

1

# Selectivity Estimation

Chuan Xiao

Osaka University and Nagoya University

# Outline

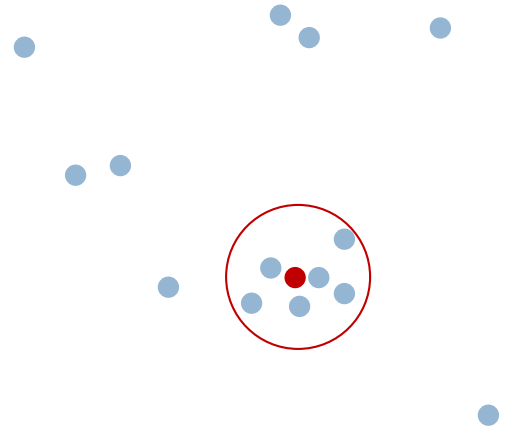
2

- Problem Definition
- Applications
- Methods
- Performance Evaluation

# Problem Definition

3

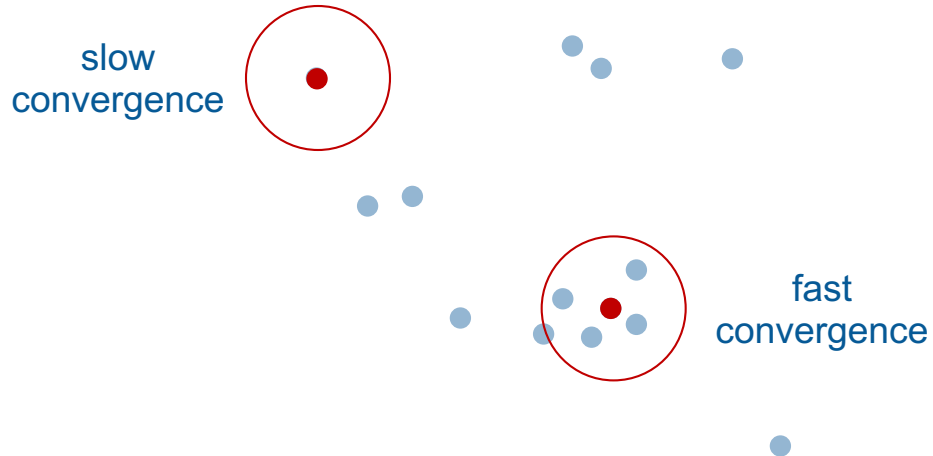
- Selectivity estimation of similarity search for high-dimensional data
  - Given:
    - a database  $X$  of high-dimensional vectors,
    - a query vector  $\mathbf{q}$ ,
    - a distance function  $dist(., .)$ ,
    - a threshold  $t$ .
  - Estimate the number of objects  $\mathbf{x}$  in  $X$  such that  $dist(\mathbf{q}, \mathbf{x}) \leq t$ .
  - a.k.a. cardinality estimation, spherical range counting
- Related problem
  - Selectivity (cardinality) estimation for relational data [KKRL+19, OBGK19, SL19, WSR19, YLKW+19, HTAK+20, PZM20]
    - Each predicate deals with a dimension.
      - `SELECT COUNT(*) FROM employee WHERE age < 30 AND salary > 50000`
    - Dimensionality is usually low.



# Application – Local Density Estimation

4

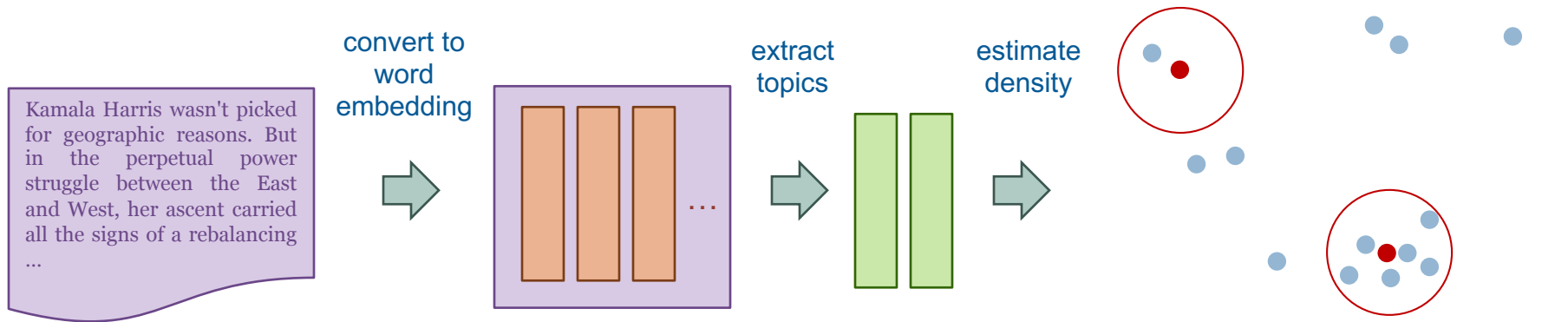
- Clustering
  - Find starting points.



# Application – Local Density Estimation

5

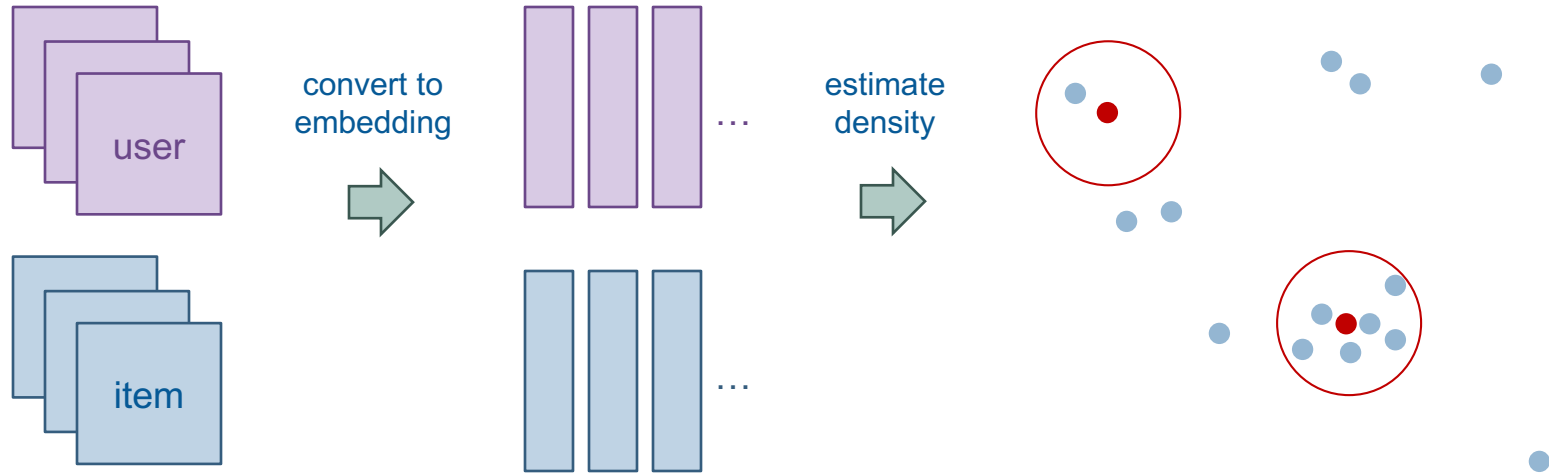
- Determine the popularities of topics.



