

Computing the Shapley Value of Facts in query Answering

SIGMOD 22 Session 22: Provenance and Uncertainty

By 钟瑞峥 林俊光

May 10, 2023

Presentation Overview

- **Background**

- Introduction to Shapley Value
- General Problems
- Problems remained in ICDT 2020

- **Solution**

- Overview
- Theoretical Analysis
- Exact Computation
- Inexact Computation

- **Conclusion**

- **References**

Background - Shapley Value

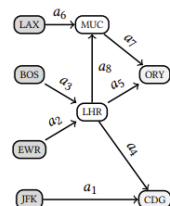
Shapley value is a game-theoretic notion for wealth distribution that is nowadays extensively used to explain complex data-intensive computation.

$$\text{Shapley}(q, D_n, D_x, f) \stackrel{\text{def}}{=} \sum_{E \subseteq D_n \setminus \{f\}} \frac{|E|! (|D_n| - |E| - 1)!}{|D_n|!} (q(D_x \cup E \cup \{f\}) - q(D_x \cup E))$$

- Research show that query evaluation over relational databases fits well in this explanation paradigm.
- Airport problem

FLIGHTS (endo)			AIRPORTS (exo)		
	Src	Dest	Name	Country	
a_1	JFK	CDG	b_1	JFK	USA
a_2	EWB	LHR	b_2	EWB	USA
a_3	BOS	LHR	b_3	BOS	USA
a_4	LHR	CDG	b_4	LAX	USA
a_5	LHR	ORY	b_5	LHR	EN
a_6	LAX	MUC	b_6	MUC	GR
a_7	MUC	ORY	b_7	ORY	FR
a_8	LHR	MUC	b_8	CDG	FR

(a) Database of flights and airports



(b) Flights in graph view. Dark and light gray depict "USA" and "FR" airports respectively

Background - Problems in Calculating Shapley Value

- Calculation of the shapley value is NP-hard in general
- The number of possible coalitions is exponential in the number of facts

ICDT 2020

- Showed mainly lower bounds on the complexity of computing the Shapley value of facts in query answering
 - gives polynomial-time algorithm for self-join-free conjunctive queries
- Need a large number of executions of the query over database subsets
- Does not provide sufficient evidence

Solution - Overview

Two Approaches

- Exact Computation
 - Capture the dependence of query answers using boolean expressions
 - Transform it to a d-DNNF circuit form in which we devise an algorithm for computing shapley value
 - May be too costly in certain circumstances

Note: Given a d-DNNF circuit, we can compute the shapley value in polynomial time with the help of c2d and Provsq1

- Inexact Computation
 - Not necessarily compute exact Shapley values, determine the order of facts according to their Shapley values
 - Faster yet inexact approach that transforms it into CNF Proxy

d-DNNF(Decomposable Deterministic Negation Normal Form)

NNF Definition

- A formula is in NNF if negation only appears in literals

$$(A \vee B) \wedge C$$

$$\neg(A \vee \neg C)$$

DNNF Definition

- A formula in NNF is in DNNF form if the decompositional property holds

$$(A \wedge B) \vee (A \wedge ((\neg B \vee E) \wedge F))$$

d-DNNF(Decomposable Deterministic Negation Normal Form)

d-DNNF Definition

- A DNNF is called deterministic if operands of a disjunction do not share models

$$(A \wedge B) \vee (A \wedge ((\neg B \vee E) \wedge F))$$

This is not in d-DNNF form as two disjunction operands share the model

$$A = 1, B = 1, E = 1, F = 1$$

The c2d Compiler



Description

c2d is a system that compiles CNF into d-DNNF (deterministic, decomposable negation normal form). d-DNNF is a tractable logical form that allows operations such as model counting, clausal entailment, and model enumeration and minimization to be performed in linear time. d-DNNF is also a strict superset of OBDD and is strictly more succinct than OBDD.

Read the [paper](#) that describes the c2d compiler.
Read the [paper](#) that describes d-DNNF and related languages.
See [results](#) of using c2d on some benchmarks.

Download

[Download](#) c2d and instructions for using it.

Related Software

[minic2d](#): A software package for knowledge compilation and model counting based on exhaustive DPLL.

[Ace](#): A system for compiling Bayesian networks which is based on c2d.

[The SDD Package](#): A system for constructing, manipulating, and optimizing Sentential Decision Diagrams (SDDs).

