

Database Learning: Toward a Database that Becomes Smarter Over Time

Yongjoo Park

Ahmad Shahab Tajik

Michael Cafarella

Barzan Mozafari

University of Michigan, Ann Arbor

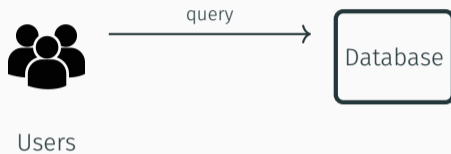
Today's databases



Users



Today's databases



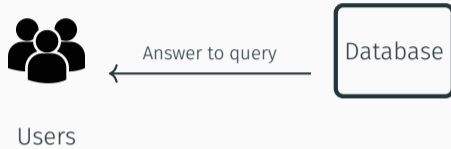
Today's databases



Users



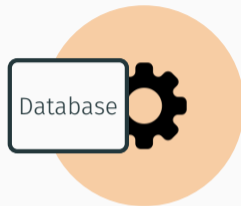
Today's databases



Today's databases



Users

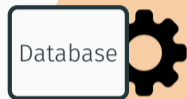


After answering queries,
THE WORK is **GONE**.

Today's databases



Users



After answering queries,
THE WORK is **GONE**.

Our Goal: reuse **the work**

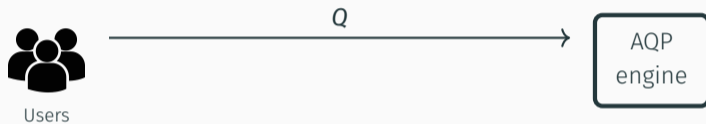
Our high-level approach



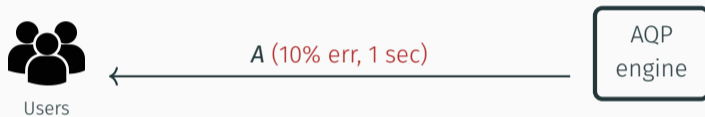