# Application – Local Density Estimation

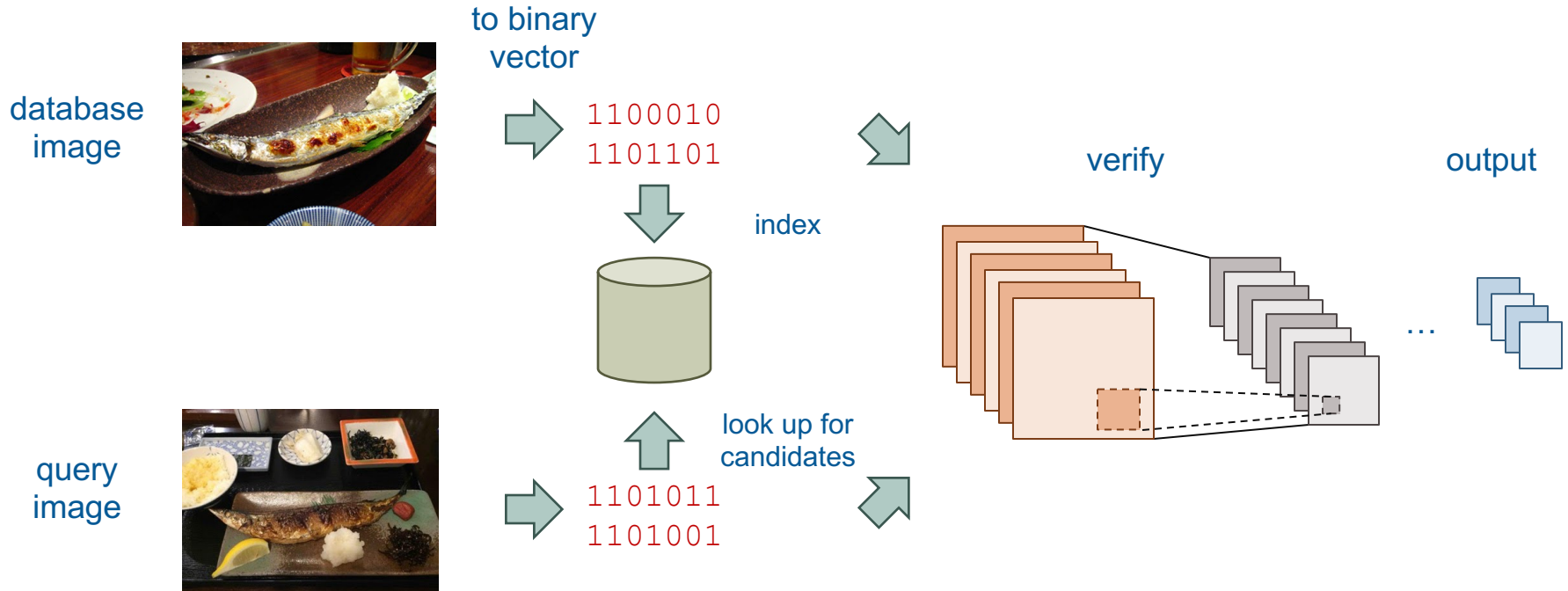
6

- Find out if a user/item is an outlier in an e-commerce application.



# Application – Image Retrieval

7



source: Wikipedia

# Application – Image Retrieval

8

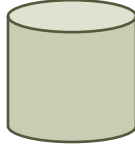
database image



to binary vector

1100010  
1101101

index



look up for candidates

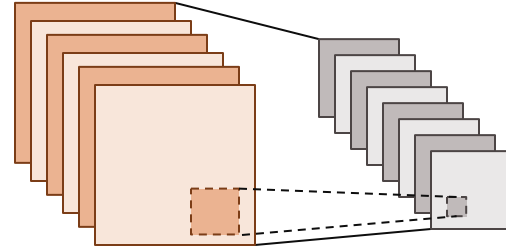
1101011  
1101001

query image



source: Wikipedia

verify



output

...



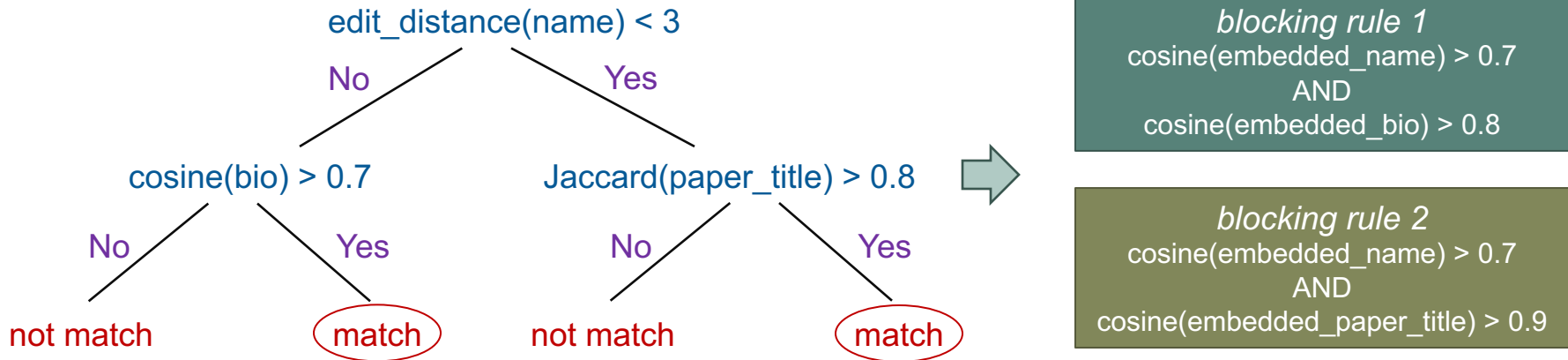
Estimate candidate size  
→ running time, SLA ...



# Application – Query Optimization

9

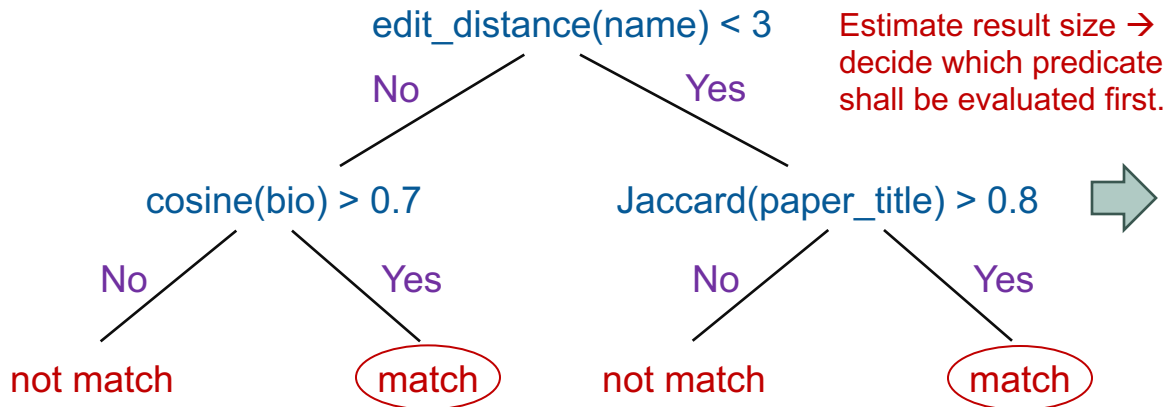
- Hands-off entity matching systems (e.g., Falcon [DCDN+17]) extract paths from random forests and take each path (a conjunction of similarity predicates) as a blocking rule.



# Application – Query Optimization

10

- Hands-off entity matching systems (e.g., Falcon [DCDN+17]) extract paths from random forests and take each path (a conjunction of similarity predicates) as a blocking rule.
- Embed textual attributes (e.g., by edit distance embedding [DYZW+20]) and process the conjunctive query.



# Evaluation Criteria of Selectivity Estimation

11

- Accuracy
  - Measured by MSE, MAPE, q-error, etc.
- Estimation speed
- Offline processing speed
  - Build an index?
  - Train a model?
- Performance guarantee
  - $\epsilon$ - $\delta$
- Consistency (monotonicity)
  - For a fixed query object, selectivity is non-decreasing in the threshold.
  - This yields more interpretability and less vulnerability.
- Updatability
  - The database may have updates.

