

# KARL: Fast Kernel Aggregation Queries

**(Edison) Tsz Nam Chan**<sup>1,2</sup>,    Man Lung Yiu<sup>2</sup>,    Leong Hou U<sup>3</sup>

<sup>1</sup>Univ. of Hong Kong

<sup>2</sup>Hong Kong Polytechnic Univ.

<sup>3</sup>Univ. of Macau



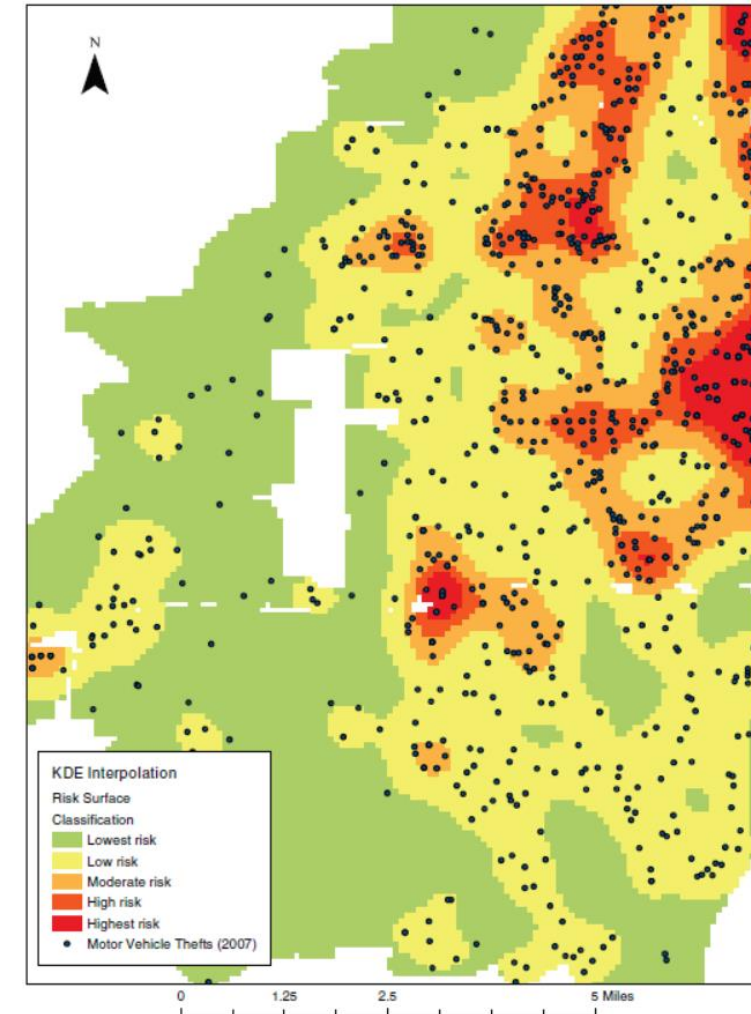
# What is Kernel Aggregation?

Compute the function:

$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} \underbrace{w_i}_{\text{weighting}} \exp\left(-\underbrace{\gamma}_{\text{constant}} \cdot \underbrace{\text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{\text{Euclidean distance}}\right)$$

query  
dataset

- *Application*: crime rate prediction
  - Each  $\mathbf{p}$  (black dot) represents the location of a crime
    - e.g., robbery,  
commercial burglary,  
motor vehicle theft
  - Predict the crime rate of a given location ( $\mathbf{q}$ ) by computing  $\mathcal{F}_P(\mathbf{q})$