Users



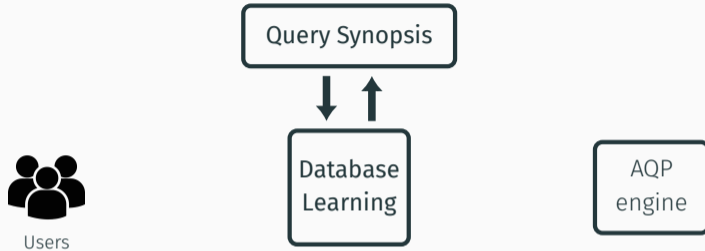
Our high-level approach



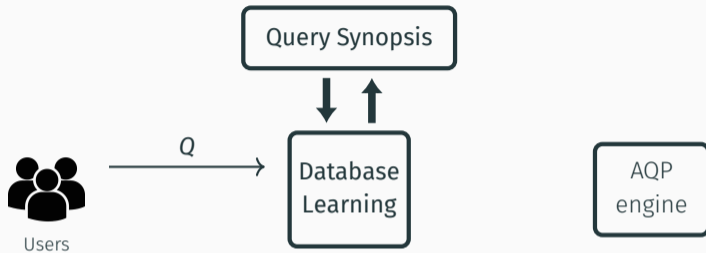
Our high-level approach



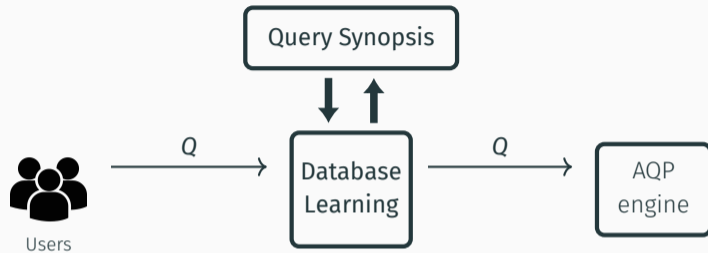
Our high-level approach



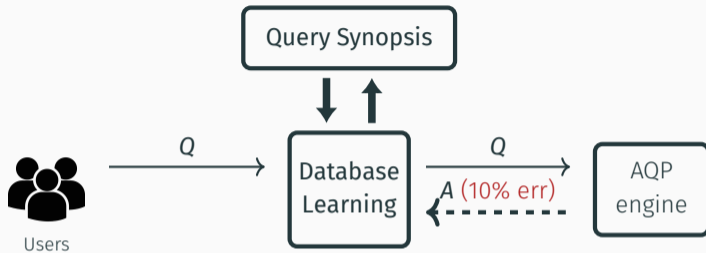
Our high-level approach



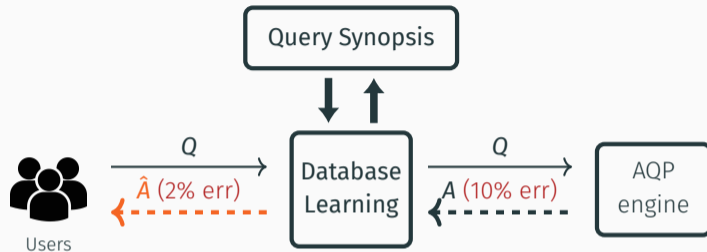
Our high-level approach



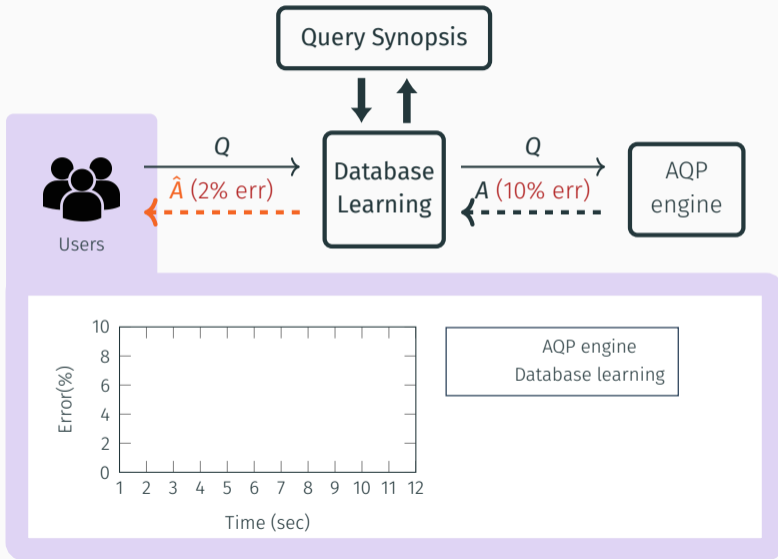
Our high-level approach



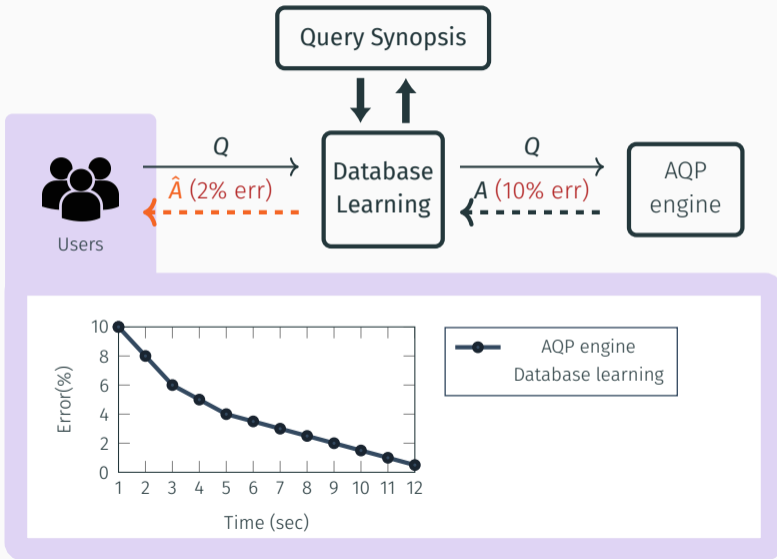
Our high-level approach



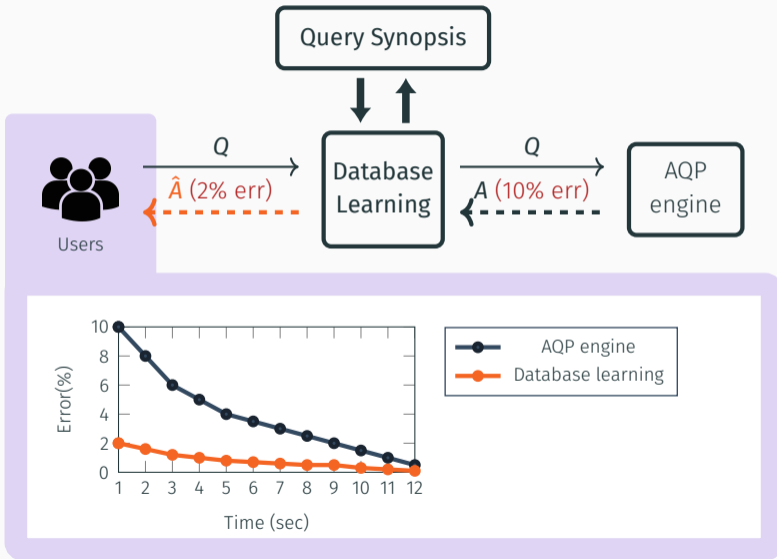
Our high-level approach



Our high-level approach



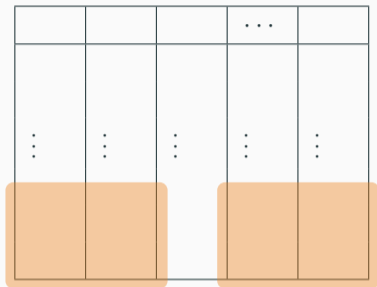
Our high-level approach



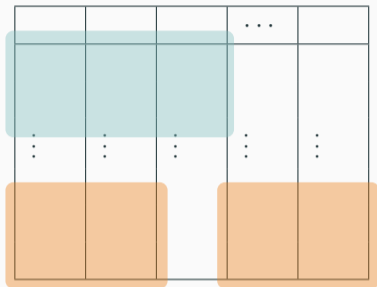
Technical challenges

