ACM SIGMOD/PODS Chicago, IL, USA

May 15, 2017

# A Relational Framework for Classifier Engineering

Benny Kimelfeld

Technion, Israel



Christopher Ré

Stanford University



# Background

- ML application constantly increasing
  - e.g., by 2020 >50% Intel servers will run ML (D. Bryant, Intel SVP)
- Rising interest in DB research for ML
  - e.g., query optimization for feature selection / evaluation [Zhang+14, Kumar+15,16], ML on factorized DB [Schleich+16]
  - DEEM workshop on Data Management for End-to-End ML
  - Dagstuhl Presp. Workshop 16151: Research Directions for PDM
- Feature Engineering (FE) critical for quality
  - Yet heavy resource consumer in ML development
  - Tooling and principles [Guyon+06 book]
  - Standard practice; here to stay!
    - Deep Learning avoids FE; applicable in certain areas / domains w/ massive training data available

### **Classic ML Classification Flow**



## **Classic ML Classification Flow**



## Framework Goal

- DB "understands" how *entities* become *features* Relational structure, constraints, queries
- Can be used for assisting FE?
  - Estimate feature quality?
  - Suggest new features?
  - Test for suitability of a feature language?
  - Detect engineering faults?
  - Implication of underlying languages on computational complexity?
  - Benefit from decades of DB theory?
- Setup for attacking questions
- Step towards DB theory for ML engineering

### Outline

- Formal Setup
- Computational Problems
- Complexity Results
- Directions



- ML task: binary classification
  - Learn a mapping entity  $\rightarrow +1/-1$
- Boolean features
  - Simplifies the framework
  - Common in practice
    - e.g., binning / bucketing
- Hence, a *classifier* has the form

 $\mathbf{C}: \{+1, -1\}^n \longrightarrow \{+1, -1\}$ 



	Txr	nInfo		(	Card		
TXN	card	country	state	id	SSN	country	state
1	100	US	GA	100	200	US	GA
2	100	US	NY	101	201	US	NY
3	101	BR	RJ	102	202	BR	RJ
4	102	US	CA				



(txn in owner's country)  $Q_2(x) \leftarrow TxnInfo(x, n, c, s), Card(n, d, s)$ (txn in NY)  $Q_3(x) \leftarrow TxnInfo(x, n, c, 'NY')$ 



+1	+1	-1	-1
-1	+1	+1	+1
-1	-1	-1	-1
-1	-1	-1	-1



# Formal Setup

• Entity schema:  $(S,\eta)$ 

- S is a relational schema (signature, constraints)
- $\eta$  is a unary relation in **S**, representing *entities*
- An instance I of S defines:
  - An entity set  $\eta^{\mathrm{I}}$  (the  $\eta$  relation of  $\mathrm{I}$ )
  - Information on the entities (all other relations)
- Feature query: unary query  $\boldsymbol{Q}$  over  $\boldsymbol{S}$
- Statistic: series  $\Pi = (Q_1, ..., Q_n)$  of feature queries
- Each  $e \in \eta^I$  has a feature vector  $\Pi(e) = (f_1, ..., f_n)$

 $\mathbf{f}_{i} = \begin{cases} +1 & \text{if } \mathbf{e} \in \mathbf{Q}_{i}(\mathbf{I}) \\ -1 & \text{if } \mathbf{e} \notin \mathbf{Q}_{i}(\mathbf{I}) \end{cases}$ 

### S



	Txr	nInfo			Card			
TXN	card	country	state	id	SSN	country	state	
1	100	US	GA	100	200	US	GA	
2	100	US	NY	101	201	US	NY	
З	101	BR	RJ	102	202	BR	RJ	
4	102	US	CA					

 $\begin{array}{l} (txn \ in \ owner's \ state) & \mathsf{Q}_1(\mathbf{x}) \leftarrow \mathsf{TxnInfo}(\mathbf{x}, n, c, s), \mathsf{Card}(n, c, s) \\ (txn \ in \ owner's \ country) & \mathsf{Q}_2(\mathbf{x}) \leftarrow \mathsf{TxnInfo}(\mathbf{x}, n, c, s), \mathsf{Card}(n, d, s) \\ (txn \ in \ NY) & \mathsf{Q}_3(\mathbf{x}) \leftarrow \mathsf{TxnInfo}(\mathbf{x}, n, c, \mathsf{'NY'}) \end{array}$ 