# Kernel Aggregation: More Applications

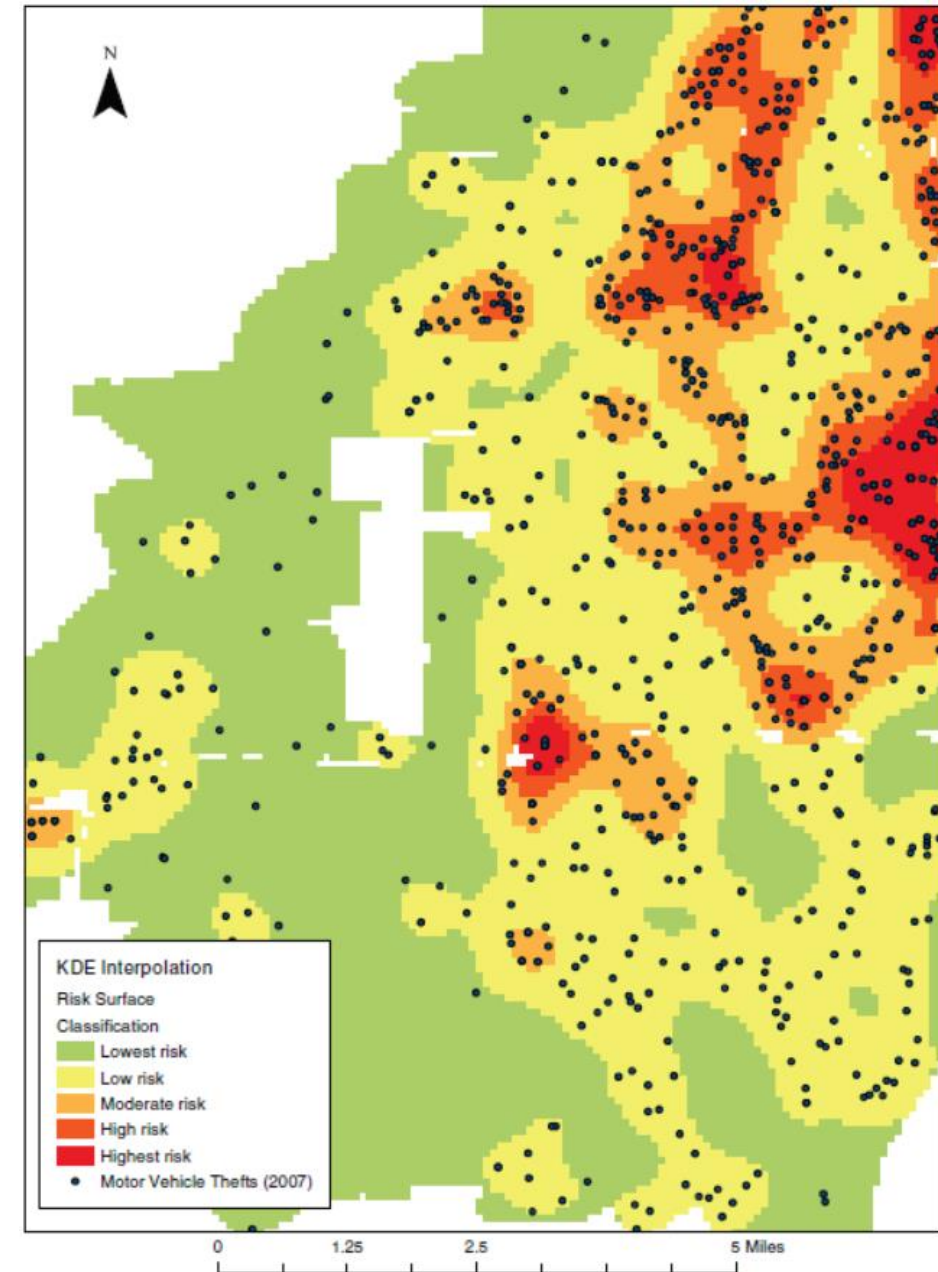
Kernel aggregation function:

$$\underbrace{\mathcal{F}_P}_{\text{dataset}}(\underbrace{\mathbf{q}}_{\text{query}}) = \sum_{\mathbf{p}_i \in P} \underbrace{w_i}_{\text{weighting}} \exp(-\underbrace{\gamma}_{\text{constant}} \cdot \underbrace{\text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{\text{Euclidean distance}})$$

Type of weighting	Used in model	Application(s)
Type I: identical, positive $w_i$ (most specific)	Kernel Density Classification	Crime rate prediction Ecological modeling Particle Searching
Type II: positive $w_i$ (subsuming Type I)	1-class SVM	Network fault detection Outlier detection
Type III: no restriction on $w_i$ (subsuming Type II)	2-class SVM	Network fault detection Tumor samples classification Image classification

# Kernel Aggregation Queries

- *Application*: crime rate prediction
- Just classify  $\mathcal{F}_P(\mathbf{q})$  into
  - **Lowest risk**:  $[0, 10)$
  - **Low risk**:  $[10, 20)$
  - **Moderate risk**:  $[20, 50)$
  - **High risk**:  $[50, 100)$
  - **Highest risk**:  $[100, \infty)$
- Threshold kernel aggregation query ( $\tau$ -KAQ)
  - Input: query vector  $\mathbf{q}$ , dataset  $P$ , threshold  $\tau$
  - Output: Boolean value
    - true (if  $\mathcal{F}_P(\mathbf{q}) \geq \tau$ ) or
    - false (if  $\mathcal{F}_P(\mathbf{q}) < \tau$ )



# Expensive to compute the exact $\mathcal{F}_P(\mathbf{q})$ !

Le et al. “*Despite their successes, what makes kernel methods difficult to use in many large scale problems is the fact that computing the decision function is typically expensive, especially at prediction time.*” [1]

$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} w_i \exp(-\gamma \cdot \text{dist}(\mathbf{q}, \mathbf{p}_i)^2)$$

O(|P| × d) time

What are typical values of |P| and d?