			...	
⋮	⋮	⋮	⋮	⋮

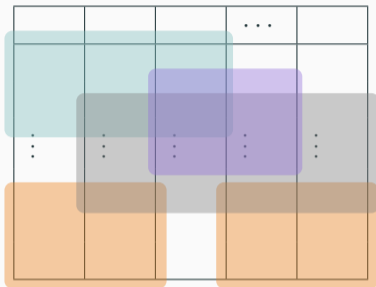
Technical challenges



Technical challenges

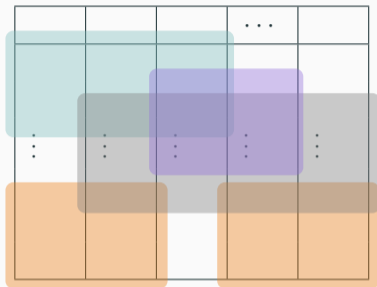


Technical challenges



Queries use the data in different columns/rows.

Technical challenges



Queries use the data in different columns/rows.

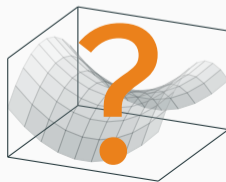
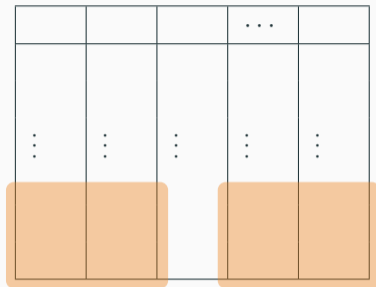
How to leverage those queries for future queries?

Our idea

			...	
⋮	⋮	⋮	⋮	⋮



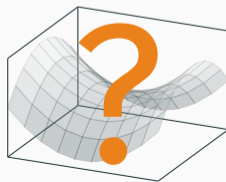
Q1



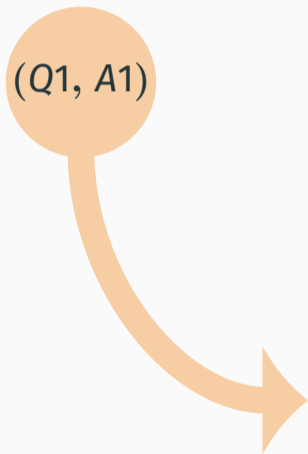
$(Q1, A1)$



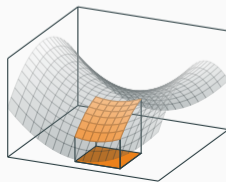
			...	
⋮	⋮	⋮	⋮	⋮



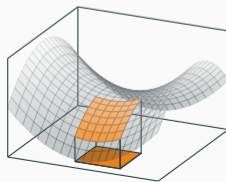
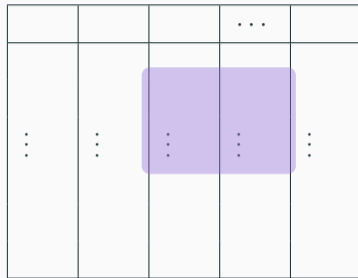
Our idea



