

Differentiating Relational Queries

PhD Workshop at VLDB21

Paul Peseux

Tutors: T.Paquet, M.Berar, V.Nicollet



x



August 16, 2021

Outline

- 1 Context
- 2 Formalization
- 3 Tables Relations
- 4 Automatic Differentiation
- 5 Implementation
- 6 Conclusion

Outline

1 Context

2 Formalization

3 Tables Relations

4 Automatic Differentiation

5 Implementation

6 Conclusion

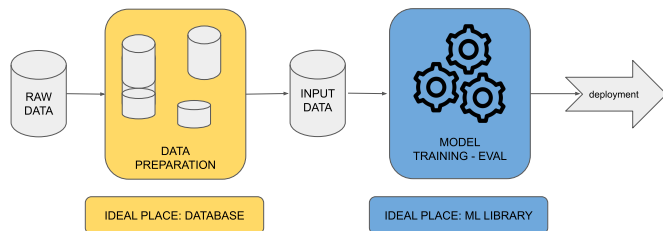


Figure: Classic Machine Learning Pipeline.

Context

- costly data transfer (Schüle 2019)
-

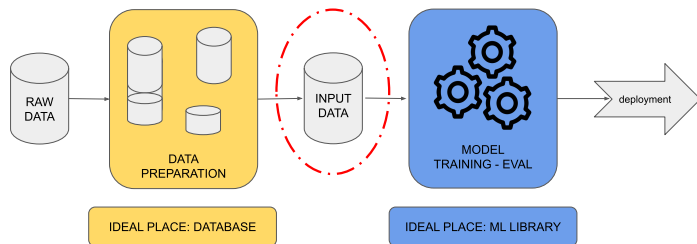


Figure: Classic Machine Learning Pipeline.

Context

- costly data transfer (Schüle 2019)
- ML libraries built for computer vision, NLP ...
/ **inadapted to relational data**

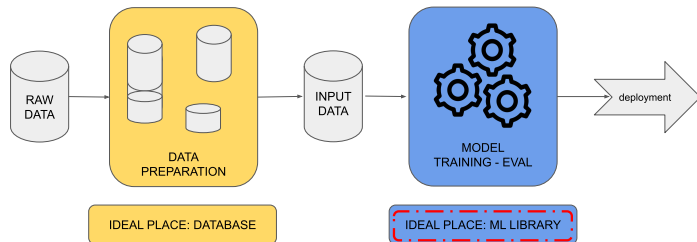


Figure: Classic Machine Learning Pipeline.

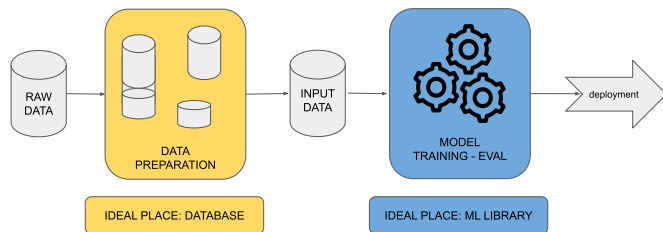


Figure: Classic Machine Learning Pipeline.

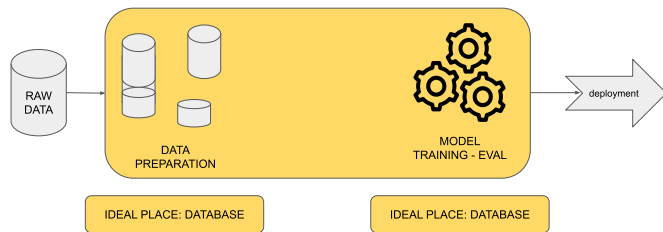


Figure: Proposed Pipeline.

Many Machine Learning methods are based on gradient methods.

Figure: Gradient Descent, source (Hutson)

Many Machine Learning methods are based on gradient methods.