[Jc2d](#): A CNF preprocessor. See results on c2d in: J.-M. Lagniez, E. Lonca, P. Marquis, "Improving Model Counting by Leveraging Definability", Proceedings of IJCAI16, pages 751-757. [paper](#)

Contact

Send questions and comments to darwiche@cs.ucla.edu.

Solution - Overview

Shapley value can be computed in polynomial time whenever the query can be evaluated in polynomial time over tuple-independent probabilistic databases.

- Algorithm 1: compiling to a deterministic and decomposable circuit
- Algorithm 2: resort to CNF proxy which is fast yet inexact if timeout reaches
- Hybrid approach: Algorithm 1 + Algorithm 2

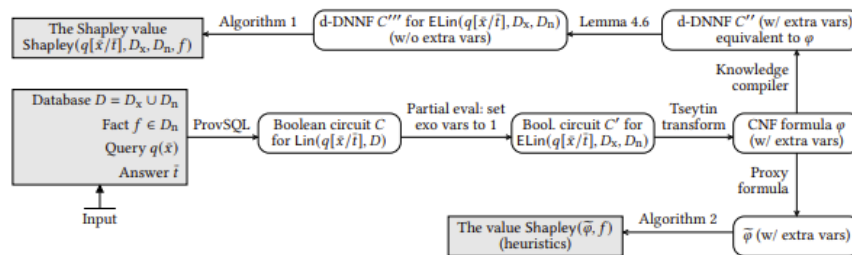


Figure 3: Our implementation architecture.

Note: For non-boolean queries, we are interested in the Shapley Value of the fact f for every individual tuple in the output. Therefore, the computational challenge reduces to that of the boolean queries.

Tseytin Transformation

- Knowledge Compiler usually takes boolean formulas in CNF form, and not arbitrary boolean circuits.
- Tseytin transformation is a method for converting a boolean circuit into an equisatisfiable CNF formula.

$$((p \vee q) \wedge r) \rightarrow (\neg s)$$

Consider all subformulas(excluding simple variables)

$$\neg s$$

$$x_1 \leftrightarrow \neg s$$

$$p \vee q$$

$$x_2 \leftrightarrow p \vee q$$

$$(p \vee q) \wedge r$$

$$x_3 \leftrightarrow x_2 \wedge r$$

$$((p \vee q) \wedge r) \rightarrow (\neg s)$$

$$x_4 \leftrightarrow x_3 \rightarrow x_1$$

All substitutions can be transformed into CNF form

$$x_2 \leftrightarrow p \vee q \equiv (\neg x_2 \vee p \vee q) \wedge (\neg p \vee x_2) \wedge (\neg q \vee x_2)$$

Theoretical Analysis - Reduction

PROBLEM: Shapley(q)
 INPUT: A database $D = D_x \cup D_n$ and an endogenous fact $f \in D_n$.
 OUTPUT: The value $\text{Shapley}(q, D_n, D_x, f)$.

PROBLEM: PQE(q)
 INPUT: A tuple-independent database (D, π) .
 OUTPUT: The value $\Pr(q, (D, \pi))$.

For every self-join free boolean queries, either q is hierarchical and $\text{shapley}(q)$ can be solved in polynomial time, or q is not hierarchical and $\text{shapley}(q)$ is intractable

Proposition: For every Boolean query q , we have that $\text{Shapley}(q) \leq_T^P \text{PQE}(q)$

Proof:

$$\# \text{ Slices}(q, D_x, D_n, k) \stackrel{\text{def}}{=} |\{E \subseteq D_n \mid |E| = k \text{ and } q(D_x \cup E) = 1\}|.$$

$$\text{Shapley}(q, D_n, D_x, f) =$$

$$\sum_{k=0}^{|D_n|-1} \frac{k! (|D_n| - k - 1)!}{|D_n|!} (\# \text{ Slices}(q, D_x \cup \{f\}, D_n \setminus \{f\}, k) - \# \text{ Slices}(q, D_x, D_n \setminus \{f\}, k)).$$

Arithmetic terms can be computed in polynomial time

Theoretical Analysis - Reduction

We want to prove that #Slices can be computed in polynomial time

$$(1 + z)^n \Pr(q, (D_z, \pi_z)) = \sum_{i=0}^n z^i \# \text{Slices}(q, D_x, D_n, i)$$

We can now call an oracle to PQE(q) on $n+1$ databases D_{z_0}, \dots, D_{z_n} for distinct values $z_0 = 0, z_1 = 1, \dots, z_n = n$.

