

Demo: Probabilistic Predicates to Accelerate ML Inference

Microsoft

Research Yao Lu, Srikanth Kandula, Surajit Chaudhuri

UNIVERSITY of WASHINGTON

Motivation



Context: ML serving + Relational big-data platforms
SQL + UDFs on images, text, videos...

Example

```
SELECT ImageID FROM UDF_YOLOv2(Images)
WHERE has_orange AND has_banana;
```

Query is **slow** because each image \rightarrow UDF.

Speedup? Images \rightarrow UDF_YOLOv2 \rightarrow $\sigma_{orange\&banana}$ \rightarrow Result

Predicate pushdown? **X**
No "orange" column here

Pre-materialize? **X**
Too many UDFs & predicates for ad-hoc queries; high storage cost

Probabilistic Predicates (PPs) [2]:

Input \rightarrow PP_{orange&banana} \rightarrow YOLOv2Detector \rightarrow σ



Can offer sizable speed-up if:

- relative cost of PP is small
- achieves large data reduction
- meets target accuracy

Technical problems addressed

- Train PPs for diverse inputs - SVMs, KDE, shallow NNs... - model selection.
- Avoid per-query training - train PP for *simple preds.*, - use QO to build PP comb.
- Sizable gains on images, video and text datasets

Interactive Demo

Probabilistic Predicates

PPs are **fast predicate-specific** filters on input; they reduce work for expensive UDFs.

Example query: retrieve images containing "oranges":

Images \rightarrow PP_{orange} \rightarrow UDF_YOLOv2 \rightarrow σ \rightarrow Result

No harm on query accuracy, but less speedup

Accuracy vs. Data reduction

Execution cost

Simple PP using linear SVM:

$$PP_{orange}: f(x) = w \cdot x + b$$

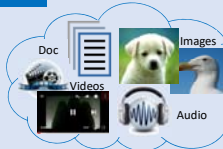


PP discards $f(x) \leq th$

Accuracy/reduction tradeoff; choose threshold (th)



Key ideas



Complex query predicates:

- Large space of possible predicates
- Costly to train/store PP per predicate
- Indiv. PPs do not generalize

Challenge 1: Build PPs for arbitrary inputs? **Dimensionality, Sparsity, Linear separability?** ...

Our solution:

Linear SVM: $f(x) = w \cdot x + b < th$

Kernel Density Estimator: $f(x) = kde^*(x) / kde^*(x) < th$

Shallow DNN: $f^*(x) = g_n(W_n \cdot f^{(n-1)}(x) + b_n) < th$

Any classifier that fits $f(x) < th$ can work for PP

- Pre-process: dim. reduction (PCA, feature hashing.) (Opt.)
- Model Selection

Challenge 2: Use PPs for complex predicates w/o per-query training (i.e., small overhead)?

5 cols x 10 values \Rightarrow $>10^5$ predicates

PP_{redASUV} will not work for σ_{SUV}

(i) Build PPs for simple predicates; QO identifies best necessary condition

(ii) Assign accuracy budget per PP

(iii) Assemble PPs into plan

Injection rule for $p \vee q$

Injection rule for $p \wedge q$

Solved using dynamic programming

Find the plan with best reduction/cost

Machine Learning + Query Processing

1. COCO Image Dataset

Input: images

UDF: YOLOv2 Object Detector

Output: [0,1] vector \Rightarrow Multi-Class labels

```
SELECT img_id FROM YOLOv2()
WHERE has_dog & has_person &
has_car | has_bicycle;
```

2. UCF101 Video Dataset

Input: Videos

UDF: CNN Video Features, Linear SVM

Output: category - Single-class label

```
SELECT video_id FROM LSVM()
WHERE cat. = "ApplyLipstick";
```

3. IMDB Movie Review Dataset

Input: movie review text

UDF: Gated-Recurrent-Unit (GRU) RNN

Output: score - [1,10] real value

```
SELECT docId FROM RNNClas.()
WHERE score < 2 | score > 8;
```

4. DETRAC Traffic Surveillance Video Analytics

Input: Video frames

UDF: YOLOv2 vehicle detector

For each detected vehicle patch: CNN image feature extractor #1,2 LinearSVM classifier #1 \Rightarrow Veh Color #2 \Rightarrow Vehicle Type

Output: veh_color, veh_type

```
SELECT * FROM $color, $type
ON $color.{fid,vid} =
$type.{fid, vid}
WHERE veh_color = red &
veh_type = SUV;
```

System & Demonstration

A cross-platform, data-parallel engine for DNN-enabled ML

User query: Standard SQL + UDF syntax. See [Lu, et al. SoCC16].

Timely Dataflow: Rust impl. of Naiad [Murray et al. SOSP13]

Query engine: Join node, filter logic, columnar store for blobs etc.

μ DL: A light-weight DNN engine in C/C++.

μ KDE: A Kernel Density Estimator in C/C++.

Demonstration Overview

(1) ML query processing with/without PPs. (2) User-specified filtering accuracy. (3) QO with available PPs

Demonstration Systems:

a) Baseline System w/o PPs

b) Constructing PPs. (Not demonstrated)

c) Full system w/ PPs

Example:

Images where has_orange & has_banana

Original ML query plan:

Input \rightarrow YOLOv2Det \rightarrow σ

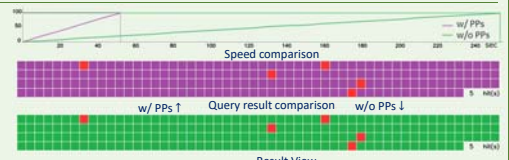
Candidates plans:

[has_orange&has_banana]: 0.974 *

[has_orange] 0.939, [has_banana]: 0.873

Query plan with PPs, target $a = 0.9$:

Input \rightarrow orange \rightarrow banana \rightarrow FileReader \rightarrow YOLOv2Det \rightarrow σ



Please try out different ML queries, with different predicates, PPs on/off, and different target accuracies.

References

- <https://www.reportlinker.com/p04519487/Machine-Learning-as-a-Service-MLaaS-Global-Market-Outlook.html>
- Y. Lu, A. Chowdhery, S. Kandula, S. Chaudhuri. Accelerating Machine Learning Inference with Probabilistic Predicates. SIGMOD 2018.
- Y. Lu, A. Chowdhery, S. Kandula. A Relational Platform for Efficient Large-Scale Video Analytics. SoCC 2016.
- D. Murray et al. Naiad: a timely dataflow system. SOSP 2013.