# (Representative) Selectivity Estimation Methods

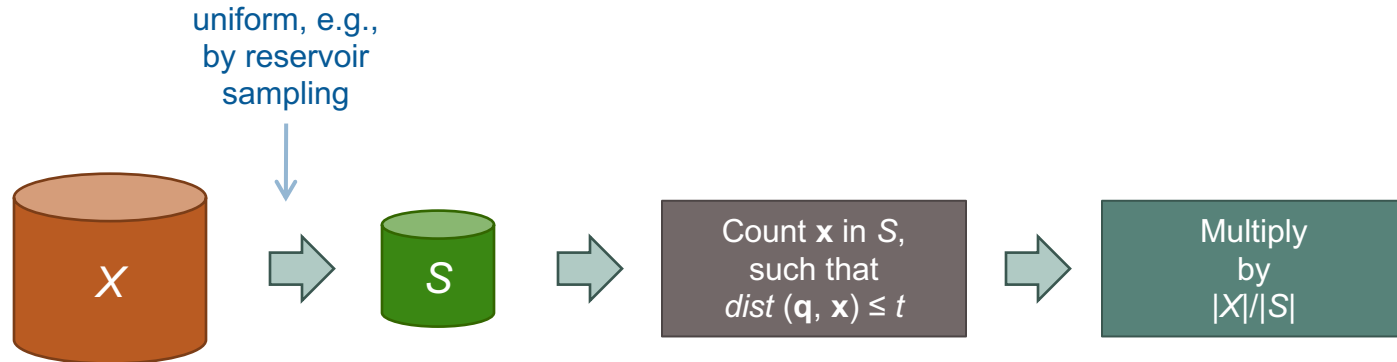
12

- Sampling
- Kernel density estimation
- Regression
  - ▣ Non-deep learning
    - XGBoost
  - ▣ Deep learning
    - Vanilla deep neural network
    - Recursive model index
    - Deep lattice network
    - CardNet (incremental prediction + deep learning)
    - SelNet (piecewise linear function + deep learning)

# Sampling – Uniform Sampling

13

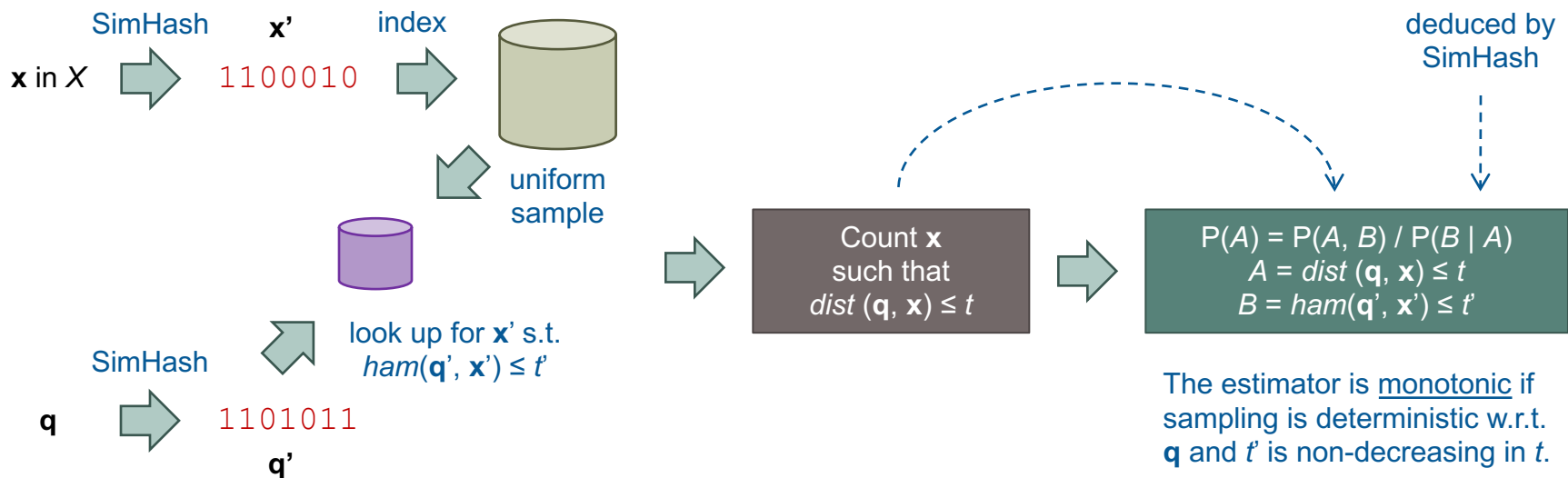
- A natural baseline
  - ▣ It is lightweight with well-understood performance, and easy to support monotonicity and handle updates.
  - ▣ Weakness: the probability that  $\text{dist}(\mathbf{q}, \mathbf{x}) \leq t$  is small, especially when  $\mathbf{q}$  is an outlier. So we need a very large sample size for accurate estimation.



# Sampling – Importance Sampling

14

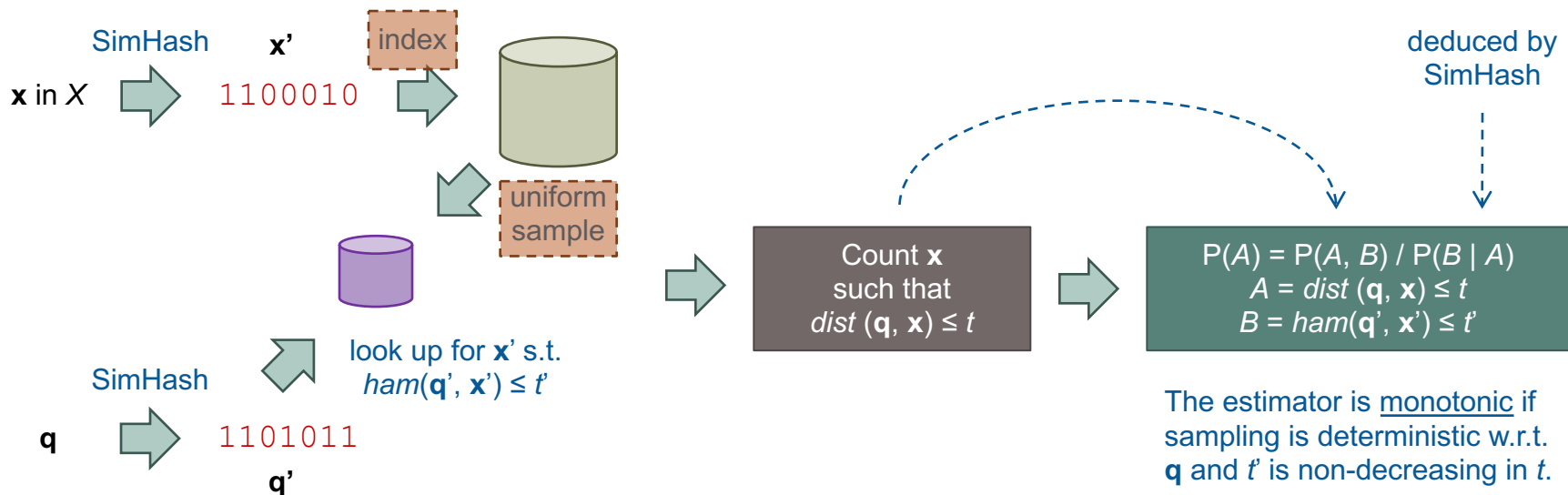
- Estimate selectivity by generating samples from another distribution.
  - ▣ SimHash for angular distance (cosine similarity) [WCN18].
    - $dist(., .)$  is captured by Hamming distance between hash values.
    - Use  $L$  independent hash tables for better accuracy.



# Sampling – Importance Sampling

15

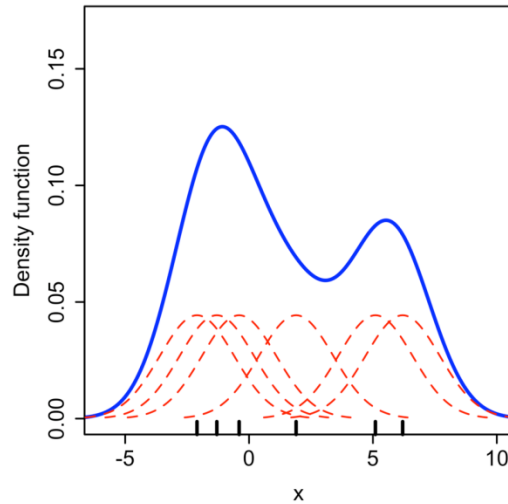
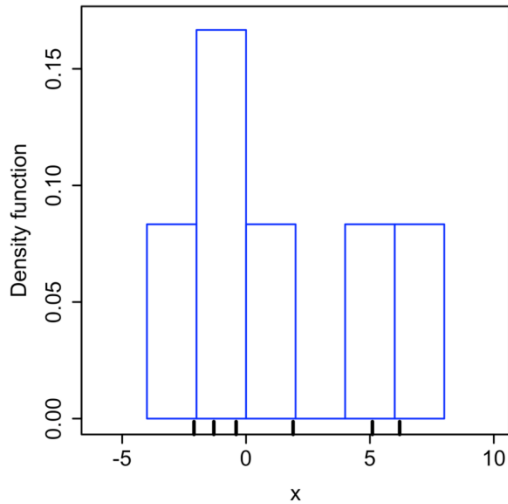
- Deal with updates.
  - Update the index and sample.