			...	
⋮	⋮	⋮	⋮	⋮



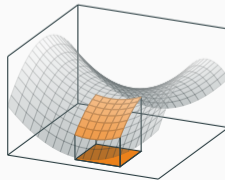
Q2



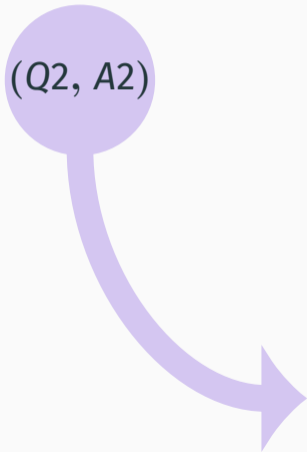
$(Q2, A2)$



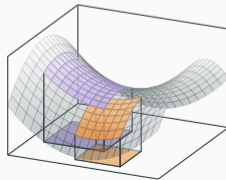
			...	
⋮	⋮	⋮	⋮	⋮



Our idea



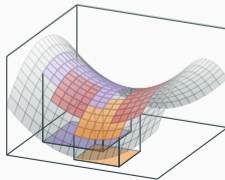
			...	
⋮	⋮	⋮	⋮	⋮



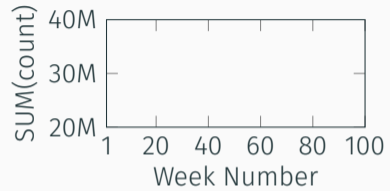
Our idea

more queries
and answers

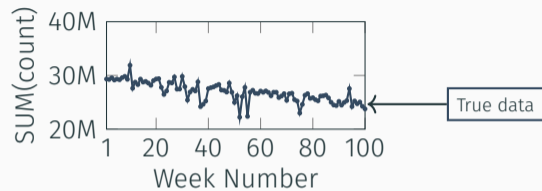
			...	
⋮	⋮	⋮	⋮	⋮



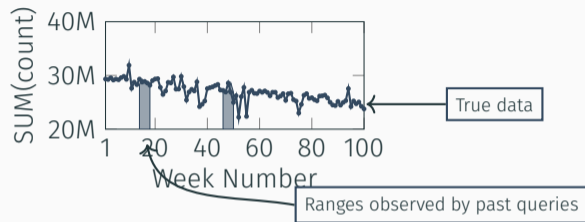
Concrete example



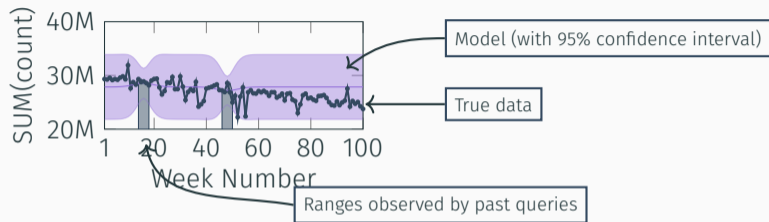
Concrete example



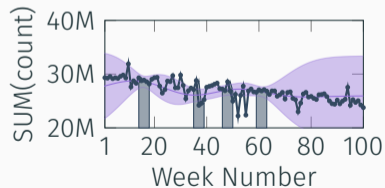
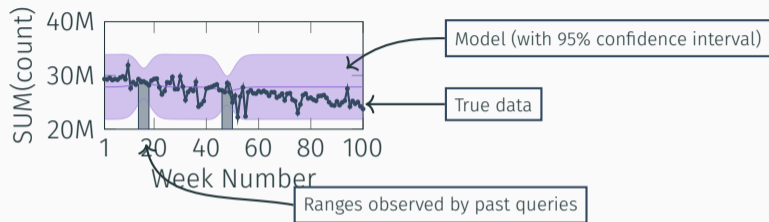
Concrete example



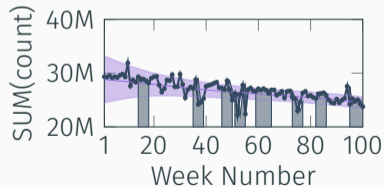
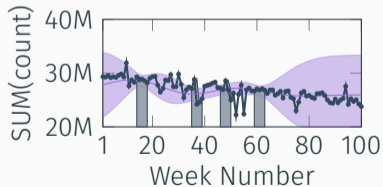
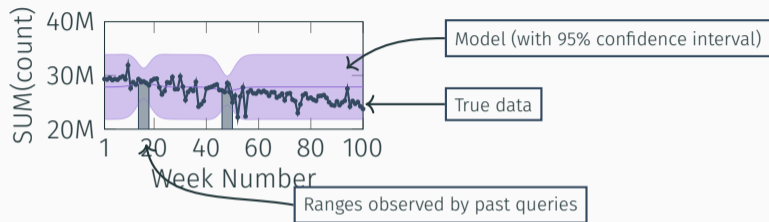
Concrete example



Concrete example



Concrete example



Design goals

```
select X3, avg(Y1)
from t
where 5 <
```

```
select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

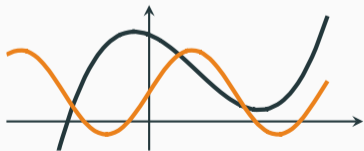
1. Support a **wide class of** SQL queries

Design goals

```
select X3, sum(Y1)
from t
where 5 <
```

```
select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

1. Support a **wide class of** SQL queries



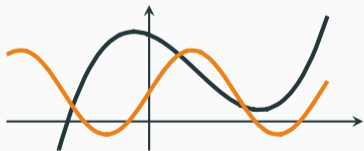
2. **No Assumptions** about Data

Design goals

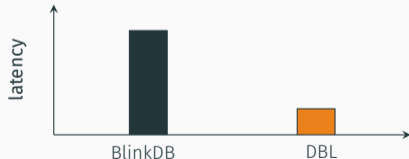
```
select X3, sum(Y1)
from t
where 5 < X2;

select sum(Y2)
from t
where X2 between Apr and May
group by X3;
```

1. Support a **wide class of** SQL queries



2. **No Assumptions** about Data



3. **Lightweight**

Our Approach

Problem statement

Problem statement

Problem:

Given past queries (q_1, \dots, q_n) , a new query (q_{n+1}) , and their approximate answers,

Find the **most likely** answer to the new query (q_{n+1}) and **its estimated error**.

Problem statement

Problem:

Given past queries (q_1, \dots, q_n) , a new query (q_{n+1}) , and their approximate answers,
Find the **most likely** answer to the new query (q_{n+1}) and **its estimated error**.

Our result:

Under a *certain model assumption*,

our answer's error bound \leq **original answer's error bound**
(in practice, much more accurate)

if the error bounds provide the same probabilistic guarantees.

