

# Locally Sensitive Hashing for the Content Based Image Retrieval

MASOUD BADIEI KHUZANI, PROF. YINYU YE, PROF. SANDY NAPEL,  
PROF. LEI XING

---

JULY, 2020

STANFORD UNIVERSITY, DEPARTMENT OF RADIATION ONCOLOGY.

# Content-based Image Retrieval

- Content-based image retrieval is the task of searching images by analyzing the contents of the image rather than the metadata such as keywords, tags, or descriptions associated with the image.

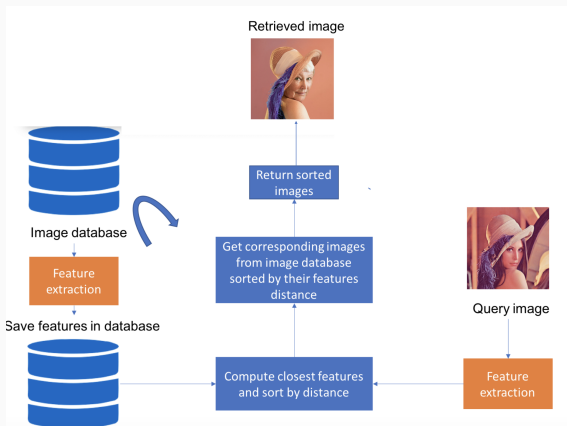


Figure 1: Content-Based Image Retrieval

# Content-based Image Retrieval for Medical Imaging

- Consider the image retrieval problem:

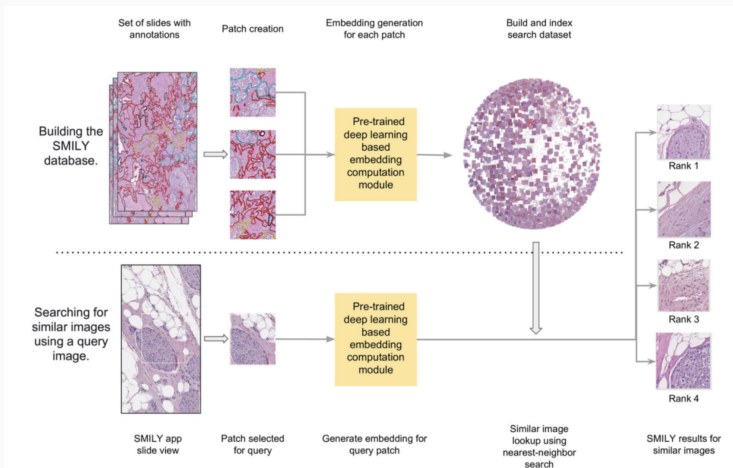
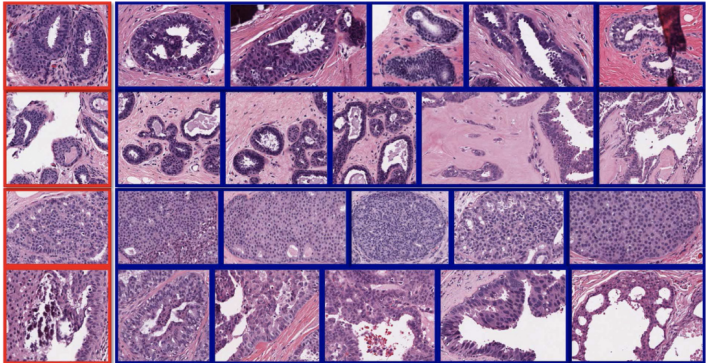


Figure 2: Image retrieval (N Hedge, et al., Nature 2019).

## Other Variants of the Retrieval Systems.

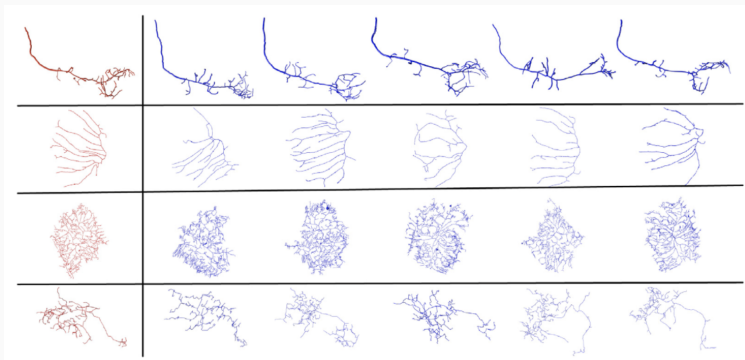
- Such retrieval systems can be applied to histopathological images



