

M2177.0043 Introduction to Deep Learning

Lecture 19: Metric Learning

Hyun Oh Song¹

¹Dept. of Computer Science and Engineering, Seoul National University

May 25, 2020

Last time

- ▶ Disentanglement

Outline

Metric learning

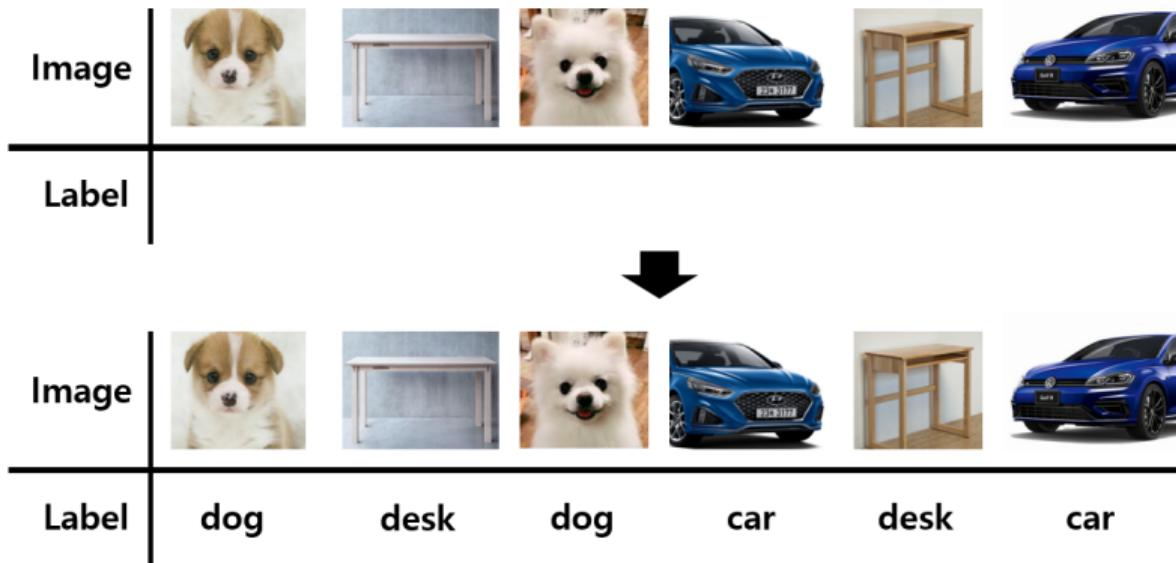
Definition

Efficient inference via hashing

Classification

- The objective of classification problem is to classify the images into one of the specified labels.

Label $\in \{dog, desk, car\}$



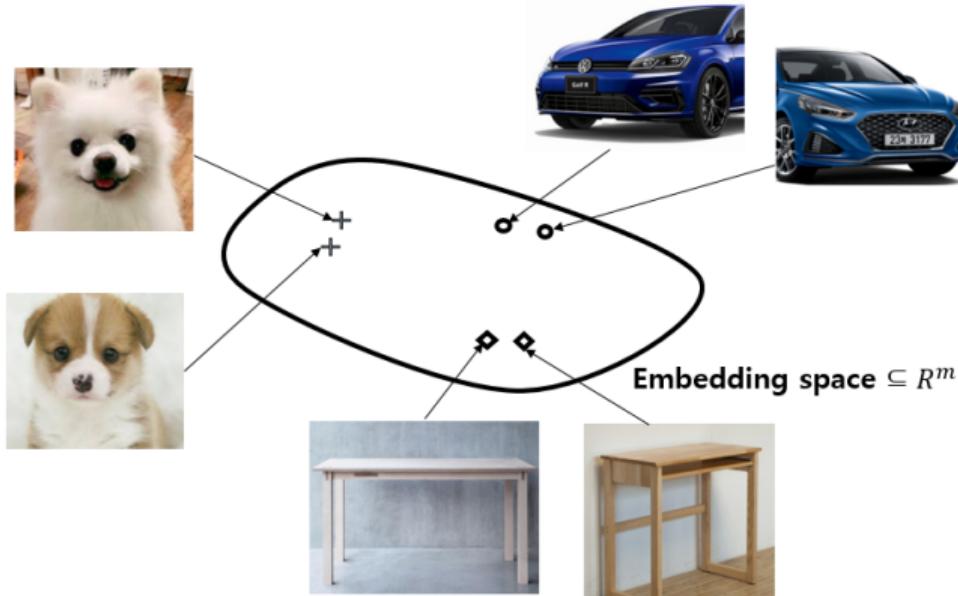
The limitation of classifications

- ▶ If there are extremely many possible labels, classification logit dimension scales along with it
- ▶ Classification is not adequate to classify the images with new labels.
- ▶ e.g. We cannot classify human faces with 7 billion labels.



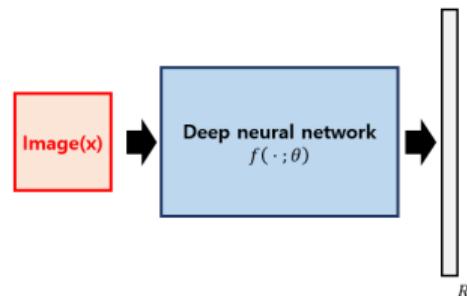
What is metric learning?