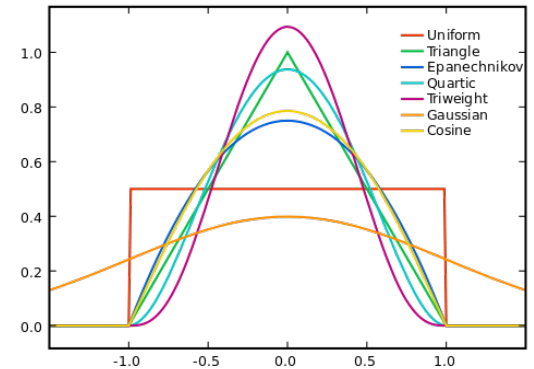
# Kernel Density Estimation (KDE)

16

- Sample and smooth by kernel



$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$



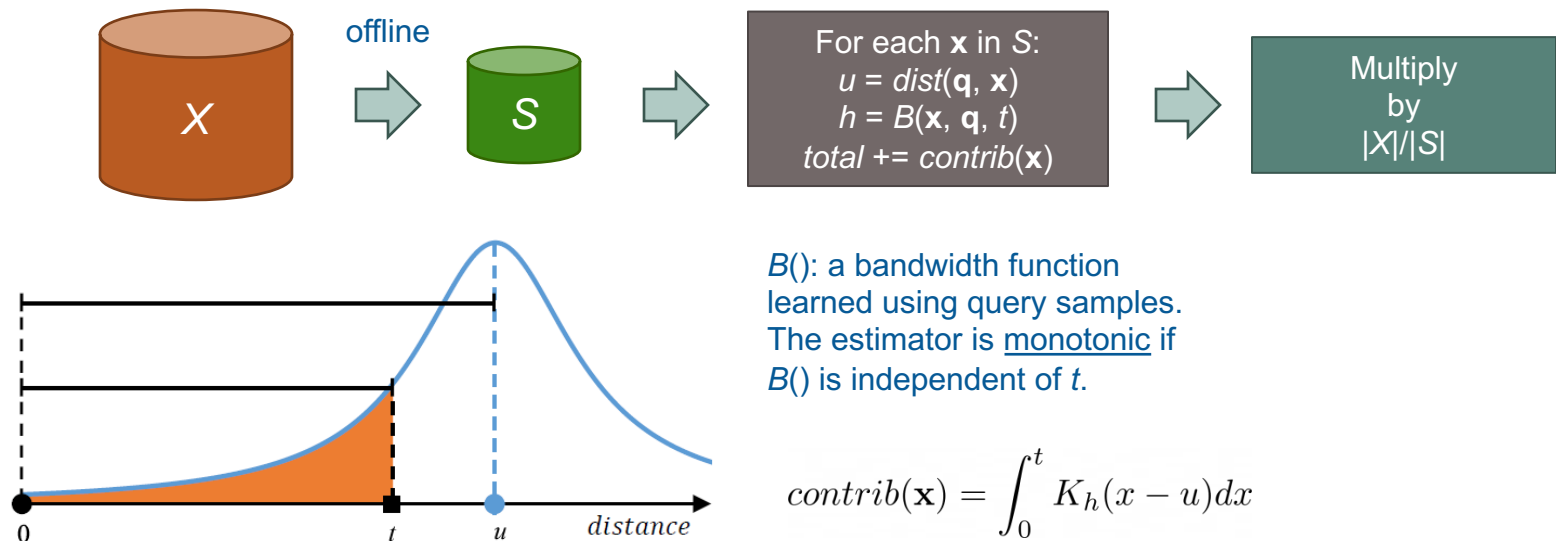
source: Wikipedia



# Kernel Density Estimation (KDE)

17

- Model the probability density function of  $dist(\mathbf{q}, \mathbf{x})$  by KDE [MFBS18].
  - Sample objects and compute their contributions to the selectivity.



$B()$ : a bandwidth function learned using query samples. The estimator is monotonic if  $B()$  is independent of  $t$ .

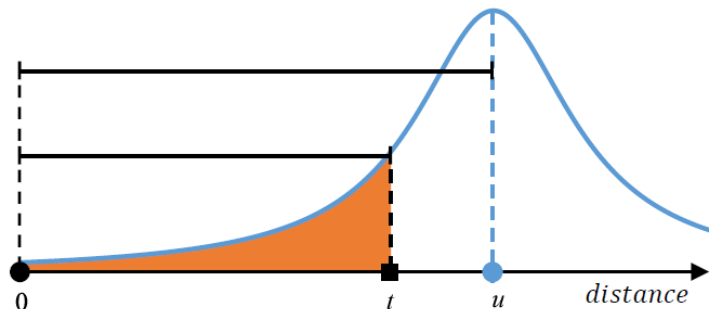
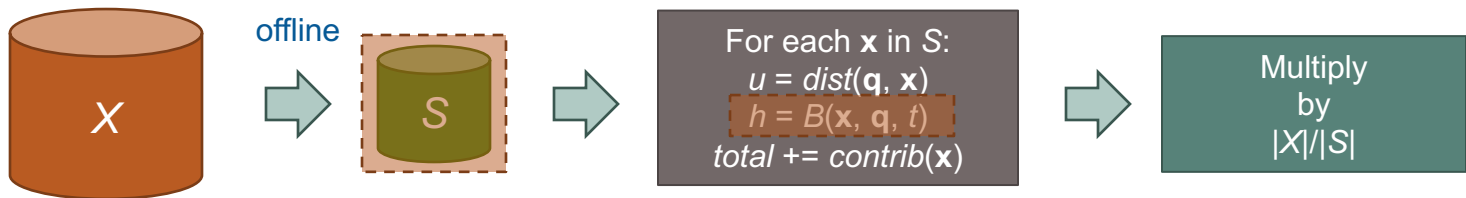
$$contrib(\mathbf{x}) = \int_0^t K_h(x - u) dx$$

source: [MFBS18]

# Kernel Density Estimation (KDE)

18

- Deal with updates.
  - Incrementally sample more objects.
  - Retrain bandwidth functions.



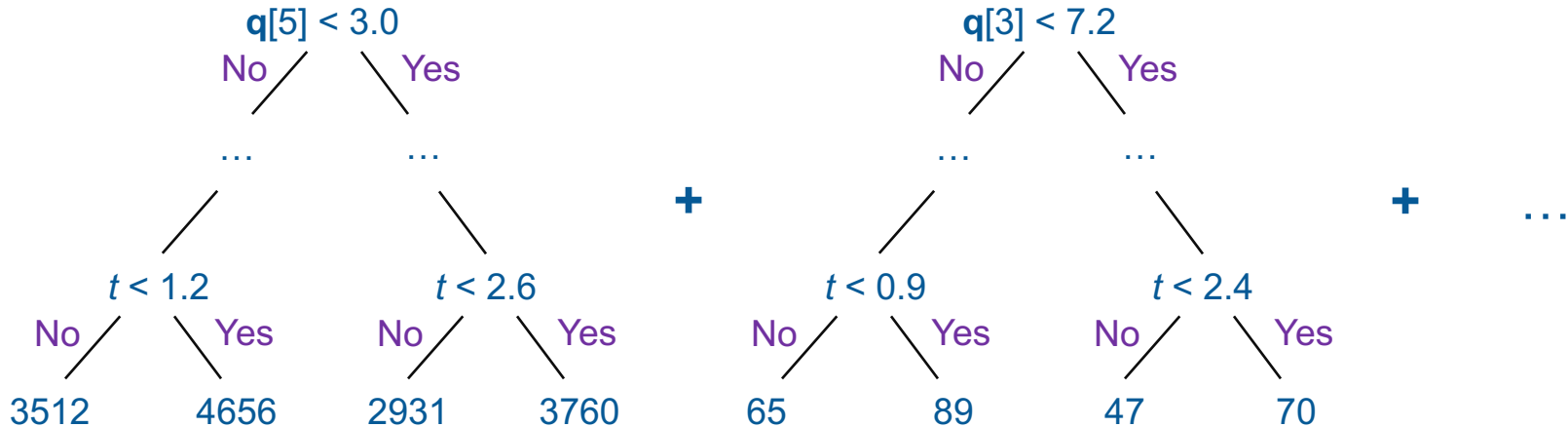
$B()$ : a bandwidth function learned using query samples. The estimator is monotonic if  $B()$  is independent of  $t$ .