**Figure 3:** Image retrieval (Zhang, et al., IEEE Trans. Med 2015).

## Other Variants of the Retrieval Systems.

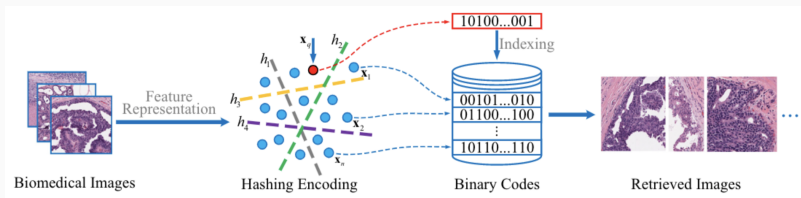
- Such Retrieval system can also be applied to morphological neuron data-base



**Figure 4:** Image retrieval (Zhongyu Li, IEEE Trans. Pattern Recogn. 2017).

## Other Variants of the Retrieval Systems.

- The crux of such retrieval system is a hash function that computes a binary code from the representations of the input



**Figure 5:** Image retrieval (Zhongyu Li, IEEE Trans. Pattern Recogn. 2017).

# What is hashing?

- A hash function is a function that takes a group of characters or numbers (called a key) and maps it into a value of a certain length.
- Hash functions and their associated hash tables are used in data storage and retrieval applications to access data in a small and nearly constant time per retrieval.
- A **collision** or **clash** occurs when two distinct pieces of data have the same hash value

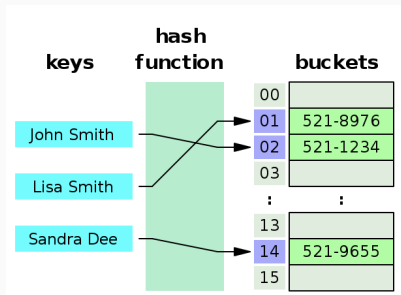
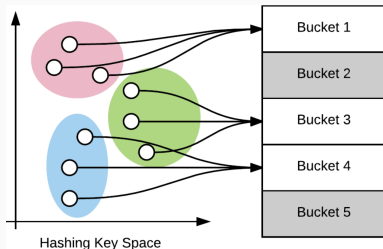


Figure 6: Hash function

# What is the locally sensitive hashing (LSH)?

- A locally sensitive hash (LSH) function is an algorithmic technique that hashes similar input items into the same “buckets” with high probability.
- In a sense, we want a *controlled* collision in LSH.
- Mathematically,

$$\mathbb{P}_{h \in \mathcal{H}}[h(\mathbf{x}) = h(\tilde{\mathbf{x}})] = \text{sim}(\mathbf{x}, \tilde{\mathbf{x}}). \quad (1)$$

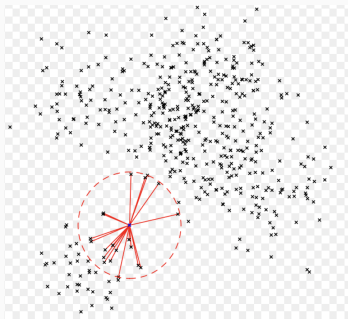


**Figure 7:** Locally Sensitive Hash function.



# What is the use case of LSH?

- Given a query, we need to do a look up via nearest neighbor search.
- The complexity of  $k$ -NN using Euclidean distance is  $\mathcal{O}(mdk)$ , where  $m$  is the size of data-base and  $d$  is the dimension of data-points.