Figure: Gradient Descent, source (Hutson)

- ! To optimize models, relational queries differentiation is missing (Schale 2019)

Differentiating Relational Queries Derivative of the Relational Queries

"

This is not differential data flow (Mcsherry 2021)

Figure: What we are looking for

Outline

- 1 Context
- 2 **Formalization**
- 3 Tables Relations
- 4 Automatic Differentiation
- 5 Implementation
- 6 Conclusion

Formalization

For the rest of the presentation, optimisation means minimisation and is allowed through gradient descent.

$$x^* = \arg \min_x f(x)$$

Figure: Gradient Descent, source (Hutson)

f is called loss

We want to minimize ϕ and thus compute the gradient of

For that we need:

- a framework

- constraints on the query

Minimization is only possible on scalar.

$$\text{Loss} = \sum_{i=1}^{\text{Obs}} X \quad \text{loss} = \sum_{i=1}^{\text{Obs}} X \quad f(\text{data}_i)$$

Minimization is only possible on scalar.

$$\text{Loss} = \sum_{i=1}^n \text{loss}_i = \sum_{i=1}^n f(\text{data}_i)$$

Constraint 1

Loss is computed line by line

Example

Let's make it concrete with the Chicago taxi trip dataset.

Figure: Chicago trips dataset, source (Chicago)

Example

Objective: explain the trip's tip with distance and company "quality"

With Linear Regression as the machine learning model.

Linear Regression on the Chicago dataset

Model

$$\text{Tip}_{\text{estimated}} = a_{\text{company}} \text{ distance} + b$$

One slope per company; Intercept is shared among all the taxis.

Comparing the matrix approach (ML Libraries) and relational one

$$(M:A) X + b$$

is the point-wise product

Figure: Matrix approach

Figure: Relational approach

Figure: Matrix approach

Figure: Relational approach

Model

$$\text{Tip}_{\text{estimated}} = a_{\text{company}} \text{ distance} + b$$

In SQL it gives

SQL query of our model.

Formalization

" Trips = Observations

$$\text{Loss} = \sum_{t \in \text{Trips}} \text{loss}_t = \sum_{t \in \text{Trips}} f(\text{data}_t) = \sum_{t \in \text{Trips}} (a_{\text{comp}_t} \text{dist}_t + b \text{tip}_t)^2$$

with

$$f(a; x; b; y) = (ax + b - y)^2$$

" Trips = Observations

$$\text{Loss} = \sum_{t \in \text{Trips}} \log \xi = \sum_{t \in \text{Trips}} f(\text{data}_t) = \sum_{t \in \text{Trips}} (a_{\text{comp}_t} \text{dist}_t + b \text{tip}_t)^2$$

with

$$f(a; x; b; y) = (ax + b y)^2$$

Then it is feasible to compute gradients!

$$\frac{\partial}{\partial a} ; \quad \frac{\partial}{\partial x} ; \quad \frac{\partial}{\partial b} ; \quad \frac{\partial}{\partial y}$$

Constraint 2

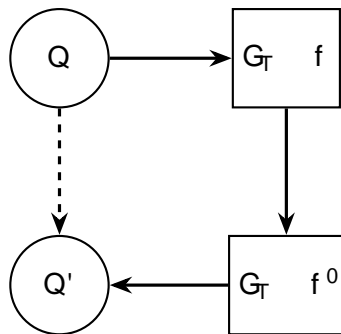
f has to be differentiable.

Figure: Inputs origin.

Figure: Inputs origin.

SQL query of our model.

Approach



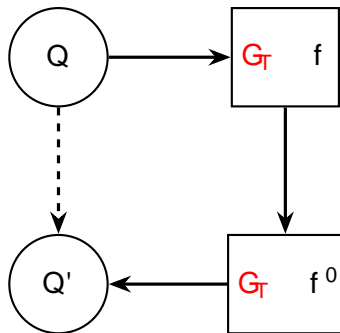
Q : query
 G_T : tables graph
 f : loss function

Figure: Path to Differentiating Relational Queries.

Outline

- 1 Context
- 2 Formalization
- 3 **Tables Relations**
- 4 Automatic Differentiation
- 5 Implementation
- 6 Conclusion

Approach



Q : query
G_T : tables graph
f : loss function

Figure: Path to Differentiating Relational Queries.

Definition 1 (Broadcast)

Let's note " $T_A \rightarrow T_B$ " when the primary key of T_A is a foreign key in T_B . It is said that T_A broadcasts into T_B .

Definition 1 (Broadcast)

Let's note " $T_A \rightarrow T_B$ " when the primary key of T_A is a foreign key in T_B . It is said that T_A broadcasts into T_B .

Tables used in the query with the relationship \rightarrow forms a graph G_T .