Overview of our technique

```
select avg(Y2)
from t
where 6 < X1 < 8;
```

Overview of our technique

```
select sum(Y2)
from t
where 5 < X1 < 8;
```

Overview of our technique

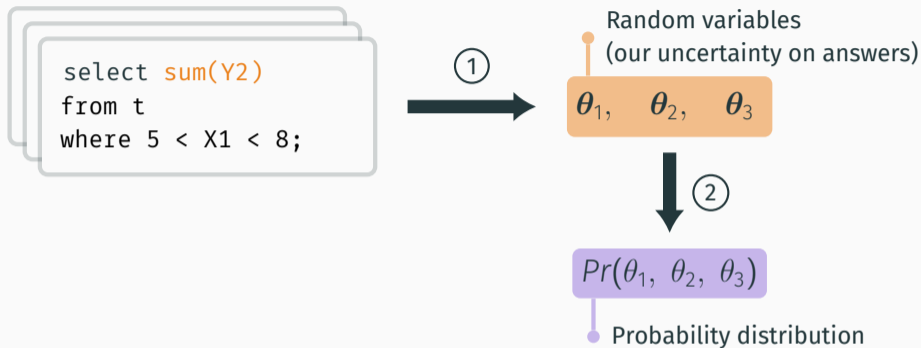
```
select sum(Y2)
from t
where 5 < X1 < 8;
```



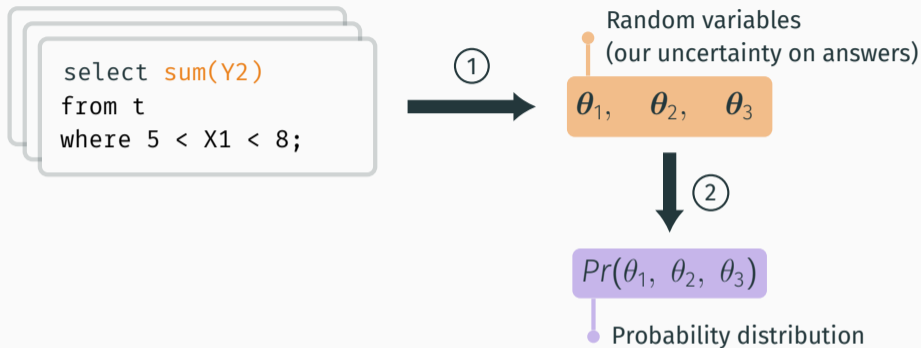
Random variables
(our uncertainty on answers)

$\theta_1, \theta_2, \theta_3$

Overview of our technique

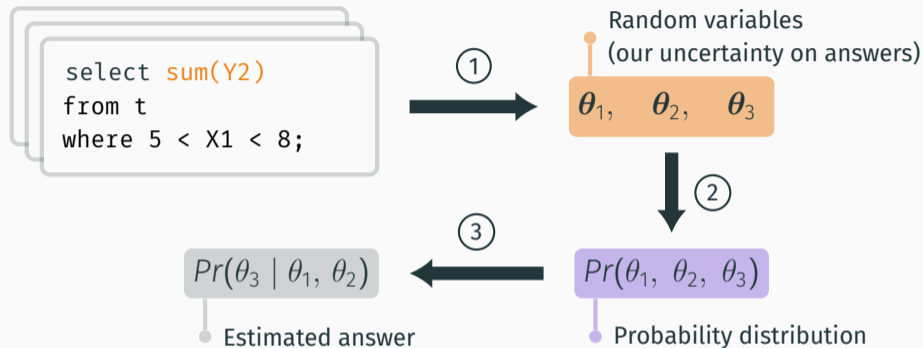


Overview of our technique



Two aggregations involve **common values**
→ **correlation** between answers

Overview of our technique



Two aggregations involve **common values**
→ **correlation** between answers

How to define random variables

```
select sum(Y2)
from t
where 5 < X1 < 8;
```

How to define random variables

```
select sum(Y2)
from t
where 5 < X1 < 8;
```

We define a random variable θ
for every combination of:

How to define random variables

We define a random variable θ
for every combination of:

```
select sum(Y2)  
from t  
where 5 < X1 < 8;
```

● Aggregate function

How to define random variables

We define a random variable θ
for every combination of:

```
select sum(Y2)  
from t  
where 5 < X1 < 8;
```

● Aggregate function

● Selection predicates

How to define random variables

We define a random variable θ
for every combination of:

```
select sum(Y2)
from t
where 5 < X1 < 8;
```

● Aggregate function

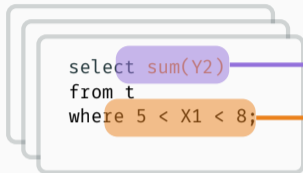
● Selection predicates

```
select X3, avg(Y1), sum(Y2)
from t
where 5 < X1 < 8
      and X2 between Apr and May
group by X3;
```

What if your query is complex?

How to define random variables

We define a random variable θ
for every combination of:



Aggregate function

Selection predicates



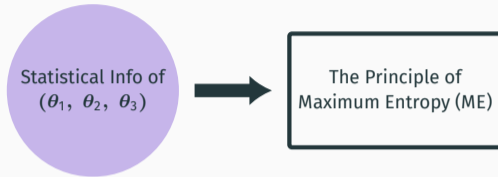
```
select X3, avg(Y1), sum(Y2)
from t
where 5 < X1 < 8
      and X2 between Apr and May
group by X3;
```

What if your query is complex?