# Our contribution: KARL

- Kernel Aggregation Rapid Library (KARL)
  - **2.5-738x** speed up over state-of-the-art in different datasets
  - The **Fastest** library in **online prediction phase**

Library	Support indexing	Response time
LibSVM	no	high
Scikit-learn	yes	high
KARL (this paper)	yes	low

**KARL:** <https://github.com/edisonchan2013928/KARL-Fast-Kernel-Aggregation-Queries>

**LibSVM:** <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

**Scikit-learn:** <https://scikit-learn.org/>

# How to speed up?

$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} w \exp(-\gamma \cdot \text{dist}(\mathbf{q}, \mathbf{p}_i)^2)$$

$$LB(\mathbf{q}) \leq \mathcal{F}_P(\mathbf{q}) \leq UB(\mathbf{q})$$

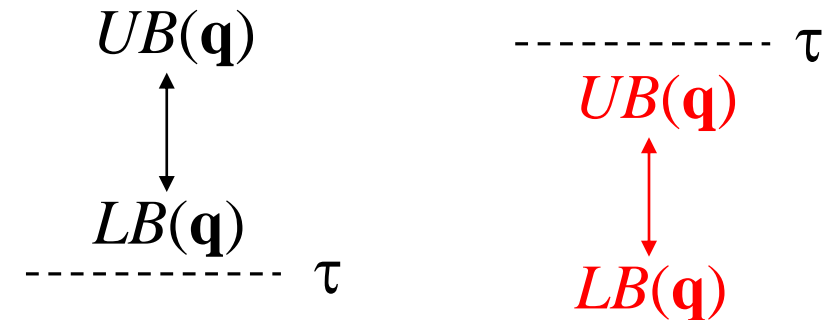


Requirements:

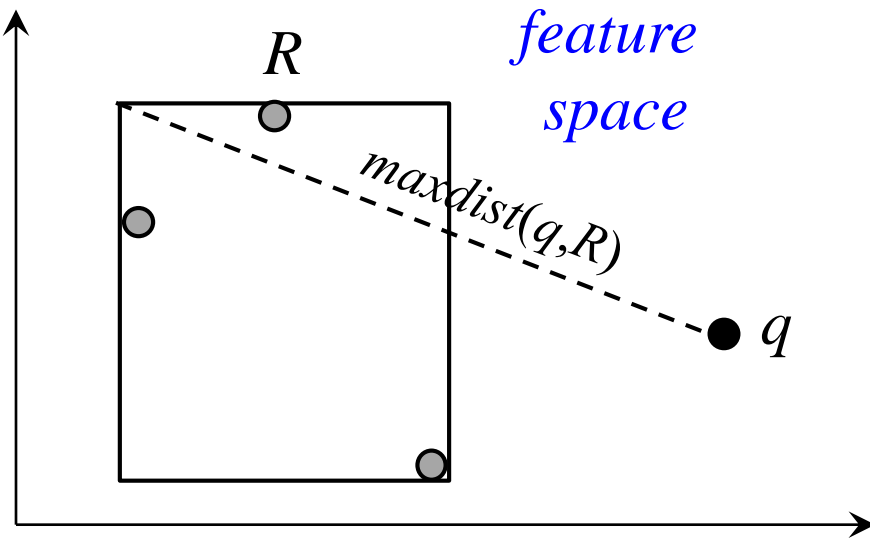
1. Fast to compute, e.g.,  $O(d)$  time
2. As tight as possible

Cases to avoid computing exact  $\mathcal{F}_P(\mathbf{q})$ :

Case 1 ( $\mathcal{F}_P(\mathbf{q}) \geq \tau$ ):      Case 2 ( $\mathcal{F}_P(\mathbf{q}) < \tau$ ):



