075028_model_1	0.223552	0.472812	0.450720	0.334355
 - USER FEEDBACK**

Actions taken and discoveries made by AutoML

timestamp	level	stage	message
0 07:50:28.176	Info	Workflow	AutoML job created: 2017.06.08 07:50:28.173
1 07:50:28.301	Info	DataImport	Automatically split the training data into training, validation and leaderboard datasets in the ratio 0.70:0.15:0.15

Tpot (python)

- (*genetic search*)

```
from tpot import TPOTClassifier
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split

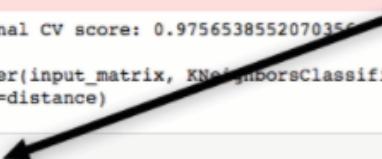
digits = load_digits()
X_train, X_test, y_train, y_test = train_test_split(digits.data, digits.target,
                                                    train_size=0.75, test_size=0.25)

tpot = TPOTClassifier(generations=5, population_size=50, verbosity=2, n_jobs=-1)
tpot.fit(X_train, y_train)

Optimization Progress: 33%|███████| 100/300 [01:02<09:07, 2.74s/pipeline]
Generation 1 - Current best internal CV score: 0.9644750872792087
Optimization Progress: 50%|███████| 150/300 [01:35<05:41, 2.27s/pipeline]
Generation 2 - Current best internal CV score: 0.9681323584103183
Optimization Progress: 67%|███████| 200/300 [01:59<01:59, 1.19s/pipeline]
Generation 3 - Current best internal CV score: 0.9718282518620386
Optimization Progress: 83%|███████| 250/300 [02:23<00:41, 1.21s/pipeline]
Generation 4 - Current best internal CV score: 0.9756538552070356
Optimization Progress: 99%|██████████| 299/300 [02:23<00:01, 0.00s/pipeline] ~99% accuracy on MNIST
Generation 5 - Current best internal CV score: 0.9756538552070356
Best pipeline: KNeighborsClassifier(input_matrix, KNeighborsClassifier__n_neighbors=10, KNeighborsClassifier__p=DEFAU
LT, KNeighborsClassifier__weights=distance)
print(tpot.score(X_test, y_test))
```

0.995555555556

~99% accuracy on MNIST
out of the box



Tpot (python)

- Export pipeline

should be exported to the `tpot_mnist_pipeline.py` file and look similar to the following:

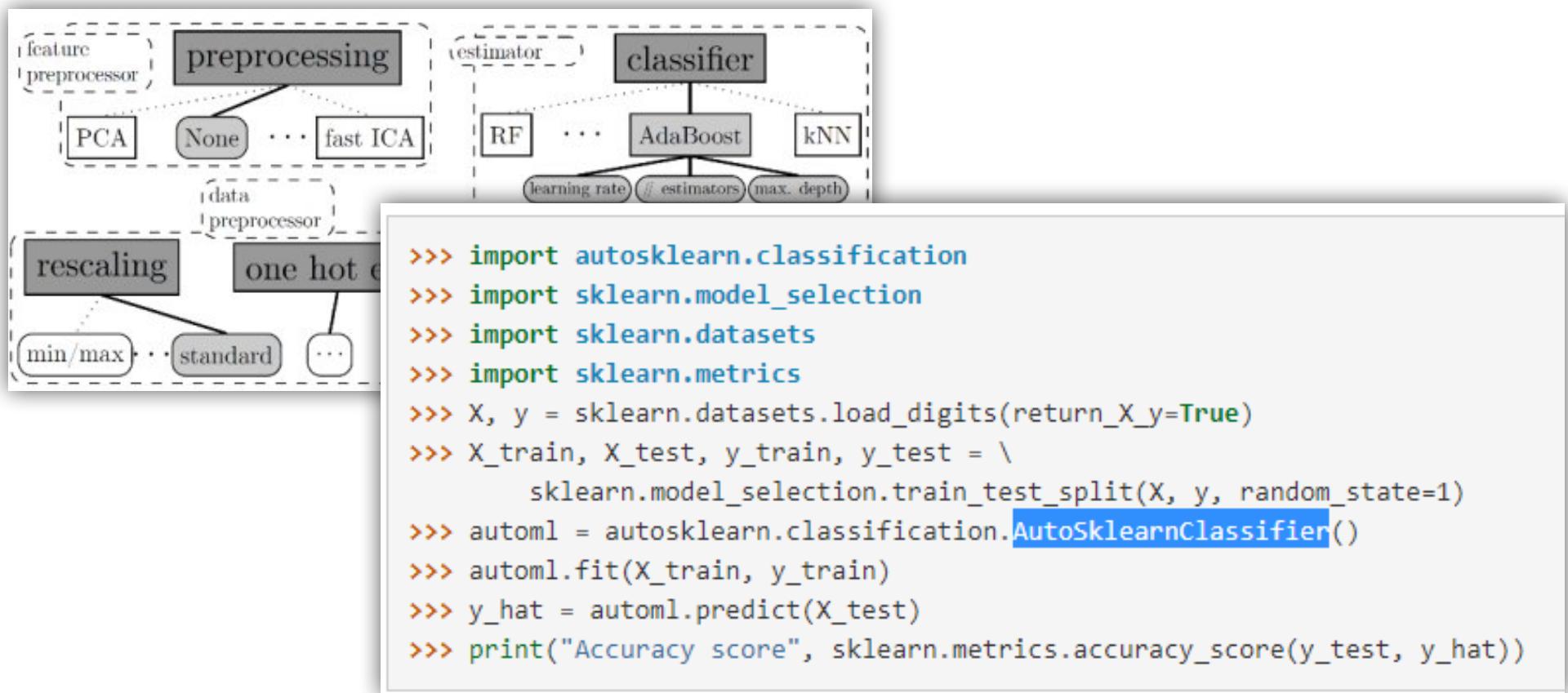
```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

# NOTE: Make sure that the class is labeled 'target' in the data file
tpot_data = pd.read_csv('PATH/TO/DATA/FILE', sep='COLUMN_SEPARATOR', dtype=np.float64)
features = tpot_data.drop('target', axis=1).values
training_features, testing_features, training_target, testing_target = \
    train_test_split(features, tpot_data['target'].values, random_state=42)

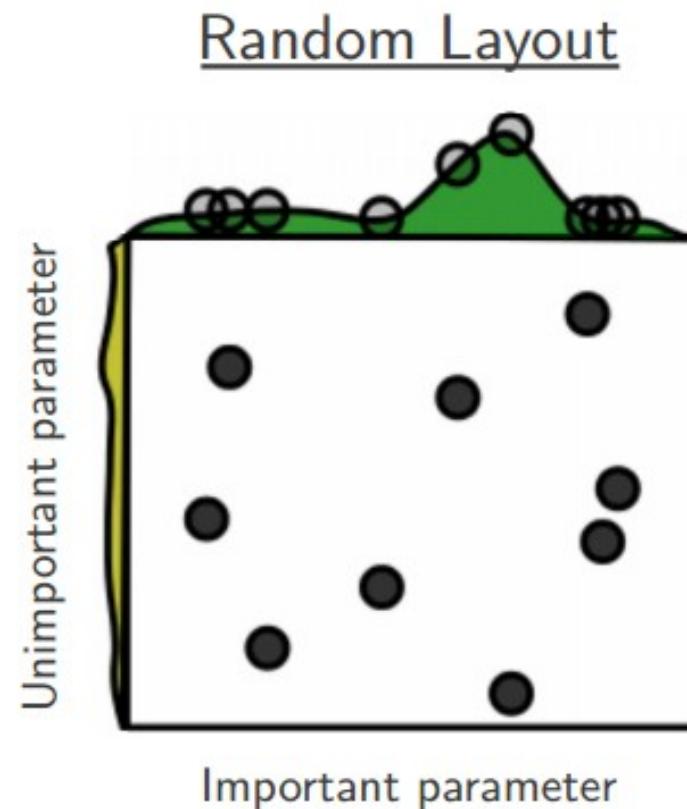
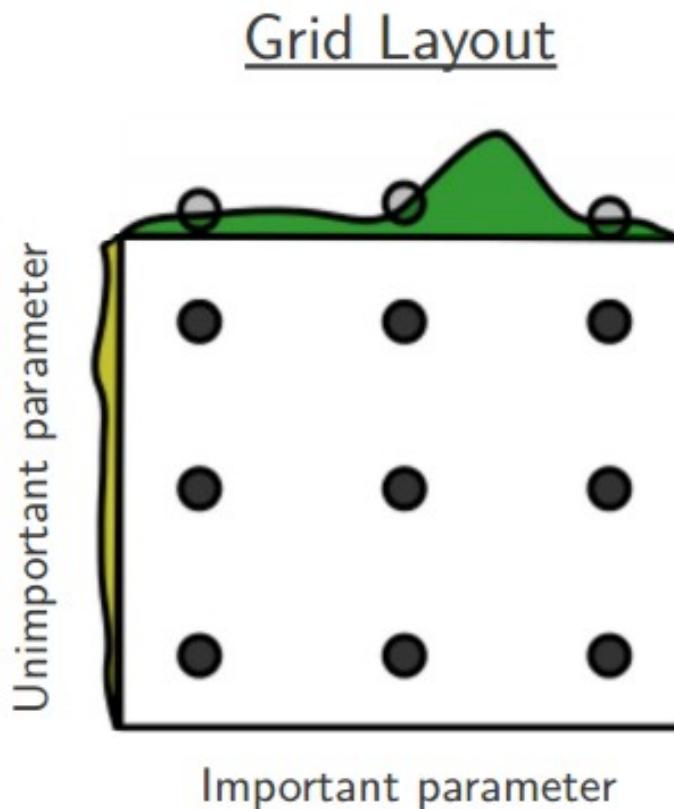
exported_pipeline = KNeighborsClassifier(n_neighbors=6, weights="distance")

exported_pipeline.fit(training_features, training_classes)
results = exported_pipeline.predict(testing_features)
```