$$\text{contrib}(\mathbf{x}) = \int_0^t K_h(x - u) dx$$

# Non-deep Regression – XGBoost

19

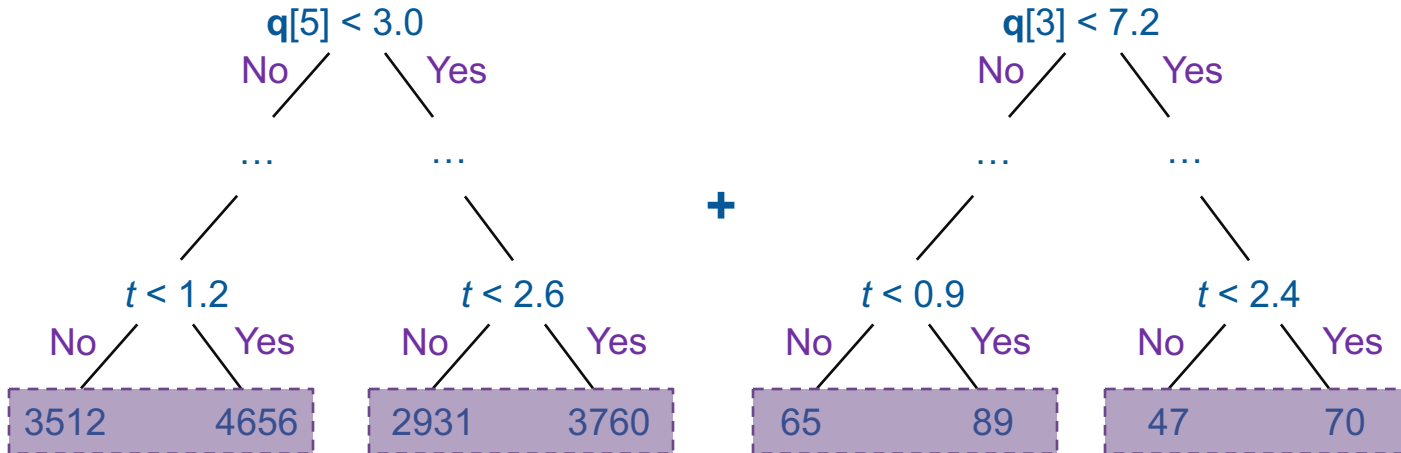
- Gradient boosting [CG16]
  - Ensemble of weak prediction models.
    - For example, decision trees, with each rule in the form of  $q[i] < \alpha$  or  $t < \beta$ .
  - Each model is learned to fit the residual of previous ones.



# Non-deep Regression – XGBoost

20

- Gradient boosting [CG16]
  - Ensemble of weak prediction models.
    - For example, decision trees, with each rule in the form of  $q[i] < \alpha$  or  $t < \beta$ .
  - Each model is learned to fit the residual of previous ones.

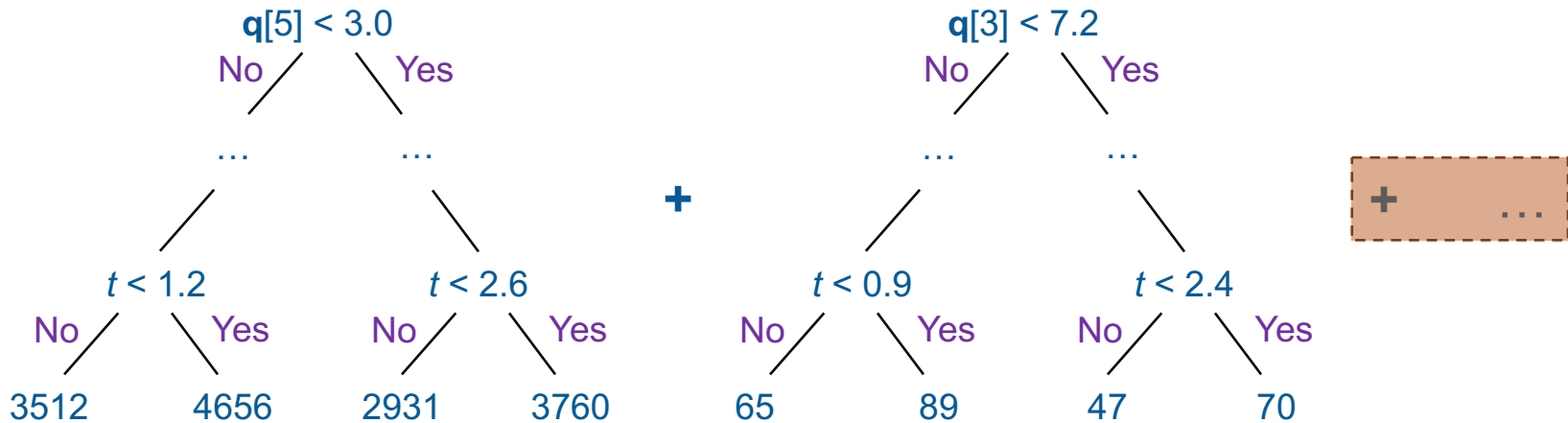


The estimator is monotonic if  $t < \beta$  is only at the bottom level and leaf node values are non-decreasing w.r.t. Yes/No.

# Non-deep Regression – XGBoost

21

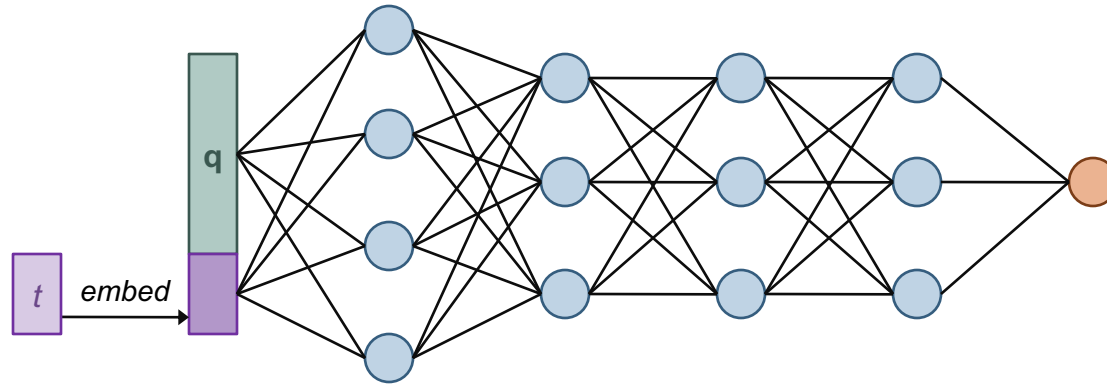
- Deal with updates.
  - It is time-consuming to retrain existing decision trees.
  - Train more decision trees to fit the residual.



# Deep Regression – Vanilla Deep Neural Network

22

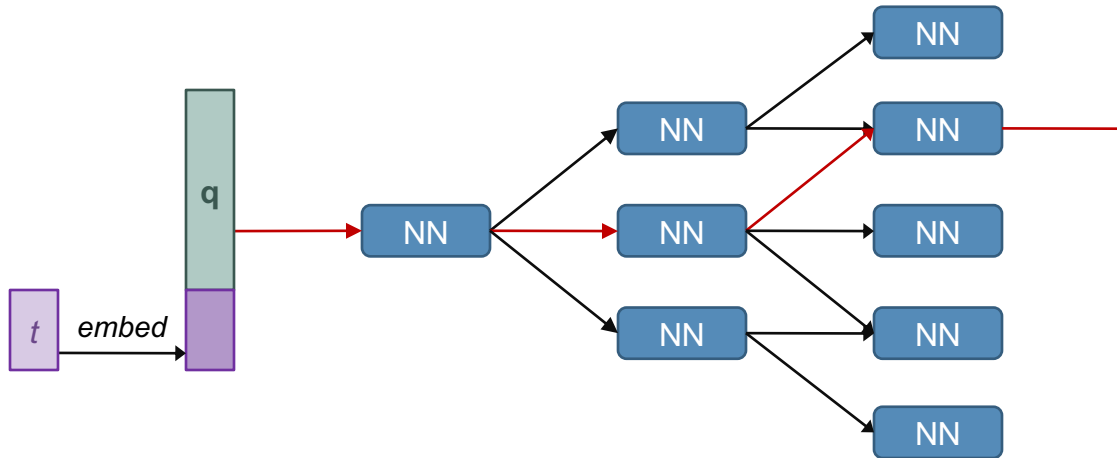
- Fully connected neural network.
  - ▣ Number of hidden layers  $\approx 4$ .
  - ▣ For higher accuracy, embed  $t$  (dim  $\approx 5$ ) and concatenate to  $\mathbf{q}$  as input.
  - ▣ Non-monotonic.



