

Deep Learning (BEV033DLE)

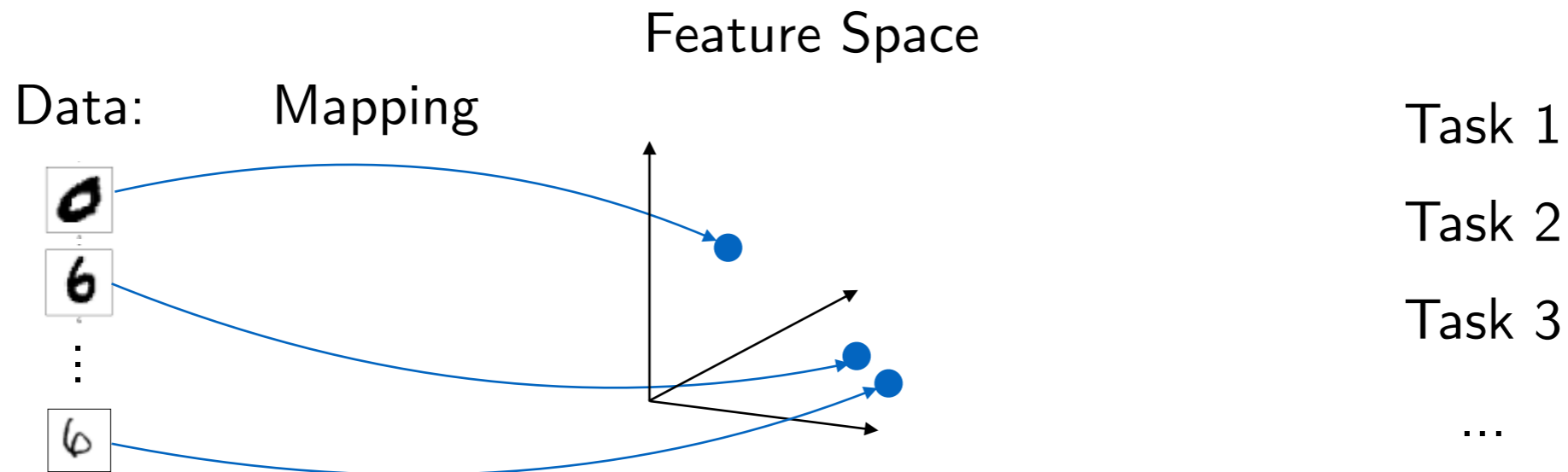
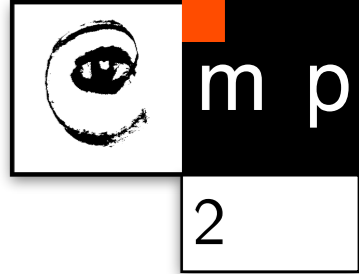
Lecture 10

Learning Representations I

Czech Technical University in Prague

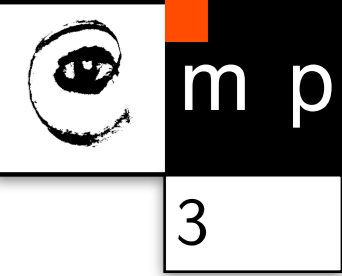
- ◆ Lecture 10: LR-1:
 - Feature Space Representations
 - Word Vectors
 - Similarity / Metric Learning
 - Cross-Modality Representations

Feature Space Representation



- ✦ With good features many tasks are easy:
 - E.g. logistic regression atop of deep features could easily classify butterflies
 - Finding similar objects can be done by nearest neighbor search
- ✦ Suppose we are interested in high-level (semantic tasks). What would we like good features to do?
 - Keep useful information (for all relevant tasks)
 - Discard unnecessary information (view point, lighting, etc.)
 - Similar representations should correspond to semantically similar objects

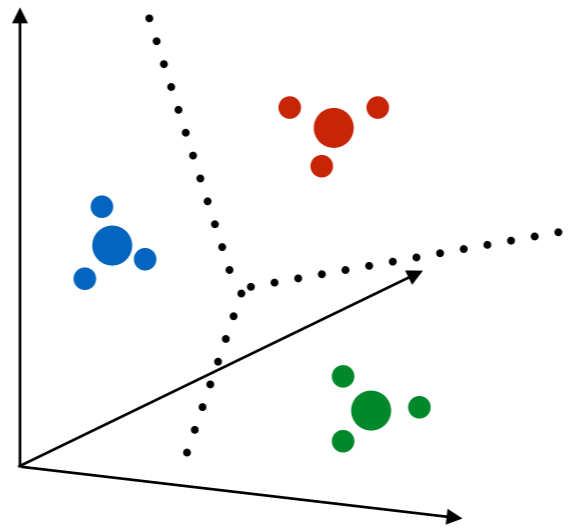
Examples



Many tasks become easier if we have good feature

✦ Classification:

- SVM
- Logistic regression
- Nearest neighbor classifier
- Fine-tune the whole model (same or different data)

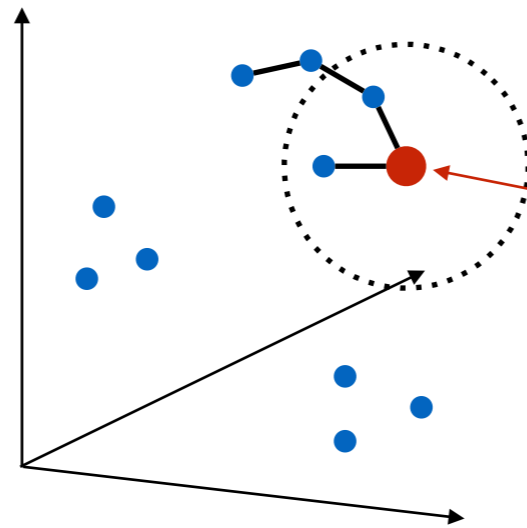


Examples

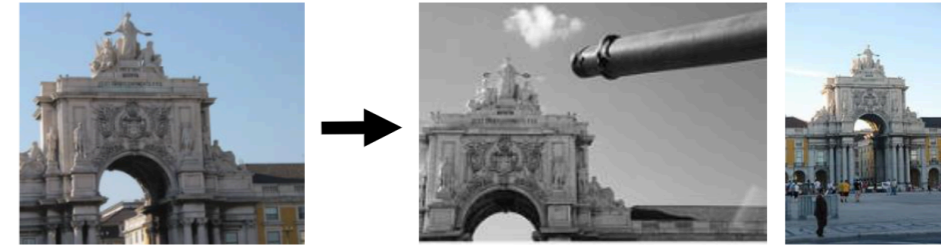
Many tasks become easier if we have good feature