**Figure 8:** Image Lookup via the Nearest Neighborhood.

## An example of LSH

- We can speed up the look up by using LSH that maps data points to a discrete set  $h : \mathbb{R}^d \mapsto \{0, 1, 2, \dots, q - 1\}$ .
- A possible hash function for  $q = 2$  (Hamming code) is proposed by Raginsky, et al. [NIPS 2009], where

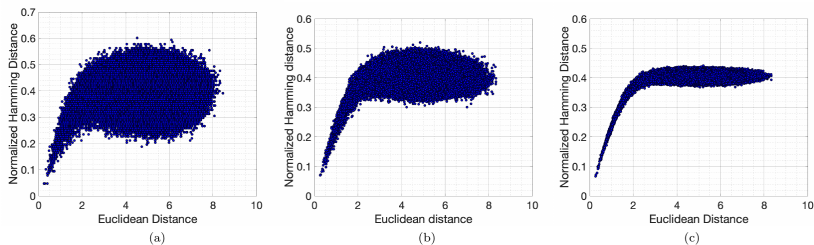
$$h_{t,b,\omega}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \left[ 1 + \text{sign}(\cos(\langle \omega, \mathbf{x} \rangle + b) + t) \right], \quad (2)$$

where

1.  $t \sim \text{Uniform}[-1, 1]$ .
  2.  $b \sim \text{Uniofmr}[-\pi, \pi]$ .
  3.  $\omega \sim \text{N}(\mathbf{0}, \mathbf{I}_{d \times d})$ .
- A hash code of length  $n$  can be constructed by sampling these random variables i.i.d., i.e.,  $h^n(\mathbf{x}) = (h_{t_1, b_1, \omega_1}(\mathbf{x}), \dots, h_{t_n, b_n, \omega_n}(\mathbf{x}))$ , where  $t_1, \dots, t_n \sim \text{Uniform}[-1, 1]$ ,  $b_1, \dots, b_n \sim \text{Uniofmr}[-\pi, \pi]$ , and  $\omega_1, \dots, \omega_n \sim \text{N}(\mathbf{0}, \mathbf{I}_{d \times d})$ .

# Simulations on MNIST Data-Set

- The hamming distance remains flat as the Euclidean distance change.



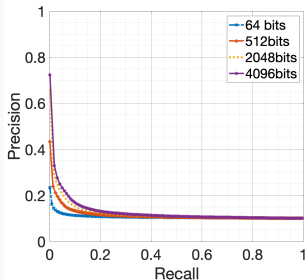
**Figure 9:** The scatter plots of normalized Hamming distance versus Euclidean distance for the hash function, where the hash functions. Panel (b): 512 bits, Panel (c): 4096 bits

# Simulations on MNIST Data-Set

- Another way to look at the performance is via the Precision-Recall curve

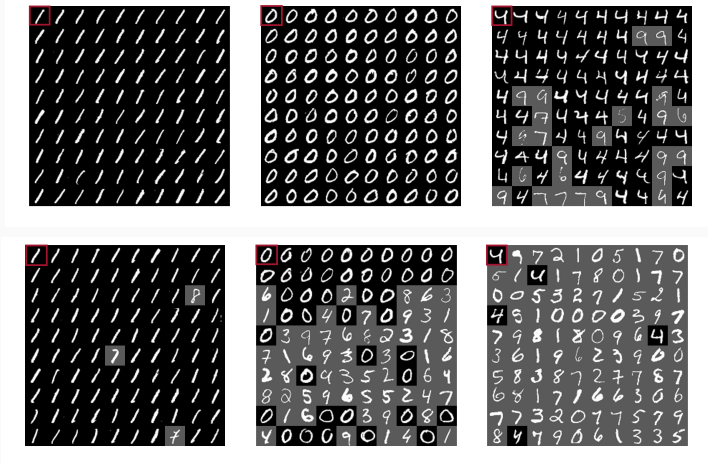
$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}.$$



**Figure 10:** The precision-recall curve for different lengths of the hash codeword.

# Simulations on MNIST Data-Set



**Figure 11:** Examples of image retrieval for three query images on the MNIST database. The images are compressed via an autoencoder with 500 latent variables. The query image is in the top left of the collage (red box), and the incorrect retrieved images are shaded with the gray color. The top row shows top 100 neighbors for each query according to Euclidean distance. The bottom row shows nearest neighbors according to normalized Hamming distance with a 4096-bit code using an untrained