- ▶ Learn an embedding representation space where similar data are close to each other and different data are far from each other.



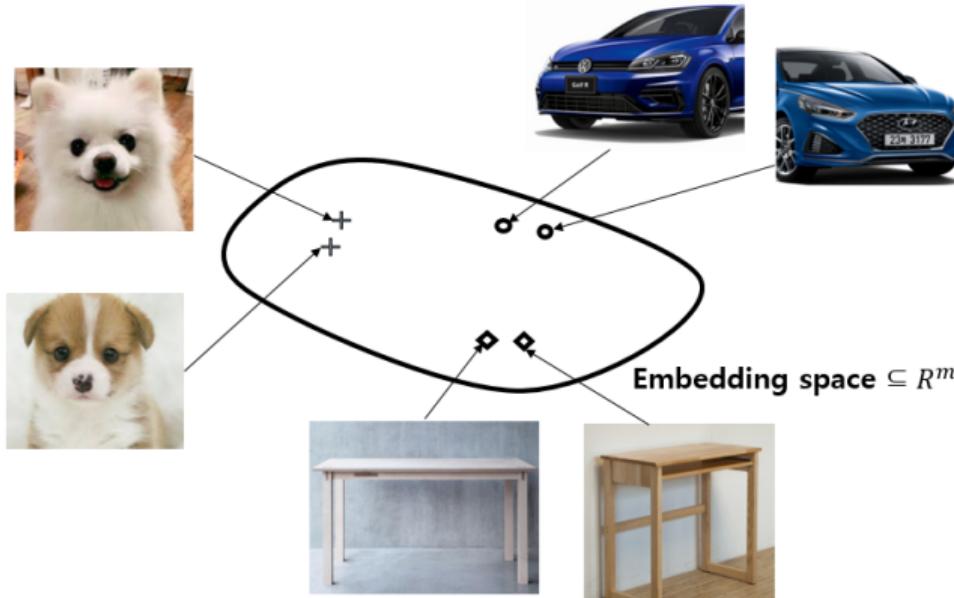
Notation

- ▶ $\mathcal{X} = \{x_1, \dots, x_n\}$: a set of images
- ▶ $x_i \in \mathcal{X}$: image
- ▶ $y_i \in \{0, 1, \dots, C - 1\}$: class label of image x_i where C is the number of classes
- ▶ $f : \mathcal{X} \rightarrow \mathbb{R}^m$ is the embedding function where m is the embedding dimension. m is a fixed small number which does not scale with C
- ▶ θ is the parameter to be learned in $f(\cdot; \theta)$