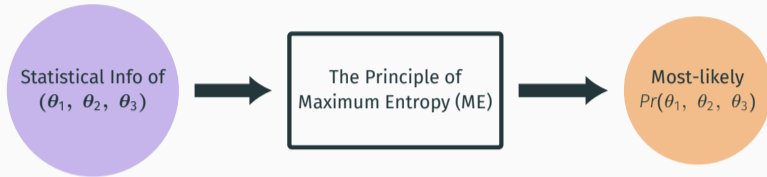
How to determine the probability distribution

The Principle of
Maximum Entropy (ME)

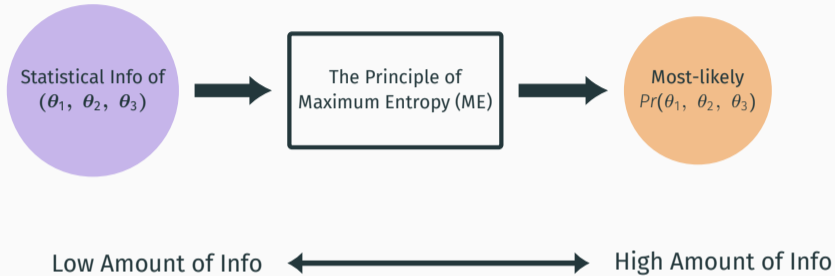
How to determine the probability distribution



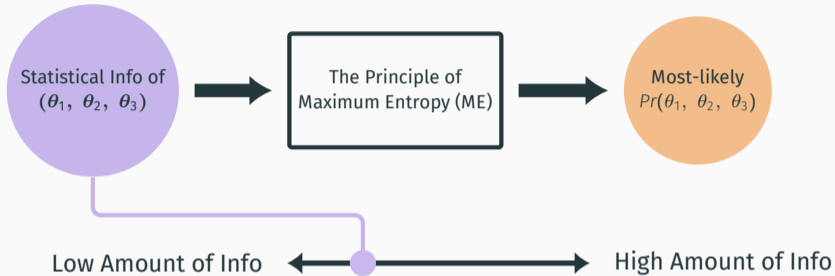
How to determine the probability distribution



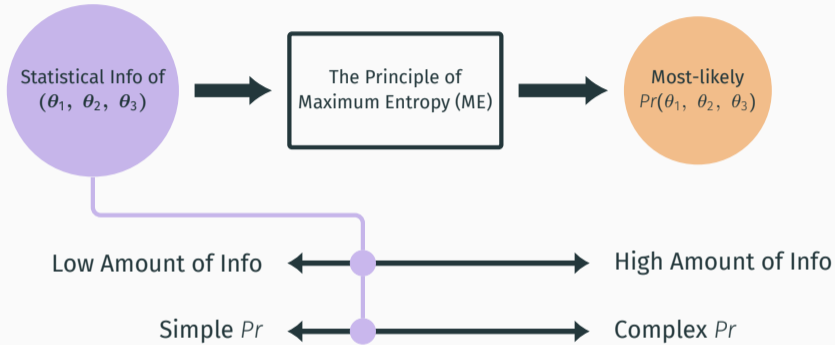
How to determine the probability distribution



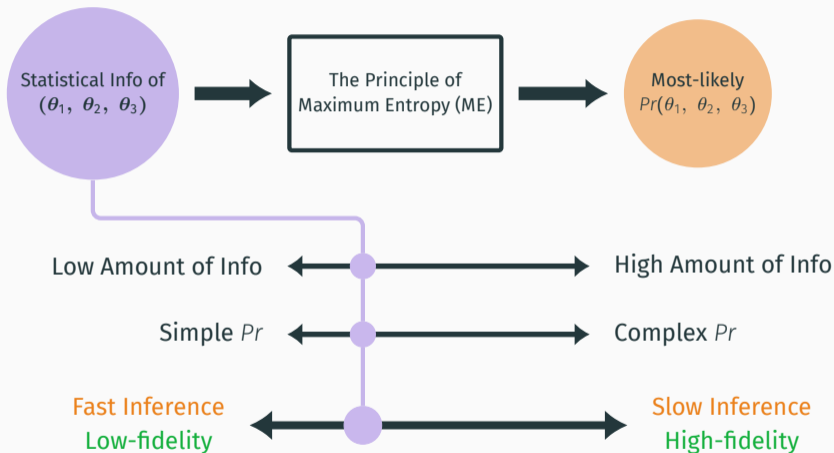
How to determine the probability distribution



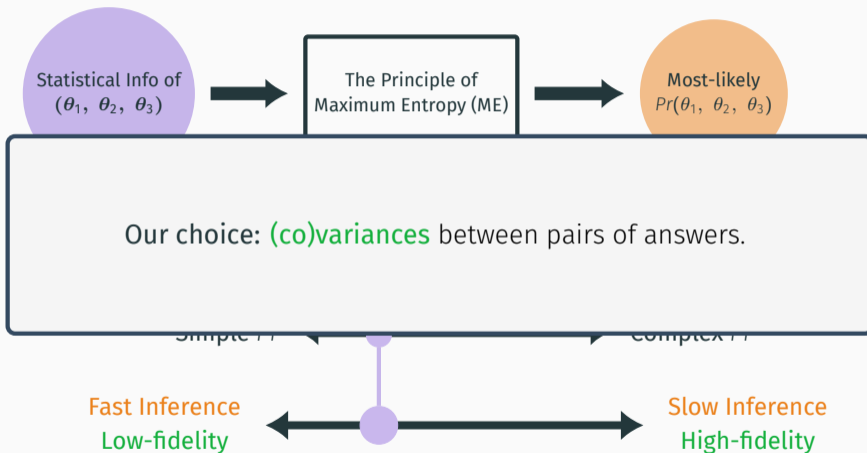
How to determine the probability distribution