◆ Classification:

- SVM
- Logistic regression
- Nearest neighbor classifier
- Fine-tune the whole model (same or different data)



◆ Visual search (retrieval):



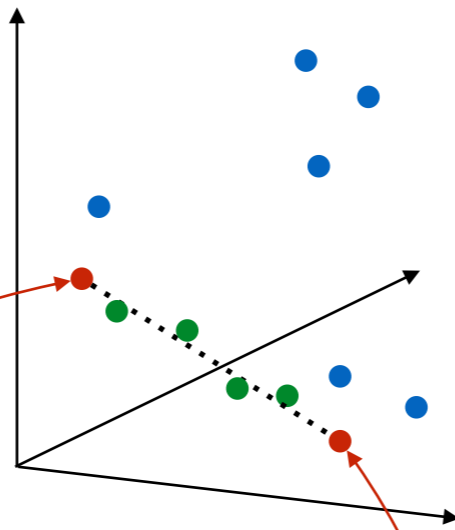
- query with an image
- retrieve similar objects or views
- Euclidean nearest neighbor
- NN graph distance

Examples

Many tasks become easier if we have good feature

◆ Classification:

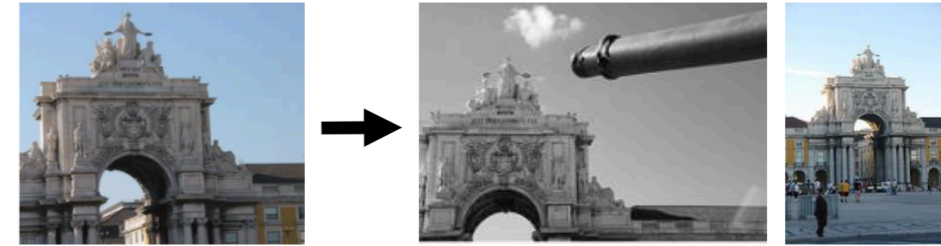
- SVM
- Logistic regression
- Nearest neighbor classifier
- Fine-tune the whole model (same or different data)



◆ Data exploration



◆ Visual search (retrieval):



- query with an image
- retrieve similar objects or views
- Euclidean nearest neighbor
- NN graph distance

[Johnson et al. (2017)]

Word Vectors

◆ Problem Formulation:

- Assume a finite vocabulary I , $|I| = n$
- Given a word x , predict nearby words y

◆ Simple Model:

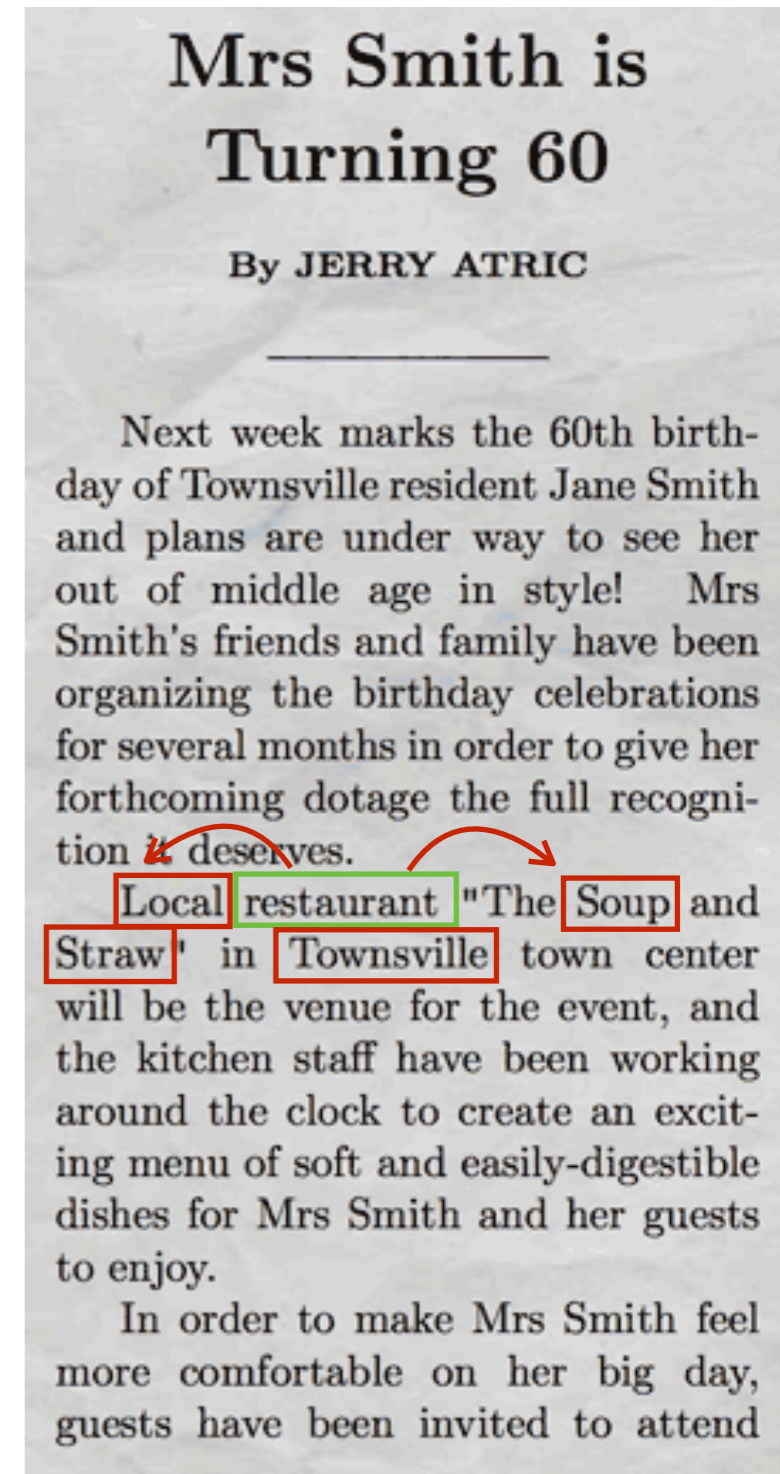
- Predict one word by categorical distribution $p(y|x)$
- Let $v(x) \in \mathbb{R}^d$ – word vector for x
- x is discrete and not structured — parameterize by a matrix V of all word vectors of size $n \times d$: $v(x) = V_{x,:}$
- Define categorical distribution:

$$p(y|x) = \frac{\exp(u(y)^\top v(x))}{\sum_{y'} \exp(u(y')^\top v(x))}$$

- Need a different embedding for context words: $u(y) = U_{y,:}$
- Learn via maximum likelihood classification:

$$\max_{U,V} \mathbb{E}_{t,t'} \left[\log p(y_{t'}|x_t) \right],$$

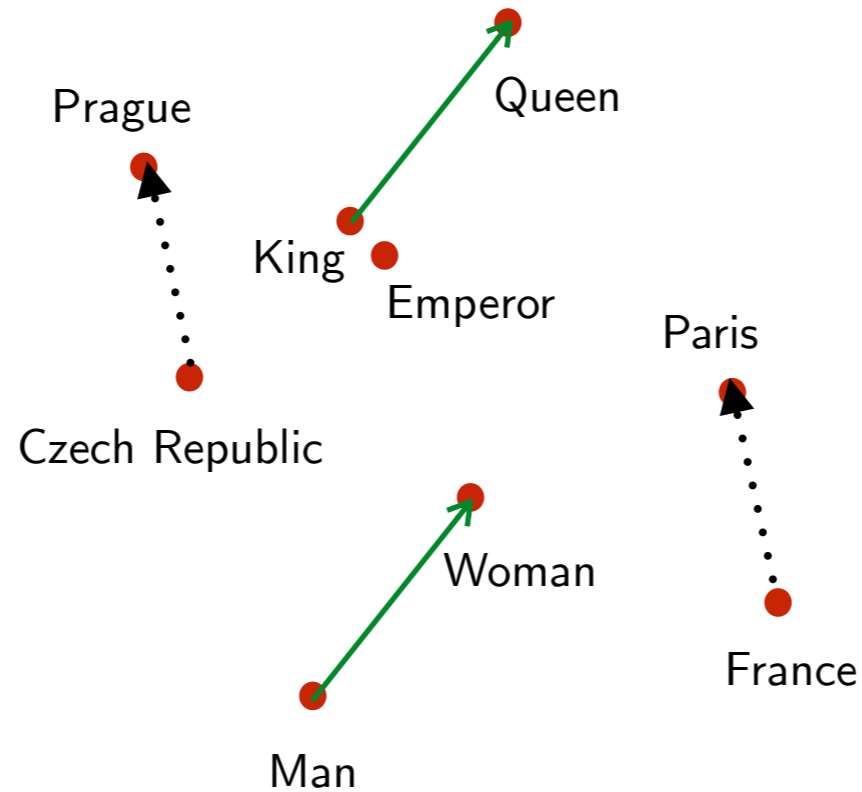
where t, t' – nearby positions in the text



Word Vectors



- ◆ Learned a **representation** of each word x as the embedding $v(x) = V_{x,:} \in \mathbb{R}^d$



- ◆ Direction of $v(x)$ appears to capture abstract relations:

- Semantic:

"King" - "Man" + "Woman" \approx "Queen"

"Prague" - "Czech Republic" + "France" \approx "Paris"

"Czech" + "currency" \approx "koruna"

- Syntactic:

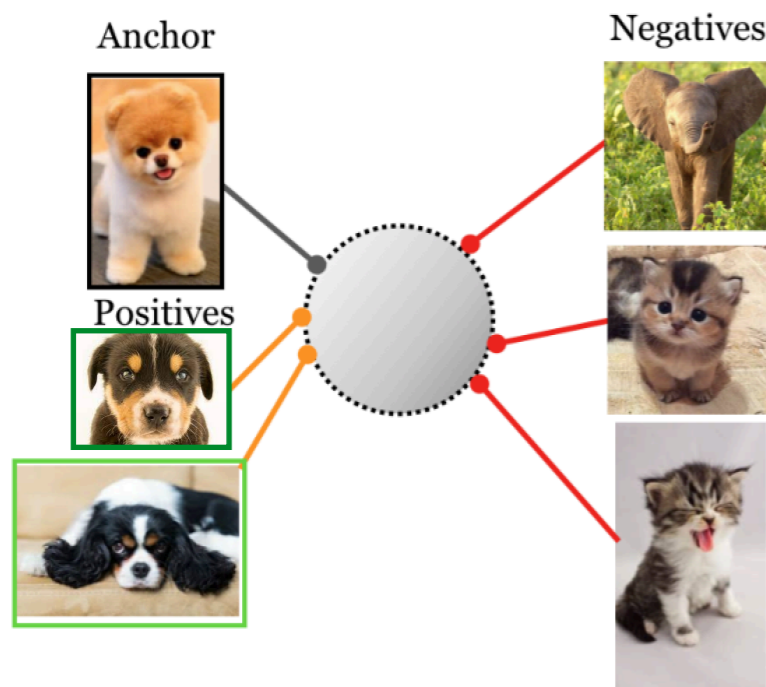
"quick" - "quickly" \approx "slow" - "slowly"