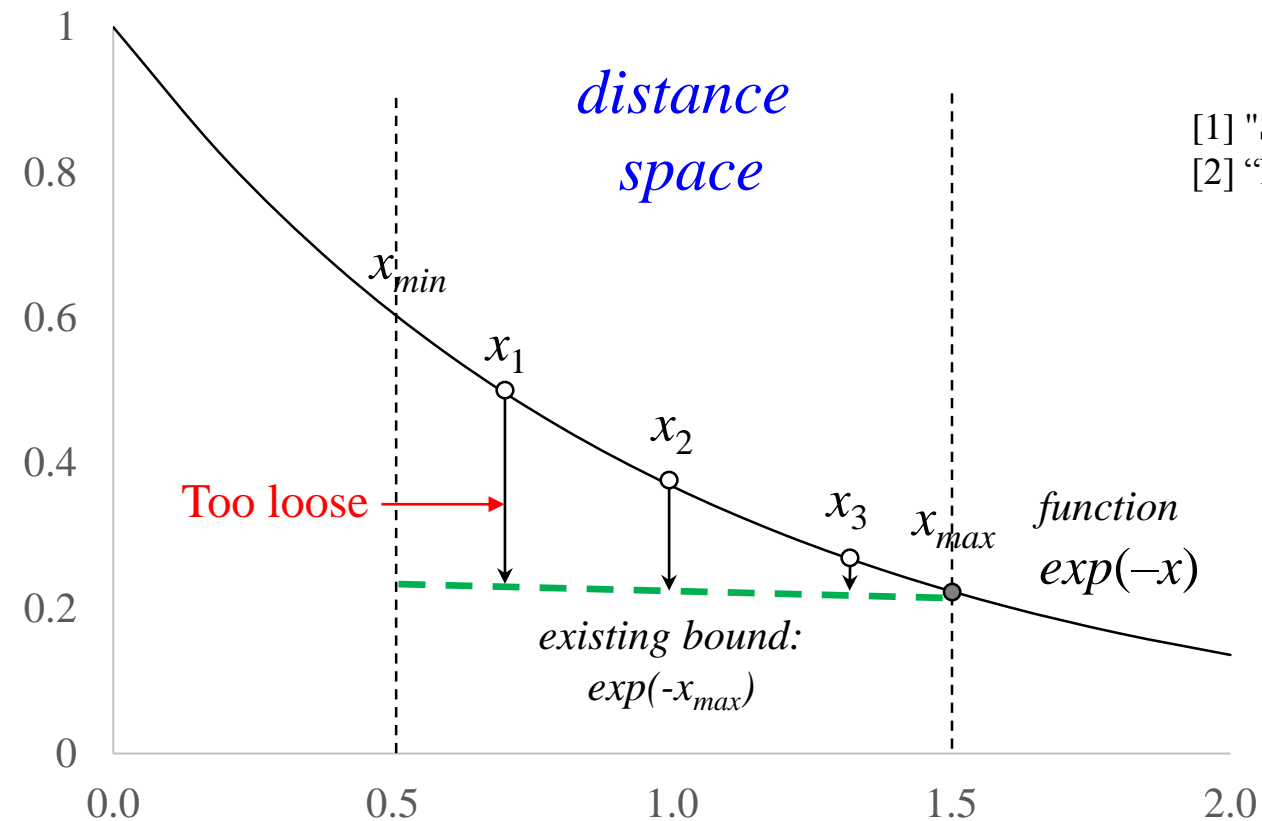
# Existing Bounding Function



$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} w \exp(-\gamma \cdot \underbrace{dist(\mathbf{q}, \mathbf{p}_i)^2}_{X_i})$$

$$LB_R(\mathbf{q}) = w \cdot R.\text{count} \cdot \exp(-\gamma \cdot \underbrace{maxdist(\mathbf{q}, R)^2}_{X_{\max}})$$

- [1] "Scalable kernel density classification via threshold based pruning" SIGMOD2017  
 [2] "Nonparametric Density Estimation: Toward Computational Tractability" SDM03



