

deep image retrieval and manifold learning

Yannis Avrithis

Inria Rennes-Bretagne Atlantique

Paris, March 2018



outline

feature pooling
manifold learning
fine-tuning
ranking on manifolds
fast spectral ranking
mining on manifolds

feature pooling

image ranking by CNN features

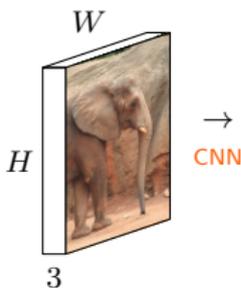
[Krizhevsky et al. 2012]



- 3-channel RGB input, 224×224
- AlexNet pre-trained on ImageNet for classification
- last fully connected layer (fc_6): global descriptor of dimension $k = 4096$

image ranking by CNN features

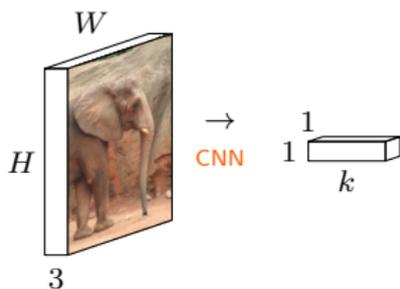
[Krizhevsky et al. 2012]



- 3-channel RGB input, 224×224
- AlexNet pre-trained on ImageNet for classification
- last fully connected layer (fc_6): global descriptor of dimension $k = 4096$

image ranking by CNN features

[Krizhevsky et al. 2012]



- 3-channel RGB input, 224×224
- AlexNet pre-trained on ImageNet for classification
- last fully connected layer (fc_6): **global descriptor** of dimension $k = 4096$

image ranking by CNN features



- query images
- nearest neighbors in ImageNet according to Euclidean distance

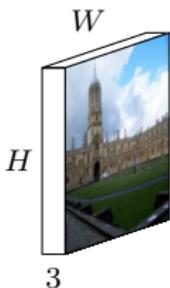
image ranking by CNN features



- query images
- nearest neighbors in ImageNet according to Euclidean distance

neural codes for image retrieval

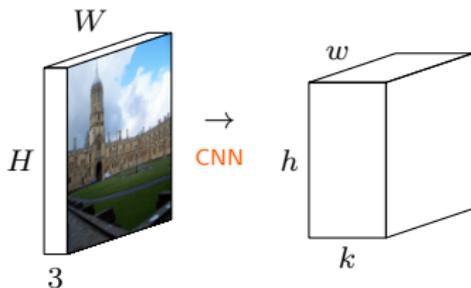
[Babenko et al. 2014]



- 3-channel RGB input, 224×224
- AlexNet last pooling layer, global descriptor of dimension $w \times h \times k = 6 \times 6 \times 256 = 9216$
- alternatively: fully connected layers fc_6, fc_7 , global descriptors of dimension $k' = 4096$ (best is fc_6)
- in each case: PCA-whitening, ℓ_2 normalization

neural codes for image retrieval

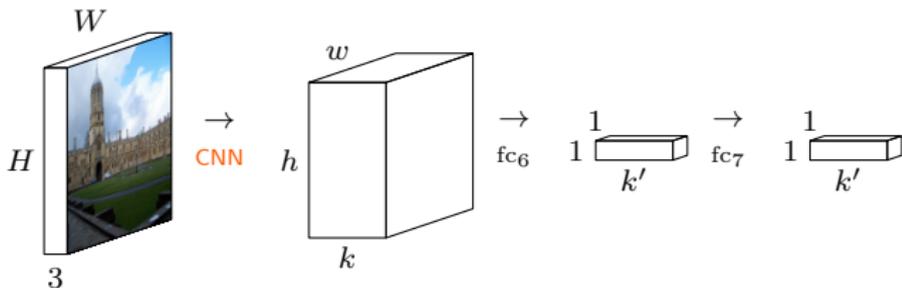
[Babenko et al. 2014]



- 3-channel RGB input, 224×224
- AlexNet last pooling layer, **global descriptor** of dimension $w \times h \times k = 6 \times 6 \times 256 = 9216$
- alternatively: fully connected layers fc_6, fc_7 , **global descriptors** of dimension $k' = 4096$ (*best is fc_6*)
- in each case: PCA-whitening, ℓ_2 normalization

neural codes for image retrieval

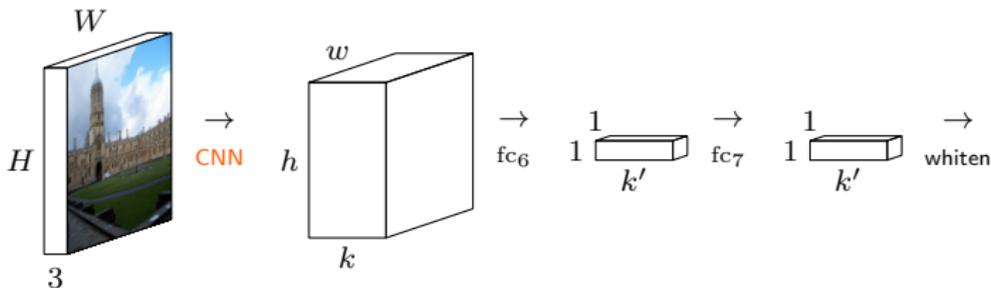
[Babenko et al. 2014]