- ◆ Evaluated on a corpus of relation prediction tasks
- ◆ More complex tasks become easier when using such vector representations

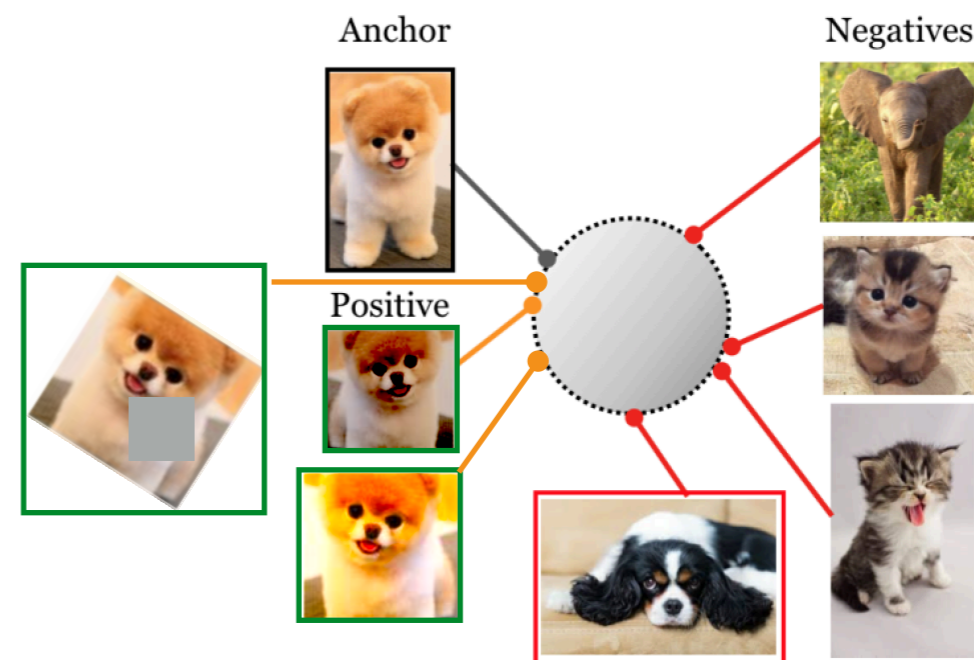
Deep Similarity/Metric Learning

Similarity Learning

- ◆ Goals:
 - learn the concept of **similarity** of two inputs
 - quantify this similarity
- ◆ Supervised learning:
 - Given examples of "similar" and "distinct" pairs learn the function $\text{sim}(x,y)$
 - no reference for the values of $\text{sim}(x,y)$
- ◆ Self-supervised learning:
 - Create "similar" pairs by identity-preserving transforms