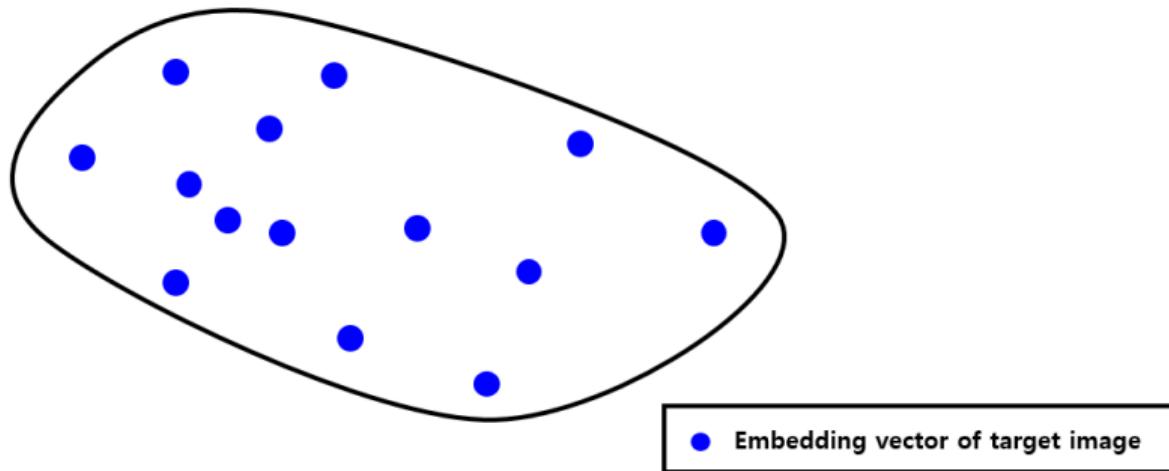


Training objective of metric learning

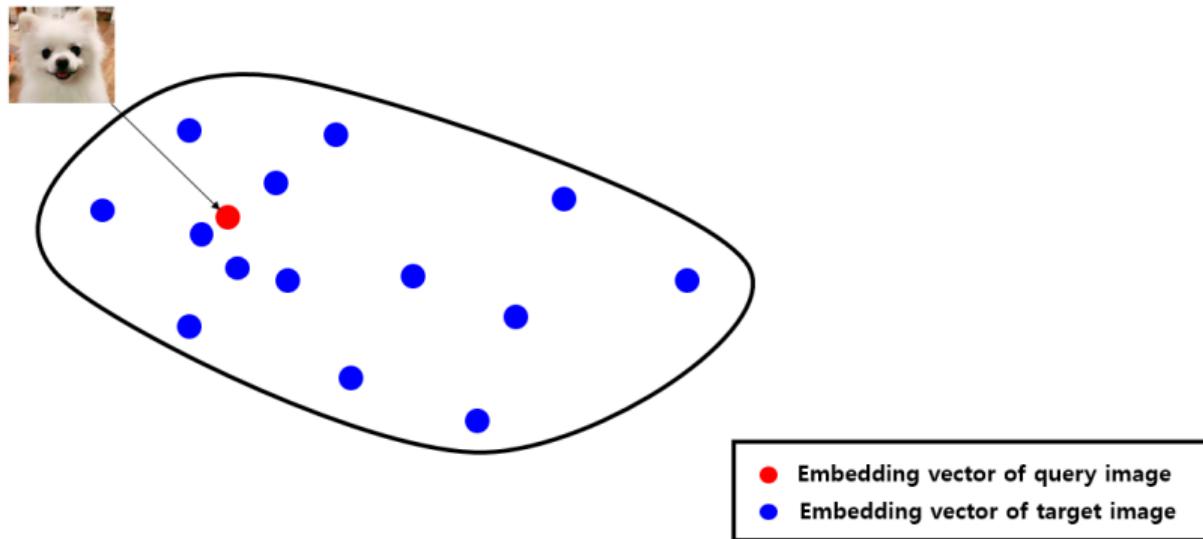
$$\triangleright D_{i,j} = \|f(x_i; \theta) - f(x_j; \theta)\|_p = \begin{cases} \text{small} & \text{if } y_i = y_j \\ \text{large} & \text{if } y_i \neq y_j \end{cases}$$



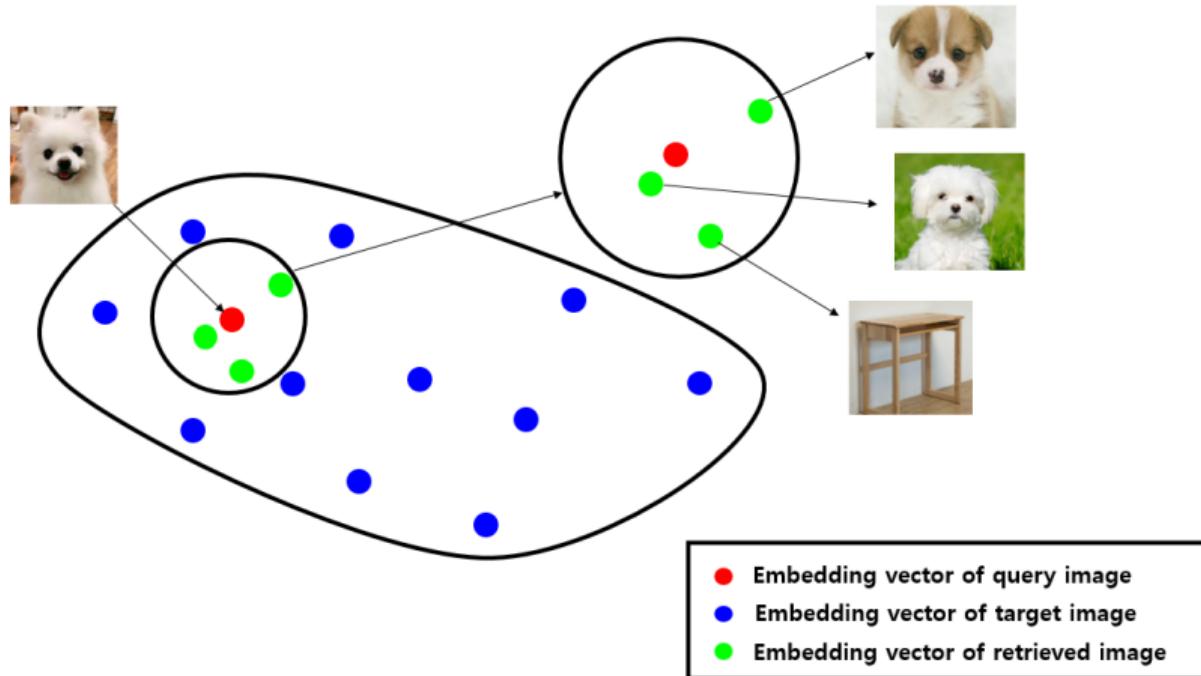
Retrieval Procedure



Retrieval Procedure

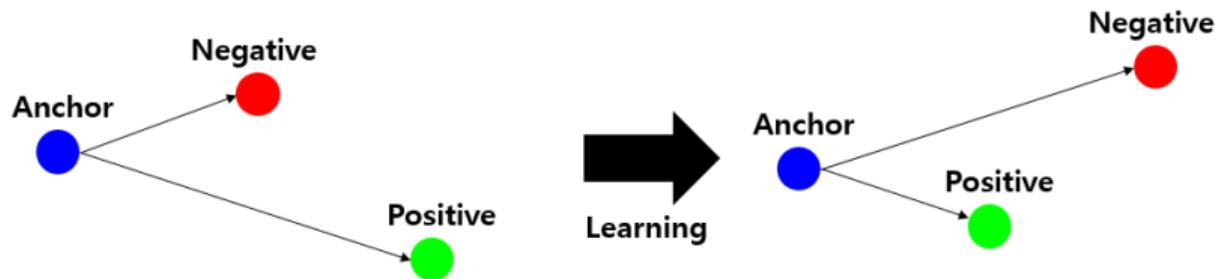


Retrieval Procedure



Triplet loss

- ▶ $\mathcal{T} = \{(a, p, n) \mid y_a = y_p \neq y_n\}$: a set of triplets(anchor, positive, negative)
- ▶ $D_{a,p}$ is trained to be reduced and $D_{a,n}$ is trained to be increased.