- 3-channel RGB input, 224×224
- AlexNet last pooling layer, **global descriptor** of dimension $w \times h \times k = 6 \times 6 \times 256 = 9216$
- alternatively: fully connected layers fc_6, fc_7 , **global descriptors** of dimension $k' = 4096$ (**best is fc_6**)
- in each case: PCA-whitening, ℓ_2 normalization

neural codes for image retrieval

[Babenko et al. 2014]



- 3-channel RGB input, 224×224
- AlexNet last pooling layer, **global descriptor** of dimension $w \times h \times k = 6 \times 6 \times 256 = 9216$
- alternatively: fully connected layers fc_6, fc_7 , **global descriptors** of dimension $k' = 4096$ (**best is fc_6**)
- in each case: PCA-whitening, ℓ_2 normalization

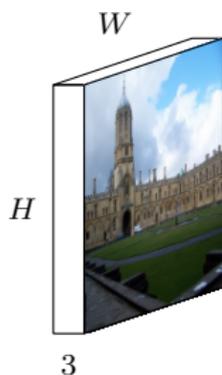
neural codes for image retrieval



- **fine-tuning** by softmax on 672 classes of 200k landmark photos
- outperforms VLAD and Fisher vectors on standard retrieval benchmarks, but still inferior to SIFT local descriptors

regional CNN features

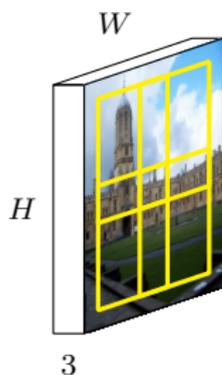
[Razavian et al. 2015]



- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, warped into $w \times h = 227 \times 227$
- each region yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial max-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

regional CNN features

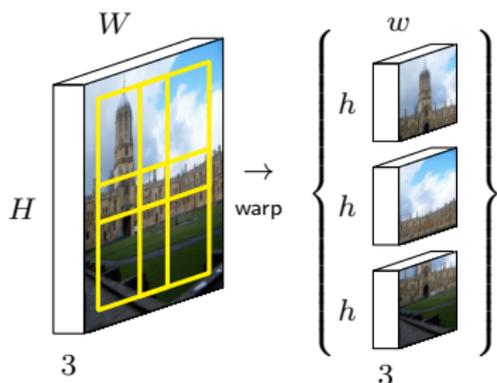
[Razavian et al. 2015]



- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, warped into $w \times h = 227 \times 227$
- each region yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial max-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

regional CNN features

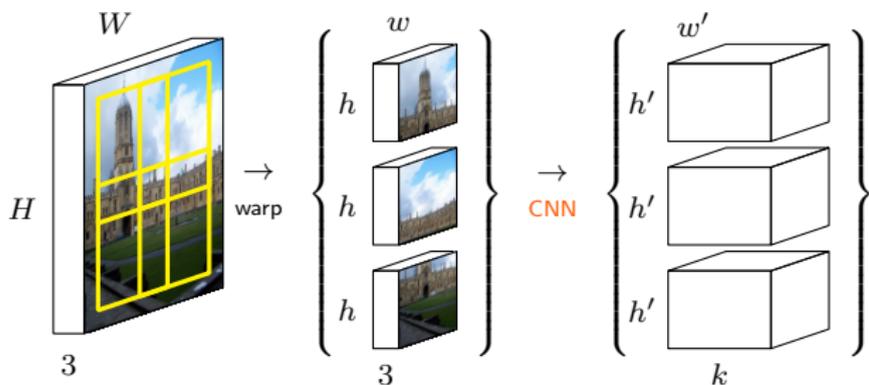
[Razavian et al. 2015]



- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, **warped** into $w \times h = 227 \times 227$
- **each region** yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial **max**-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

regional CNN features

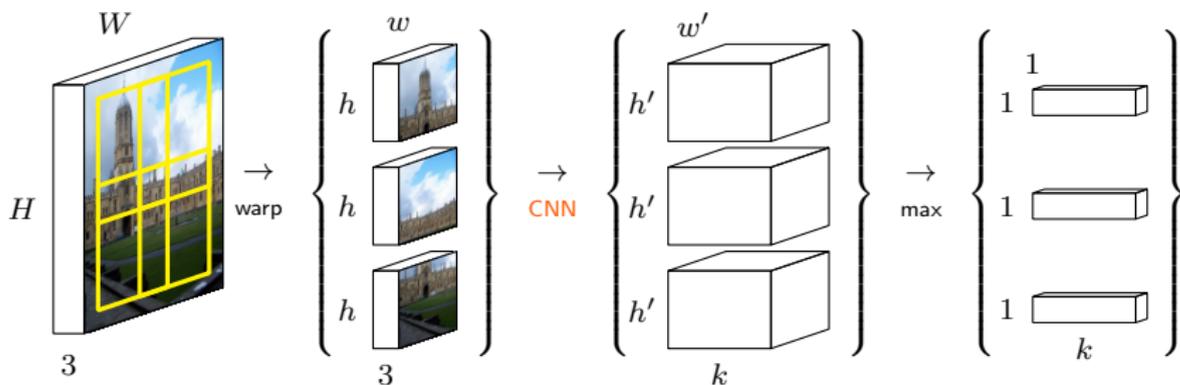
[Razavian et al. 2015]



- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, **warped** into $w \times h = 227 \times 227$
- **each region** yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial **max**-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

regional CNN features

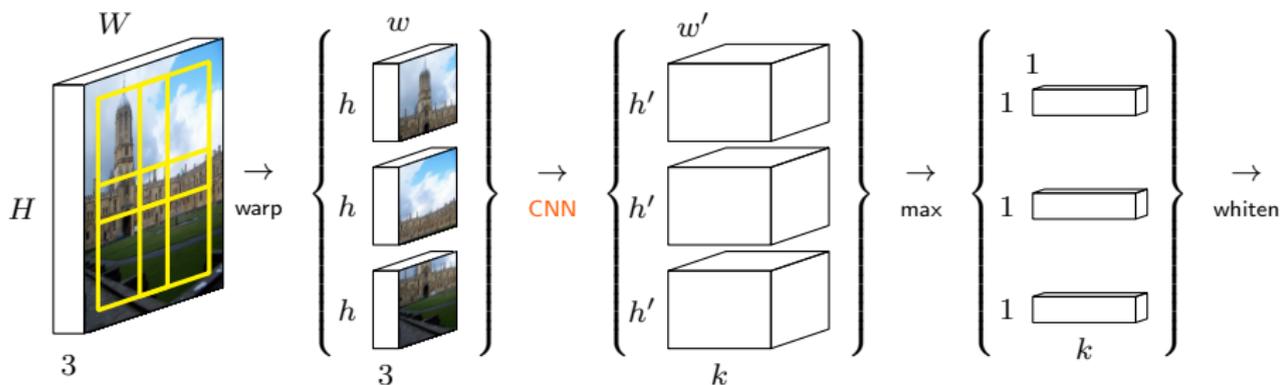
[Razavian et al. 2015]



- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, **warped** into $w \times h = 227 \times 227$
- **each region** yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial **max**-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

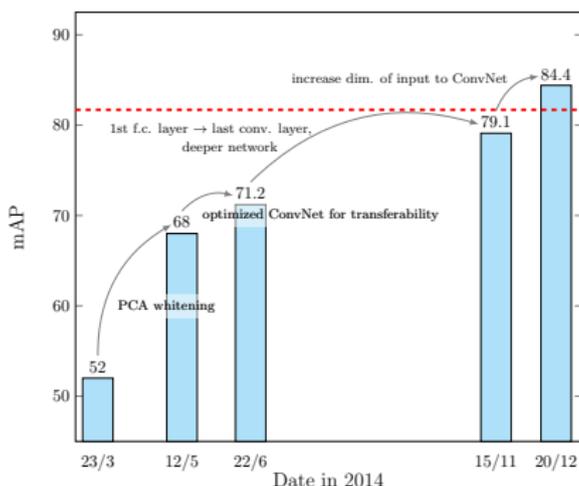
regional CNN features

[Razavian et al. 2015]



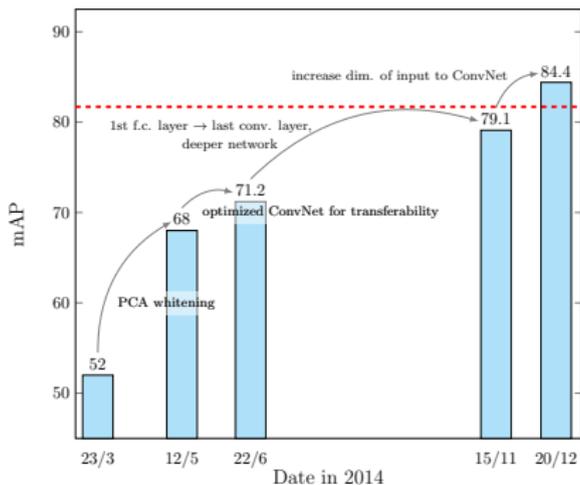
- 3-channel RGB input, largest square region extracted
- fixed multiscale overlapping regions, **warped** into $w \times h = 227 \times 227$
- **each region** yields a $w' \times h' \times k = 36 \times 36 \times 256$ dimensional feature at the last convolutional layer of AlexNet
- global spatial **max**-pooling
- ℓ_2 -normalization, PCA-whitening of each descriptor

regional CNN features



- CNN visual representation jumps by more than 30% mAP to outperform standard SIFT pipeline in a few months
- however, this is based on **multiple** regional descriptors per image and **exhaustive** pairwise matching of all descriptors of query and all dataset images, which is not practical