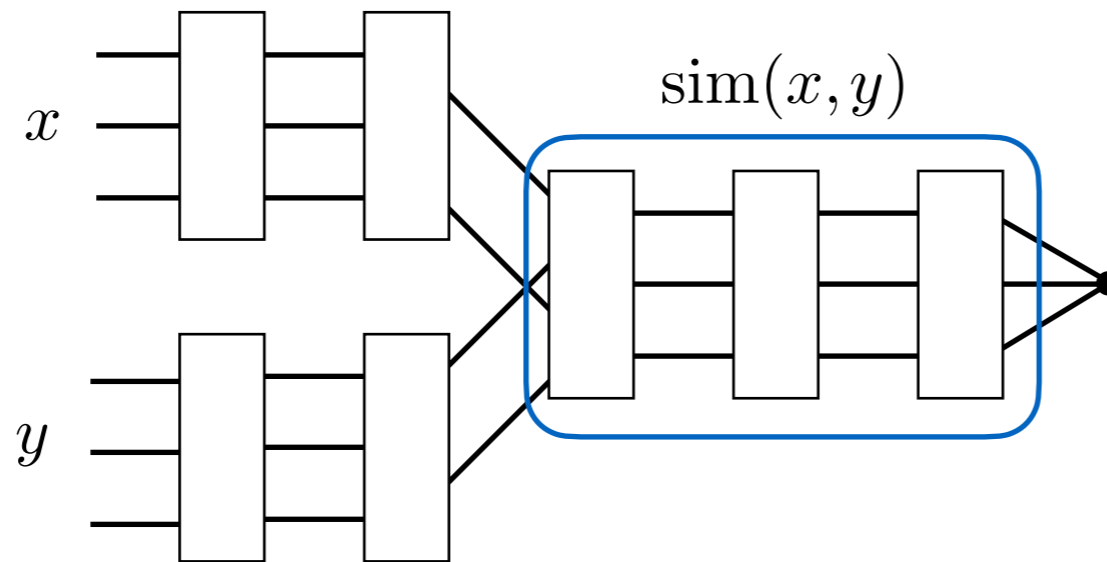


Supervised Contrastive



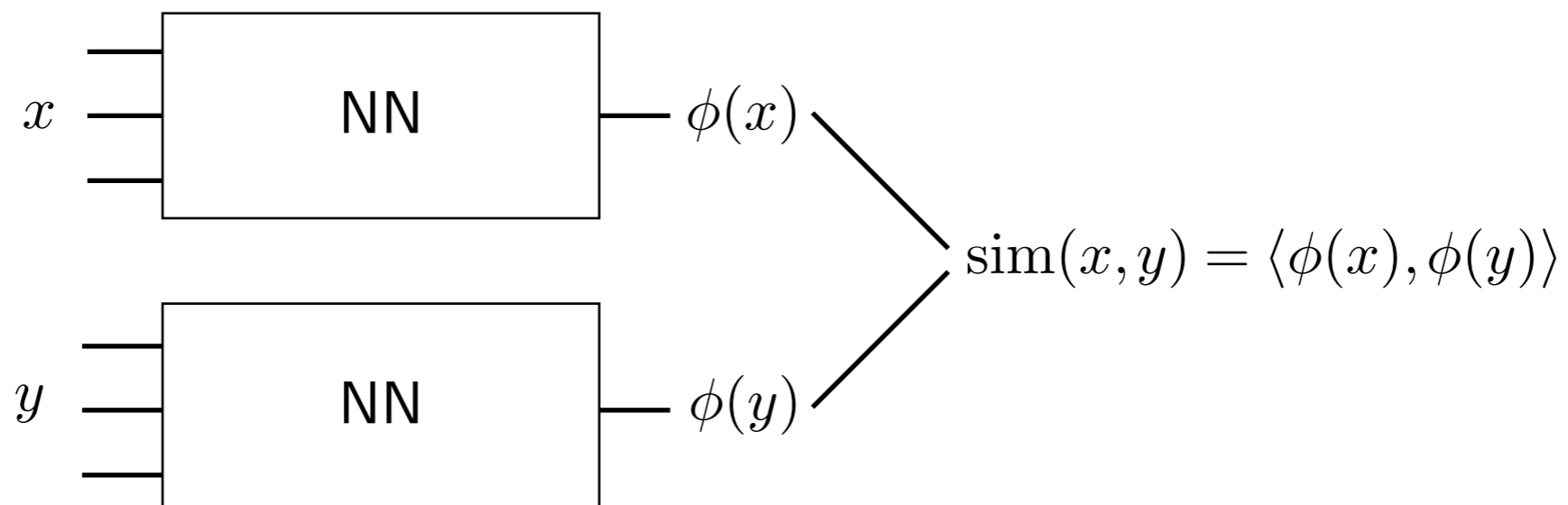
Self Supervised Contrastive

- ◆ Approach 1: generic network with two inputs



first layers extract generic features -- can be shared

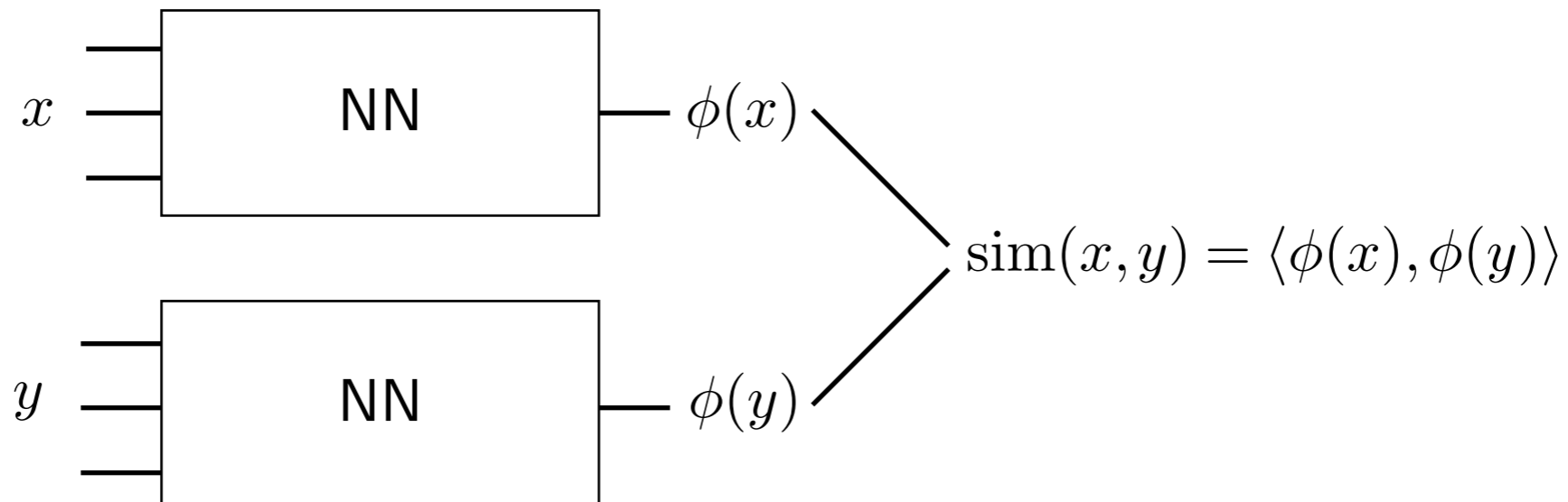
- ◆ Approach 2: network creates representations (embeddings / features)



- Inner product: $\langle \phi(x), \phi(y) \rangle$ can approximate any kernel $K(x, y)$