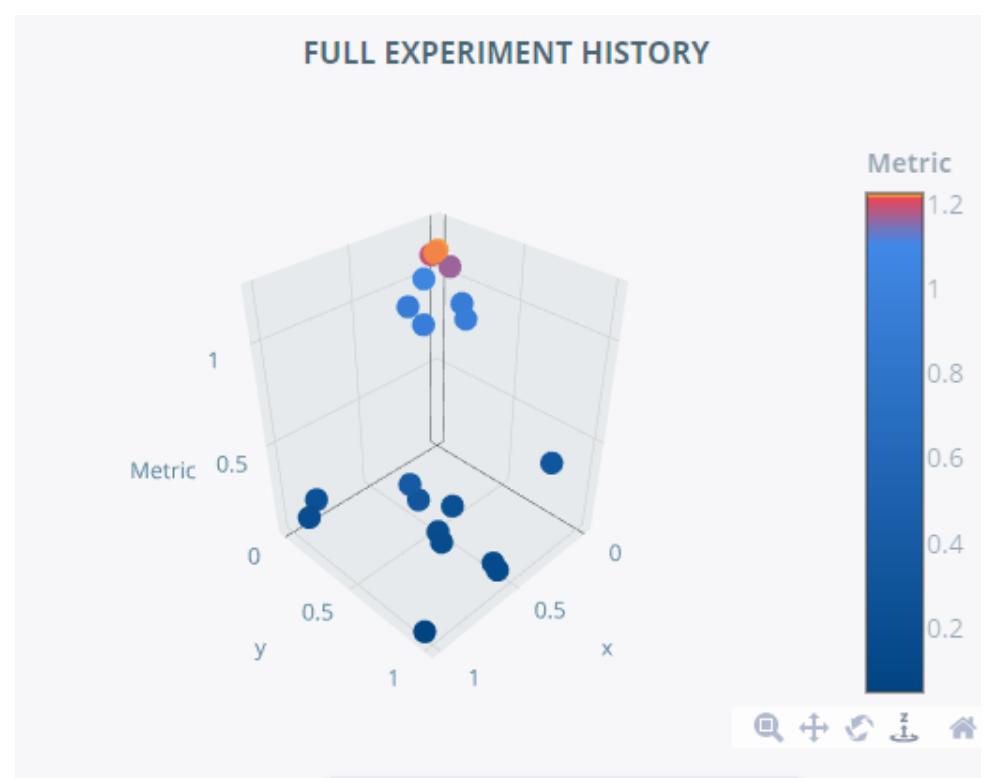
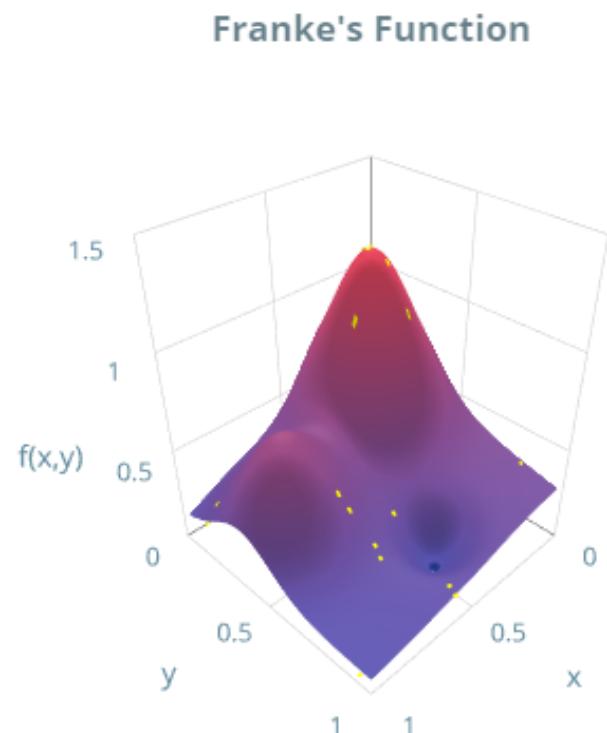
Auto-sklearn



