

Deep Feature Synthesis: Towards Automating Data Science Endeavors



Xue Huanran 1701214297

2018.12.10



Introduction

In this paper, the Data Science Machine is developed, which is able to derive predictive models from raw data automatically.

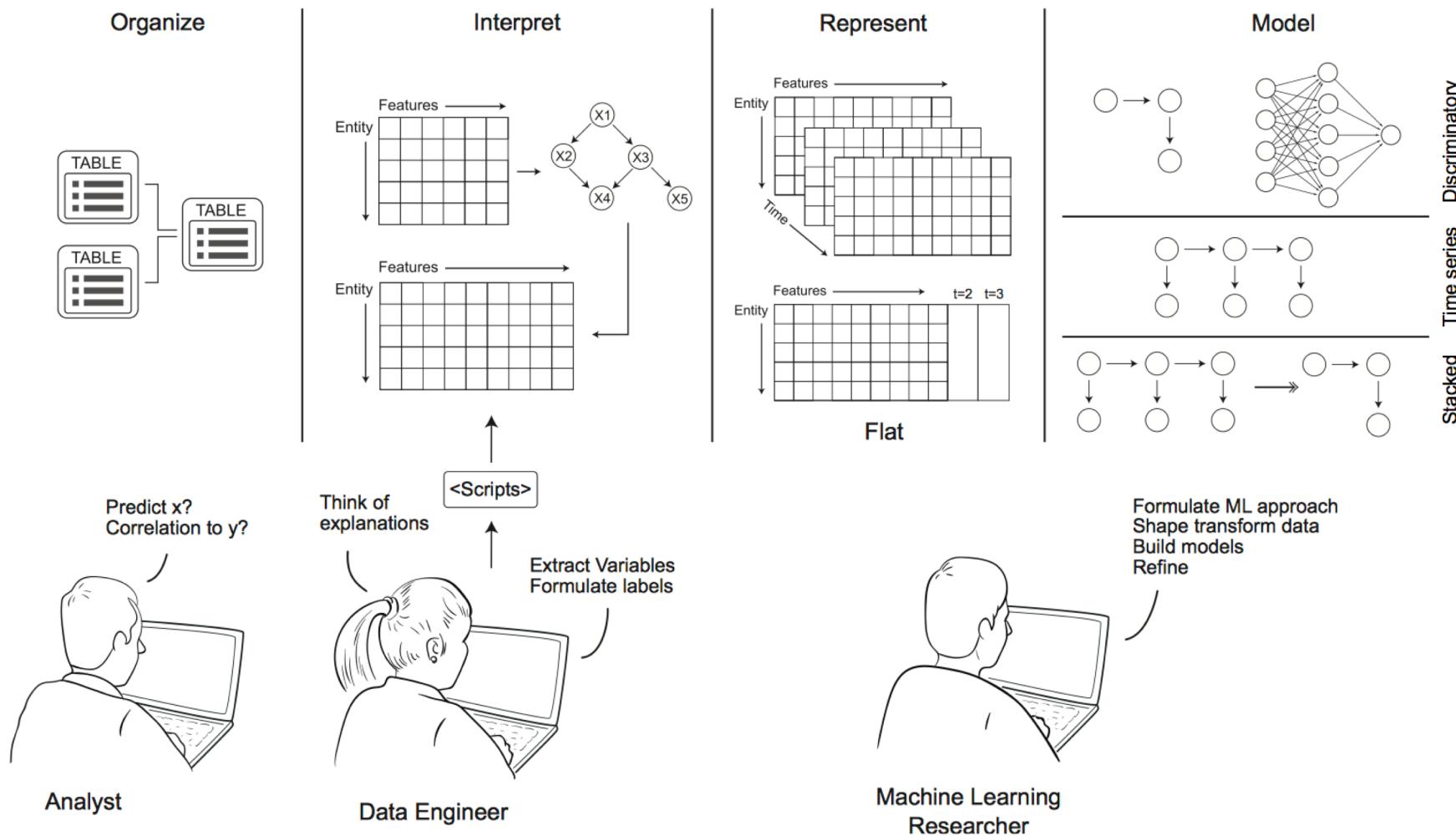
The pipeline of data science:

- Feature Engineering ---- Deep Feature Synthesis
- Select the predictive methodology
- Tuning the hyperparameters



Data Science Machine

A typical data science endeavor





Peking University

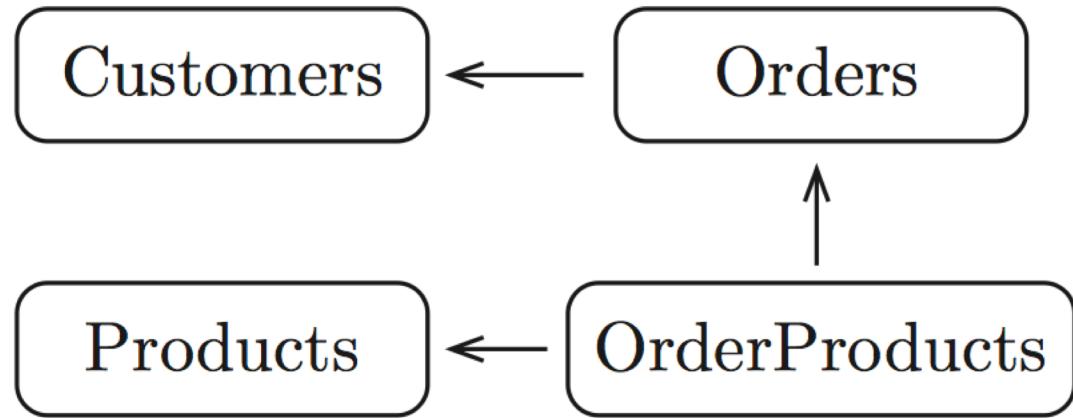
Contributions

- Automatically generate features via **Deep Feature Synthesis**.
- Autotune a machine learning pipeline to extract the most value out of the synthesized features.
- Produce submissions for online data science competitions.

Deep Feature Synthesis



Prototypical Problems and Motivation



A simplified schema for an e-commerce website. There are 4 entities. An arrow from one entity to another signifies that the first entity references the second in the databases.

Feature Synthesis Abstractions



Peking University

Given a dataset, we have entities given by $E^{1 \dots K}$, where each entity table has $1 \dots J$ features. Denote a specific entry as $x_{i,j}^k$, which is the value for feature j for the i^{th} instance of the k^{th} entity.

Next, given entities and their relationships, define a number of mathematical functions that are applied at two levels: at the **entity** level and at the **relational** level.



Two types of features

Peking University

Entity features (efeat): Entity features derive features by computing a value for each entry $x_{i,j}$.

LENGTH(), YEAR(), MONTH(), DAY(), WEEKDAY(),...

The second set of features is derived by jointly analyzing two related entities, E^l and E^k . These two entities relate to each other in one of two ways: **forward** and **backward**.