## Why do we get a poor performance?

- The main theorem of Raginsky, et al. [NIPS 2009] computes the probability of collision in terms of a kernel function

$$\mathbb{P}[h_{t,b,\omega}(\mathbf{x}) = h_{t,b,\omega}(\tilde{\mathbf{x}})] = \frac{8}{\pi^2} \sum_{m=0}^{\infty} \frac{1 - K(m\mathbf{x}, m\tilde{\mathbf{x}})}{4m^2 - 1},$$

where we recall the hash function

$$h_{t,b,\omega}(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \left[ 1 + \text{sign}(\cos(\langle \omega, \mathbf{x} \rangle + b) + t) \right], \quad (3)$$

- The main cause of the poor performance is due to the kernel  $K(\mathbf{x}, \tilde{\mathbf{x}}) = \exp(-\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2)$  that is not refined for the hashing task.
- **Task 1:** We need to improve the performance of hash function by optimizing the kernel.
- **Task 2:** We need to extend the hamming codes to general non-binary alphabets.

# A novel algorithm for the optimization of the radial kernels

- Since we would like to capture the Euclidean structure in image retrieval problems, we focus on the class of the *radial kernels*

$$\mathcal{K} \stackrel{\text{def}}{=} \left\{ K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} : K(\mathbf{x}_i, \mathbf{x}_j) = \psi(\|\mathbf{x}_i - \mathbf{x}_j\|_2), \psi \in C^1(\mathbb{R}) \right\}.$$

**Theorem:** (I. J. SCHÖENBERG) A continuous function  $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}$  is positive semi-definite if and only if it admits the following integral representation

$$\psi(r) = \int_0^\infty e^{-tr^2} d\mu(t),$$

for a finite positive Borel measure  $\nu$  on  $\mathbb{R}_+$ . Moreover, if  $\text{supp}(\nu) \neq \{0\}$ , then  $\psi$  is positive definite.

# A novel algorithm for optimization of the kernel

- From the Schönberg's representation theorem, the following integral representation of the radial kernels follows

$$K(\mathbf{x}, \tilde{\mathbf{x}}) = \int_0^\infty e^{-\xi \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2} \mu(d\xi), \quad \forall \mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{R}^d, \mu \in \mathcal{M}_+(\mathbb{R}_+),$$

- We optimize the kernel-target alignment over the distribution  $\mu$  instead of the radial kernels. In particular,

$$\sup_{\mu \in \mathcal{P}} \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} y_i y_j \int_0^\infty e^{-\xi \|\mathbf{x}_i - \mathbf{x}_j\|_2^2} \mu(d\xi),$$

where  $(\mathbf{x}_i, y_i)_{1 \leq i \leq n} \sim P_{\mathbf{x}, y}$  are labeled data.

- **Questions:**

- How do we obtain labeled data  $(y_i, \mathbf{x}_i)_{1 \leq i \leq n}$  in unsupervised tasks such as the image retrieval?
- How do we optimize a distribution?



# $k$ -mean clustering for generating labeled data in unsupervised tasks

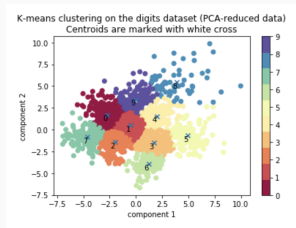
- We perform  $k$ -mean clustering on a small *subset* of the data-base.

$$\arg \min_{S_1, \dots, S_k} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2,$$

where  $\boldsymbol{\mu}_\ell$  is the mean of points in the cluster  $S_\ell$ .

- Choose an index  $\ell \in \{1, 2, \dots, k\}$ , and assign the class labels

$$y_i = \begin{cases} +1 & \mathbf{x}_i \in S_\ell \\ -1 & \mathbf{x}_i \notin S_\ell. \end{cases} \quad (4)$$



# A novel algorithm for the optimization of the radial kernels

- The main idea is to optimize the samples (particles) from the distribution  $\xi = (\xi^1, \dots, \xi^N) \sim_{\text{i.i.d.}} \mu$ .
- Let  $\mathbf{z}_m = (y_m, \mathbf{x}_m)$  and  $\tilde{\mathbf{z}}_m = (\tilde{y}_m, \tilde{\mathbf{x}}_m)$  denote two samples randomly drawn from the training data-set.
- We update the position of the particles using projected Langevin dynamics