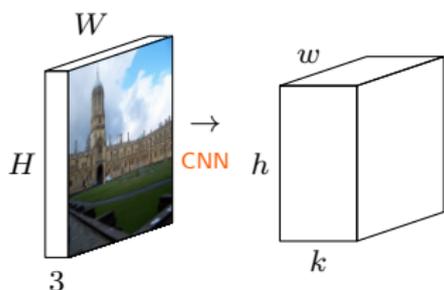
regional CNN features



- CNN visual representation jumps by more than 30% mAP to outperform standard SIFT pipeline in a few months
- however, this is based on **multiple** regional descriptors per image and **exhaustive** pairwise matching of all descriptors of query and all dataset images, which is not practical

regional max-pooling (R-MAC)

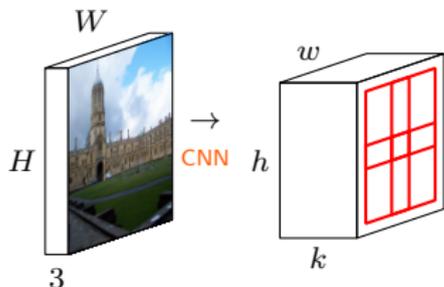
[Tolias et al. 2016]



- VGG-16 last convolutional layer, $k = 512$
- fixed multiscale overlapping regions, spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- sum-pooling over all descriptors, ℓ_2 -normalization

regional max-pooling (R-MAC)

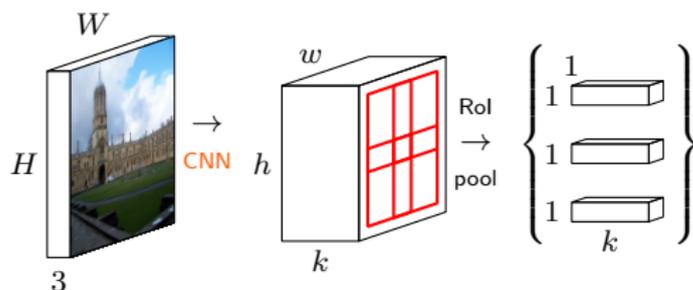
[Tolias et al. 2016]



- VGG-16 last convolutional layer, $k = 512$
- fixed multiscale overlapping regions, spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- sum-pooling over all descriptors, ℓ_2 -normalization

regional max-pooling (R-MAC)

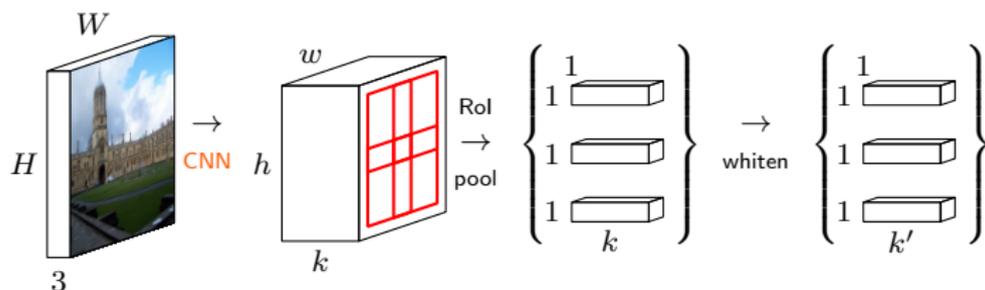
[Tolias et al. 2016]



- VGG-16 last convolutional layer, $k = 512$
- fixed multiscale overlapping regions, spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- sum-pooling over all descriptors, ℓ_2 -normalization

regional max-pooling (R-MAC)

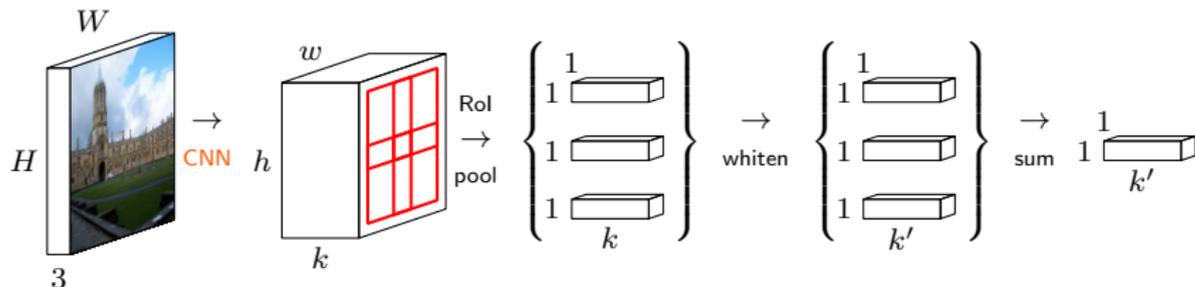
[Tolias et al. 2016]



- VGG-16 last convolutional layer, $k = 512$
- fixed multiscale overlapping regions, spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- sum-pooling over all descriptors, ℓ_2 -normalization

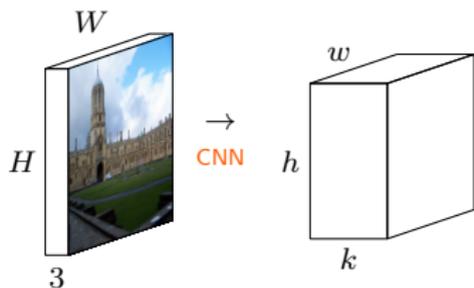
regional max-pooling (R-MAC)