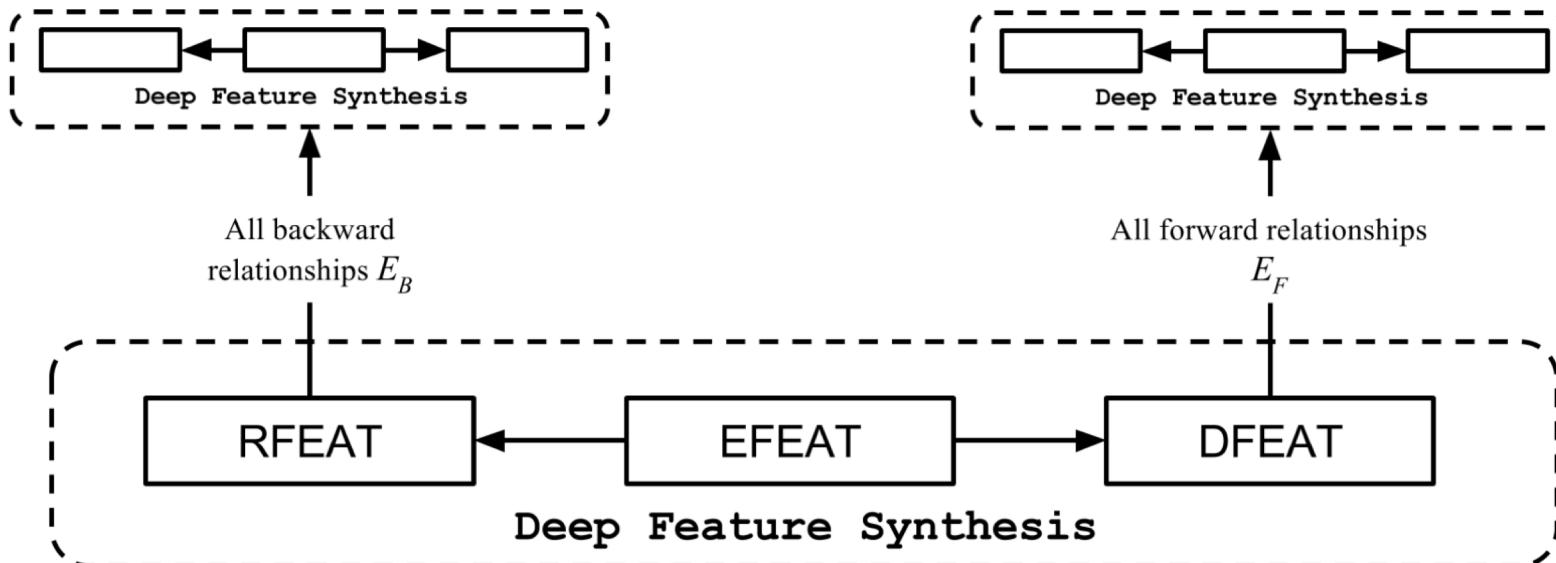
Forward (dfeat): Instance i has an explicit dependence on instance m . E.g. the **Orders** entity has a forward relationship with the **Customers**.

Backward (rfeat): The backward relation is the relationship from an instance i in E^k to all the instances $m = \{1 \dots M\}$ in E^l . E.g. the **Customers** entity has backward relationship with **Orders**.

Avg(), MAX(), MIN(), SUM(), STD(), COUNT()

How to generate features

Consider a dataset of K entities, the goal is to extract **rfeat**, **dfeat** and **efeat** features for a target E^k . Additionally, all the entities with which E^k has **forward** and **backward** relationships. These are denoted by sets E_F and E_B .





The recursive scenario of DFS

Algorithm 1 Generating features for target entity

```
1: function MAKE_FEATURES( $E^i, E^{1:M}, E_V$ )
2:    $E_V = E_V \cup E^i$ 
3:    $E_B = \text{BACKWARD}(E^i, E^{1:M})$ 
4:    $E_F = \text{FORWARD}(E^i, E^{1:M})$ 
5:   for  $E^j \in E_B$  do
6:     MAKE_FEATURES( $E^j, E^{1:M}, E_V$ )
7:      $F^j = F^j \cup \text{RFEAT}(E^i, E^j)$ 
8:   for  $E^j \in E_F$  do
9:     if  $E^j \in E_V$  then
10:       CONTINUE
11:     MAKE_FEATURES( $E^j, E^{1:M}, E_V$ )
12:      $F^i = F^i \cup \text{DFEAT}(E^i, E^j)$ 
13:    $F^i = F^i \cup \text{EFEAT}(E^i)$ 
```