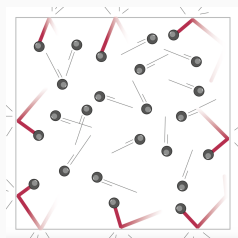
$$\xi_m = P_{\Xi^N} \left( \xi_{m-1} - \eta \nabla \hat{J}^N(\xi_{m-1}; \mathbf{z}_m, \tilde{\mathbf{z}}_m) + \sqrt{\frac{2\eta}{\beta}} \zeta_m \right), \quad (5)$$

- $\xi_m \sim \mathcal{N}(0, \mathbf{I}_{N \times N})$  is isotropic Gaussian noise,  $\beta$  is the inverse temperature, and  $\hat{J}^N(\xi_{m-1}; \mathbf{z}_m, \tilde{\mathbf{z}}_m)$  is a cost function whose  $k$ -th partial derivative is defined as follows

$$\frac{\partial \hat{J}_{\varepsilon, h}^N(\xi_m; \mathbf{z}_m, \tilde{\mathbf{z}}_m)}{\partial \xi_m^k} \stackrel{\text{def}}{=} \frac{\partial E_\gamma(\xi_m; \mathbf{z}_m, \tilde{\mathbf{z}}_m)}{\partial \xi_m^k} + \frac{h}{2} \frac{\partial W_{\varepsilon, 2}^2(\hat{\mu}_m^N, \hat{\nu}^N)}{\partial \xi_m^k}. \quad (6)$$

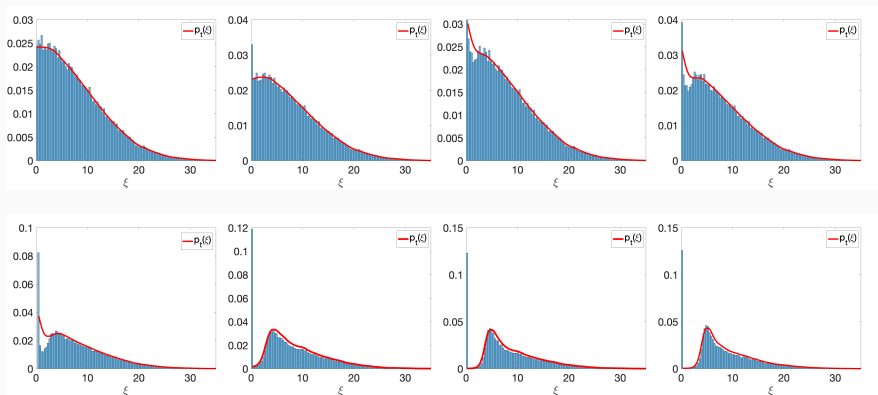
# A novel algorithm for the optimization of the radial kernels

- The intuition for this algorithm is from **the kinetic theory of non-ideal gases** in statistical physics.
- The gas has the inverse temperature of  $T = \frac{1}{\beta}$ , and the particles are interactive.
- Upon hitting the boundaries (wall), the confined particles are reflected **elastically**.



# Histogram of the particles

- We can plot the histogram of the particles at each given iteration of the algorithm



**Figure 12:** The evolution of the histogram of Langevin particles  $\xi^1, \dots, \xi^N \in \Xi = \mathbb{R}_+$

- Define the continuous-time embedding of the empirical measure  $(\widehat{\mu}_k^N)_{k \in \mathbb{N}}$  associated with the projected Langevin particles in Eq. (??), i.e.,

$$\mu_t^N(\xi) \stackrel{\text{def}}{=} \widehat{\mu}_{\lfloor \frac{t}{\eta} \rfloor}^N = \frac{1}{N} \sum_{k=1}^N \delta(\xi - \xi_{\lfloor \frac{t}{\eta} \rfloor}^k), \quad 0 \leq t \leq T, \quad (7)$$

- Suppose the step size  $\eta = \eta_N$  satisfies  $\eta_N \rightarrow 0$ ,  $N/\log(\eta_N/N) \rightarrow \infty$ , and  $\eta_N/\log(\eta_N/N) \rightarrow 0$  as  $N \rightarrow \infty$ .
- Furthermore, suppose the Lebesgue density of the initial particles  $q_0(\xi) = d\mu_0/d\xi$  exists.
- Then, for any fixed  $t \in [0, T]$ ,  $\mu_t^N \xrightarrow{\text{weakly}} \mu_t$  as  $N \rightarrow \infty$ .