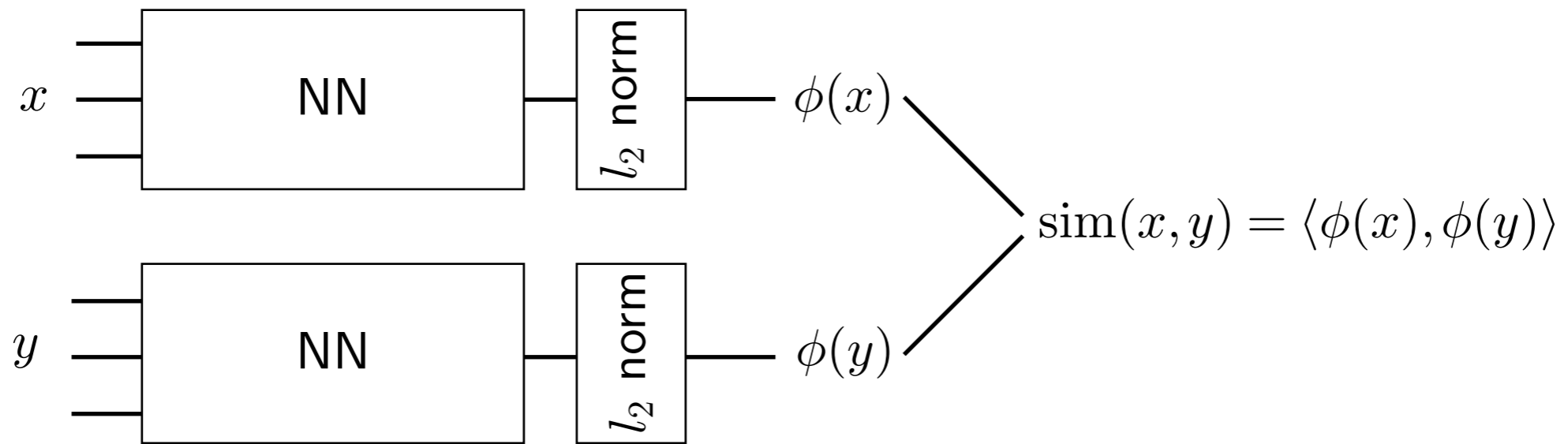
Similarity Function

- ◆ Network creates representations (embeddings / features)

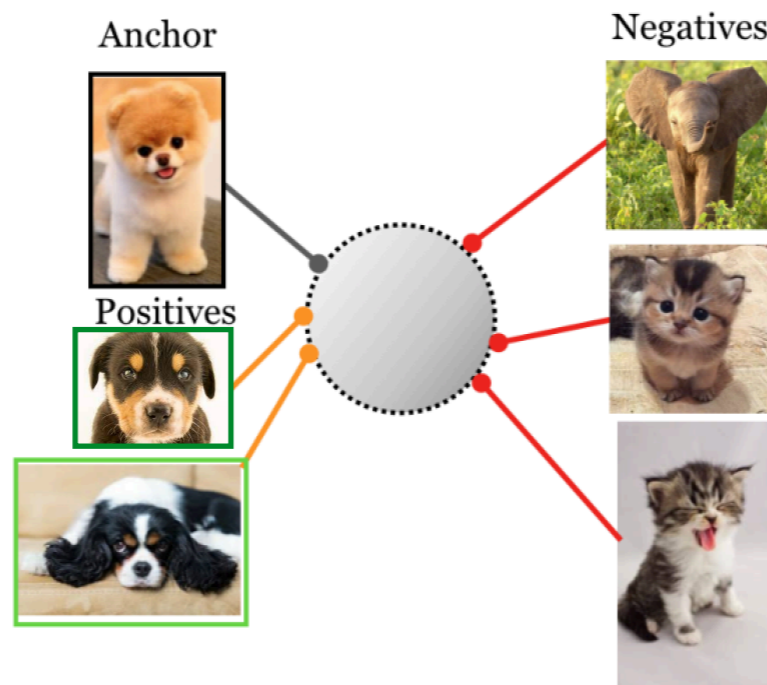


- Inner product: $\text{sim}(x, y) = \langle \phi(x), \phi(y) \rangle$
Retrieval: Maximum Inner Product Search (MIPS) is more difficult (no triangle inequality)
 - Euclidean: $\text{sim}(x, y) = -\|\phi(x) - \phi(y)\|^2$
nearest neighbor search (NNS): sub-linear approximate methods
 - Correlation: $\text{sim}(x, y) = \frac{\langle \phi(x), \phi(y) \rangle}{\|\phi(x)\| \|\phi(y)\|}$ correlation-NNS: sub-linear approximate methods
- ◆ All equivalent if $\|\phi(x)\| = 1$ for all x
 - ◆ There are known mappings to approximate $\langle u, v \rangle$ with $\|P(u) - Q(v)\|^2$ or $\frac{\langle P(u), Q(v) \rangle}{\|P(u)\| \|Q(v)\|}$

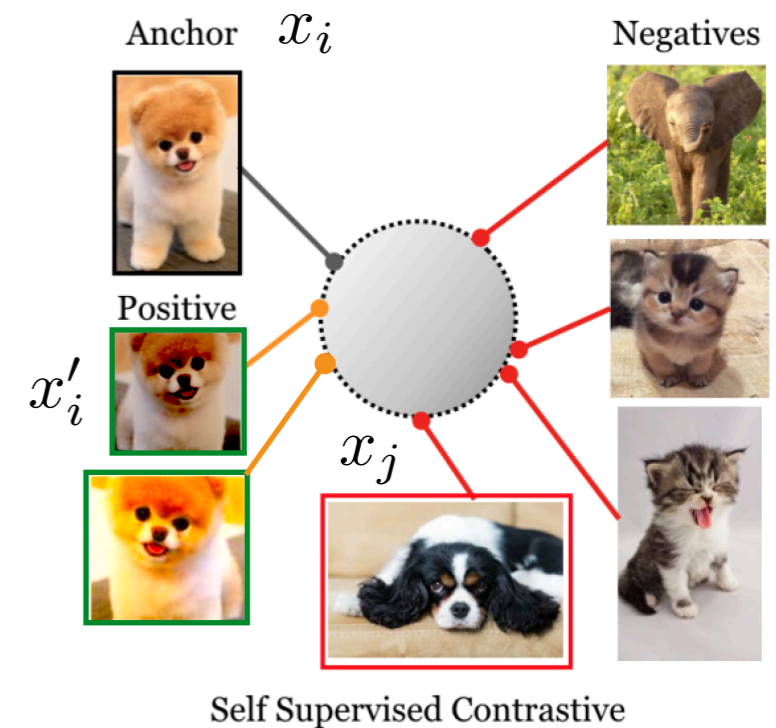
- ◆ Convenient common model: normalize representations (equivalent to using correlation for similarity)