[Tolias et al. 2016]



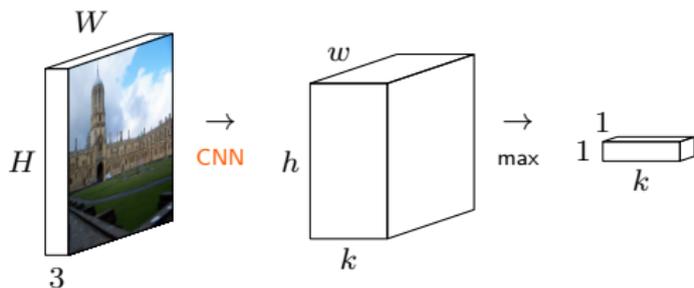
- VGG-16 last convolutional layer, $k = 512$
- fixed multiscale overlapping regions, spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- sum-pooling over all descriptors, ℓ_2 -normalization

global max-pooling (MAC)



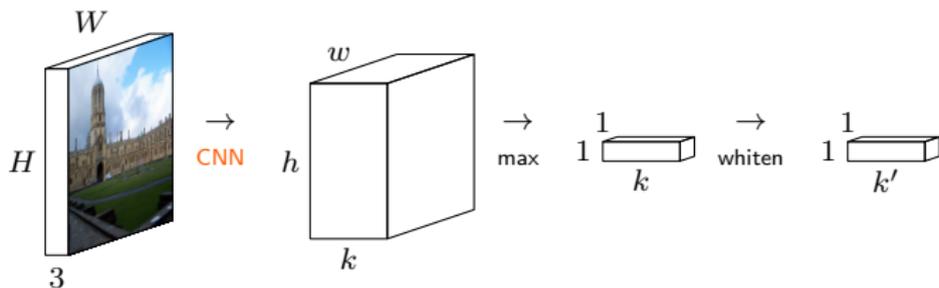
- VGG-16 last convolutional layer, $k = 512$
- global spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- MAC: maximum activation of convolutions

global max-pooling (MAC)



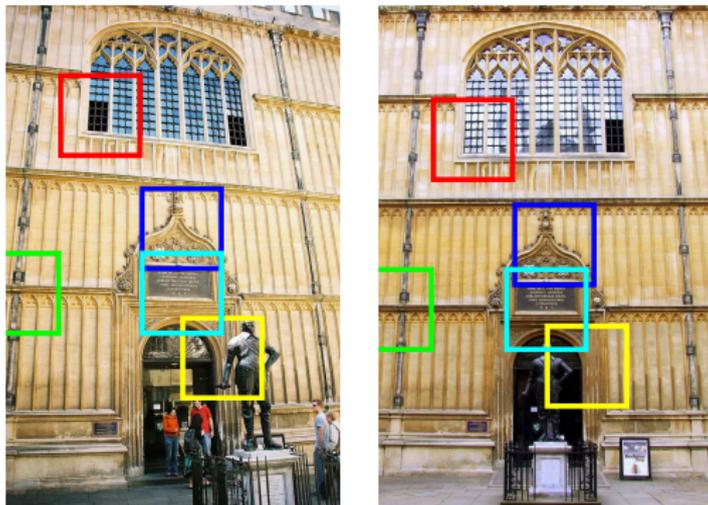
- VGG-16 last convolutional layer, $k = 512$
- global spatial max-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- MAC: maximum activation of convolutions

global max-pooling (MAC)



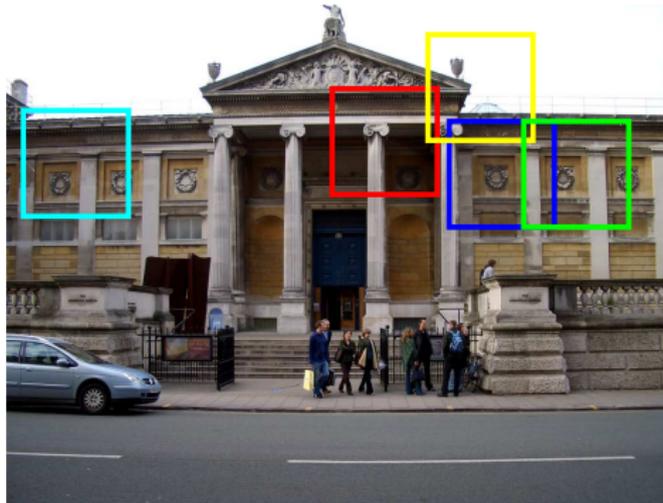
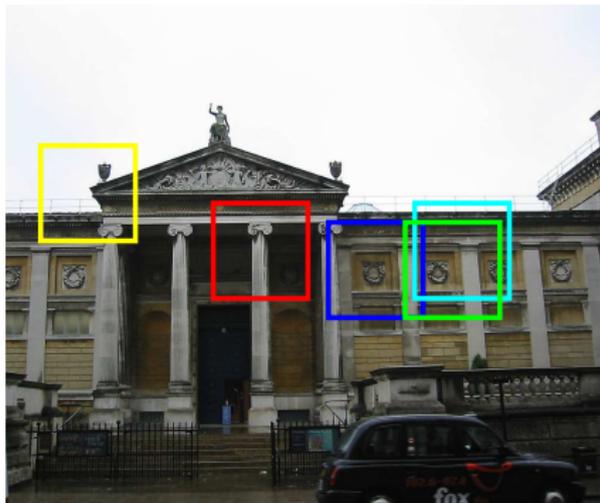
- VGG-16 last convolutional layer, $k = 512$
- global spatial **max**-pooling
- ℓ_2 -normalization, PCA-whitening, ℓ_2 -normalization
- **MAC**: maximum activation of convolutions

