

Introduction

- Statistical-Relational Learning:** Learn a joint statistical model for *all* tables in the input database.
- New approach to *SRL system building*.
- The RDBMS stores structured objects for statistical analysis as *first-class citizens* in the database.
- SQL is used to build and transform statistical objects:
 - Structured Model (*Bayesian network, Markov Logic Network*).
 - Parameter Estimates.
 - Sufficient Statistics.
- Empirical evaluation: leveraging the RDBMS capabilities achieves scalable learning and fast model testing.
- All code and datasets are available online [1].

Contributions

- Identifying new system requirements for multi-relational machine learning that go beyond single table machine learning.
- An integrated set of SQL-based solutions for providing these system capabilities.

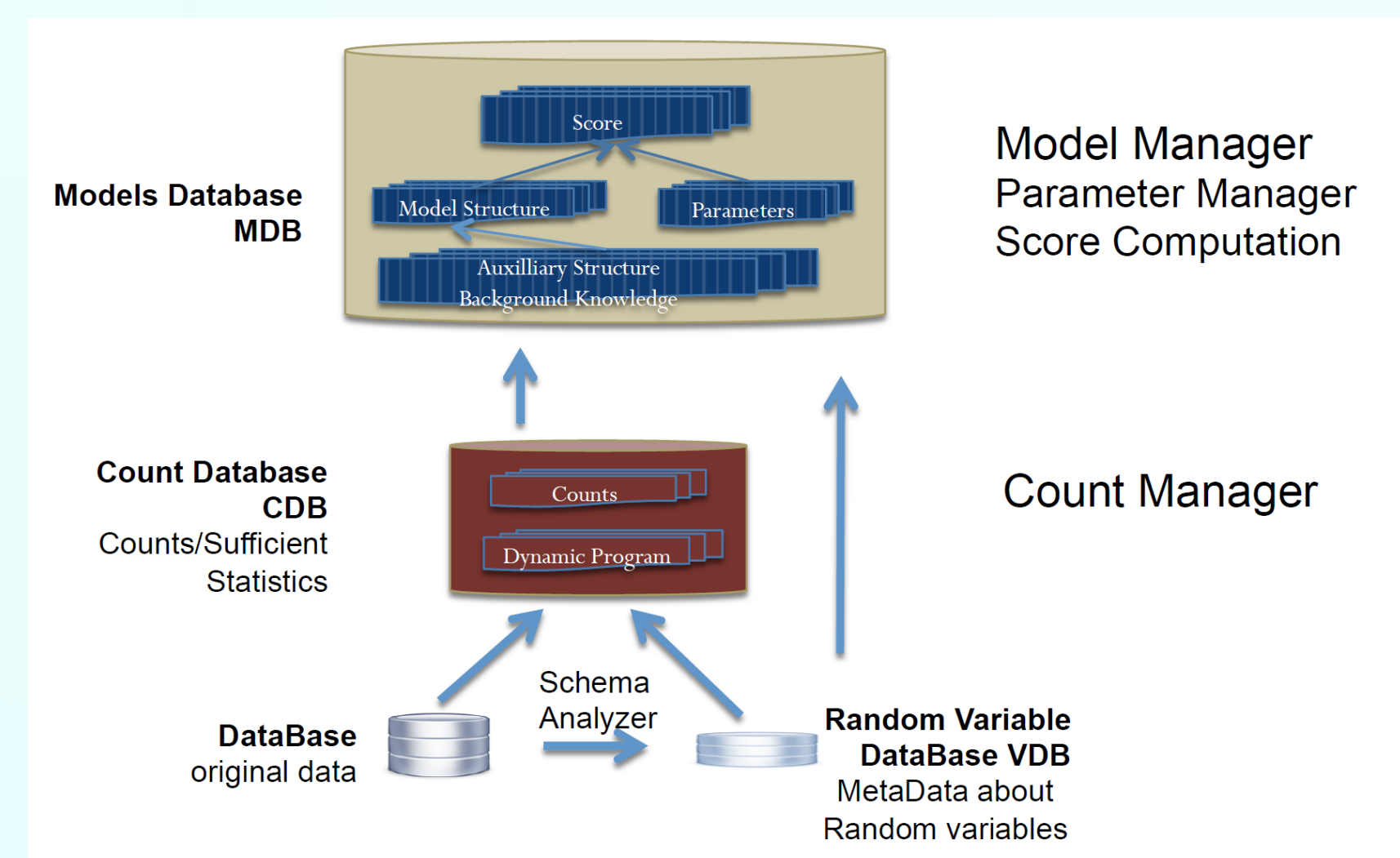
Related Works

- BayesStore [3]: all statistical objects are first-class citizens in a relational database. Inference, no learning.
- MadLib [5]: leverages SQL for *single-relational* data table analysis.