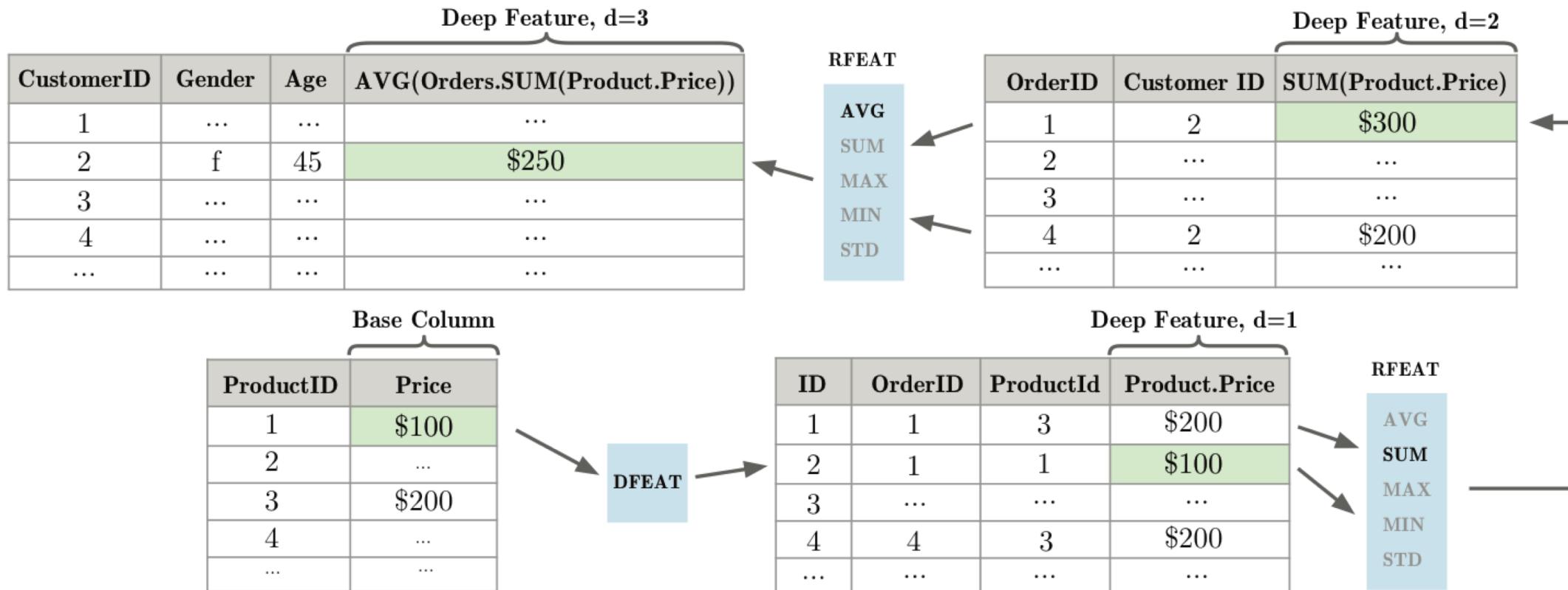
The MAKE_FEATURES is presented to make features, F^i , for the i^{th} entity.

E_V keeps track of the entities that have been visited.

An Example of DFS



Peking University



Data Science Machine



Machine Learning Pipeline

Peking University

Data Science Machine implements a parametrized pipeline for data preprocessing, feature selection, dimensionality reduction, modeling and evaluation.

Data preprocessing: Remove the null variables, convert the categorical variables using one-hot encoding, normalizing the features.

Feature selection and dimensionality reduction: First, use truncated SVD transformation. Then, rank each feature by f-value with the target value and select the $r\%$ highest ranking features.

Modeling: Random Forest is used by constructing n decision trees

Hyperparameter Optimization: Gaussian Process

Experimental Results

Competitions



Peking University

KDD CUP 2014: Using past projects histories on DonrosChoose.org, predict if a crowd-funded project is exciting.

IJCAI: Using past merchant and customer shopping data, predict if a customer making a purchase during a promotional period will turn into a repeat buyer.

KDD CUP 2015: Using student interaction with resources on an online course, predict if the student will dropout in the next 10 days.