global max-pooling: matching



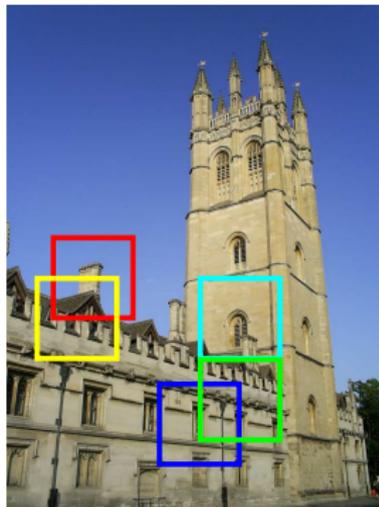
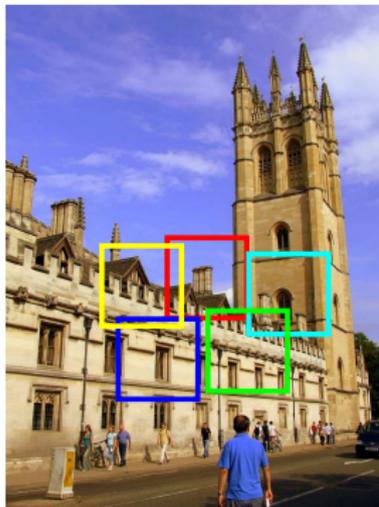
- receptive fields of 5 components of MAC vectors that contribute most to image similarity

global max-pooling: matching



- receptive fields of 5 components of MAC vectors that contribute most to image similarity

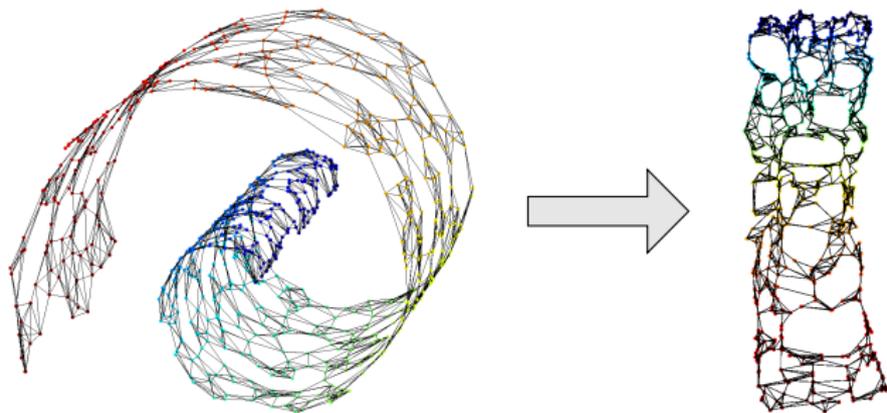
global max-pooling: matching



- receptive fields of 5 components of MAC vectors that contribute most to image similarity

manifold learning

manifold learning



- e.g. Isomap: apply PCA to the geodesic (graph) distance matrix
- e.g. kernel PCA: apply PCA to the Gram matrix of a nonlinear kernel
- other **topology-preserving** methods are only focusing on distances to nearest neighbors
- many classic methods use eigenvalue decomposition and most do **not** learn and **explicit mapping** from the input to the embedding space

siamese architecture

[Chopra et al. 2005]

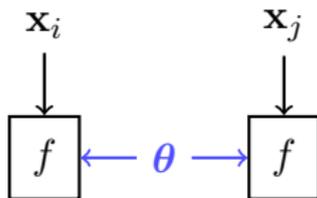
\mathbf{x}_i

\mathbf{x}_j

- an input sample is a pair $(\mathbf{x}_i, \mathbf{x}_j)$
- both $\mathbf{x}_i, \mathbf{x}_j$ go through the same function f with shared parameters θ
- loss ℓ_{ij} is measured on output pair $(\mathbf{y}_i, \mathbf{y}_j)$ and target t_{ij}

siamese architecture

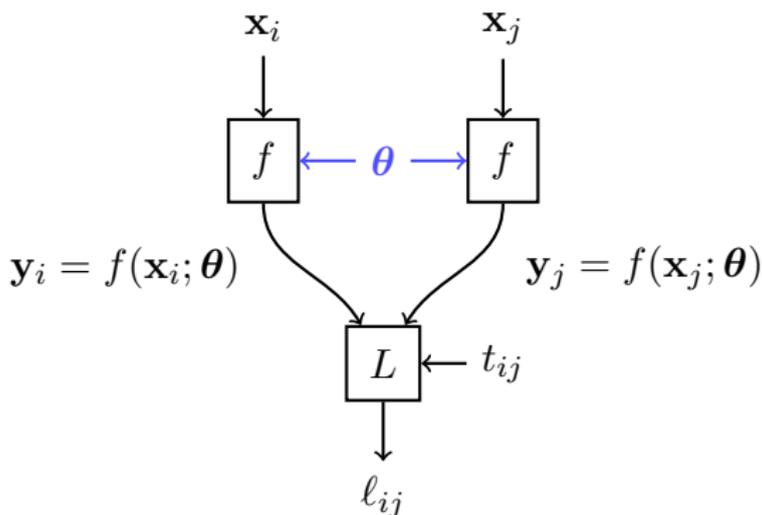
[Chopra et al. 2005]