How to determine the probability distribution



How to determine the probability distribution



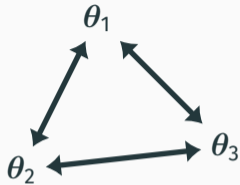
Most-likely probability distribution

θ_1

θ_2

θ_3

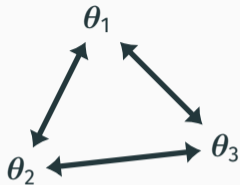
Most-likely probability distribution



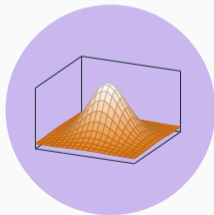
Statistical Information:

Mean, variances, covariances

Most-likely probability distribution



↓ MaxEnt

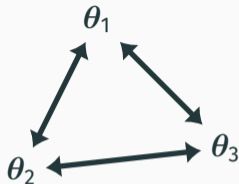


Statistical Information:

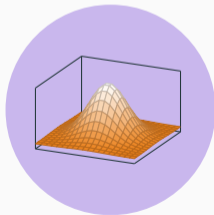
Mean, variances, covariances

Multivariate normal distribution

Most-likely probability distribution



MaxEnt



Statistical Information:

Mean, variances, covariances

Multivariate normal distribution

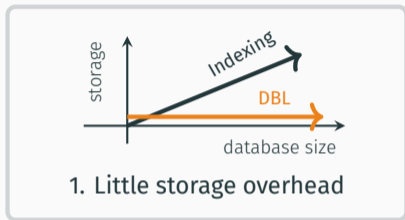
Fast inference using a **closed form**

Benefits of database learning

Database learning vs. indexing

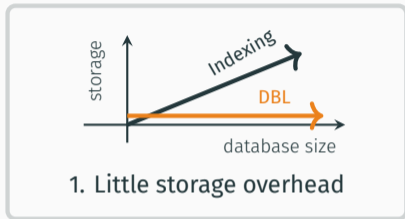
Benefits of database learning

Database learning vs. indexing



Benefits of database learning

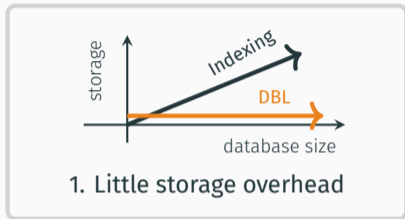
Database learning vs. indexing



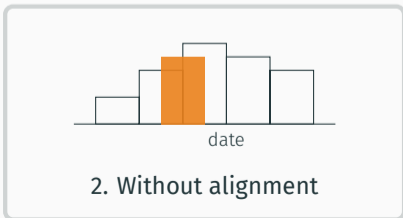
Database learning vs. materialized view

Benefits of database learning

Database learning vs. indexing

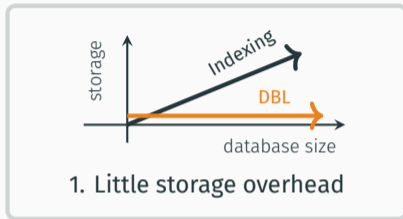


Database learning vs. materialized view

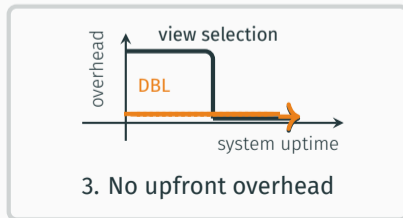
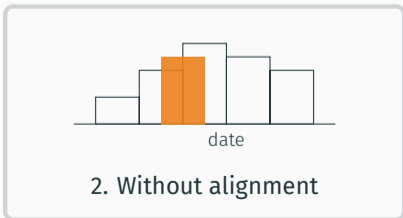


Benefits of database learning

Database learning vs. indexing



Database learning vs. materialized view



Experiment

Experiment setup

1. Two systems:

- NOLEARN: Approximate query processing engine (The longer runtime, the more accurate answer)
- VERDICT: Our database learning system (on top of NOLEARN)

Experiment setup

1. Two systems:

- NOLEARN: Approximate query processing engine (The longer runtime, the more accurate answer)
- VERDICT: Our database learning system (on top of NOLEARN)

2. Datasets:

- **Customer1**: 536GB data and query log from a customer
- TPC-H: 100GB TPC-H dataset

Experiment setup

1. Two systems:

- NOLEARN: Approximate query processing engine (The longer runtime, the more accurate answer)
- VERDICT: Our database learning system (on top of NOLEARN)

2. Datasets:

- **Customer1**: 536GB data and query log from a customer
- **TPC-H**: 100GB TPC-H dataset

3. Environment:

- 5 Amazon EC2 workers (**m4.2xlarge**) + 1 master
- SSD-backed HDFS for Spark's data loading

