Warper



Efficiently Adapting Learned Cardinality Estimators to Data and Workload Drifts



Beibin Li^{1,2}, Yao Lu², Srikanth Kandula² ¹University of Washington, Seattle, WA ²Microsoft Research, Redmond, WA 2022



Machine Learning (ML) for Database Systems



Challenge in AI Applications: Training and Testing Distribution Shifts (Drifts)

Cardinality Estimation and Its Machine Learning Process



Data Shifts (e.g., insert, delete, update) f(weight <=3 and price <= 3)



Workload Shifts

P(x): distribution of input data features

Querie s	Weight Bound	Price Bound	
Previo us	< 1.0	< 1.0	
	< 2.1	< 2.0	
	< 4.0	< 3.8	
New	> 5.0	< 1.0	
	< 1.5	> 5.2	

Performance After Data Shifts



Adaptaion Progress (min)

State-of-the-art for Adaption in CE

Scheme	How		
LM ^[1] Naru ^[3]	Re-train		
MSCN ^[2]	Completely re-train or fine-tune model		
DeepDB ^[4]	Partial re-train		

- Other DB tasks (e.g., ^[5]) use apriori training
 - Ad-hoc
 - Domain insights
 - Overly general.

^[1] Dutt, Anshuman, et al. "Selectivity estimation for range predicates using lightweight models." *Proceedings of the VLDB Endowment* 12.9 (2019): 1044-1057.

^[2] Kipf, Andreas, et al. "Learned cardinalities: Estimating correlated joins with deep learning." CIDR (2019).

^[3] Yang, Zongheng, et al. "Deep unsupervised cardinality estimation." *Proceedings of the VLDB Endowment, 13.3 (2019)*

^[4] Hilprecht, Benjamin, et al. "DeepDB: learn from data, not from queries!." *Proceedings of the VLDB Endowment* 13.7 (2020): 992-1005.

^[5] Ma, Lin, et al. "Active Learning for ML Enhanced Database Systems." *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020.

Goals

- Agnostic to ML Models
- Small Computation Overhead
- Quick Adaptation



Drift Case 1: Workload Distribution Drift



- Table stays the same
- Workload changed



Drift Case 1: Warper with GAN (Generative Adversarial Network)



Drift Case 1: Experiments

- Black Box Models: LM (VLDB 19'), MSCN (CIDR 19')
- Three datasets
- 12 new queries arrive per minute.
 - Fine-Tune
 - Hard Example Mining

- Mixture
 - Mixture Wa Augmentation (Add Noise)
- Warper (Ours)







Drift Case 1: Summary

- Workload Drift
- Scenario: number of new query predicates is small
- Goal: adapt the black-box ML model quickly
- Solution: synthesize additional queries for the model

• However, what if the number of new queries is large?



Drift Case 2: Too Many to Label

• Table stays the same

Workload changed

New Distribution



Previous Distribution

	Weight Bound	Price Bound	Card. GT
Many Queries	< 1.0	< 1.0	10
	< 2.1	< 2.0	20
	> 2.9	> 3.0	20

Drift Case 2: Warper with Picker



Drift Case 2: Summary

- Workload Drift
- Scenario: number of new queries is large
- Solution: select novel queries with higher priority





Drift Case 3: Data Shifted

Data Table



Queries



< 2.0

...

?

?

< 2.1

Re-calculate Ground Truth

Drift Case 3: Solution (for Data Shift)



3.0x 1.3x 1.5x

CE

Model

Summary of These Drift Cases



- Workload and Data shifts can happen at the same time.
- Other different drift examples (scenarios), and End-to-End experiment with continuous drifts are shown in the paper.

Related Work

- Active Learning
 - HAL (SIGMOD 19'), ADCP (SIGMOD 20'), Wilds (PMLR 21'), ...
- Generative Adversarial Network (GAN)
 - Deep Learning: GAN (NeurIPS 14'), InfoGAN (NeurIPS, 16'), CycleGAN (ICCV 17'),
 StyleGAN2 (ECCV 20'), TransGAN (NeurIPS 21')
 - DB: Relative Data Synthesis (VLDB 20')

Conclusion

- ML for System also Suffers from Data and Workload Shifts
- Create *Warper* to Adapt for Cardinality Estimation
 - Low Computation Overhead
 - 3x 6x Faster Adaptation in Slow.
 - 1-2x Faster Adaptation in Fast.
- Future: Examine More Real Workloads in End-to-End Setting







Thank You!

Warper: Efficiently Adapting Learned Cardinality Estimators to Data and Workload Drifts

Beibin Li, Yao Lu, Srikanth Kandula *beibin.li@microsoft.com*