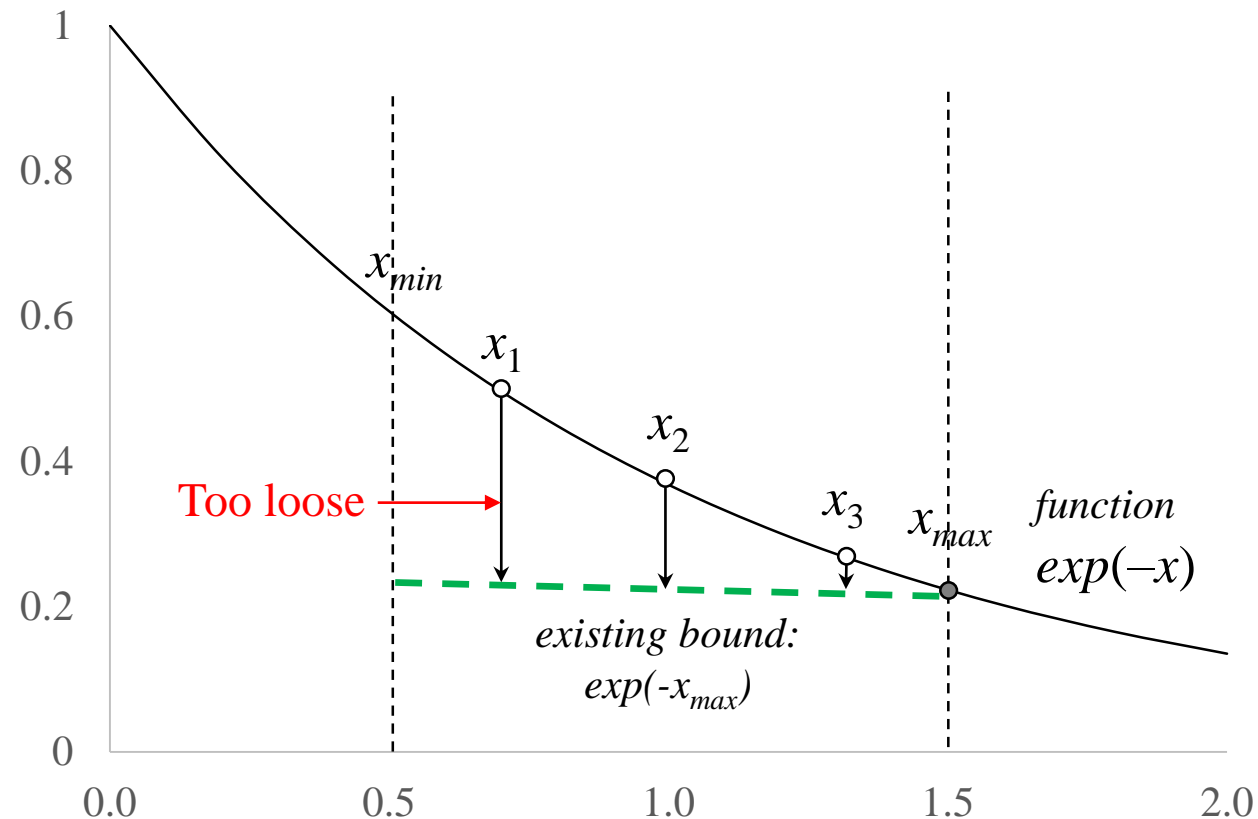
*Summary:*

- ☑ Fast to compute, e.g.,  $O(d)$  time
- ☒ Too loose

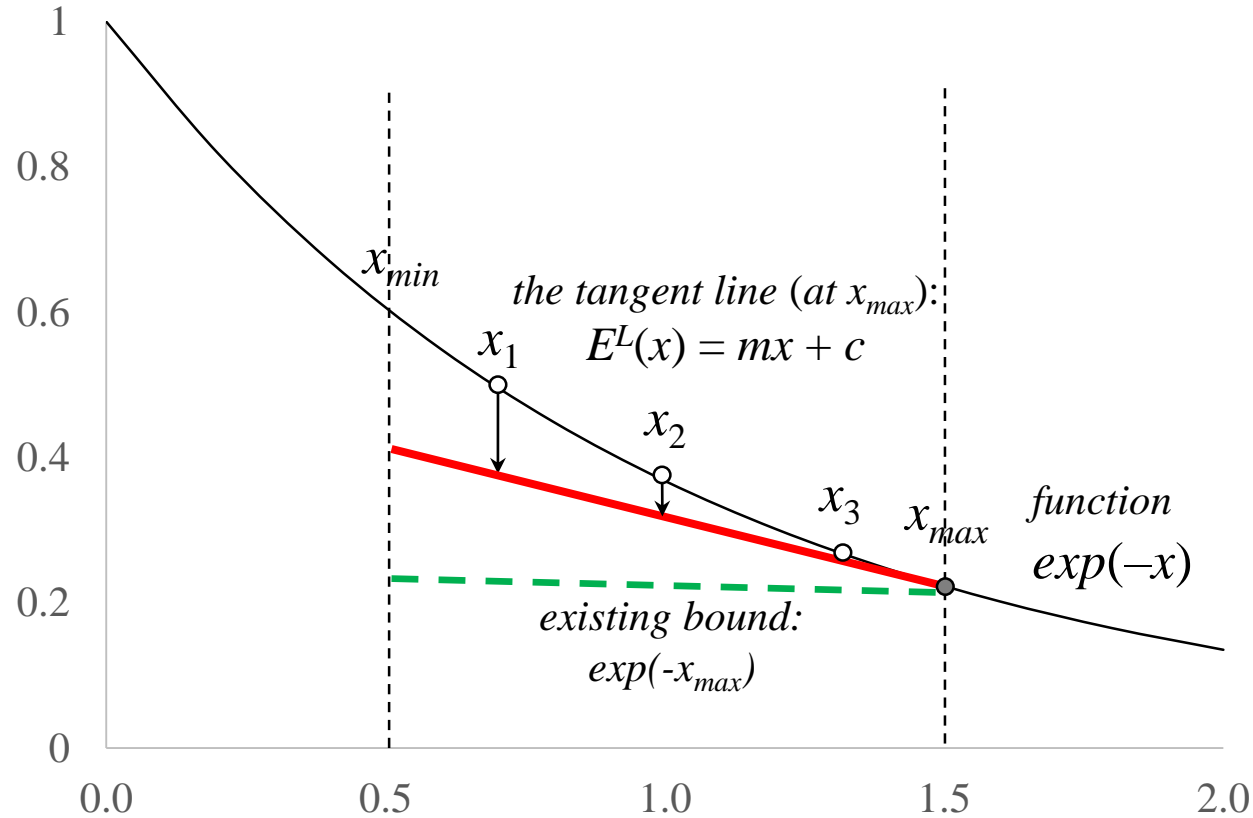


## Research question 1:

How to develop a **tighter** lower bound function (than the existing one)?



# Our Idea: Tangent Bound



$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} w \exp(-\underbrace{\gamma \cdot \text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{X_i})$$

Tangent Bound:

$$\mathcal{FL}_P(\mathbf{q}, \text{Lin}_{m,c}) = \sum_{\mathbf{p}_i \in P} w (m(\underbrace{\gamma \text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{X_i}) + c)$$