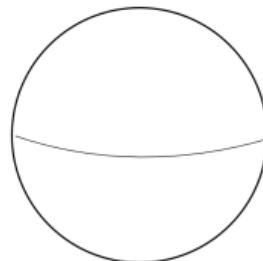
Triplet loss

$$\ell(\mathcal{X}, \mathbf{y}) = \frac{1}{|\mathcal{T}|} \sum_{(a,p,n) \in \mathcal{T}} [D_{a,p}^2 + \alpha - D_{a,n}^2]_+$$

- ▶ \mathcal{X} : a set of images
- ▶ \mathbf{y} : a set of classes
- ▶ \mathcal{T} : a set of triplets(anchor, positive, negative)
- ▶ $[\cdot]_+ = \max(\cdot, 0)$
- ▶ Training process reduces $\ell(\mathcal{X}, \mathbf{y})$ which means $D_{a,p}$ decreases and $D_{a,n}$ increases.

Facenet¹

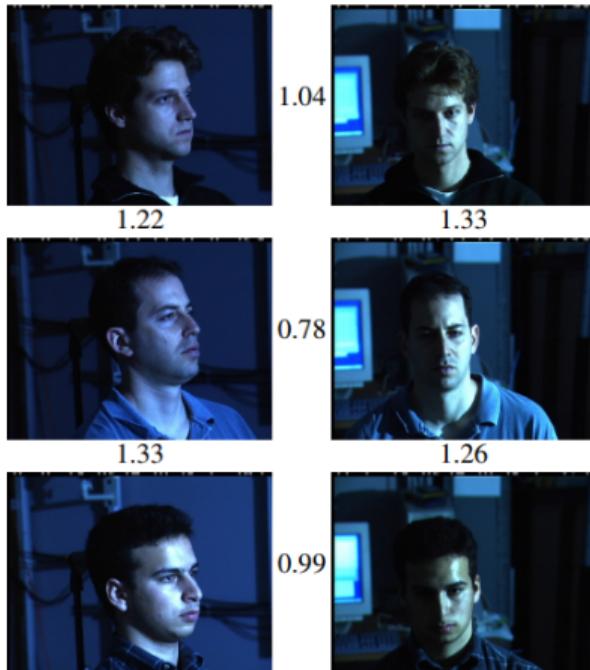
- ▶ This paper applied triplet loss for the metric learning on a face dataset, Labeled Faces in the Wild(LFW).
- ▶ The numbers in the figure mean the distances between images.
- ▶ Facenet restricted the embedding space of triplet loss on the hypersphere. $f(x; \theta) \in \{y \in \mathbb{R}^m \mid \|y\|_2^2 = 1\}$



$$\text{Hypersphere } \|f(x; \theta)\|_2^2 = 1$$

¹Schroff et al. "FaceNet: A unified embedding for face recognition and clustering" CVPR2015

Facenet



- ▶ The distances between same identities with different poses are shorter than that between different identities.

Hard negative mining

- ▶ In practice, the performance of triplet loss depends highly on the triplet sampling strategy.
- ▶ Facenet paper suggested the following online hard negative mining strategy.
- ▶ The idea is to construct triplets by associating with each positive pair in the minibatch a “semi-hard” negative example.
- ▶ This is an example which is further away from the anchor i than the positive exemplar j is, but still hard because the distance is close to the $i - j$ distance.

Hard negative mining

- We can define triplet loss with this hard negative mining.

$$\ell(X, \mathbf{y}) = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \left[D_{i,j}^2 + \alpha - D_{i,k^*(i,j)}^2 \right]_+$$

$$\text{where } k^*(i,j) = \operatorname{argmin}_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} D_{i,k}^2 \text{ s.t. } D_{i,k}^2 > D_{i,j}^2$$

Hard negative mining

- We can define triplet loss with this hard negative mining.

$$\ell(X, \mathbf{y}) = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \left[D_{i,j}^2 + \alpha - D_{i,k^*(i,j)}^2 \right]_+$$