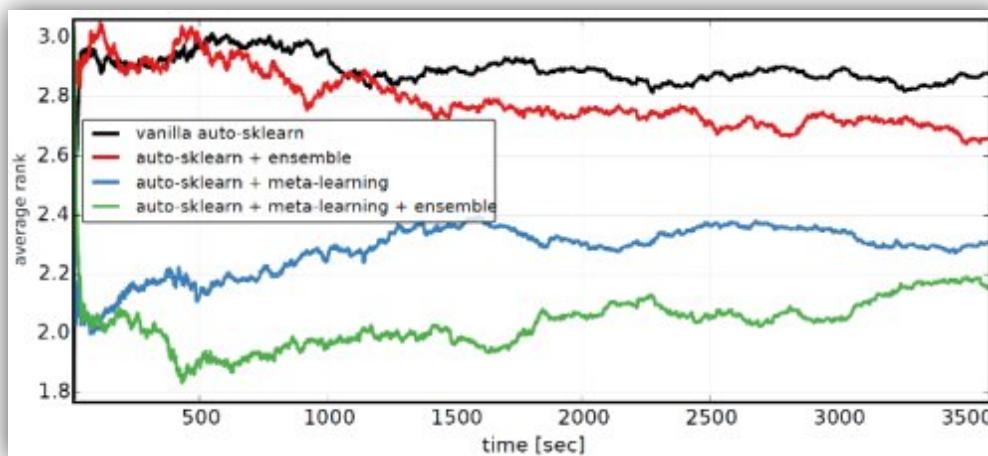
Search Strategies: Random & Grid



Search Strategies: Bayesian



Bayesian Search & Meta Learning – Fast Start!



The OpenML homepage features a central logo with the text "Machine learning, better, together" and "beta". Below the logo are four circular icons representing different data types and counts:

- 19999 data sets (green circle)
- 67122 tasks (orange circle)
- 5744 flows (blue circle)
- 9009219 runs (red circle)

Below each icon is a corresponding call-to-action text:

- Find or add **data** to analyse
- Download or create scientific **tasks**
- Find or add data analysis **flows**
- Upload and explore all **results** online.

Commercial AutoML Kaggle “Grandmasters”



H2O Driverless AI (Enterprise, Not O.S.)

- *Algorithms/Hyper Params*
- *Ensembles*
- *Feature Engineering*
- *Interpretability*
- *Fast (GPUs)*



TRAINING DATA

DATASET
creditcard.csv

ROWS COLUMNS DROPPED IGNORED
24K 25 -- --

TARGET COLUMN

default payment next month

TYPE COUNT MEAN STD DEV
int32 23999 0.2237 0.4167

EXPERIMENT SETTINGS



SCORER
R2
AUC
RMSE
MSE
MAE
LOGLOSS

CLASSIFICATION ONE SPLIT ONLY AUTO SEL. COLS

LAUNCH EXPERIMENT

H2O.ai

0.0.1

TRAINING DATA

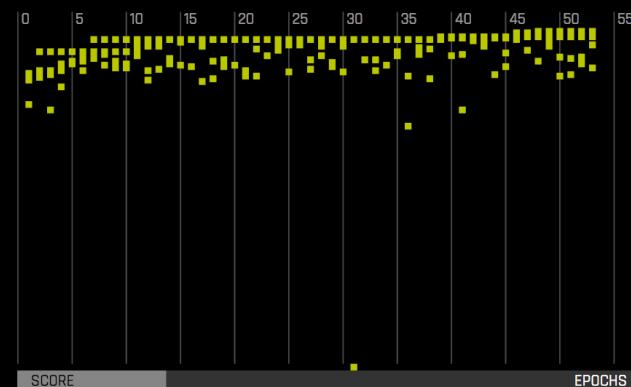
DATASET
creditcard.csvROWS
24KCOLUMNS
25DROPPED
--IGNORED
--

TARGET COLUMN

default payment next month

TYPE
int32COUNT
23999MEAN
0.2237STD DEV
0.4167

ITERATION SCORES



STATUS: COMPLETE

EXPERIMENT SETTINGS

1 WORKERS

53 ITERATIONS

3 CV FOLDS

4 POPULATION

INTERPRET THIS MODEL

CLASSIFICATION

DETECT IDS

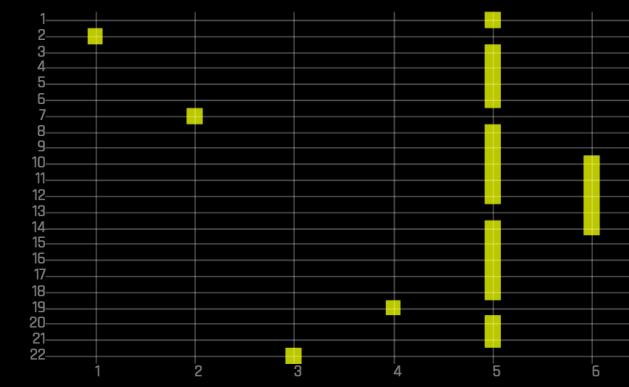
DROP DUPS.