Tables Relations

Definition 1 (Broadcast)

Let's note " $T_A \rightarrow T_B$ " when the primary key of T_A is a foreign key in T_B . It is said that T_A broadcasts into T_B .

Tables used in the query with the relationship forms a graph G_T .

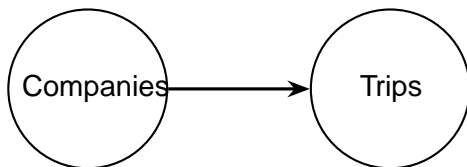


Figure: Graph from our linear regression model.

Figure: Inputs origin.

Tables Relations

Let be

T a table used in the query

$T:A$ be a column of T

a the input of f representing $T:A$

If T (transitively) broadcasts into Observation then a the input of f representing $T:A$ is a scalar.

Tables Relations

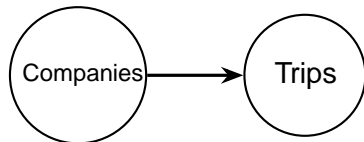
Let be

T a table used in the query

T:A be a column of T

a the input of f representing T:A

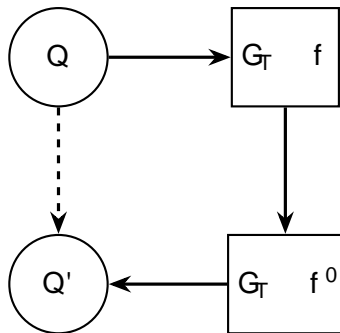
If T (transitively) broadcasts into Observations then a the input of f representing T:A is a scalar.



$$\text{Tip}_{\text{estimated}} = a_{\text{company}} \text{ distance} + b$$

Figure: Graph from our linear regression model.

Approach



Q : query
 G_T : tables graph
 f : loss function

Figure: Path to Differentiating Relational Queries.

Outline

- 1 Context
- 2 Formalization
- 3 Tables Relations
- 4 Automatic Differentiation**
- 5 Implementation
- 6 Conclusion

Approach

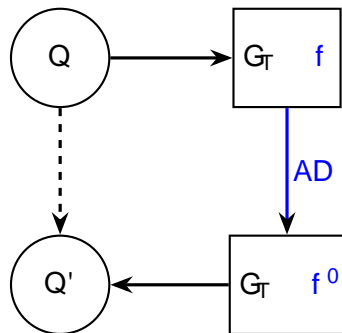


Figure: Path to Differentiating Relational Queries.

Q : query

G_T : tables graph

f : loss function

AD : Automatic Differentiation

Automatic Differentiation

P a program that apply the mathematical function to its inputs.

Automatic Differentiation constructs program the program P^0 that apply f^0 to its inputs.

Automatic Differentiation

P a program that apply the mathematical function to its inputs.

Automatic Differentiation constructs program the program P^0 that apply f^0 to its inputs.

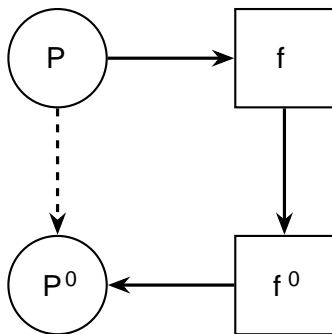


Figure: Automatic Differentiation.

Automatic Differentiation

Fortran, C: Tapenade

Python: Tangent, Myia

Julia: Zygote

F#: Di Sharp

...

Fortran, C: Tapenade
Python: Tangent, Myia

Julia: Zygote
F#: Di Sharp
...

not differentiating a specific programming language.
define a narrowed programming language **ADSL**. Similar to
(Abadi 2019) (Hu 2020) (Mak 2020).

ADSL is closed by differentiation

Automatic Differentiation compilation

We can use this pipeline to differentiate a function written in any programming language. You just need to pay the price of compilation.

Outline

- 1 Context
- 2 Formalization
- 3 Tables Relations
- 4 Automatic Differentiation
- 5 **Implementation**
- 6 Conclusion

Implementation

This work has been implemented at Lokad:
on the DSLEnvision
live in production

Optimization through gradient descent is used daily and triggers orders of millions of SKUs.

Outline

- 1 Context
- 2 Formalization
- 3 Tables Relations
- 4 Automatic Differentiation
- 5 Implementation
- 6 Conclusion

Conclusion

In this work we've presented a framework on automatic differentiation on relational queries.