Feature queries

Statistic:  $\Pi = (Q_1, Q_2, Q_3)$ 

# Training

- Let  $(\mathbf{S},\eta)$  be an entity schema
- A training instance is a pair (I, $\lambda$ ) where
  - I is an instance over  ${\boldsymbol{S}}$
  - $\lambda: \eta^{I} \longrightarrow \{+1,-1\}$  is a labeling function
- $(I,\lambda)$  + statistic  $\Pi$  define the training collection

#### $T = \{ \langle \Pi(e), \lambda(e) \rangle \mid e \in \eta^{I} \}$

- Training finds a classifier from a hypothesis class  ${\bf H}$  by minimizing a risk function over  ${\bf T}$ 



Classifier (model)  $Q_1(x) \leftarrow TxnInfo(x, n, c, s), Card(n, c, s)$  $Q_2(x) \leftarrow TxnInfo(x, n, c, s), Card(n, d, s)$  $Q_3(x) \leftarrow TxnInfo(x, n, c, 'NY')$ 

 $\mathbf{\Pi}=(\mathbf{Q}_1,\mathbf{Q}_2,\mathbf{Q}_3)$ 



Txn

id

1

2

З

4

Txn			Txr	nInfo		Card				
	id	λ(e)	TXN	card	country	state	id	SSN	country	state
	1	-1	1	100	US	GA	100	200	US	GA
335	2	+1	2	100	US	NY	101	201	US	NY
	3	-1	3	101	BR	RJ	102	202	BR	RJ
	4	-1	4	102	US	CA				
			Q]	ĹЭ	$ \left\{\begin{array}{c} Q_1(\mathbf{x}) \\ Q_2(\mathbf{x}) \\ Q_3(\mathbf{x}) \end{array}\right. $	$\leftarrow Tx \\ \leftarrow Tx \\ \leftarrow Tx$	nInfo(2 nInfo(2	x, n, c, x, n, c, x, n, c,	s), Card( s), Card( 'NY')	n, d, s)
-	+1 +1	-1 -1	]		1	П=	=(Q <sub>1</sub> ,	Q <sub>2</sub> ,Q	3)	
	-1 +1	+1 +1		Clas (mo						
	-1 -1	-1 -1		sifier del)						
	-1 -1	-1 -1		·						
	Т									

### Outline

- Formal Setup
- Computational Problems
  - Complexity Results
  - Directions



# Problem 1: Separability

The naïve "noise-free" training from ML textbooks: Is full separation possible?

#### (H,QL)-separablity

Given a training instance  $(I,\lambda)$  over a schema  $(S,\eta)$ , is there any statistic  $\Pi$  in **QL** such that  $(I,\lambda)$  can be perfectly realized by a classifier in **H**?





# Redundancy / Identifiablity



- Linear column dependence in the feature matrix often means redundant features

   e.g., linear/logistic classification/regression
- ML libraries often require full column rank
  - For their optimization solution to be "identifiable"
  - c.f. "Theory of Point Estimation" [LehmannCasella83]





TXN in owner's US state	TXN in different US state	TXN in East Coast	TXN US but not East Coast
+1	-1	+1	-1
-1	+1	-1	+1
-1	-1	-1	-1



TXN in owner's US state	TXN in different US state	TXN in East Coast	TXN US but not East Coast		
+1	-1	+1	-1		
-1	+1	-1	+1		
-1	-1	-1	-1		
SL	ım =	sum			

# Problem 2: Identifiability

#### **QL**-identifiability

Given a statistic  $\Pi$  in **QL** over entity schema (**S**, $\eta$ ), is there any instance **I** with a column-independent feature matrix?