GPU STATS

1

2

FEATURE TRANSFORMATIONS



VARIABLE IMPORTANCE

1_PAY_0	7497.51
2_PAY_2	1020.18
5_BILL_AMT1	888.59
0_LIMIT_BAL	826.45
11_PAY_AMT2	773.38
10_PAY_AMT1	684.03
20_CV_TE_SEX_EDUCATION_0	634.89
12_PAY_AMT3	503.78
13_PAY_AMT4	499.84
6_BILL_AMT2	428.59
9_BILL_AMT6	425.90
14_PAY_AMT5	418.09
15_PAY_AMT6	409.82
3_PAY_4	326.66

Single Row Lookup

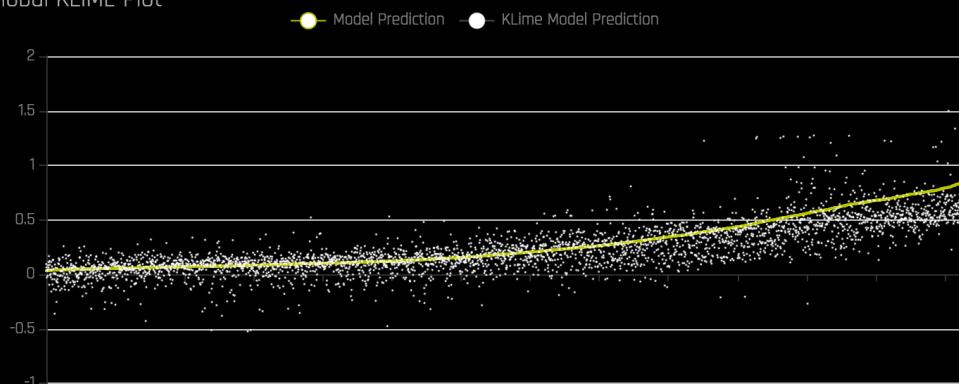
Column: H2O Frame Row # Value:

SEARCH

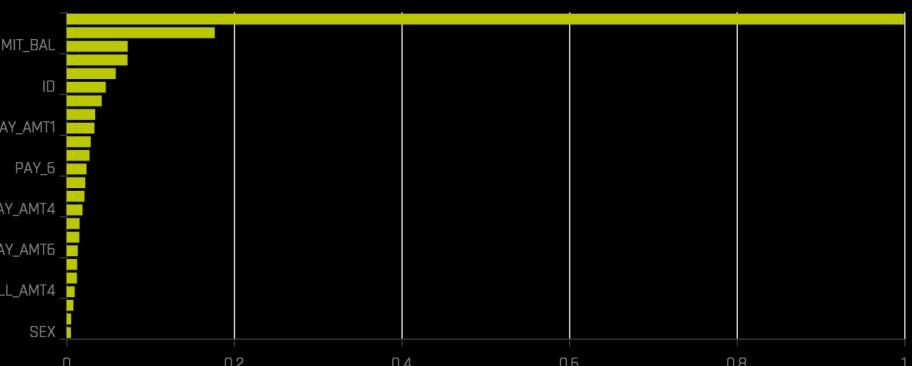
Plot: Global

EXPLANATIONS

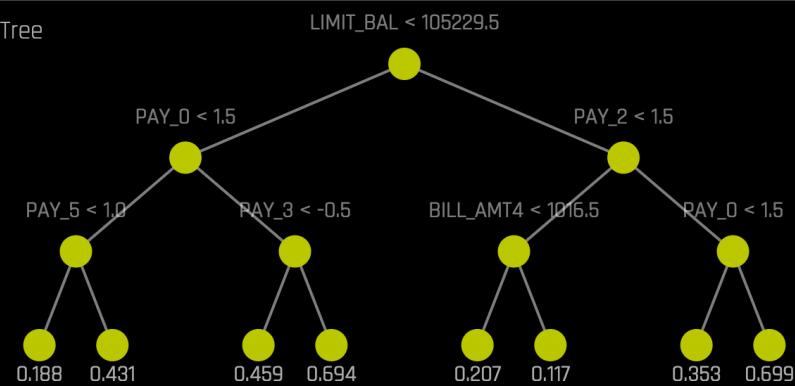
Global KLIME Plot



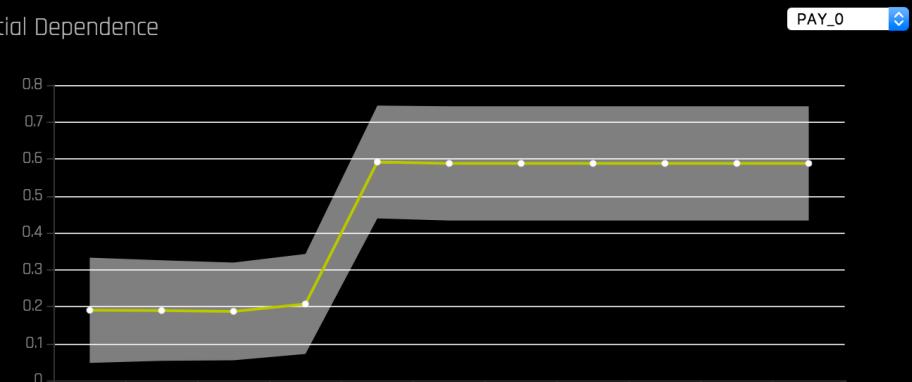
Variable Importance



Decision Tree

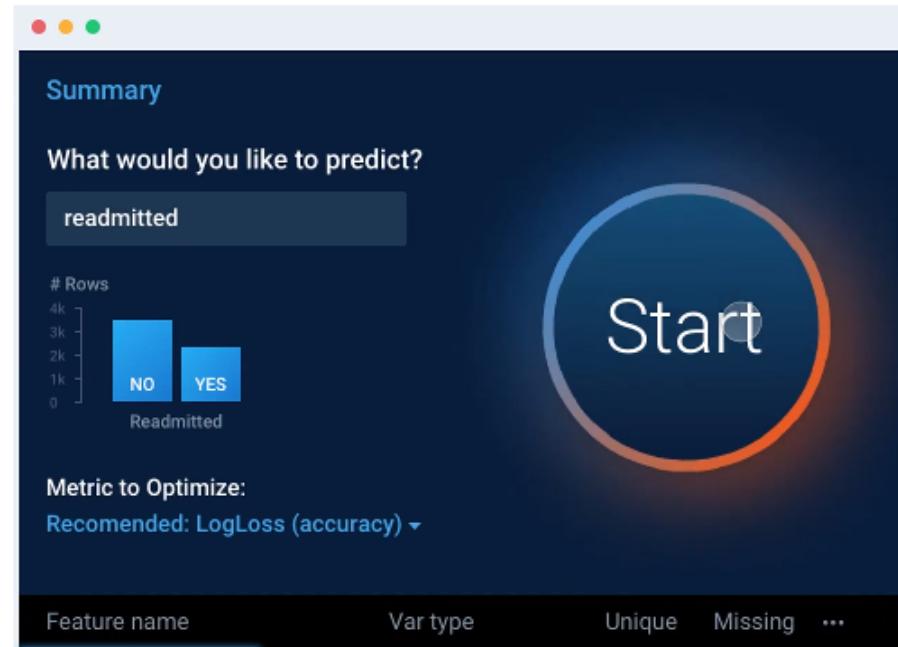


Partial Dependence



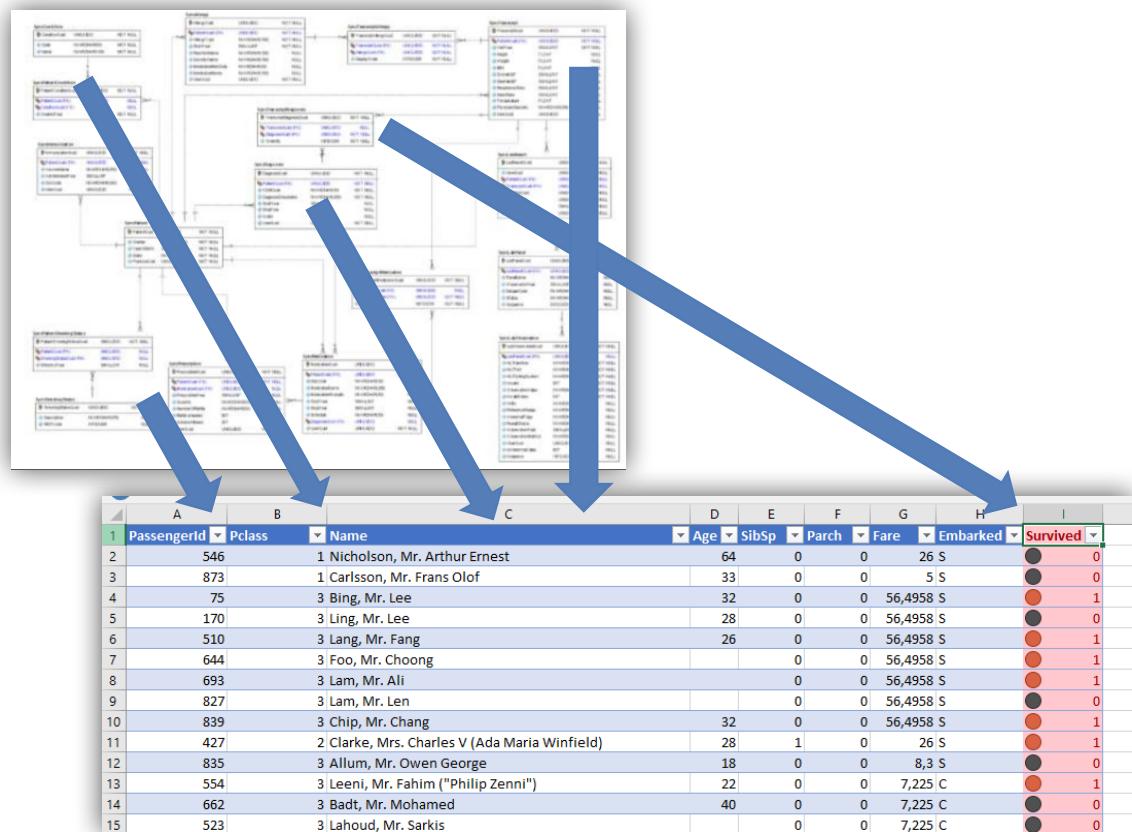
DataRobot

- ① Ingest your data
- ② Select the target variable
- ③ Build 100s of models in one click
- ④ Explore top models and get insights
- ⑤ Deploy best model and make predictions



What about Feature Discovery/Extraction?

- Machine Learning models only use a single table (with features & goal/label/target)

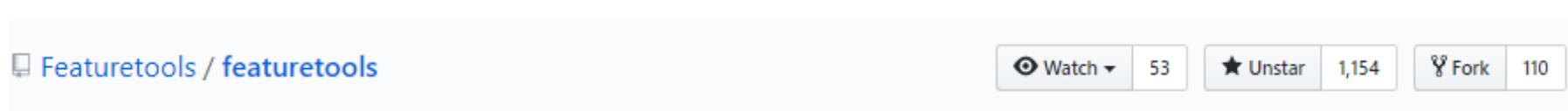


Getting Value from Machine Learning Isn't About Fancier Algorithms — It's About Making It Easier to Use

by [Ben Schreck](#), [Max Kanter](#), [Kalyan Veeramachaneni](#), [Sanjeev Vohra](#), and [Rajendra Prasad](#)

MARCH 06, 2018

Feature Tools & Deep Feature Synthesis (MIT/DARPA)



"One of the holy grails of machine learning is to automate more and more of the feature engineering process." — Pedro Domingos, [A Few Useful Things to Know about Machine Learning](#)

```
Entityset: transactions
  Entities:
    customers (shape = [5, 3])
    sessions (shape = [35, 4])
    products (shape = [5, 2])
    transactions (shape = [500, 5])
  Relationships:
    transactions.product_id -> products.product_id
    transactions.session_id -> sessions.session_id
    sessions.customer_id -> customers.customer_id
```

Feature Tools & Deep Feature Synthesis (MIT/DARPA)

```
>> feature_matrix, features_defs = ft.dfs(entityset=es, target_entity="customers")
>> feature_matrix.head(5)

t)) STD(transactions.MAX(amount)) NUM_UNIQUE(DAY(session_start)) MIN(transactions.am
150 5.857976 1 -0.
350 7.420480 1 -0.
976 12.537259 1 -0.
969 12.738488
385 5.599228
```

Cutoff times

We can specify the time for each instance of the `target_entity` to calculate features. The timestamp represents the last time data can be used for calculating features. This is specified using a dataframe of cutoff times. Below we show an example of this dataframe for our customers example.