Our experimental claims

1. VERDICT supports a large portion of real-world queries

Our experimental claims

1. VERDICT supports a large portion of real-world queries
2. VERDICT achieves speedup compared to NOLEARN

Our experimental claims

1. VERDICT supports a large portion of real-world queries
2. VERDICT achieves speedup compared to NOLEARN
3. VERDICT works with small memory and computational overhead

Generality of VERDICT

Dataset	# Analyzed	# Supported	Percentage
Customer1	3,342	2,463	73.7%
TPC-H	21	14	66.7%

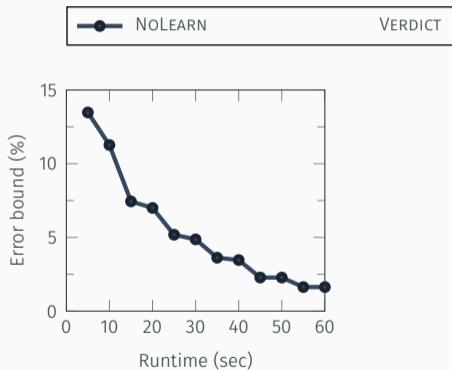
Unsupported queries:

1. Nested queries (that cannot be flattened)
2. Textual filters: `city like '%arbor%'`

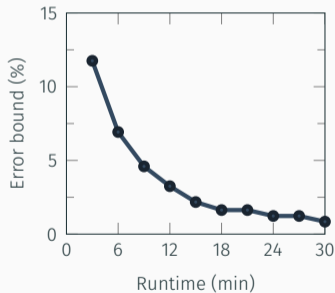
Runtime-error trade-off

Results on the TPC-H dataset (the paper has the `Customer1` results)

Number of past queries fixed to 50



(a) Data in Memory

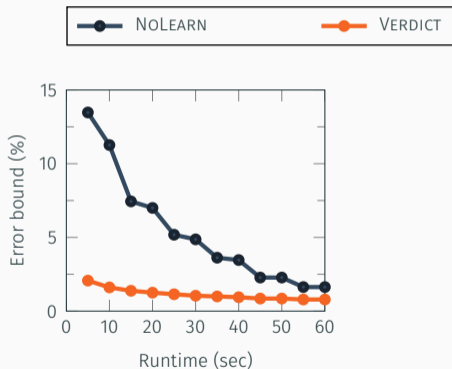


(b) Data on SSD

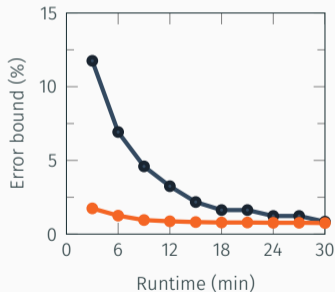
Runtime-error trade-off

Results on the TPC-H dataset (the paper has the `Customer1` results)

Number of past queries fixed to 50



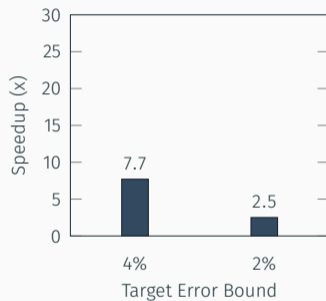
(a) Data in Memory



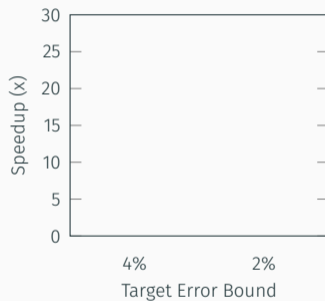
(b) Data on SSD

Speedup

The results on the `Customer1` dataset (the paper has the TPC-H results)



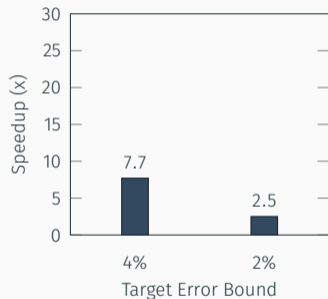
(a) Data in memory



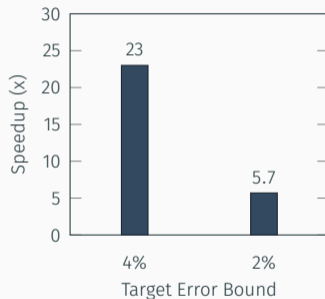
(b) Data on SSD

Speedup

The results on the `Customer1` dataset (the paper has the TPC-H results)



(a) Data in memory



(b) Data on SSD

Memory and computational overhead

1. Memory overhead:

Memory and computational overhead

1. Memory overhead:

- Queries and their answer, some matrices and their inverses

Memory and computational overhead

1. Memory overhead:

- Queries and their answer, some matrices and their inverses
- 23.2 KB per query for the **Customer1** dataset
- 15.8 KB per query for the **TPC-H** dataset

Memory and computational overhead

1. Memory overhead:

- Queries and their answer, some matrices and their inverses
- 23.2 KB per query for the **Customer1** dataset
- 15.8 KB per query for the **TPC-H** dataset

2. Computational overhead:

	Latency for memory	Latency for SSD
NOLEARN	2.083 sec	52.50 sec
VERDICT	2.093 sec	52.51 sec
Overhead	0.010 sec (0.48%)	0.010 sec (0.02%)

Thank You!