Exact Computation

Recently, research proved that when the models are given as circuits from knowledge compilation can be computed in polynomial time (SHAP scores).

The same approach can be applied to the computation of Shapley value.

Proposition: Given as input a deterministic and decomposable circuit C representing $\text{ELin}(q, D_n, D_x)$ for a database $D = D_x \cup D_n$ and Boolean query q , and an endogenous fact $f \in D_n$, we can compute in polynomial time (in $|C|$) the value $\text{Shapley}(q, D_n, D_x, f)$.

We can rewrite the equation as following:

$$\text{Shapley}(q, D_n, D_x, f) = \sum_{k=0}^{|D_n|-1} \frac{k! (|D_n| - k - 1)}{|D_n|} (\#SAT_k(C_1) - \#SAT_k(C_2)).$$

Lemma: Given as input a d-DNNF Boolean Circuit C , and an integer k , we can compute in polynomial time the quantity $SAT_k(C)$.

Exact Computation

Proof: Let $X = \text{Vars}(C)$ and $n = |X|$ We denote by

ϕ_g the boolean function over the variables $\text{Var}(g)$ that is represented by this gate

$\alpha_l^g = \#SAT_l(\phi_g)$ the number of assignments of size l to $\text{Vars}(g)$ that satisfy ϕ_g

Variable gate.

- α_g^0 is 0 and α_g^1 is 1 .

\neg gate

- If g is a \neg -gate with input gate g' , then $\alpha_g^\ell = \binom{|\text{Vars}(g)|}{\ell} - \alpha_{g'}^\ell$ for every $\ell \in \{0, \dots, |\text{Vars}(g)|\}$.

Exact Computation

Deterministic \vee -gate

Define $S_1 \stackrel{\text{def}}{=} \text{Vars}(g_2) \setminus \text{Vars}(g_1)$ and similarly $S_2 \stackrel{\text{def}}{=} \text{Vars}(g_1) \setminus \text{Vars}(g_2)$. Since g is deterministic, we have:

$$\text{SAT}(\varphi_g) = (\text{SAT}(\varphi_{g_1}) \otimes 2^{S_1}) \cup (\text{SAT}(\varphi_{g_2}) \otimes 2^{S_2})$$

with the union being disjoint. By intersecting with the assignments of $\text{Vars}(g)$ of size ℓ , we obtain:

$$\begin{aligned} \text{SAT}_\ell(\varphi_g) = & [(\text{SAT}(\varphi_{g_1}) \otimes 2^{S_1}) \cap \{v \subseteq \text{Vars}(g) \mid |v| = \ell\}] \\ & \cup [(\text{SAT}(\varphi_{g_2}) \otimes 2^{S_2}) \cap \{v \subseteq \text{Vars}(g) \mid |v| = \ell\}] \end{aligned}$$

$$\begin{aligned} \# \text{SAT}_\ell(\varphi_g) = & |(\text{SAT}(\varphi_{g_1}) \otimes 2^{S_1}) \cap \{v \subseteq \text{Vars}(g) \mid |v| = \ell\}| \\ & + |(\text{SAT}(\varphi_{g_2}) \otimes 2^{S_2}) \cap \{v \subseteq \text{Vars}(g) \mid |v| = \ell\}| \end{aligned}$$

$$|(\text{SAT}(\varphi_{g_1}) \otimes 2^{S_1}) \cap \{v \subseteq \text{Vars}(g) \mid |v| = \ell\}| = \sum_{i=\max(0, \ell-|S_1|)}^{\min(\ell, |\text{Vars}(g_1)|)} \alpha_{g_1}^i \times \binom{|S_1|}{\ell-i}.$$

Exact Computation

Decomposable \wedge -gate.

$$\text{SAT}(\varphi_g) = \text{SAT}(\varphi_{g_1}) \otimes \text{SAT}(\varphi_{g_2})$$

We now intersect with the set of assignments of $\text{Vars}(g)$ of size ℓ to obtain

$$\alpha_g^\ell = \# \text{SAT}_\ell(\varphi_g) = \sum_{i=\max(0, \ell-|\text{Vars}(g_2)|)}^{\min(\ell, |\text{Vars}(g_1)|)} \alpha_{g_1}^i \times \alpha_{g_2}^{\ell-i}$$

Exact Computation

Algorithm 1: Shapley values from deterministic and decomposable Boolean circuits

Input : Deterministic and decomposable Boolean circuit C with output gate g_{output} representing $\text{ELin}(q, D_X, D_n)$ and an endogenous fact $f \in D_n$.