# Deep Regression – Recursive Model Index

23

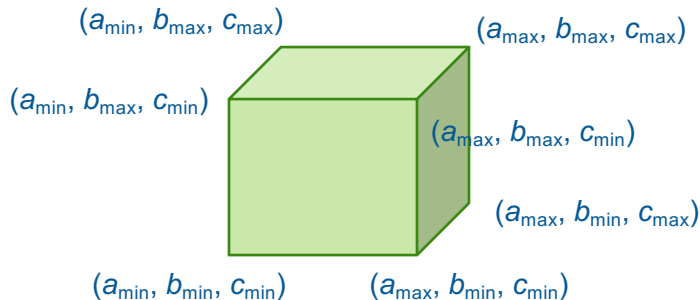
- Originally developed for range indexing in RDBMS [KBCD+18].
  - ▣ Inspired by the mixture-of-experts model.
  - ▣ Each model (e.g., a neural network) picks another one in the next stage.
  - ▣ Non-monotonic.



# Deep Regression – Deep Lattice Network

24

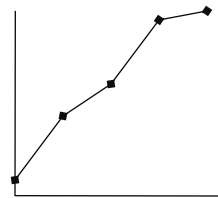
- Developed for monotonic regression tasks [YDCP+17].
  - ▣ Input: monotonic features + non-monotonic features.
- Components
  - ▣ Lattice: regression for a  $d$ -dimensional input.



Only learn for vertex values.  
Others are processed by  
multilinear interpolation.

$lat(a, b, c) \leq lat(a', b, c)$  for  
monotonic feature  $a \leq a'$ .

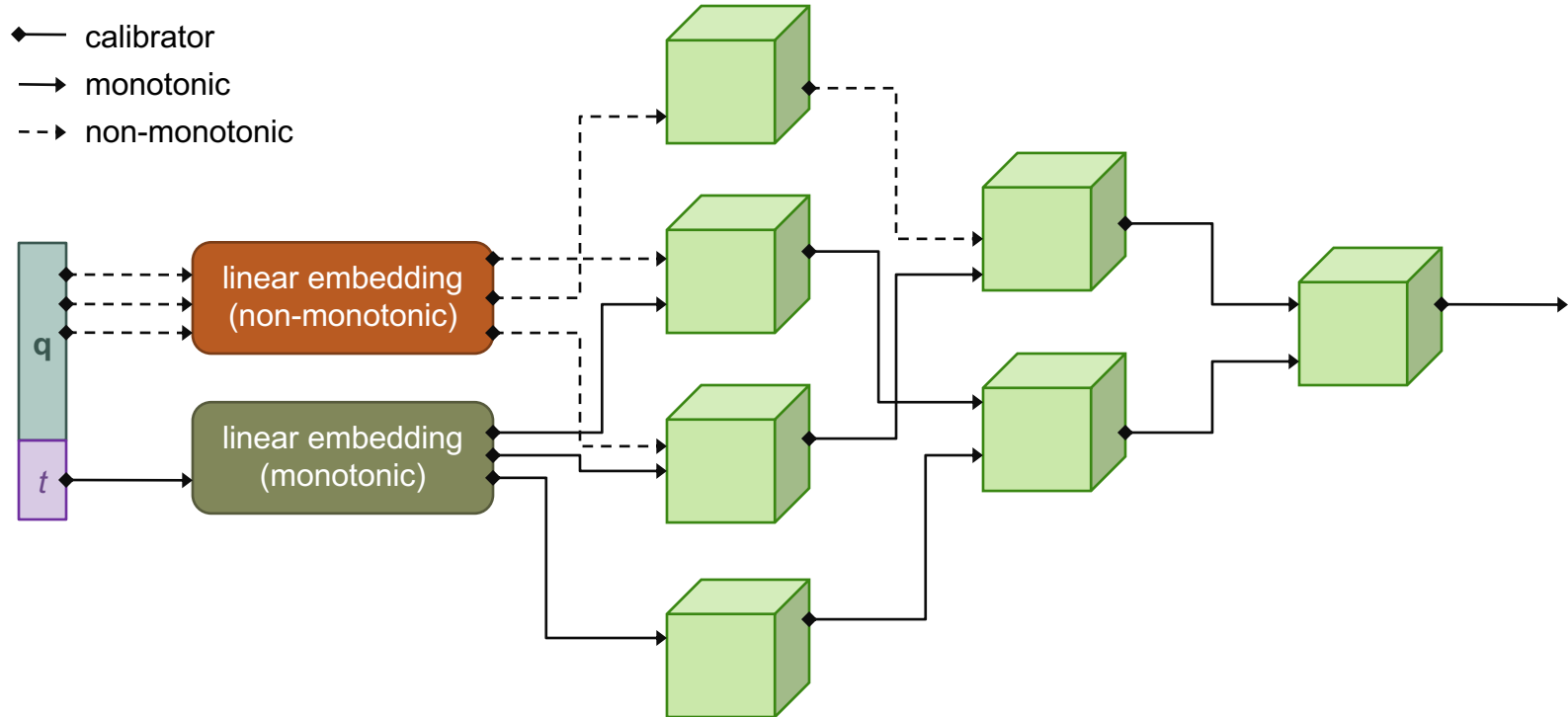
- ▣ Calibrator: 1-dimensional non-decreasing piecewise linear function. →
- ▣ Linear embedding: a matrix (all elements  $\geq 0$  for monotonicity).





# Deep Regression – Deep Lattice Network

25



# Deep Regression – Incremental Prediction (CardNet)

26

- Idea: use multiple regressors, each dealing with an interval  $[WXQC+20]$ .
- Procedure

- Feature extraction

- query vector  $\mathbf{q}$   $\rightarrow$  binary vector  $\mathbf{r}$

- Map  $\mathbf{q}$  to an integer  $B$  by LSH (e.g., random projection) and then set the  $B$ -th bit to 1
- Repeat  $L$  rounds using  $L$  hash functions and concatenate the resulting bit vectors.



- threshold  $t \rightarrow$  integer  $\tau$

- $0 \rightarrow 0$ .  $t_{\max} \rightarrow \tau_{\max}$ . Other values are mapped in a non-decreasing manner  $\rightarrow$  monotonicity.