AUC Scores of Data Science Machine



Peking University

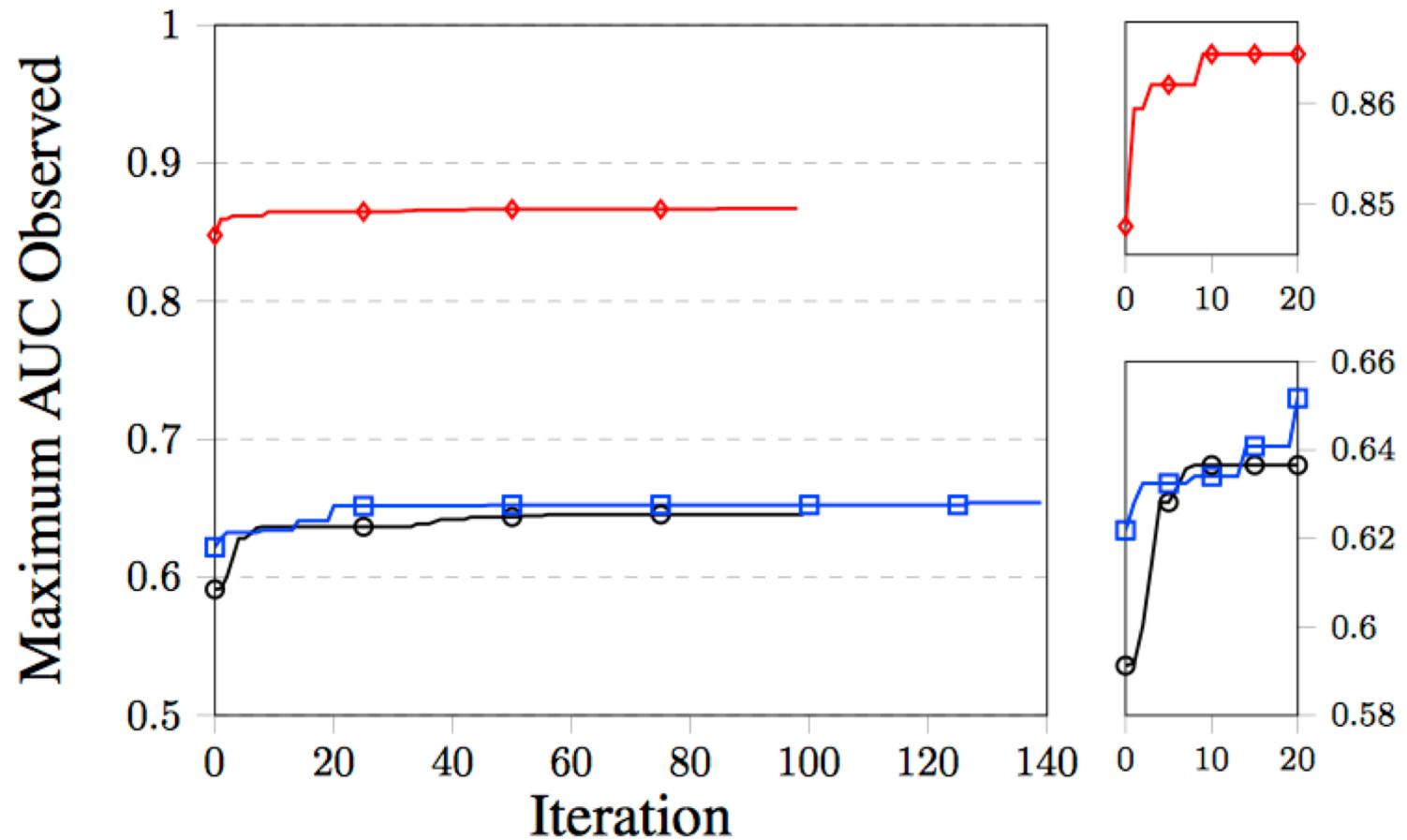
Parameter Selection	KDD cup 2014		IJCAI		KDD cup 2015	
	Local	Online	Local	Online	Local	Online
Default	.6059	.55481	.6439	.6313	.8444	.5000
GCP Optimized	.6321	.5863	.6540	.6606	.8672	.8621

The local score is the result of running 3-folds cross validation, while online score is based on submitting predictions to the competition.