☺ always tighter than the existing bound

Research question 2:

How to compute this tighter bound quickly?

# O(d)-time Tighter Linear Bound

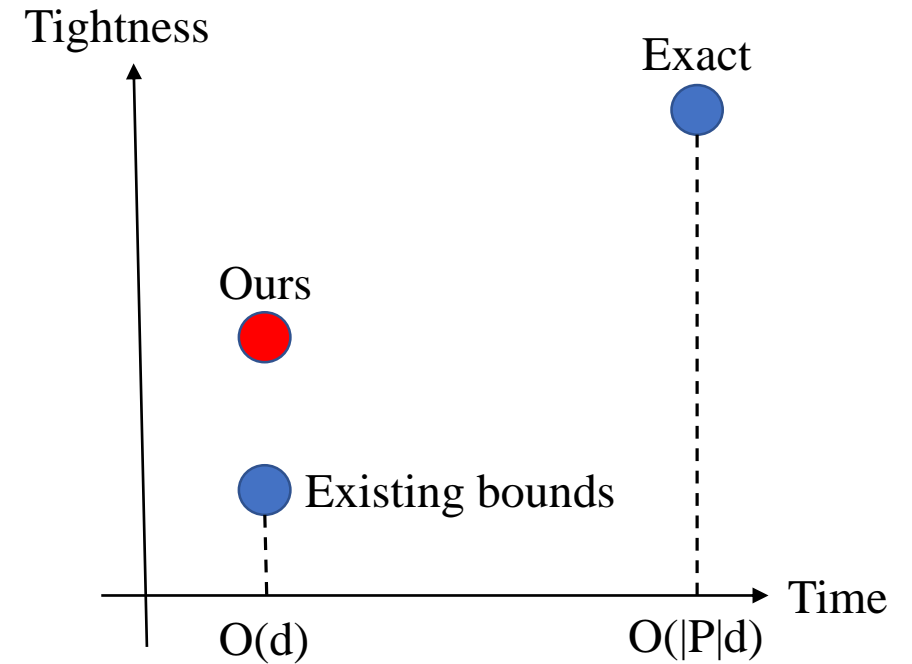
$$\mathcal{F}_P(\mathbf{q}) = \sum_{\mathbf{p}_i \in P} w \exp(-\underbrace{\gamma \cdot \text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{X_i})$$

Linear bound of  $\mathcal{F}_P(\mathbf{q})$ :

$$\mathcal{FL}_P(\mathbf{q}, \text{Lin}_{m,c}) = \sum_{\mathbf{p}_i \in P} w (m(\underbrace{\gamma \text{dist}(\mathbf{q}, \mathbf{p}_i)^2}_{X_i}) + c)$$