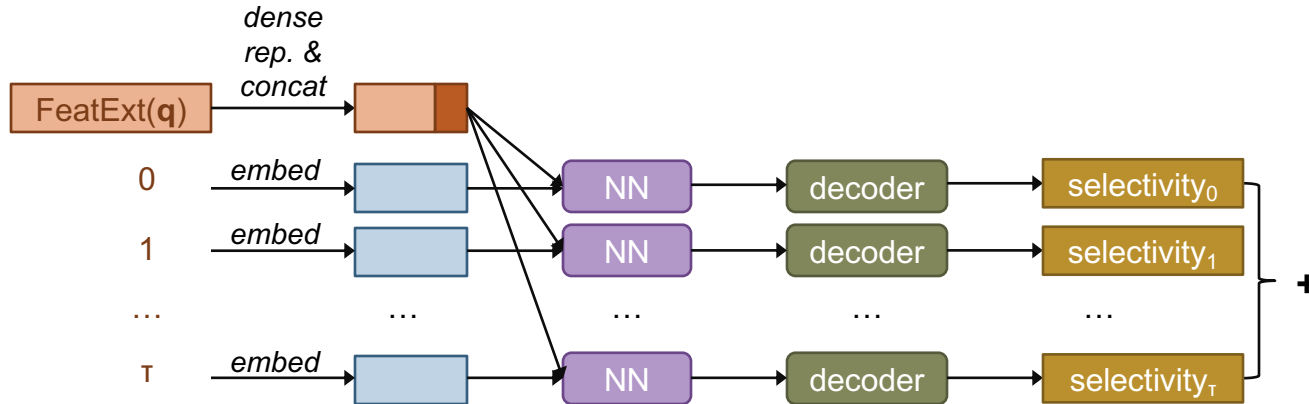


- Regression

# Deep Regression – Incremental Prediction (CardNet)

27

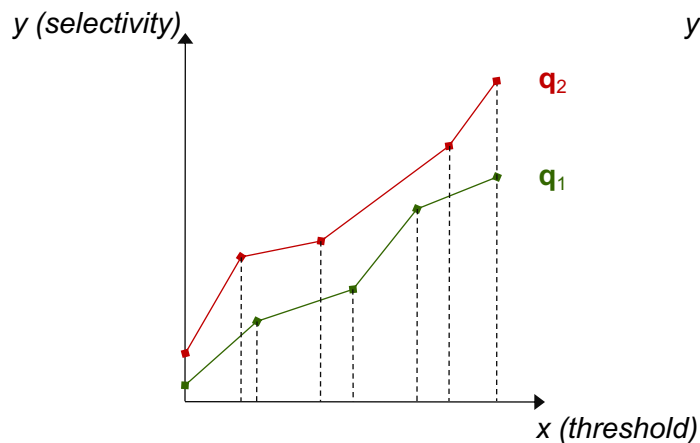
- Idea: use multiple regressors, each dealing with an interval.
- Procedure
  - ▣ Feature extraction
  - ▣ Regression
    - Use  $(\tau + 1)$  regressors, each for a distance in  $0, 1, \dots, \tau$ .



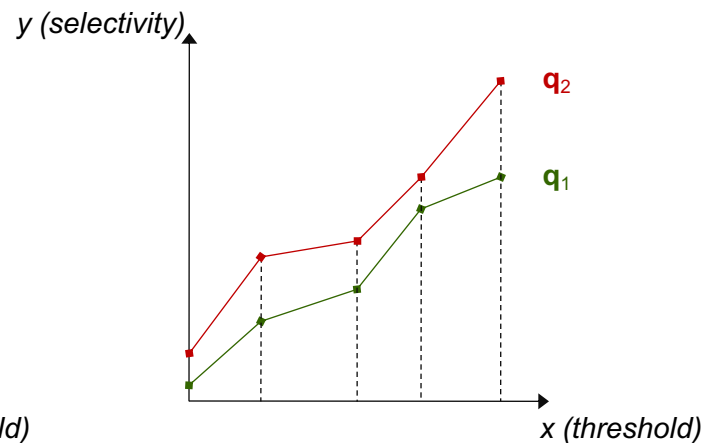
# Deep Regression – Piecewise Linear Function (SelNet)

28

- Query-dependent PLF v.s. query-independent PLF



query-dependent:  
variable (adaptive)  
control points

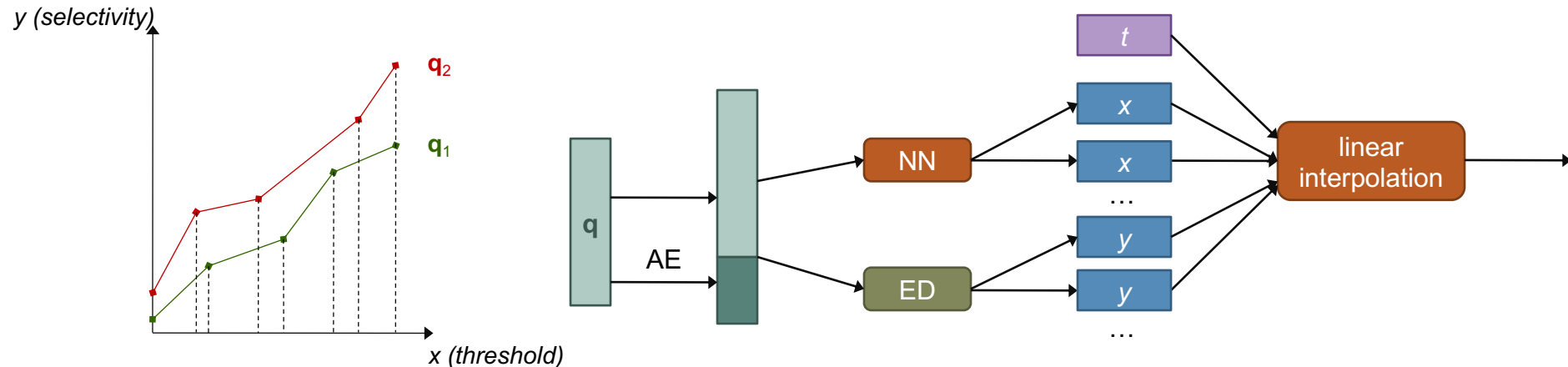


query-independent  
(deep lattice network):  
fixed control points

# Deep Regression – Piecewise Linear Function (SelNet)

29

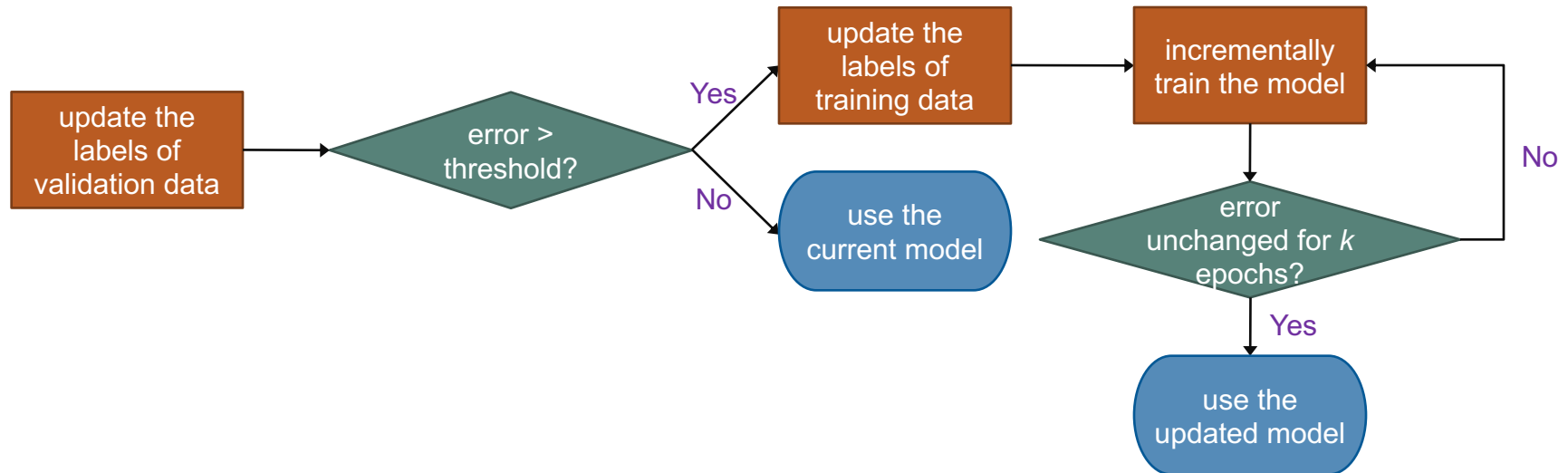
- Learn query-dependent PLFs [WXQM+20].
  - ▣ NN: a neural network that outputs the  $x$  values of control points.
  - ▣ ED: an encoder-decoder model that outputs the  $y$  values of control points.
  - ▣ Monotonic if  $y$  is non-decreasing in  $x$ .



# Deep Regression – Dealing with Updates

30

- Many deep regression models are trained through gradient descent.
- We adopt incremental learning for these models.



# Benchmarks

31

- So far there are specific benchmarks for this problem.
- Datasets used in existing work (benchmarks for other uses)
  - ▣ Text
    - GloVe: 1.9M 300-dimensional word embedding
      - <https://nlp.stanford.edu/projects/glove/>
    - fastText: 1M 300-dimensional word embedding
      - <https://fasttext.cc/docs/en/english-vectors.html>
  - ▣ Image
    - MS-Celeb-1M: 10M celebrity images for face recognition
      - <https://msceleb.org/> (terminated in 2019)
      - Pre-processed by faceNet [SKP15] to 128-dimensional vectors.
  - ▣ Video
    - YouTube: 3.4K videos of 1.6K people
      - <http://www.cs.tau.ac.il/~wolf/ytfaces/index.html>
      - 0.35M 1770-dimensional (-feature) vectors extracted from the frames.