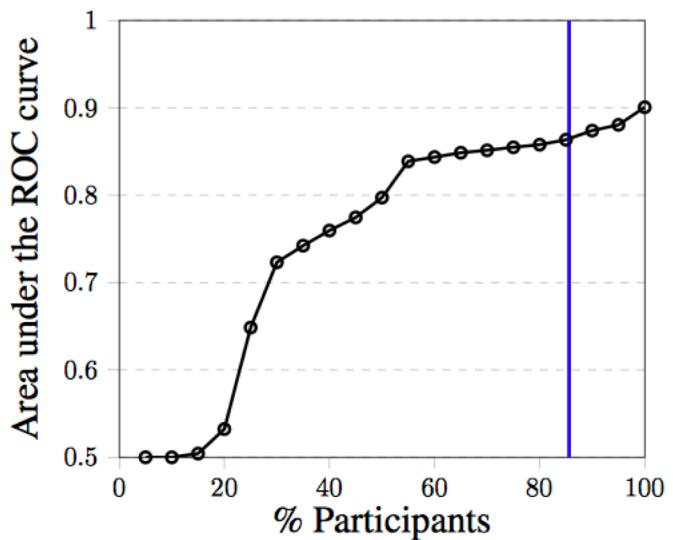
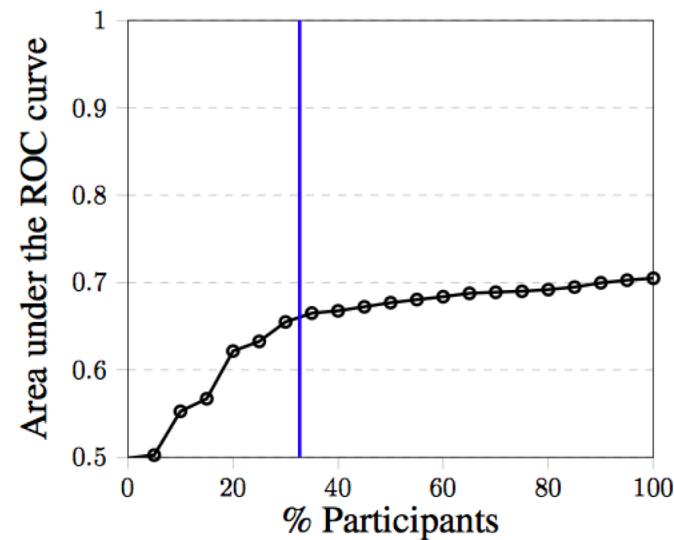
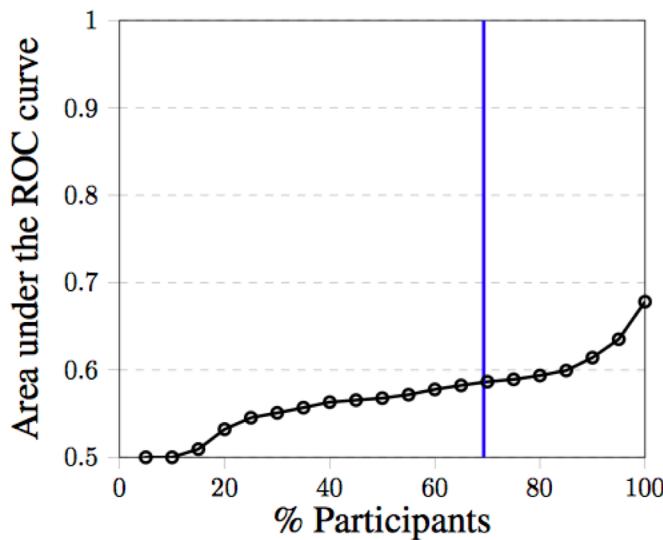
Peking University