- Tuffy [7]: reliable and scalable inference and parameter learning for Markov Logic Networks with an RDBMS. No structure learning.

References

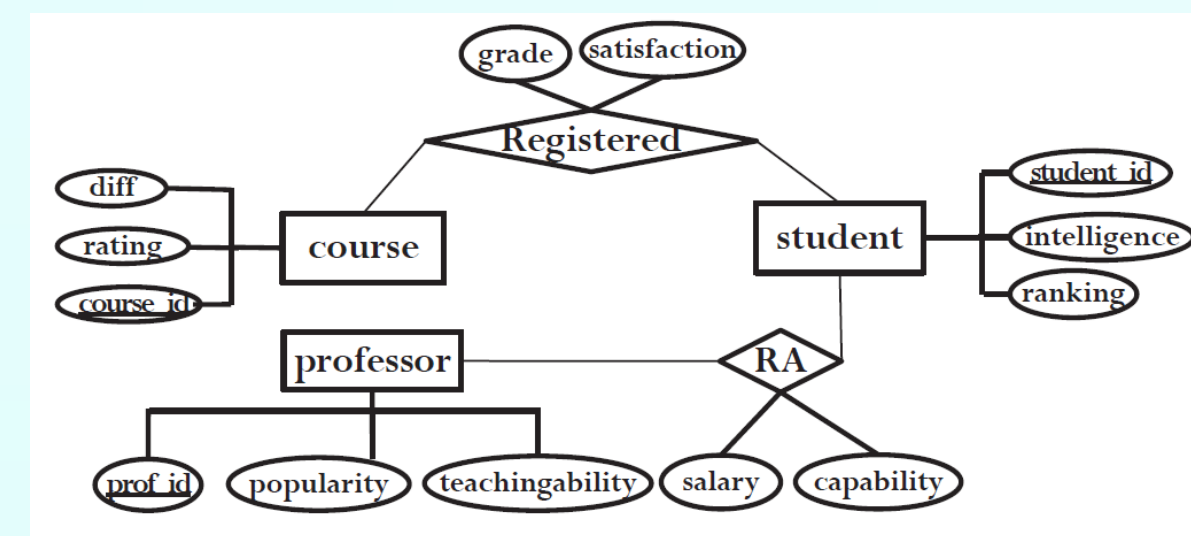
- Qian, Z.; Schulte, O. & et al. MRLBase Code. <http://www.cs.sfu.ca/~oschulte/BayesBase>
- Russell, S. & Norvig, P. Artificial Intelligence: A Modern Approach Prentice Hall, 2010.
- Wang, D. Z.; Michelakis, E.; & et al. BayesStore: managing large, uncertain data repositories with probabilistic graphical models, PVLDB, 2008, 1, 340-351.
- Qian, Z.; Schulte, O. & Sun, Y. Computing Multi-Relational Sufficient Statistics for Large Databases, CIKM 2014, 1249-1258.
- Hellerstein, J. M.; Ré, C.; Schoppmann, F.; & et al, The MADlib Analytics Library: Or MAD Skills, the SQL, PVLDB, 2012, 5, 1700-1711
- Schulte, O. & Khosravi, H. Learning graphical models for relational data via lattice search Machine Learning, 2012, 88, 331-368
- Niu, F.; Ré, C.; Doan, A. & Shavlik, J. W. Tuffy: Scaling up Statistical Inference in Markov Logic Networks using an RDBMS PVLDB, 2011, 4, 373-384