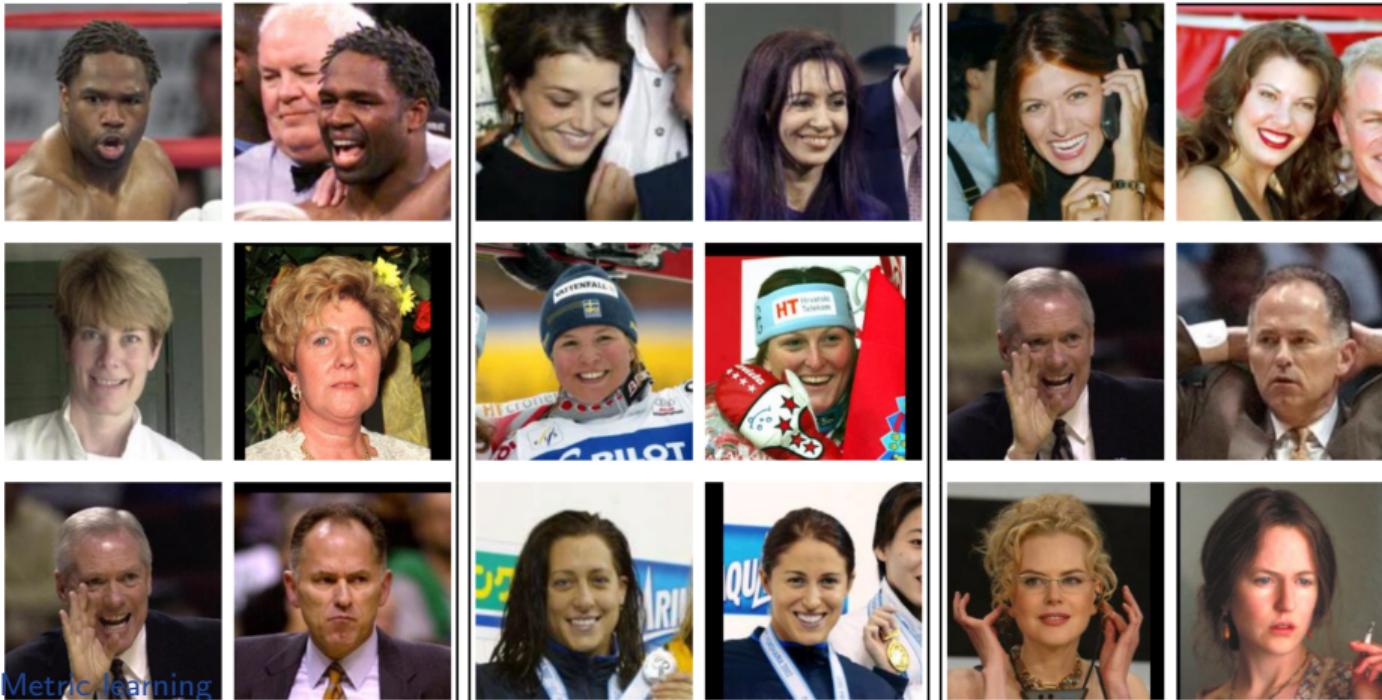
$$\text{where } k^*(i,j) = \operatorname{argmin}_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} D_{i,k}^2 \text{ s.t. } D_{i,k}^2 > D_{i,j}^2$$

- If there is no such negative example satisfying the constraint, we just pick the furthest negative example in the minibatch, as follows:

$$k^*(i,j) = \operatorname{argmax}_{k: \mathbf{y}[k] \neq \mathbf{y}[i]} D_{i,k}^2$$

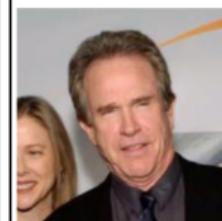
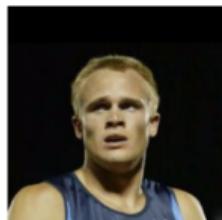
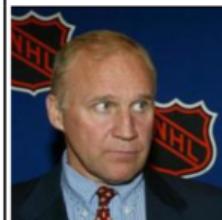
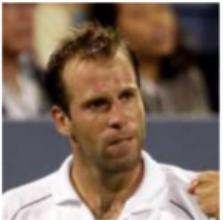
Failed k-nearest neighbor result

False reject



Failed k-nearest neighbor result

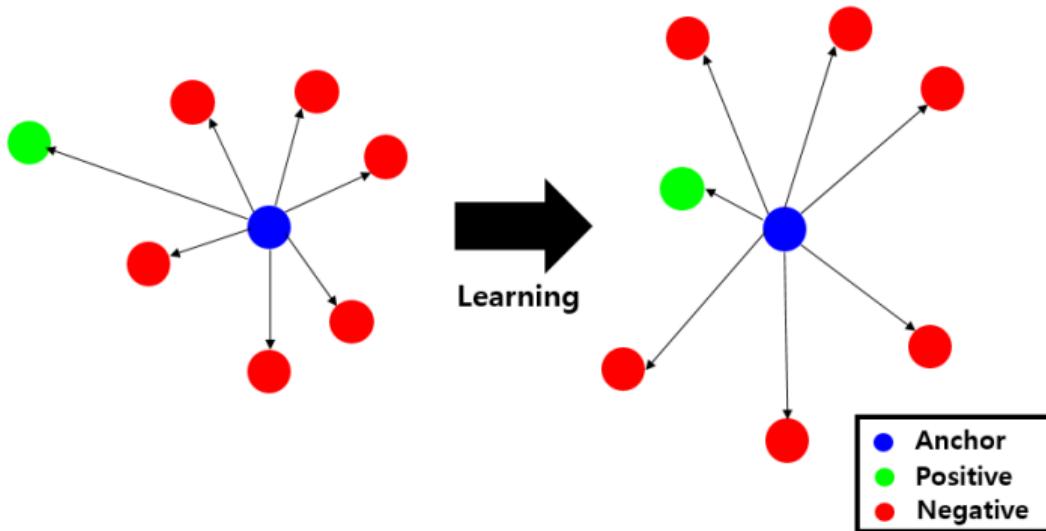
False accept



Metric learning

Npairs loss

- ▶ There are multiple negatives with an anchor and a positive.
- ▶ The objective is to reduce the distance between the **anchor** and the **positive** increasing the distance between the **anchor** and **negatives**.