IBM R&D: One Button Machine, Cognito & more

COGNITO: Automated Feature Construction
Cognito: Automated Feature Engineering

Tree config: Feature Selection: No Model: Default Metric: Default Budget (mins): 2

Transforms: None Transforms: Select All Selected Name: Auto Detect

Status: Running...

• Employ black-box optimization to build neural network structures (see here).

(train = orange, val = green)

learning Rate

dropout

1 - convpool

1 - filterSize

1 - noFilters

One Button Machine : Automated Feature Engineering

ICS workflow and automation

Model selection & tuning

Feature engineering

cleaning & curation

Video player icon

Bike Sharing Demand

Input: Supervised prediction problem

Data Scientist can interact with System

User Interface (or REST-API directly)

Submit to CADS

AI technology automatically determines best analytics pipeline

Learning Controller

Analytic Monitoring and Adaptation

Science of Analytics Repository

Tactical Planner, Orchestrator and Scheduler

Deployed Analytic

Analytic Platforms

Cross-Platform Deployment and Evaluation

Knowledge Acquisition External Knowledge about Analytics

1, 2, 3, 4, 5, 6

The diagram illustrates the One Button Machine architecture. It starts with a 'User Interface (or REST-API directly)' which interacts with a 'Science of Analytics Repository'. This repository feeds into a 'Tactical Planner' and 'Orchestrator and Scheduler' (labeled 2). These components interact with 'Deployed Analytic' units (labeled 3) running on various 'Analytic Platforms' (labeled 4). The 'Deployed Analytic' units provide feedback to the 'Science of Analytics Repository' (labeled 5). An 'Learning Controller' (labeled 1) monitors the system and provides adaptive feedback. 'Knowledge Acquisition' (labeled 6) feeds into the 'Science of Analytics Repository'. Finally, the system submits results to 'CADS'.