System Overview



- Schema Analyzer:** examines the information in the DB system catalog to define a default set of random variables.
- Count Manager:** uses the meta data in the VDB database to compute multi-relational sufficient statistics for a set of random variables [4].
- Model Manager:** supports the construction and querying of large structured statistical models.

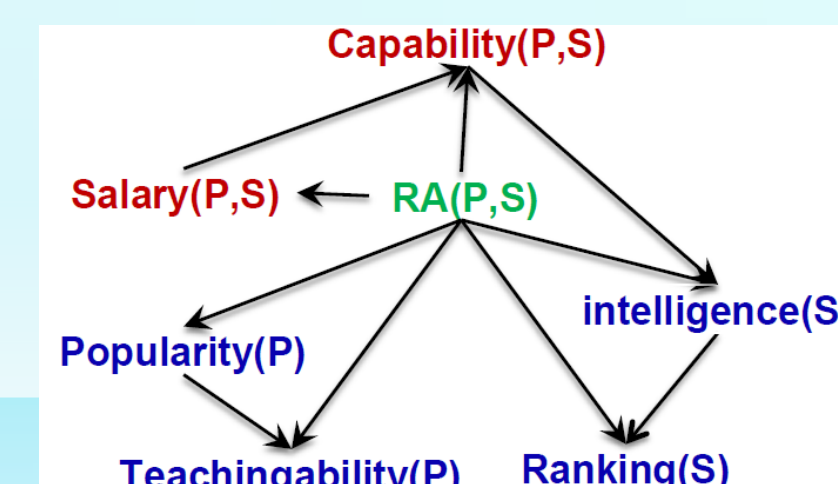
ER-Design for University Domain



The Model Manager

Goal: Learn First-Order Bayesian Network [2].

- Bayesian Network Structure Learning [6].
- Nodes = Random Variables
- Edges are stored in Database tables
- Model selection scores are also stored, not shown (BIC, AIC, BDeu)



Child	Parent
Capability(P,S)	RA(P,S)
Capability(P,S)	Salary(P,S)
Teachingability(P)	Popularity(P)
Teachingability(P)	RA(P,S)
...	...

The Parameter Manager

Goal: Learn Bayesian Network Parameters

- Stored in Conditional Probability (CP) table.
- Maximum Likelihood Estimate are easy to compute from database counts.

Capa(P,S)	RA(P,S)	Salary(P,S)	CP
4	T	high	0.45
5	T	high	0.36
3	T	high	0.18
3	T	low	0.20
2	T	low	0.40
1	T	low	0.40
2	T	med	0.22
3	T	med	0.44
1	T	med	0.33

CP table

Count	Capa(P,S)	RA(P,S)	Salary(P,S)
5	4	T	high
4	5	T	high
2	3	T	high
1	3	T	low
2	2	T	low
2	1	T	low
2	2	T	med
4	3	T	med
3	1	T	med

Contingency Table

SELECT COUNT(*) AS Count, Capability as 'Capa(P,S)', 'T' as 'RA(P,S)', Salary as 'Salary(P,S)' FROM 'RA';

Specific SQL Query

The Count Manager

Goal: for a conjunctive query, compute the instantiation count = result set size.

- Stored in Contingency (CT) Table [4].
- Main computational cost in learning.

Problem: need to generate SQL queries for **arbitrary variable lists**.