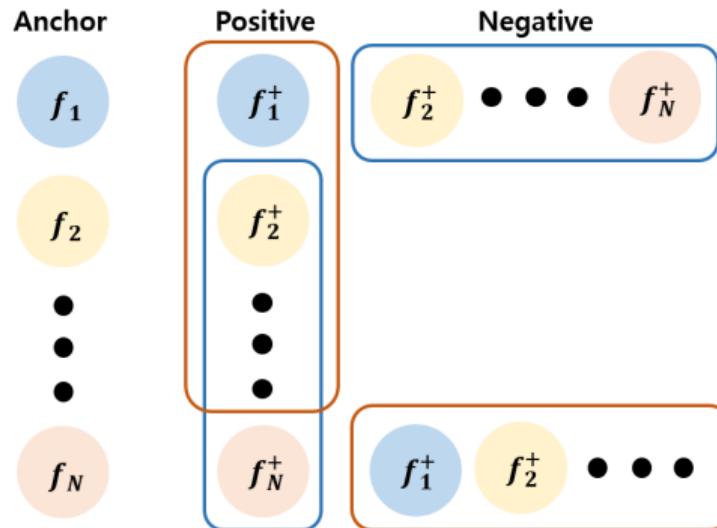
Npairs loss

$$\ell(\mathcal{X}, \mathbf{y}) = -\frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \log \frac{\exp(-D_{i,j})}{\exp(-D_{i,j}) + \sum_{k:y_k \neq y_i} \exp(-D_{i,k})} + \lambda \sum_i \|f(x_i; \boldsymbol{\theta})\|_2^2$$

- ▶ $\mathcal{P} = \{(i, j) \mid y_i = y_j\}$: a set of pairs with the same label
- ▶ $\sum_i \|f(x_i; \boldsymbol{\theta})\|_2^2$: the regularizer term
- ▶ Training process reduces $\ell(\mathcal{X}, \mathbf{y})$ which means $D_{i,j}$ with $y_i = y_j$ decreases and $D_{i,k}$ with $y_i \neq y_k$ increases.

Batch Construction

- ▶ Every batch contain N different labels with 2 images each for efficient batch construction.
- ▶ N anchor images have their positive images respectively and positive images of the other anchor images are used as negative images.



Outline

Metric learning
Definition

Efficient inference via hashing

Efficient inference via hashing

- ▶ Once we learn the embedding space, what's the most efficient way to utilize the embeddings for fast data retrieval?
- ▶ Naive approach of running nearest neighbor per each query is not scalable. Why?
- ▶ $r(\cdot) : \mathcal{X} \rightarrow \{0, 1\}^d, \quad \|r(\cdot)\|_1 = k$
- ▶ $f(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^d$ differentiable transformation with respect to parameter θ

Example hashing functions - Thresholding²

- ▶ Let $f(\cdot)$ be the embedding network trained with metric learning methods (*i.e.* triplet, npairs loss).
- ▶ Then the hash function $r(\mathbf{x})$ computes k largest dimensions from the embedding representation $f(\mathbf{x})$

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$

subject to $\|\mathbf{h}\|_1 = k,$ (1)

²Zhai et al. "Visual discovery at pinterest" ICWWW2017

Example hashing functions - Thresholding²

- ▶ Let $f(\cdot)$ be the embedding network trained with metric learning methods (*i.e.* triplet, npairs loss).
- ▶ Then the hash function $r(\mathbf{x})$ computes k largest dimensions from the embedding representation $f(\mathbf{x})$

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$

subject to $\|\mathbf{h}\|_1 = k,$ (1)

- ▶ Pros: Computing the hash code per query $r(\mathbf{x})$ is fast and simple.
 $O(1)$ in the number of data
- ▶ Cons: The training objective for $f(\cdot)$ does not take the hashing procedure into account. In other words, no end-to-end training
⇒ drop in search accuracy

²Zhai et al. “Visual discovery at pinterest” ICWWW2017

Example hashing functions - Quantization³

- ▶ Run dictionary learning or clustering (*i.e.* kmeans) on the training data in the embedding space
- ▶ Denote $\mathbf{c}_1, \dots, \mathbf{c}_d$ as the precomputed cluster centers in the embedding space
- ▶ Let $f(\mathbf{x}) \triangleq [-\|\mathbf{x} - \mathbf{c}_1\|, \dots, -\|\mathbf{x} - \mathbf{c}_d\|]$, where the distance is measured in the embedding space. Then the hash function $r(\mathbf{x})$ computes the k nearest cluster indices.

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$

subject to $\|\mathbf{h}\|_1 = k,$ (2)

Example hashing functions - Quantization³

- ▶ Run dictionary learning or clustering (*i.e.* kmeans) on the training data in the embedding space
- ▶ Denote $\mathbf{c}_1, \dots, \mathbf{c}_d$ as the precomputed cluster centers in the embedding space
- ▶ Let $f(\mathbf{x}) \triangleq [-\|\mathbf{x} - \mathbf{c}_1\|, \dots, -\|\mathbf{x} - \mathbf{c}_d\|]$, where the distance is measured in the embedding space. Then the hash function $r(\mathbf{x})$ computes the k nearest cluster indices.