- an input sample is a pair (x_i, x_j)
- both x_i, x_j go through the same function f with shared parameters θ
- loss ℓ_{ij} is measured on output pair (y_i, y_j) and target t_{ij}

siamese architecture

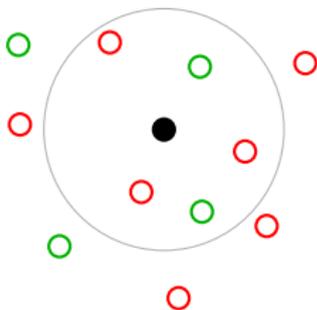
[Chopra et al. 2005]



- an input sample is a **pair** (x_i, x_j)
- both x_i, x_j go through the **same** function f with **shared** parameters θ
- loss l_{ij} is measured on output pair (y_i, y_j) and target t_{ij}

contrastive loss

[Hadsell et al. 2006]



- input samples \mathbf{x}_i , output vectors $\mathbf{y}_i = f(\mathbf{x}_i; \boldsymbol{\theta})$
- target variables $t_{ij} = \mathbb{1}[\text{sim}(\mathbf{x}_i, \mathbf{x}_j)]$
- **contrastive loss** is a function of distance $\|\mathbf{y}_i - \mathbf{y}_j\|$ only

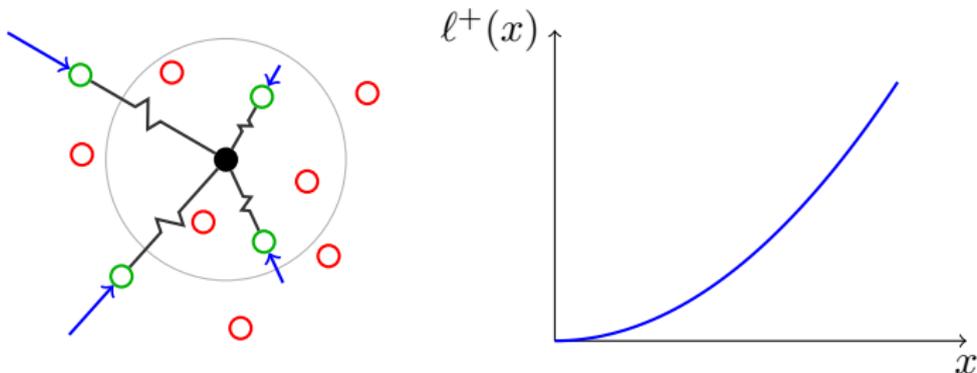
$$\ell_{ij} = L((\mathbf{y}_i, \mathbf{y}_j), t_{ij}) = \ell(\|\mathbf{y}_i - \mathbf{y}_j\|, t_{ij})$$

- **similar** samples are **attracted**

$$\ell(x, t) = t\ell^+(x) + (1 - t)\ell^-(x) = tx^2 + (1 - t)[m - x]_+^2$$

contrastive loss

[Hadsell et al. 2006]



- input samples \mathbf{x}_i , output vectors $\mathbf{y}_i = f(\mathbf{x}_i; \theta)$
- target variables $t_{ij} = \mathbb{1}[\text{sim}(\mathbf{x}_i, \mathbf{x}_j)]$
- **contrastive loss** is a function of distance $\|\mathbf{y}_i - \mathbf{y}_j\|$ only

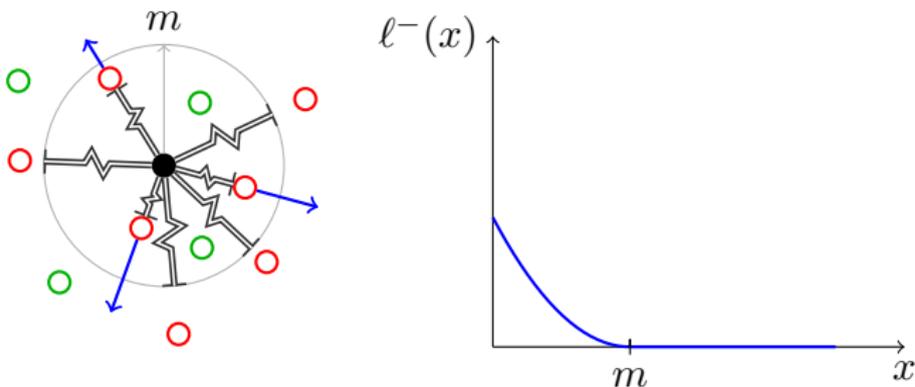
$$\ell_{ij} = L((\mathbf{y}_i, \mathbf{y}_j), t_{ij}) = \ell(\|\mathbf{y}_i - \mathbf{y}_j\|, t_{ij})$$