≡

Project: Patient Conditions

**Done**

All jobs are finished

**6/125/72**

Tables/Features/Models

**Azure/16/256**

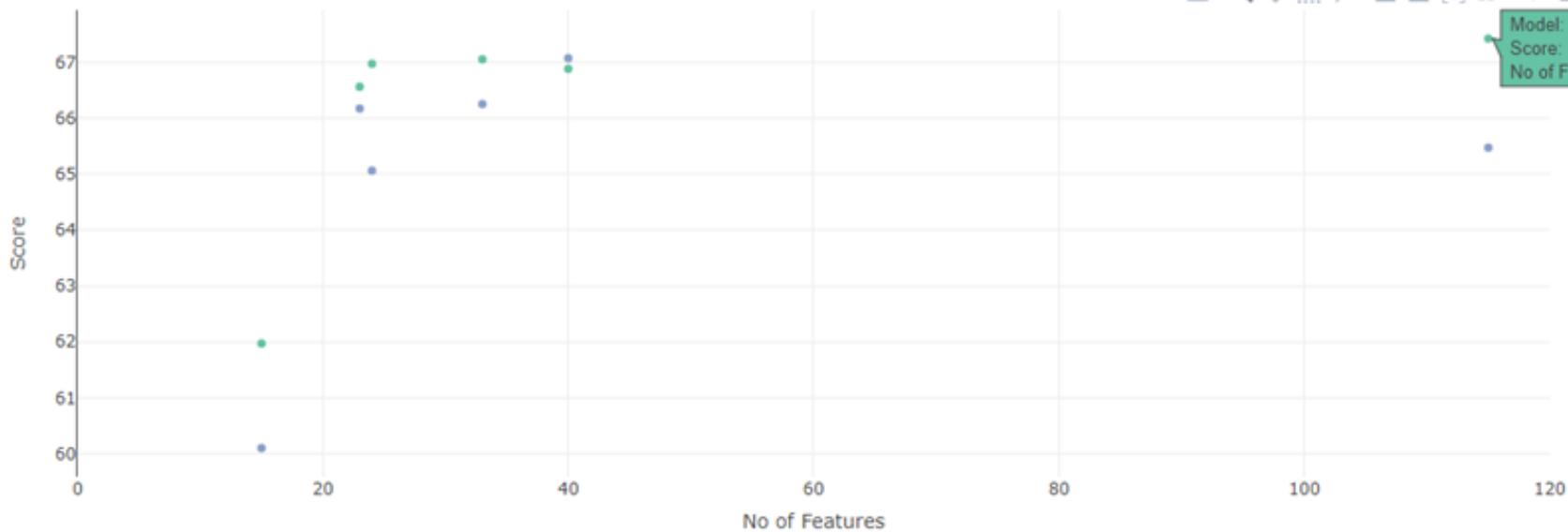
Cluster Mode/Nodes/Cores



Jobs

Evaluation

Score vs No of Features



AutoML Challenges

- “Overfitting” (+inf features), AUC/RMSE/... “hacking”
- Customization/Deployment/Commercial “Lock-in”
- Excessive focus on best performance/Kaggle style
- ROI & Costs (commercial vs in-house)
- Misuse? (from non ML Experts)
- Complexity (open source tools)

Some Best Practices



Emmanuel Ameisen [Follow](#)
AI Lead at Insight AI @EmmanuelAmeisen
Mar 6 · 9 min read

Always start with a stupid model, no exceptions.

How to efficiently build Machine Learning powered products.

Thank you!

Q&A

Should we use GUIs for Machine Learning?

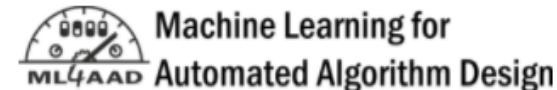
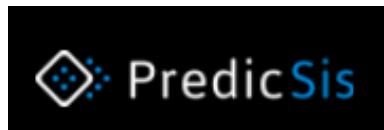
When to use Model/Hyper Params search tools ?

Will we be out of Jobs?

What cannot be automated? Where should we (humans) focus?

Thoughts? Let me know!

rui.quintino@devscope.net



Auto-WEKA

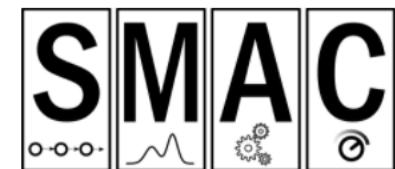
Hyperopt

Advisor

RoBO

automlk

mIrMBO



hyperband



ATM - Auto Tune Models



devscope



Rua Passos Manuel Nº 223 – 4º Andar
4000-385 Porto
Av. Sidónio Pais, Nº 2 – 3º Andar
1050-2145 Lisboa



T. +351 223 751 350/51
F. +351 223 751 352



info@devscope.net
www.devscope.net