- The Lebesgue density of the limiting measure  $\rho_t^*(\xi) = d\mu_t^*/d\xi$  is a solution to the following distributional dynamics with Robin boundary conditions as well as an initial condition

$$\frac{\partial \rho_t(\xi)}{\partial t} = \frac{\partial}{\partial \xi} \left( \rho_t(\xi) \frac{\partial}{\partial \xi} J(\xi, \rho_t(\xi)) \right) + \frac{1}{\beta} \frac{\partial^2}{\partial \xi^2} (\rho_t(\xi)),$$

$$\frac{\partial \rho_t(\xi)}{\partial \xi} + \beta \rho_t(\xi) \frac{\partial}{\partial \xi} J(\xi, \rho_t(\xi)) \Big|_{\xi=\xi_l} = 0, \quad \forall t \in [0, T],$$

$$\frac{\partial \rho_t(\xi)}{\partial \xi} + \beta \rho_t(\xi) \frac{\partial}{\partial \xi} J(\xi, \rho_t(\xi)) \Big|_{\xi=\xi_u} = 0, \quad \forall t \in [0, T],$$

$$\rho_0(\xi) = q_0(\xi), \quad \forall \xi \in [\xi_l, \xi_u]$$

$$\rho_t(\xi) \geq 0, \quad \forall \xi \in [\xi_l, \xi_u], \quad \int_{\Xi} \rho_t(\xi) d\xi = 1, \quad \forall t \in [0, T],$$

where the functional  $J(\xi, \rho_t(\xi))$  is defined

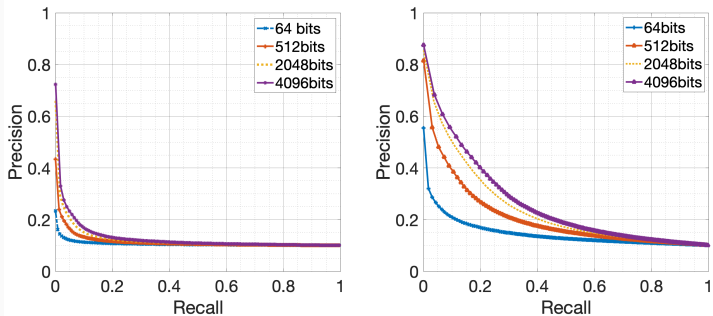
$$J(\xi, \rho_t(\xi)) = \mathbb{E} \left[ y \hat{y} \exp(-\xi \| \mathbf{x} - \hat{\mathbf{x}} \|_2^2) \right] + \int_0^\infty \mathbb{E} \left[ \exp(-(\xi + \xi') \| \mathbf{x} - \hat{\mathbf{x}} \|_2^2) \right] \rho_t(\xi') d\xi'.$$

# Simulations on MNIST Data-Set

- Another way to look at the performance is via the Precision-Recall curve

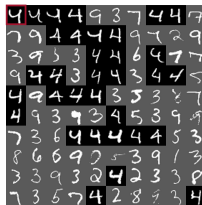
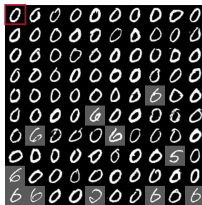
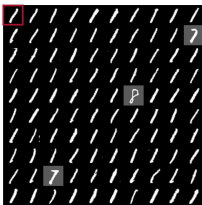
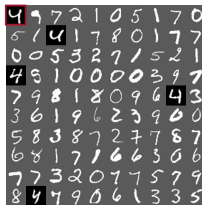
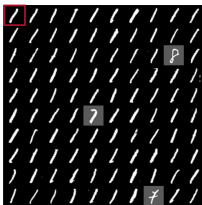
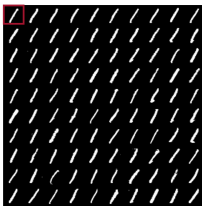
$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$$