Solution: use Meta Data + Meta Queries

General Form of SQL Count Query:

SELECT COUNT(*) AS Count, <VARIABLE-LIST> FROM TABLE-LIST GROUP BY <VARIABLE-LIST> WHERE <Join-Conditions>

Metaqueries	Entries
CREATE TABLE Select_List AS SELECT RVarID, CONCAT('COUNT(*) as "count"' AS Entries FROM Relationship UNION DISTINCT SELECT RVarID, IVarID AS Entries FROM Relationship_1Variables;	COUNT(*) as "count" 'popularity(P)' 'teachingability(P)' 'intelligence(S)' 'ranking(S)'
CREATE TABLE From_List AS SELECT RVarID, CONCAT('@database@;',TABLE_NAME) AS Entries FROM Relationship_PVariables UNION DISTINCT SELECT RVarID, CONCAT('@database@;',TABLE_NAME) AS Entries FROM Relationship;	@database@.prof AS P @database@.student AS S @database@.RA AS 'RA'
CREATE TABLE Where_List AS SELECT RVarID, CONCAT(RVarID,',',COLUMN_NAME,'=', Prid,',',REFERENCED_COLUMN_NAME) AS Entries FROM Relationship_PVariables;	'RA'.p_id = P.p_id 'RA'.s_id = S.s_id

Variable List

Meta Query

Count(*) Query

The Random Variable Database

Table Name	Column Headers in Random Variable Database			
Pvariables	Prid	TABLE_NAME		
	C	course		
	P	prof		
	S	student		
1Variables	1VarID	COLUMN_NAME		Prid
	diff(C)	diff		C
	intelligence(S)	intelligence		S
	popularity(P)	popularity		P
	ranking(S)	ranking		S
	rating(C)	rating		C
2Variables	teachingability(P)	teachingability		P
	2VarID	COLUMN_NAME1	COLUMN_NAME2	Prid1 Prid2
	capability(P,S)	p_id	s_id	P S
	grade(C,S)	c_id	s_id	C S
	salary(P,S)	p_id	s_id	P S
	sat(C,S)	c_id	s_id	C S
Relationship	RVarID	TABLE_NAME	COLUMN_NAME1 COLUMN_NAME2	Prid1 Prid2
	RA(P,S)	RA	p_id s_id	P S
	Registered(C,S)	Registered	c_id s_id	C S

Meta data about random variables stored in database tables.

- Domain of possible values.
- Pointer to corresponding data table/column.
- ...

Results

Task: learning a multi-relational Bayesian network

Dataset	# Database Tuples	# Sufficient Statistics (SS)	SS Computing Time (s)	#BN Parameters
MovieLens	1,010,051	252	2.7	292
Mutagenesis	14,540	1,631	1.67	721
UW-CSE	712	2,828	3.84	241
Mondial	870	1,740,870	1,112.84	339
Hepatitis	12,927	12,374,892	3,536.76	569
IMDB	1,354,134	15,538,430	7,467.85	60,059

Database and performance statistics for MRLBase

Comparison with other statistical-relational learning (Markov Logic Networks)

Dataset	RDN_Boost	MLN_Boost	MRLBase	MRLBase-CT
MovieLens	92.7min	N/T	1.12	0.39
Mutagenesis	118	49	1	0.15
UW-CSE	15	19	1	0.27
Mondial	27	42	102	61.82
Hepatitis	251	230	286	186.15
IMDB	N/T	N/T	524.25	439.29

The RDBMS support for multi-relational learning translates into orders of magnitude improvements in speed and scalability.

Speedup on other tasks: compute model selection score, test models, cross-validation. Not shown.

Conclusions

- Multi-relational learning requires new system capabilities.
 - leverage SQL, RDBMS.
- Fast system development through high-level SQL constructs.
- Manage large statistical objects: parameters, sufficient statistics.
- Fast native support for counting (count(*)).
- Future Directions:
 - distributed processing, in-memory computing (SparkSQL)
 - Integrate with inference systems (BayesStore, Tuffy)