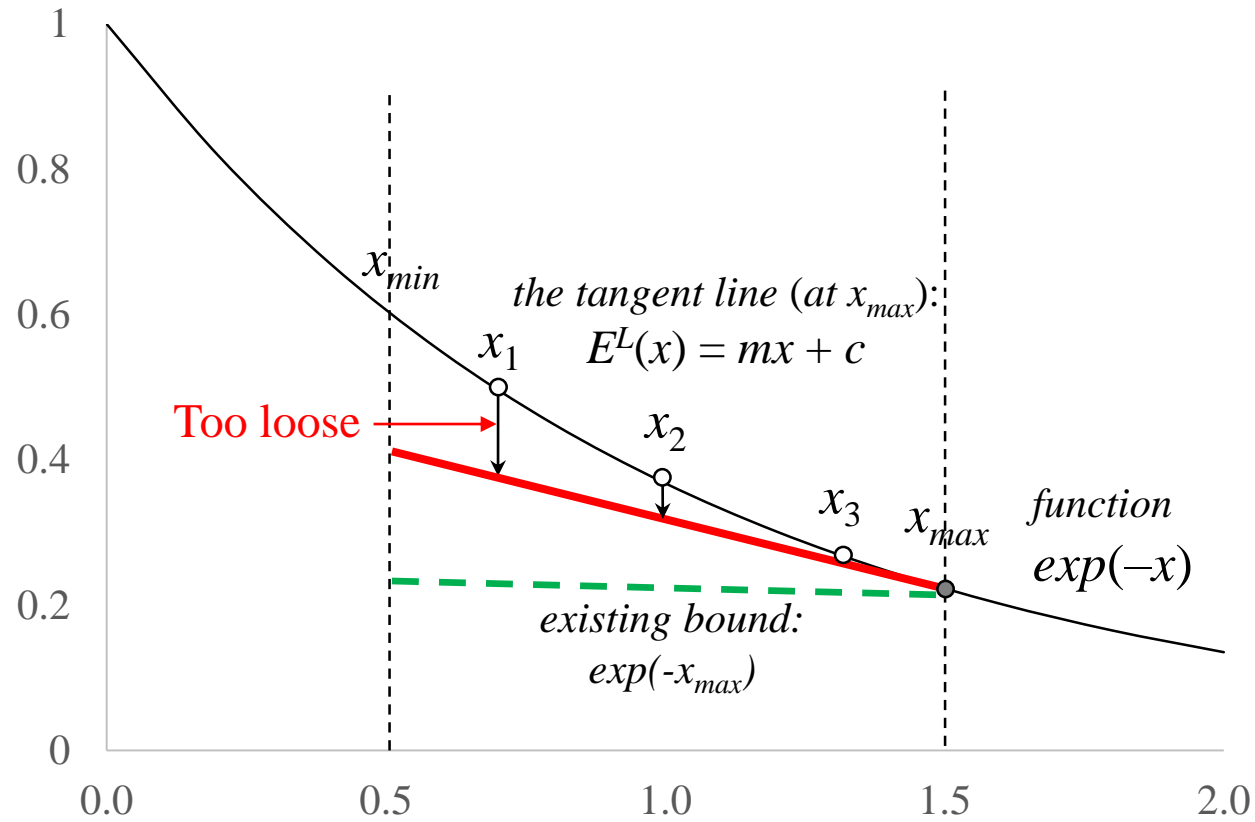
$$\mathcal{FL}_P(\mathbf{q}, \text{Lin}_{m,c}) = wm\gamma(|P| \underbrace{\|\mathbf{q}\|^2}_{O(d)} - 2\mathbf{q} \cdot \underbrace{\mathbf{a}_P}_{O(d)} + b_P) + wc|P|$$

where  $\mathbf{a}_P = \sum_{\mathbf{p}_i \in P} \mathbf{p}_i$  and  $b_P = \sum_{\mathbf{p}_i \in P} \|\mathbf{p}_i\|^2$