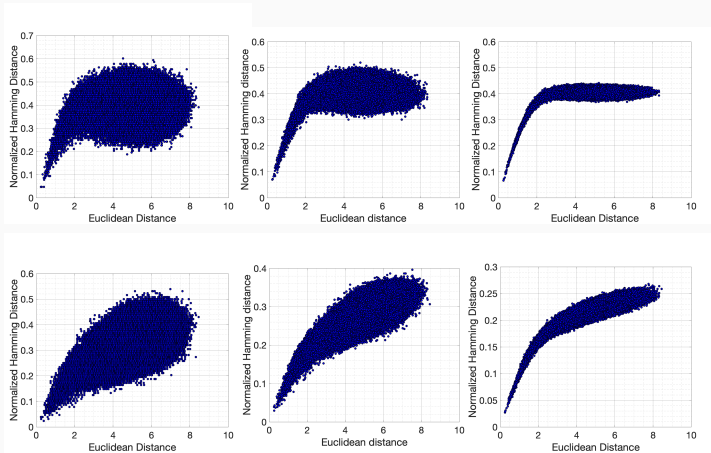
**Figure 13:** The precision-recall curve for different lengths of the hash codeword.

# Simulations on MNIST Data-Set with 4096 bits Hash codes





# Simulations on MNIST Data-Set



**Figure 14:** The scatter plots of normalized Hamming distance versus Euclidean distance, where the hash functions are generated with the random feature using untrained Gaussian kernel with the bandwidth parameter of  $\gamma = 1$  (top row), and trained kernel with our algorithm in conjunction with a  $k$ -mean clustering (bottom row). Panel (a): 128 bits, Panel (b): 512 bits, Panel (c): 4096 bits.

## Can we generalize hash function to other alphabets?

- We propose the following hash function

$$h_{t,\mathbf{w}}(\mathbf{x}) \stackrel{\text{def}}{=} \lceil q(\langle \mathbf{w}, \varphi_{N_0}(\mathbf{x}) \rangle + t) \rceil \bmod q,$$

where  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{N_0 \times N_0})$ ,  $t \sim \text{Uniform}[0, 1]$ , and

$$\varphi_{N_0}(\mathbf{x}) \stackrel{\text{def}}{=} (\cos(\langle \boldsymbol{\omega}_k, \mathbf{x} \rangle + b_k))_{1 \leq k \leq N_0},$$

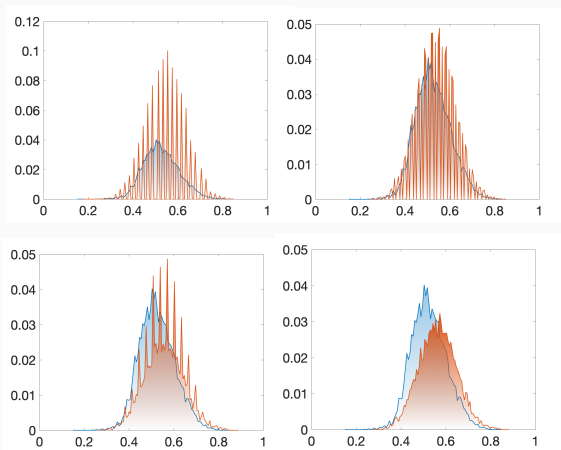
- To compute the distance between two hash codes, we leverage the Lee distance between two code-words

$\mathbf{x} = (x_1, \dots, x_n)$ ,  $\mathbf{y} = (y_1, \dots, y_n) \in \mathbb{Z}_q^n$  of length  $n$  as follows

$$\begin{aligned} d_{\text{Lee}}(\mathbf{x}, \mathbf{y}) &\stackrel{\text{def}}{=} \sum_{i=1}^n \min\{(x_i - y_i) \bmod q, (y_i - x_i) \bmod q\} \\ &= \sum_{i=1}^n \min\{|y_i - x_i|, q - |y_i - x_i|\}. \end{aligned}$$

In the special case of  $q = 2$  and  $q = 3$ , the Lee distance corresponds to the Hamming distance.

# Histograms of Lee distance vs Euclidean distance



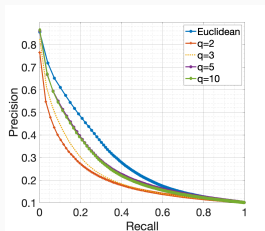
**Figure 15:** The (normalized) histogram of the (rescaled) Lee distance (red color) and the (rescaled) Euclidean distance (blue color) of a sample query from each point of the MNIST data-base for different values of  $q$ , Panel (a):  $q = 2$ , Panel (b):  $q = 20$ , Panel (c):  $q = 50$ . Increasing  $q$  in the Lee distance yields a better approximation for the Euclidean distance.