Representations live on a hypersphere



- ◆ Training data: $x_1 \dots x_N$
 - *Anchor*: x_i
 - *Positive*: $x' = T(x_i)$ – random transform
- ◆ Model as classification:
 - As many classes as there are instances (data points)
 - Score of instance i : $s_i = \phi(x_i)^\top \phi(x')$
 - $p(y=i|x') = \frac{e^{s_i}}{\sum_j e^{s_j}}$
 - Learning formulation: likelihood of classifying correctly
 - Large sum in the denominator \rightarrow common solution is to restrict to a min-batch
- ◆ Properties:
 - Ensures instances can be discriminated
 - Enforces invariance to transformations



Dosovitskiy et al. (2014): Discriminative unsupervised feature learning with convolutional neural networks

Wu et al. (2018): Unsupervised Feature Learning via Non-Parametric Instance Discrimination]

Chen et al. (2020): A Simple Framework for **Contrastive Learning** of Visual Representations

Contrastive learning: "contrasting positive pairs against negative pairs"

Triplet Loss

- ◆ Training data: $x_1 \dots x_N$
 - Positive pairs: $x_i \sim x_j$ for $(i, j) \in \mathcal{P}$
 - Negative pairs: $x_i \not\sim x_j$ for $(i, j) \in \mathcal{N}$
 - Example: known class label \rightarrow similar if same class

- ◆ Desired separation property:

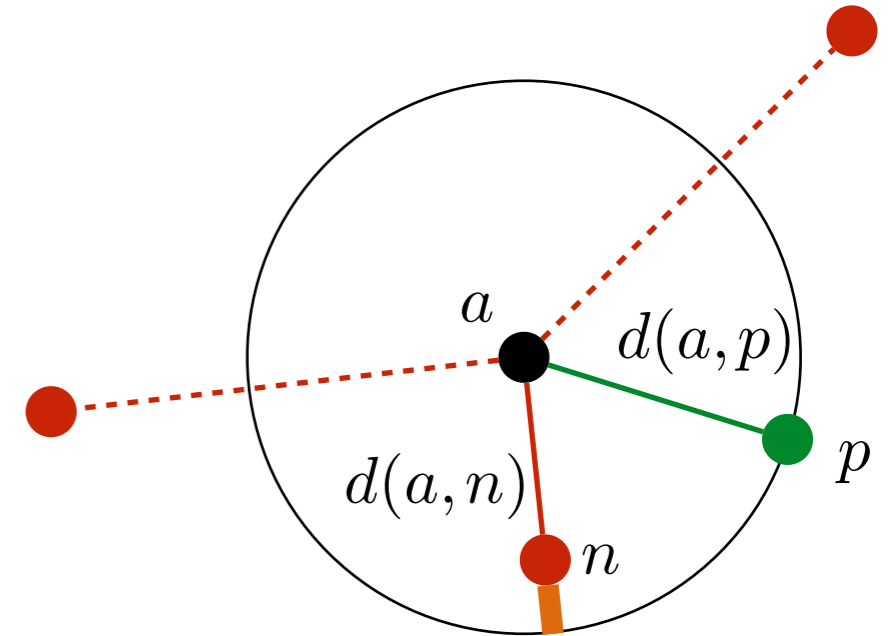
- a – anchor (any data point)
- p – positive sample for a
- n – negative
- let $D_p = d(a, p)$ and $D_n = d(a, n)$
- Want:

$$D_p < D_n \quad \forall p, n$$

$$D_p - D_n < 0$$

- Hinge loss 1: $l = \sum_{n \in \mathcal{N}(a)} \max(0, D_p - D_n)$
- Hinge loss 2: $l' = \max(0, \max_{n \in \mathcal{N}(a)} (D_p - D_n)) = \max_{x \in \mathcal{N}(a) \cup \{p\}} (D_p - D_x)$

-- **hard negative mining**



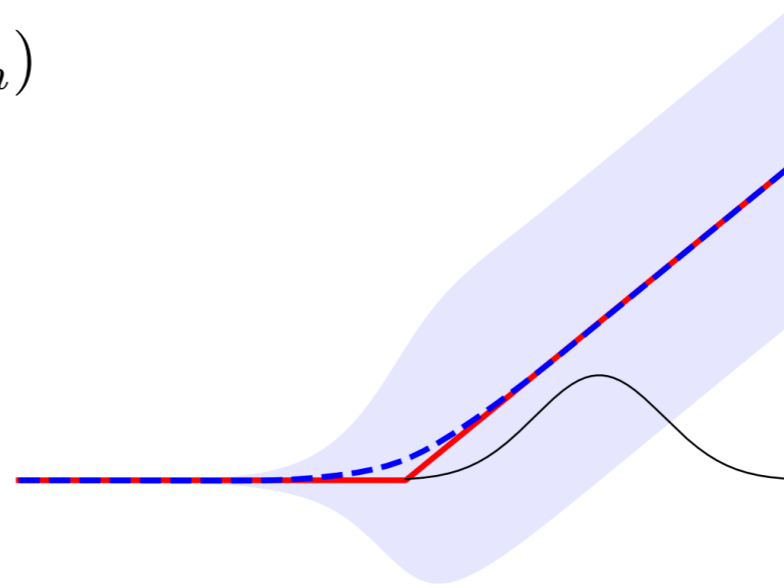
Smooth Triplet Loss = Log Likelihood Classification