The Effect of Gaussian Optimization



Performance Comparison

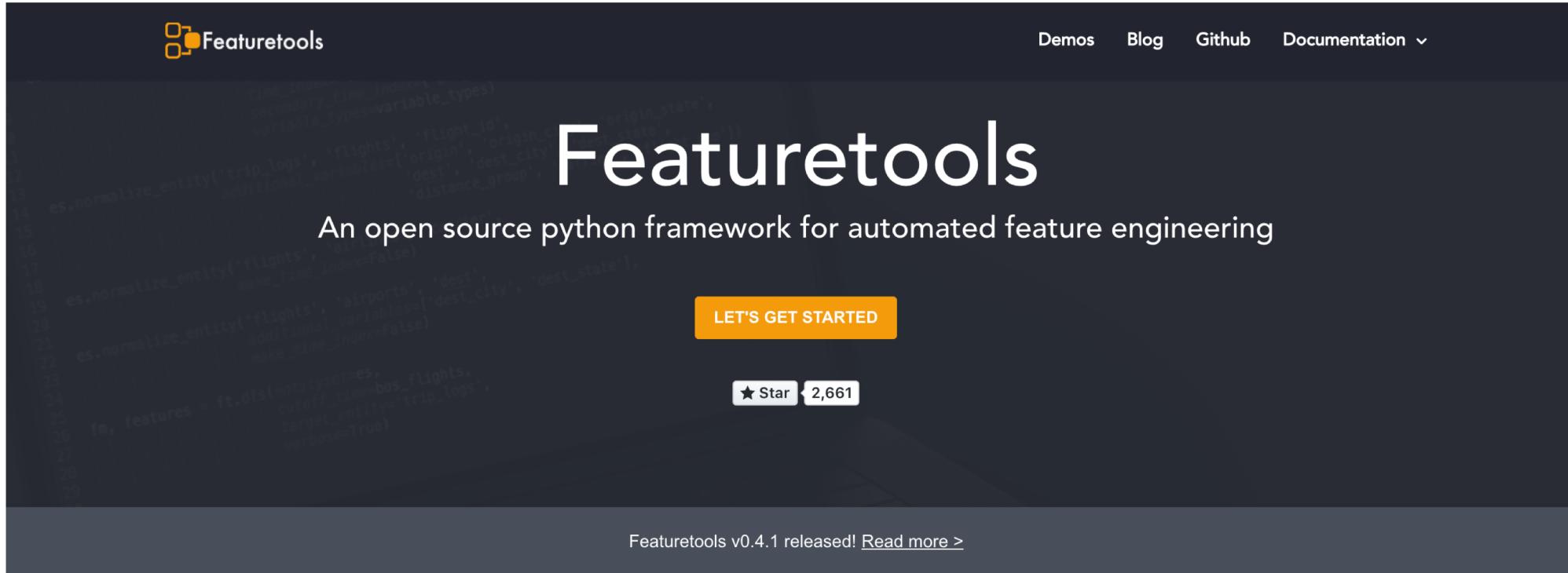
Dataset	# Teams	% of Top Submission's Score	% of Teams Worse	# Submissions worse	# Days Spent on Worse Submissions
KDD cup 2014	473	86.5%	69.3%	3873	929
IJCAI	156	93.7%	32.7%	-	-
KDD cup 2015	277	95.7	85.6%	1319	377





Peking University

Featuretools



The screenshot shows the Featuretools homepage. At the top, there is a navigation bar with the Featuretools logo, Demos, Blog, Github, and Documentation links. The main title "Featuretools" is prominently displayed in large white font. Below it, a subtitle reads "An open source python framework for automated feature engineering". A large orange button labeled "LET'S GET STARTED" is centered. To the right of the button, there is a star icon and the number "2,661". In the background, there is a faint watermark of a Python code snippet related to feature engineering.

Featuretools v0.4.1 released! [Read more >](#)

PREPARE DATA FOR MACHINE LEARNING

Featuretools automatically creates features from
temporal and relational datasets



Peking University

Feature Labs

A photograph showing a person's hands typing on a laptop keyboard. The laptop screen displays the Feature Labs website. The main headline on the site reads: "Getting Value from Machine Learning Isn't About Fancier Algorithms — It's About Making It Easier to Use". Below the headline, a subtext states: "Feature Labs accelerates the error-prone, time-intensive, and costly process of intelligently transforming raw data for machine learning algorithms". A "Learn more" button is visible at the bottom left of the laptop screen. The background of the image is a blurred indoor setting.

Technology to remove the bottlenecks in the data science process

Conclusion



Peking University

Consider the scene with multi-tables which is suitable for industry.

The novelty is only focus on the part on auto feature engineering.

This method can be applied in the system in Tencent.



Peking University

Q & A