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$
$$\text{subject to } \|\mathbf{h}\|_1 = k, \quad (2)$$

- ▶ Pros: Better search accuracy than simple thresholding
- ▶ Cons: 1) The training objective for $f(\cdot)$ does not take the hashing procedure into account. 2) running k-means on a large dataset of size n is $O(n^2)$

Example hashing functions - End-to-End learning⁴

- ▶ Can we re-design the training objective jointly optimize for both the metric learning objective of $f(\cdot)$ **and** the hashing performance of $r(\cdot)$?

⁴ Jeong & Song “Efficient end-to-end learning for quantizable representation” ICML2018

Intuition

- ▶ Finding the optimal set of embedding representations and the corresponding hash codes is a chicken and egg problem.
- ▶ Embedding representations are required to infer which k activation dimensions to set in the corresponding binary hash code.
- ▶ Binary hash codes are needed to adjust the embedding representations indexed at the activated bits so that similar items get hashed to the same buckets and vice versa.
- ▶ This notion leads to alternating minimization scheme.

Alternating minimization scheme

$$\begin{aligned} & \underset{\theta}{\text{minimize}} \quad \underbrace{\ell_{\text{metric}}(\{f(\mathbf{x}_i; \theta)\}_{i=1}^n; \mathbf{h}_1, \dots, \mathbf{h}_n)}_{\text{embedding representation quality}} + \\ & \quad \underbrace{\gamma \left(\sum_i^n -f(\mathbf{x}_i; \theta)^\top \mathbf{h}_i + \sum_i^n \sum_{j:y_j \neq y_i} \mathbf{h}_i^\top P \mathbf{h}_j \right)}_{\text{hash code performance}} \\ & \text{subject to } \mathbf{h}_i \in \{0, 1\}^d, \|\mathbf{h}_i\|_1 = k, \forall i, \end{aligned} \tag{3}$$

- ▶ Solving for binary hash codes $\mathbf{h}_{1:n}$ with sparsity k
- ▶ Updating the parameter in deep neural network θ .

Observation

$$\begin{aligned} \underset{\mathbf{h}_1, \dots, \mathbf{h}_n}{\text{minimize}} \quad & \sum_i^n -f(\mathbf{x}_i; \boldsymbol{\theta})^\top \mathbf{h}_i + \underbrace{\sum_i^n \sum_{j: y_j \neq y_i} \mathbf{h}_i^\top P \mathbf{h}_j}_{:= g(\mathbf{h}_1, \dots, \mathbf{h}_n; \boldsymbol{\theta})} \\ \text{subject to} \quad & \mathbf{h}_i \in \{0, 1\}^d, \|\mathbf{h}_i\|_1 = k, \forall i, \end{aligned} \tag{4}$$

- ▶ Unary term encourages to select k large elements in embedding vector($f(\cdot; \boldsymbol{\theta})$)
- ▶ Pairwise term selects as orthogonal elements as possible across between different classes
- ▶ NP-hard problem even in simple case $k = 1, d > 2$
- ▶ Refer to the solution to the paper (Jeong & Song ICML18) if interested. Omitted because it's outside the scope of the class.

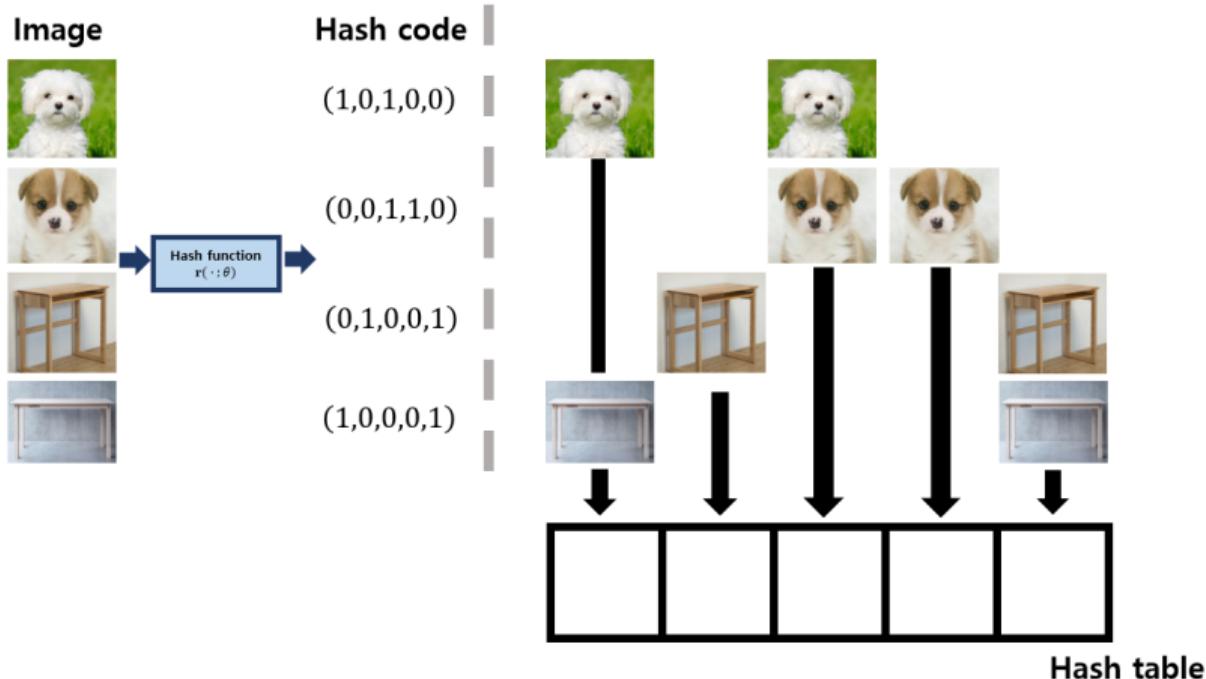
Inference

- ▶ So given the hash function $r(\cdot)$, how do we actually do inference (nearest neighbor search in the embedding space)?