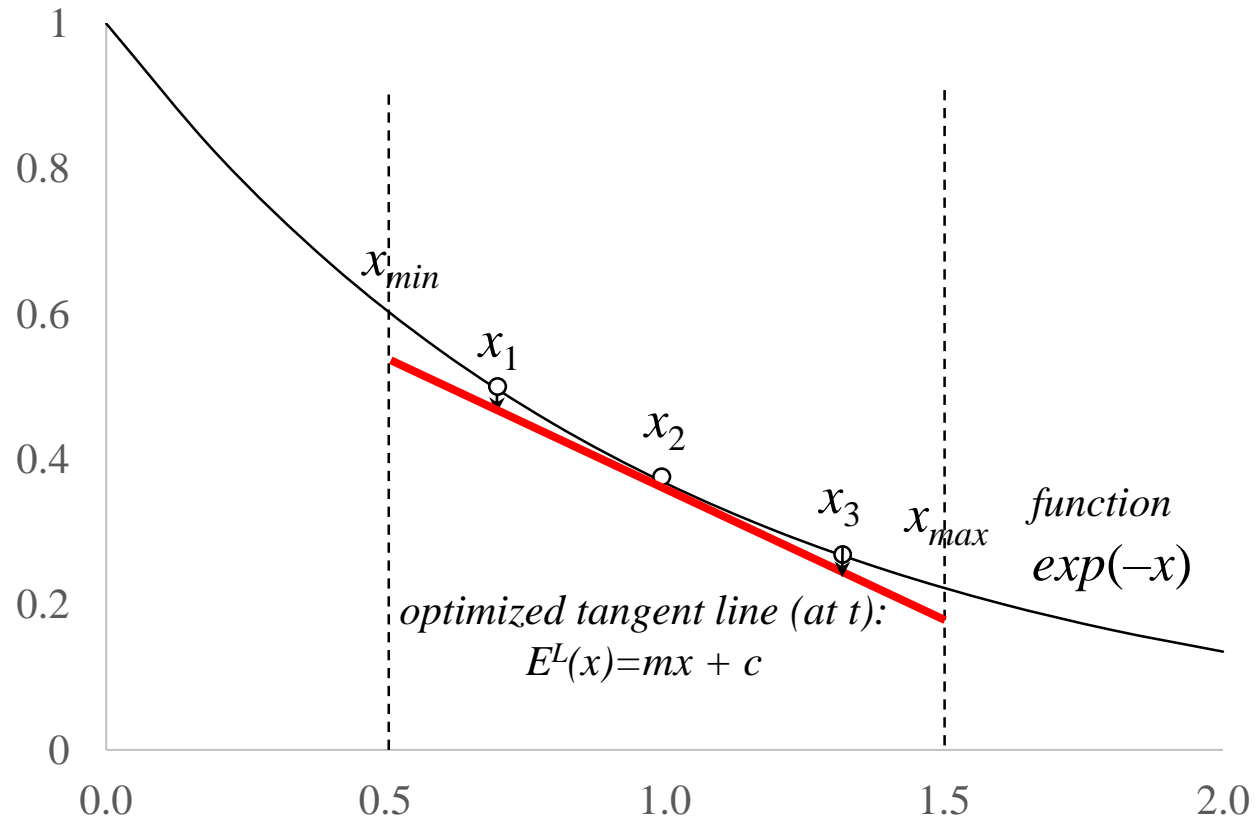
### Research question 3:

How to further tighten our bound (tangent bound)?



*Observation:*

The optimal tangent line depends on the distance values  
(e.g.,  $x_1, x_2, x_3$ )



*Challenge:*

How to compute the optimal tangent **efficiently** (e.g.,  $O(d)$  time)  
**without** knowing the exact distance values (e.g.,  $x_1, x_2, x_3$ )?

*Solved in  
our paper*

# Experimental Results

Methods: SCAN (baseline), LIBSVM, SOTA and KARL

Datasets: UCI Machine Learning Repository/ LibSVM website

## Throughput (Queries/sec)

Type	Datasets	SCAN	LIBSVM	SOTA	KARL
I- $\tau$	miniboone	36.1	34	102	<b>510</b>
	home	15.2	14.1	93.2	<b>258</b>
	susy	2.02	1.86	3.58	<b>83.4</b>
II- $\tau$	nsl-kdd	283	481	748	<b>20668</b>
	kdd99	260	520	1269	<b>11324</b>
	covtype	158	462	448	<b>6022</b>
III- $\tau$	ijcnn1	903	1170	1119	<b>826928</b>
	a9a	162	610	546	<b>6885</b>
	covtype-b	13	38.4	33.9	<b>274</b>

# Conclusion and Future Work

- What we have done:
  1. Develop Kernel Aggregation Rapid Library (KARL)
    1. Support a wide range of models (in prediction stage)
      1. Kernel Density Classification
      2. One Class SVM
      3. Two Class SVM
    2. Achieve higher throughput than the state-of-the-art by 2.5-738x times
- Future Work:
  1. Integrate the implementation into the library (Ongoing)
  2. Support different types of machine learning models (Ongoing)
  3. Support different types of kernel functions (Ongoing)
  4. Support interesting applications, e.g., Kernel Density Visualization (Ongoing)

# For more details...

- The Github page is in **Page 2** of our paper

Our proposal is Kernel Aggregation Rapid Library (KARL)<sup>1</sup>, a comprehensive solution for addressing all the issues mentioned above. It utilizes a novel bounding technique and index tuning in order to achieve excellent efficiency. Experimental studies on many real datasets reveal that our proposed method achieves speedups of 2.5-738 over the state-of-the-art.

Two widely-used libraries, namely LibSVM [8] and Scikit-learn [30], provide convenient programming support for practitioners to handle kernel aggregation queries. Implementation-wise, LibSVM is based on the sequential scan method, and Scikit-learn is based on the algorithm in [16] for query type I. We compare them with our proposal (KARL) in Table II. As a remark, since Scikit-learn supports query types II and III via the wrapper of LibSVM [30], we remove those two query types from the row of Scikit-learn in Table II. The features of KARL are: (i) it supports all three types of weightings

<sup>1</sup><https://github.com/edisonchan2013928/KARL-Fast-Kernel-Aggregation-Queries>