Two variants:

- Linear independence (arises in, e.g., least-square minimization)
- Affine independence (arises in, e.g., entropy minimization)

Txn				Txn	Info			(	Card	
	id	λ(e)	TXN	card	country	state	id	SSN	country	state
	1	-1	1	100	US	GA	100	200	US	GA
	2	+1	2	100	US	NY	101	201	US	NY
	3	-1	3	101	BR	RJ	102	202	BR	RJ
	4	-1	4	102	US	CA				
	-1 +1 1 +1 1 -1 1 -1	(txi (txn ir -1 -1 +1 -1 -1 -1 -1 -1	n in own	er's state s countr (txn in N (model)	$\begin{array}{c} \Theta \\ \Theta $	← Txı ← Txı ← Txı ← Txı	nInfo(x nInfo(x nInfo(x PCh tra	, n, c, : , n, c, : , n, c, ' ainin, r feat	s), Card(1 s), Card(1 'NY') g to tures?	n, c, s) n, d, s)

# Vapnik-Chervonenkis (VC) Dimension

What is the max #entities that can be shattered
 That is, perfectly classified on every possible labeling?



- Complexity measure for learnability
  - (not the only one)
- Estimate training amount to avoid overfitting

# Problem 3: Dimensionality

#### (H,QL)-dimensionality

Given a statistic  $\Pi$  in **QL** over an entity schema (S, $\eta$ ), what is the max m such that some instance with m entities can be shattered by **H**?



### Outline

- Formal Setup
- Computational Problems
- Complexity Results
- Directions

# **Scope of Results**

- Complexity analysis in a specific setting:
  - Hypothesis class  $\mathbf{H} = \mathbf{Lin}$ : linear classifiers
  - Query language  $\mathbf{QL} = \mathbf{CQ}$ : conjunctive queries
    - Without constants
  - No schema constraints
- Mostly intractable complexity classes (expected)
- Baseline & justification for future assumptions
- Next, a few highlights

# (Lin,CQ)-Separability



Given a training instance  $(I,\lambda)$  over a schema  $(S,\eta)$ , is there any statistic  $\Pi$  in **CQ** such that  $(I,\lambda)$  can be perfectly realized by a classifier in **Lin**?

- Every training instance is separable, unless entities with different labels are indistinguishable by CQs
  - That is, there are e and e' with  $\lambda(e) \neq \lambda(e')$  and endomorphism that maps e and e' and vice versa
  - Relationship to CQ-query-by-example
    - [Willard10,tenCateDalmau15,BarcelóRomero16]
  - coNP-complete
- Avoiding self joins  $\rightarrow$  harder:  $\Sigma_2^P$ -complete!

# **CQ-Identifiability**



Given a statistic  $\Pi$  in **CQ** over entity schema (**S**, $\eta$ ), is there any instance **I** with a column-independent feature matrix?

- The following are equivalent if CQs are connected:
  - $\Pi$  is linearly identifiable
  - $\Pi$  is affinely identifiable
  - $\Pi$  is non-redundant (no equivalent feature queries)
- Pairwise equivalences break if:
  - CQs can be disconnected
  - CQs can have negation
- Generalized characterization for disconnected CQs
- coNP-complete

# (Lin,CQ)-Dimensionality

Given a statistic  $\Pi$  in **CQ** over entity schema  $(S,\eta)$ , what is the max m such that some instance with m entities can be shattered by **Lin**?

- For connected CQs VC dim w.r.t.  $\Pi$  is  $d{+}1$ 
  - $\mathbf{d} = \textit{\texttt{#}}(\textit{equivalence classes among CQs in } \boldsymbol{\Pi})$
  - In particular, containment among CQs does not reduce the VC dimension compared to vanilla linear classification
- Can go down if we allow:
  - Disconnected CQs
  - Negation

### Outline

- Formal Setup
- Computational Problems
- Complexity Results
- Directions

# **Directions for Future Research**

- Schema constraints
- Generalized features / tasks
   Numeric, aggregate, multi-label, regression
- Realistic variants of separability
   Approximate/noisy, incremental
- Restrict model complexity
  - Small/shallow feature queries, low statistic dimension
- Connection to prob. DBs (statistical guarantees?)
- Context of text analysis
  - Doc. spanners [Fagin+2014], DeepDive [Shin+2015]



## Summary

- Framework for classifier engineering over DBs
   Entity schema, feature query, statistic, training instance
- Goal: DB smartness (schema, constraints, queries) to aid feature engineering
- Illustrated on several computational problems
  - Separability, dimensionality, identifiability
  - Preliminary results for linear classifiers and CQs
- Plethora of problems / directions to pursue

#### Thank you! Questions?