Hash table construction

- ▶ For each image $x \in \mathcal{X}$, $\|r(x)\|_1 = k$.
- ▶ There are d buckets in the hashtable.
- ▶ Insert image x in every i th bucket only if $r(x)[i] = 1$.

Inference

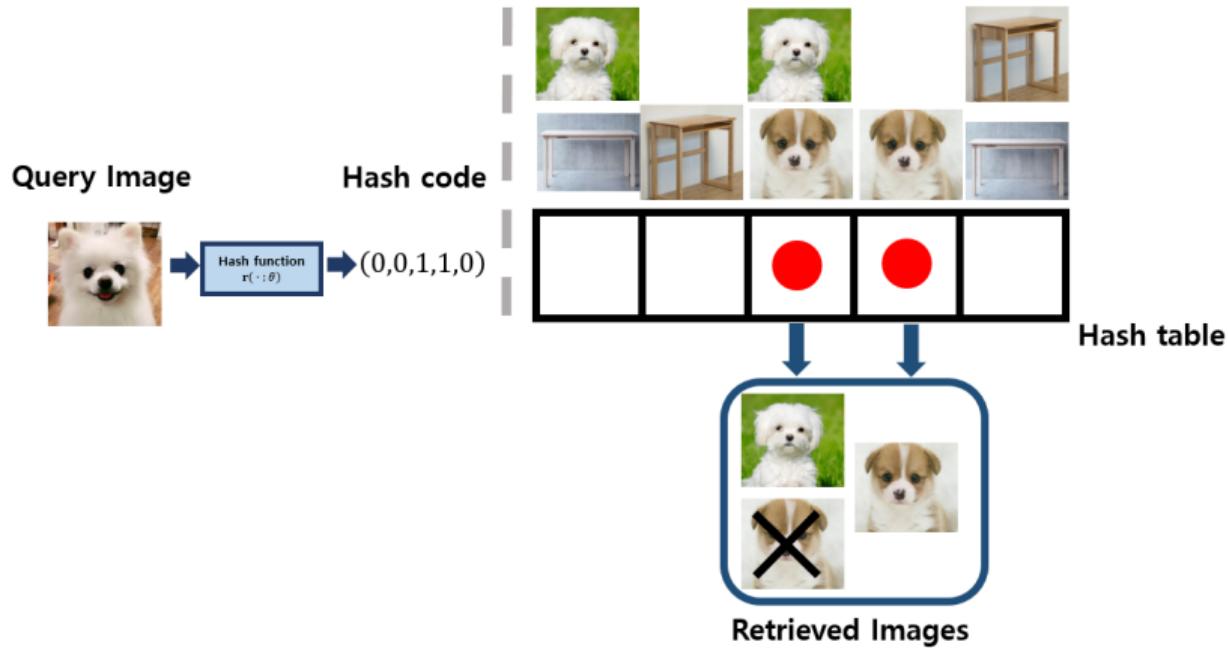


Inference

Query on the hash table

- ▶ After query image q enters, it converts to the hash codes $r(q)$.
- ▶ Retrieve every images in i th buckets where $r(q)[i] = 1$.
- ▶ Remove duplicated images.

Inference



Theoretical search speedup

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$

subject to $\|\mathbf{h}\|_1 = k,$

- ▶ Under this hashing scheme, what is the expected speedup versus the exhaustive linear search? Assume the hash table is evenly distributed

Theoretical search speedup

$$r(\mathbf{x}) = \underset{\mathbf{h} \in \{0,1\}^d}{\operatorname{argmin}} -f(\mathbf{x}; \boldsymbol{\theta})^\top \mathbf{h}$$

subject to $\|\mathbf{h}\|_1 = k$,

- ▶ Under this hashing scheme, what is the expected speedup versus the exhaustive linear search? Assume the hash table is evenly distributed
- ▶ Given the hash code for the query \mathbf{h}_q , the expected number of retrieved items is $\sum_{i=1}^n \Pr(\mathbf{h}_i^\top \mathbf{h}_q \neq 0)$
- ▶ The expected speedup is the ratio between the total number of items (n) and the expected number of retrieved items:

$$(\Pr(\mathbf{h}_i^\top \mathbf{h}_q \neq 0))^{-1} = \left(1 - \frac{\binom{d-k}{k}}{\binom{d}{k}}\right)^{-1}$$

- ▶ In case $d \gg k$, this ratio approaches $\frac{d}{k^2}$
- ▶ Dense hash code \implies less speedup but more accurate search. Vice versa for sparse hash code