# Comparison of Selectivity Estimation Methods

32

Method	Accuracy	Estimation Speed	Offline Proc. Speed	Performance Guarantee	Consistency (Monotonicity)
Uniform Sampling	Adjustable / Very Low	Adjustable	<b>None</b>	<b>Yes</b>	Possible
Importance Sampling	Adjustable / Low	Adjustable	<b>Fast</b>	<b>Yes</b>	Possible
KDE	Medium	Slow	Medium	No	Possible
XGBoost	Low	Medium	Medium	No	Possible
Vanilla DNN	Low	<b>Fast</b>	Medium	No	No
RMI	Medium	Medium	Slow	No	No
DLN	Low	Slow	Slow	No	<b>Yes</b>
CardNet	<b>High</b>	<b>Fast</b>	Slow	No	<b>Yes</b>
SelNet	<b>High</b>	Medium	Slow	No	<b>Yes</b>



# Performance in a Query Optimizer

33

- Datasets
  - AMiner (author names & publications)
  - IMDB (cast & movie titles)
  - Attributes are pre-processed by Sentence-BERT [RG19].
- Queries
  - Conjunctive queries of 2 – 5 Euclidean distance predicates.
    - For example,  $dist(name) \leq 0.25$  AND  $dist(affiliations) \leq 0.4$  AND  $dist(research\ interests) \leq 0.45$ .
  - For each query, we estimate the selectivity of each predicate.
  - The predicate with the smallest selectivity is evaluated first by index lookup. Others are checked on the fly.
- Methods
  - Uniform sampling
  - XGBoost
  - RMI
  - CardNet-A: CardNet with acceleration for estimation
  - Exact: an oracle that instantly returns the true selectivity.
  - Mean: an estimator that returns the same selectivity (mean of 10,000 random queries) for a given threshold.

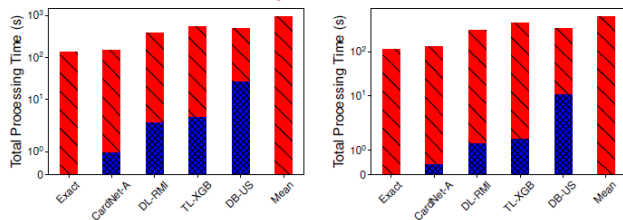
# Performance in a Query Optimizer

34

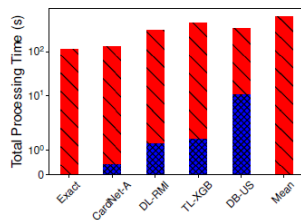
- Exact: oracle (true selectivity)
- Cardnet-A
- DL-RMI: RMI

- TL-XGB: XGBoost
- DB-US: uniform sampling
- Mean: same (mean selectivity)

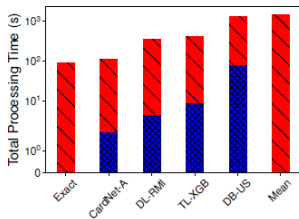
Query processing time:  
estimation / lookup+check



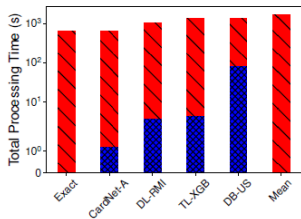
(a) Time, AMiner-Publication



(b) Time, AMiner-Author

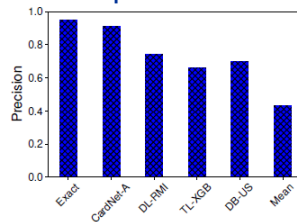


(c) Time, IMDB-Movie

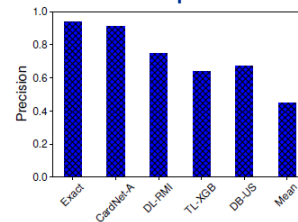


(d) Time, IMDB-Actor

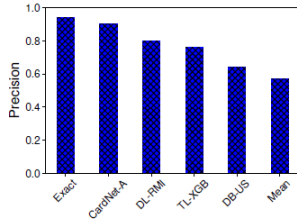
Precision of query planning:  
% of queries on which a method picks the fastest plan



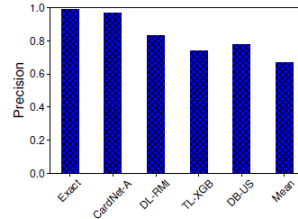
(a) Time, AMiner-Publication



(b) Time, AMiner-Author



(c) Time, IMDB-Movie



(d) Time, IMDB-Actor

# References

- T. Chen and C. Guestrin. XGBoost: A scalable tree boosting system. In KDD, pages 785–794, 2016.
- X. Dai, X. Yan, K. Zhou, Y. Wang, H. Yang, and J. Cheng. Convolutional Embedding for Edit Distance. In SIGIR, pages 599–608, 2020.
- S. Das, P. S. G. C., A. Doan, J. F. Naughton, G. Krishnan, R. Deep, E. Arcaute, V. Raghavendra, and Y. Park. Falcon: Scaling up hands-off crowdsourced entity matching to build cloud services. In SIGMOD, pages 1431–1446, 2017.
- S. Hasan, S. Thirumuruganathan, J. Augustine, N. Koudas, and G. Das. Deep learning models for selectivity estimation of multi-attribute queries. In SIGMOD, pages 1035–1050, 2020.
- A. Kipf, T. Kipf, B. Radke, V. Leis, P. A. Boncz, and A. Kemper. Learned cardinalities: Estimating correlated joins with deep learning. In CIDR, 2019.
- T. Kraska, A. Beutel, E. H. Chi, J. Dean, and N. Polyzotis. The case for learned index structures. In SIGMOD, pages 489–504, 2018.
- M. Mattig, T. Fober, C. Beilschmidt, and B. Seeger. Kernel-based cardinality estimation on metric data. In EDBT, pages 349–360, 2018.
- J. Ortiz, M. Balazinska, J. Gehrke, and S. S. Keerthi. An empirical analysis of deep learning for cardinality estimation. arXiv preprint arXiv:1905.06425, 2019.
- Y. Park, S. Zhong, and B. Mozafari. Quicksel: Quick selectivity learning with mixture models. In SIGMOD, pages 1017–1033, 2020.

# References

- N. Reimers and I. Gurevych. Sentence-BERT: Sentence embeddings using siamese bert-networks. In EMNLP-IJCNLP, pages 3980–3990, 2019.
- F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In CVPR, pages 815–823, 2015.
- N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. arXiv preprint arXiv:1701.06538, 2017.
- J. Sun and G. Li. An end-to-end learning-based cost estimator. PVLDB, 13(3):307–319, 2019.
- B. Walenz, S. Sintos, S. Roy, and J. Yang. Learning to sample: Counting with complex queries. PVLDB, 13(3):390–402, 2019.
- Y. Wang, C. Xiao, J. Qin, X. Cao, Y. Sun, W. Wang, and M. Onizuka. Monotonic cardinality estimation of similarity selection: A deep learning approach. In SIGMOD, pages 1197–1212, 2020.
- Y. Wang, C. Xiao, J. Qin, R. Mao, M. Onizuka, W. Wang, R. Zhang. Consistent and flexible selectivity estimation for high-dimensional data. arXiv preprint arXiv:2005.09908 (2020).
- A. Wehenkel and G. Louppe. Unconstrained monotonic neural networks. In NeurIPS, pages 1543–1553, 2019.
- X. Wu, M. Charikar, and V. Natchu. Local density estimation in high dimensions. In ICML, pages 5293–5301, 2018.
- Z. Yang, E. Liang, A. Kamsetty, C. Wu, Y. Duan, P. Chen, P. Abbeel, J. M. Hellerstein, S. Krishnan, and I. Stoica. Deep unsupervised cardinality estimation. PVLDB, 13(3):279–292, 2019.
- S. You, D. Ding, K. Canini, J. Pfeifer, and M. Gupta. Deep lattice networks and partial monotonic functions. In NIPS, pages 2981–2989, 2017.