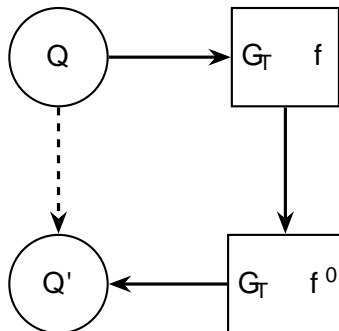


Figure: Path to Differentiating Relational Queries.

Conclusion

This will unlock ML model construction and optimisation in databases.

Figure: Proposed Pipeline.

Thanks for listening!

References

- [Abadi 2019] Martn Abadi et Gordon Plotkin. A simple differentiable programming language. Proceedings of the ACM on Programming Languages, vol. 4, pages 1{28, 12 2019.
- [Chicago] City Of Chicago. image. <https://data.cityofchicago.org/Transportation/Taxi-Trips/wrvz-psew> . Accessed: 2021-07-13.
- [Hu 2020] Y. Hu, L. Anderson, Tzu-Mao Li, Q. Sun, N. Carr, Jonathan Ragan-Kelley et F. Durand. Di Taichi: Differentiable Programming for Physical Simulation. ArXiv, vol. abs/1910.00935, 2020.
- [Hutson] Matthew Hutson. image. <https://www.sciencemag.org/news/2018/05/ai-researchers-allege-machine-learning-alchemy> . Accessed: 2021-07-15.
- [Mak 2020] Carol Mak et C. Ong. A Differentiable Pullback Programming Language for Higher-order Reverse-mode Automatic Differentiation . ArXiv, vol. abs/2002.08241, 2020.
- [Mcsherry 2021] Frank Mcsherry, Derek Murray, Rebecca Isaacs et Michael Isard. Differentiable data flow . 08 2021.
- [Schöle 2019] Maximilian E. Schöle, Frédéric Simonis, Thomas Heyenbrock, A. Kemper, Stephan Gannemann et T. Neumann. In-Database Machine Learning: Gradient Descent and Tensor Algebra for Main Memory Database Systems In BTW, 2019.

[Link to an example](#)

$\langle \textit{Pred} \rangle ::= .$

- | $\langle \textit{And } v \ w \rangle$
- | $\langle \textit{Or } v \ w \rangle$
- | $\langle \textit{Not } v \rangle$
- | $\langle v < w \rangle$
- | $\langle v \leq w \rangle$

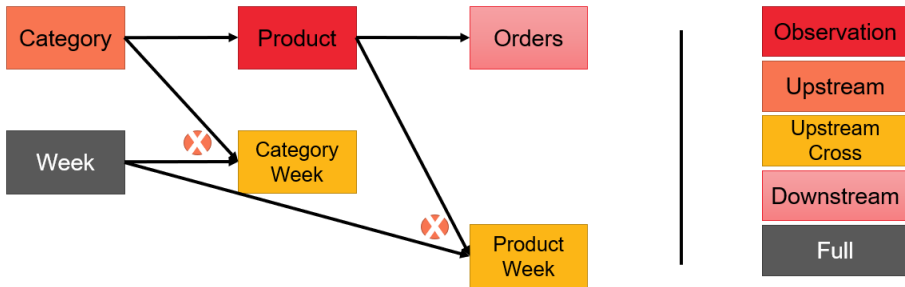


Figure: PolyStar

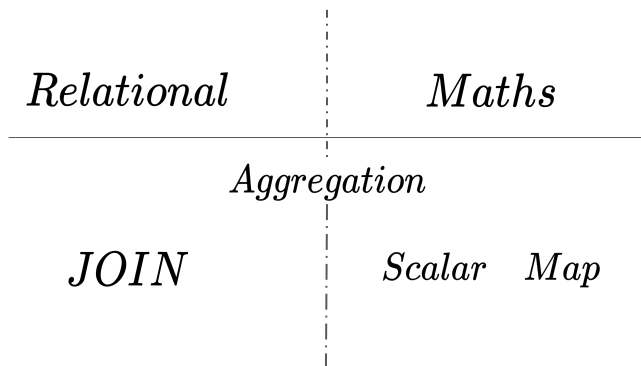


Figure: Relational - Math decomposition