# Simulations on MNIST Data-Set

- Another way to look at the performance is via the Precision-Recall curve

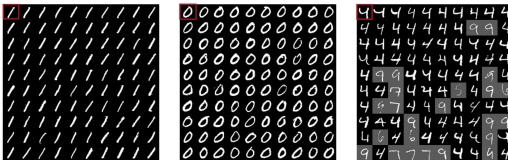
$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}.$$

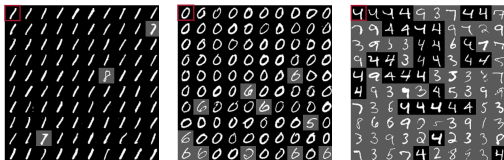


**Figure 16:** The precision-recall curve for different lengths of the hash codeword.

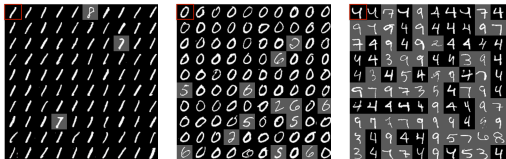
# Simulations on MNIST Data-Set



(Euclidean distance)



(4096 Hamming  $Z_2^n$  codes)



(1024 Hamming  $Z_{10}^n$  codes)

## Application to Classification Tasks

- The proposed method for optimizing RBF kernels can also be used for other machine learning tasks such as classification using kernel SVMs.
- Given the training samples  $(y_i, \mathbf{x}_i)_{1 \leq i \leq n}$ , the kernel SVMs solve the following optimization problem

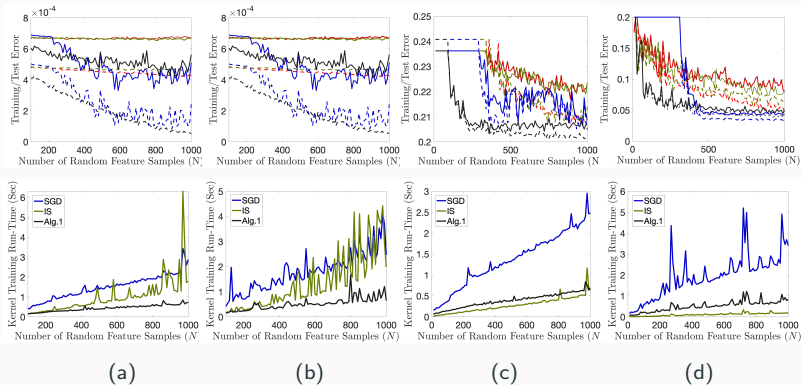
$$\min_{\alpha \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - \sum_{k=1}^D y_i \alpha_i \varphi(\mathbf{x}_i; \omega_k) \right\} + \frac{\lambda}{2} \|\alpha\|_2^2.$$

- Here,  $\varphi(\mathbf{x}_i; \omega_k)$  is the *random feature* associated with a kernel in Rahimi and Recht integral representation of the translation invariant kernels

$$K(\mathbf{x}, \tilde{\mathbf{x}}) = \int_{\Omega} \varphi(\mathbf{x}; \omega) \varphi(\tilde{\mathbf{x}}; \omega) \nu(d\omega), \quad (9)$$

where  $K(\mathbf{x}, \tilde{\mathbf{x}}) = \phi(\mathbf{x} - \tilde{\mathbf{x}})$ , and  $\nu$  is the probability distribution of random variable  $\omega$ .

# Application to Classification Tasks



**Figure 17:** The training and test errors (first row) and the run-times (second row) of kernel optimization algorithms using proposed method, the SGD optimization, and Importance Sampling (IS). Panel (a): BUZZ, Panel (b): ONLINE NEWS POPULARITY, Panel (c): ADULT, Panel (d): SEIZURE

# Application to Classification Tasks

- Many ideas and techniques in statistical physics are relevant to machine learning problems.
- The proposed method for optimizing RBF kernels in our work is based on idea of non-ideal gases.
- Currently, we are trying to apply the methods we developed to medical imaging and the isodose distribution for radiation therapy treatment planning.