- ◆ Hinge loss variant 1: $l = \sum_{n \in \mathcal{N}(a)} \max(0, D_p - D_n)$

Smooth maximum (SoftPlus):

$$Z \sim \text{Logistic}(0, \tau)$$

$$\mathbb{E}[\max(0, x + Z)] = \tau \log(1 + e^{x/\tau})$$



Let \tilde{D}_p, \tilde{D}_n be noisy descriptors with additive noise $\text{Gumbel}(0, \tau)$

$$\mathbb{E}[\max(0, \tilde{D}_p - \tilde{D}_n)] = \mathbb{E}[\max(0, D_p - D_n + Z)] = \tau \log(1 + e^{(D_p - D_n)/\tau})$$

Binary classification of positive vs negative pairs, e.g. [Sohn 2016]

- ◆ Hinge loss variant 2: $l' = \max_{x \in \mathcal{N}(a) \cup \{p\}} (D_p - D_x)$

- Let \tilde{D}_p, \tilde{D}_x be noisy descriptors with additive noises $\text{Gumbel}(0, \tau)$

- $\mathbb{E}[l] = -\tau \log \frac{e^{-D_p/\tau}}{\sum_x e^{-D_x/\tau}}$

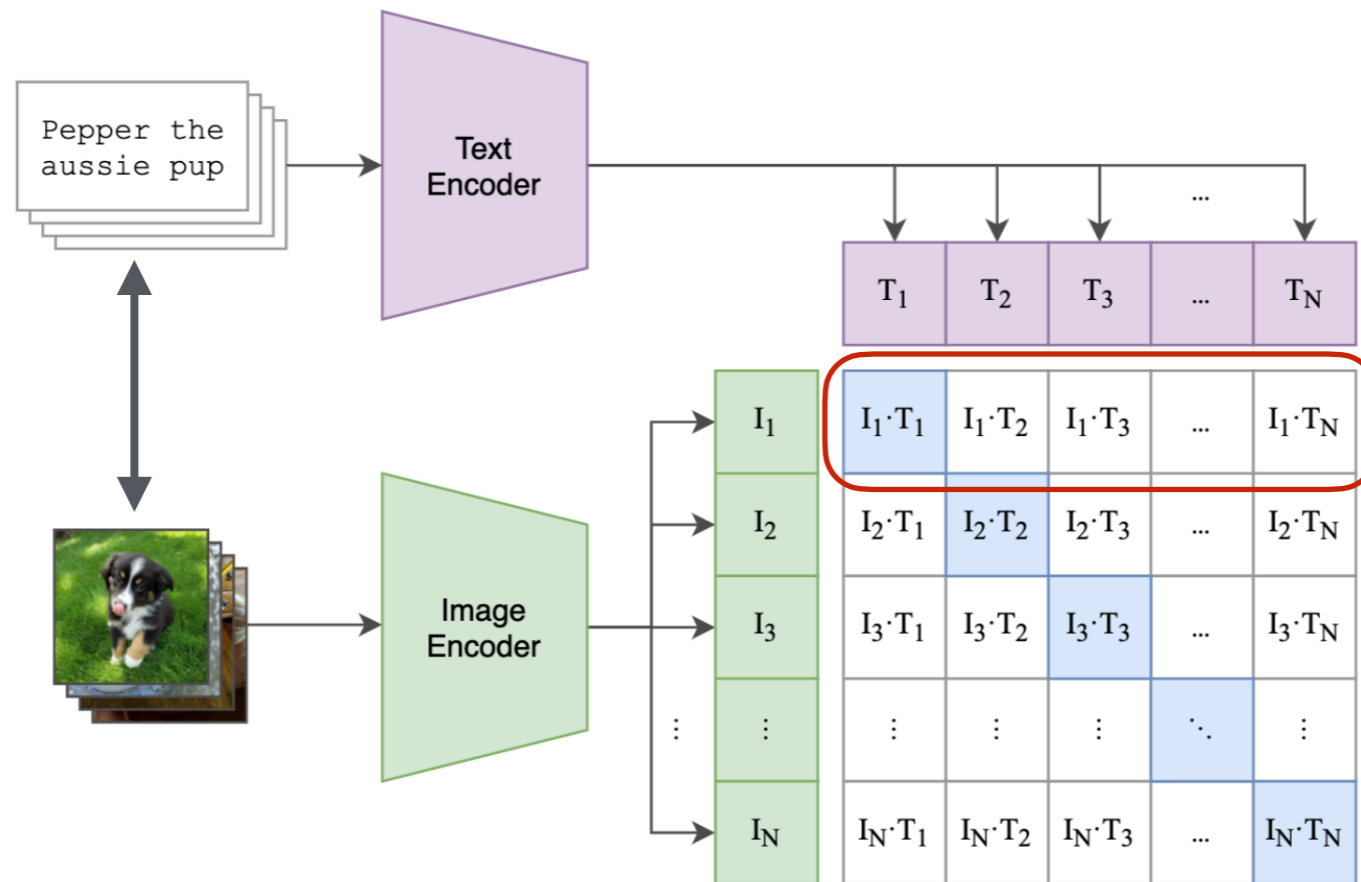
- Log softmax loss back again

- For the inner product similarity: $\mathbb{E}[l] = -\tau \log \frac{e^{\phi(a)^\top \phi(p)/\tau}}{\sum_x e^{\phi(a)^\top \phi(x)/\tau}}$

Classification loss of predicting the positive out of all candidates

Cross-Modality Representations

CLIP: Connecting Text and Images



- ◆ In each batch of image-text pairs with embeddings I_i, T_j :
 - Predict which text corresponds to image i , model: $p_1(j|i) = \frac{e^{\langle I_i, T_j \rangle / \tau}}{\sum_{j'} e^{\langle I_i, T_{j'} \rangle / \tau}}$
 - Predict which image corresponds to text j , model: $p_2(i|j) = \frac{e^{\langle I_i, T_j \rangle / \tau}}{\sum_{i'} e^{\langle I_{i'}, T_j \rangle / \tau}}$
 - Learning: symmetric cross-entropy loss:

$$-\sum_i \left(\log p_1(i|i) + \log p_2(i|i) \right)$$