Output: The value $\text{Shapley}(q, D_X, D_n, f)$.

```

1 Complete  $C$  so that  $\text{Vars}(g_{\text{output}}) = D_n$ ;
2 Compute  $C_1 = C[f \rightarrow 1]$  and  $C_2 = C[f \rightarrow 0]$ ;
   // Partial evaluations of  $C$  by setting  $f$  to 1
   and to 0
3  $\Gamma = \text{ComputeAll}\#\text{SAT}_k(C_1)$ ;           // As an array
4  $\Delta = \text{ComputeAll}\#\text{SAT}_k(C_2)$ ;           // As an array
5 return  $\sum_{k=0}^{|D_n|-1} \frac{k! (|D_n| - k - 1)!}{|D_n|!} \cdot (\Gamma[k] - \Delta[k])$ ;

6 Def  $\text{ComputeAll}\#\text{SAT}_k(C)$ :
7   Preprocess  $C$  so that each  $\vee$ -gate and  $\wedge$ -gate has
   fan-in exactly 0 or 2;
8   Compute the set  $\text{Vars}(g)$  for every gate  $g$  in  $C$ ;
9   Compute values  $\alpha_g^\ell$  for every gate  $g$  in  $C$ 
   and  $\ell \in \{0, \dots, |\text{Vars}(g)|\}$  by bottom-up induction
   on  $C$  using the inductive relations from the proof of
   Lemma 4.5;
10  return  $[\alpha_{g_{\text{output}}}^0, \dots, \alpha_{g_{\text{output}}}^{|\text{Vars}(C)|}]$ ;
```


Inexact Computation

Based on the observation that having a high shapley score is correlated with

- Appearing many times of the provenance
- Having few alternatives

$$\text{Shapley}(h, x) \stackrel{\text{def}}{=} \sum_{S \subseteq X \setminus \{x\}} \frac{|S|!(|X|-|S|-1)!}{|X|!} (h(S \cup \{x\}) - h(S)).$$

$$\Phi(\psi_i, x) = \begin{cases} \frac{1}{(a_i+b_i) \cdot \binom{a_i+b_i-1}{b_i}} & \text{if } x \text{ appears in } \psi_i \text{ in positive form;} \\ \frac{-1}{(a_i+b_i) \cdot \binom{a_i+b_i-1}{a_i}} & \text{if } x \text{ appears in } \psi_i \text{ in negative form;} \\ 0 & \text{otherwise.} \end{cases}$$

Inexact Computation

Algorithm 2: CNF Proxy

Input : CNF φ and a set of endogenous facts D_n .

Output: The value $\text{Shapley}(\tilde{\varphi}, x)$ for each $x \in D_n$.

```
1  $n \leftarrow |\varphi.\text{clauses}();$   
2  $v \leftarrow 0^{|D_n|};$  // As an array  
3 for  $\psi \in \varphi.\text{clauses}()$  do  
4    $\mathcal{L} \leftarrow \psi.\text{literals}();$   
5    $m \leftarrow |\mathcal{L}|;$   
6    $\text{pos} \leftarrow \{\ell \in \mathcal{L} \mid \ell \text{ is positive}\};$   
7    $\text{neg} \leftarrow \{\ell \in \mathcal{L} \mid \ell \text{ is negative}\};$   
8   for  $\ell \in \text{pos} \cap D_n$  do  
9      $v[\ell.\text{var}()] \leftarrow v[\ell.\text{var}()] + \frac{1}{nm \binom{m-1}{|\text{neg}|}};$   
10  end  
11  for  $\ell \in \text{neg} \cap D_n$  do  
12     $v[\ell.\text{var}()] \leftarrow v[\ell.\text{var}()] - \frac{1}{nm \binom{m-1}{|\text{pos}|}};$   
13  end  
14 end  
15 return  $v$ 
```

Conclusion - Experimental Results

- knowledge compilation using the c2d compiler
- datasets
 - TPC-H (removing nested queries and queries with aggregation)
 - IMDB
- Algorithm1 yields an overall 84.43% success rate for TPC-H and 99.96% for IMDB,
- Algorithm2 is very efficient and that the ranking based on CNF proxies is very close to the ranking based on the exact Shapley value