- **similar** samples are **attracted**

$$\ell(x, t) = \boxed{t\ell^+(x)} + (1 - t)\ell^-(x) = \boxed{tx^2} + (1 - t)[m - x]_+^2$$

contrastive loss

[Hadsell et al. 2006]



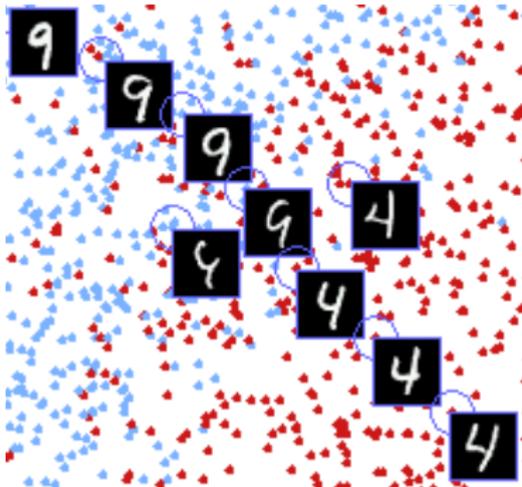
- input samples \mathbf{x}_i , output vectors $\mathbf{y}_i = f(\mathbf{x}_i; \theta)$
- target variables $t_{ij} = \mathbb{1}[\text{sim}(\mathbf{x}_i, \mathbf{x}_j)]$
- **contrastive loss** is a function of distance $\|\mathbf{y}_i - \mathbf{y}_j\|$ only

$$\ell_{ij} = L((\mathbf{y}_i, \mathbf{y}_j), t_{ij}) = \ell(\|\mathbf{y}_i - \mathbf{y}_j\|, t_{ij})$$

- **dissimilar** samples are **repelled** if closer than **margin** m

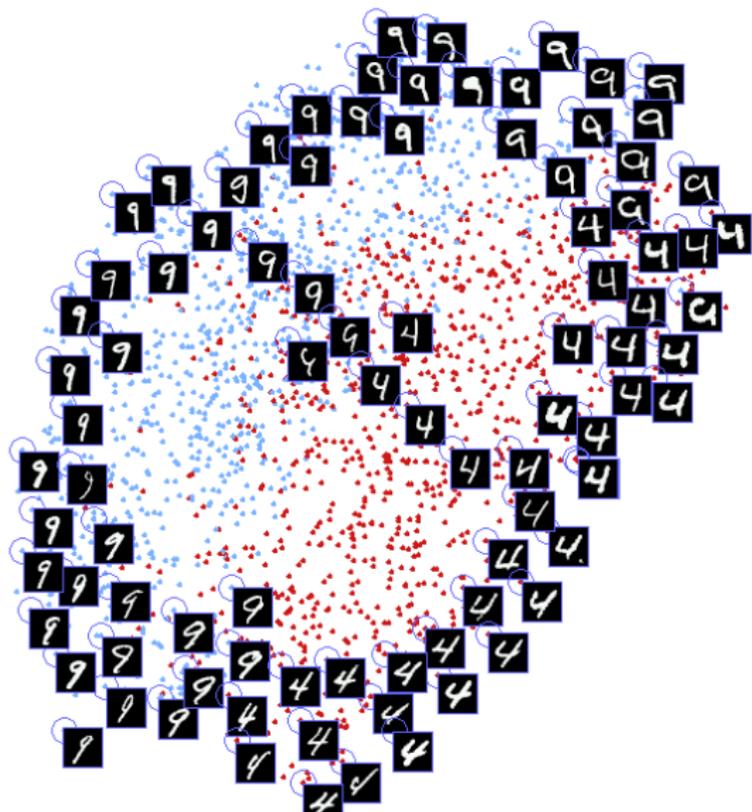
$$\ell(x, t) = t\ell^+(x) + (1 - t)\ell^-(x) = tx^2 + (1 - t)[m - x]_+^2$$

manifold learning: MNIST

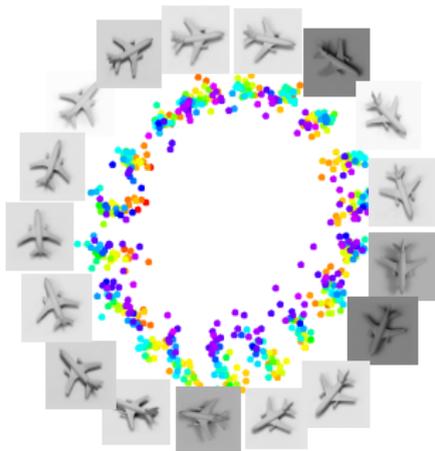


- 3k samples of each of digits 4, 9
- each sample similar to its 5 Euclidean nearest neighbors, and dissimilar to all other points
- 30k similar pairs, 18M dissimilar pairs

manifold learning: MNIST

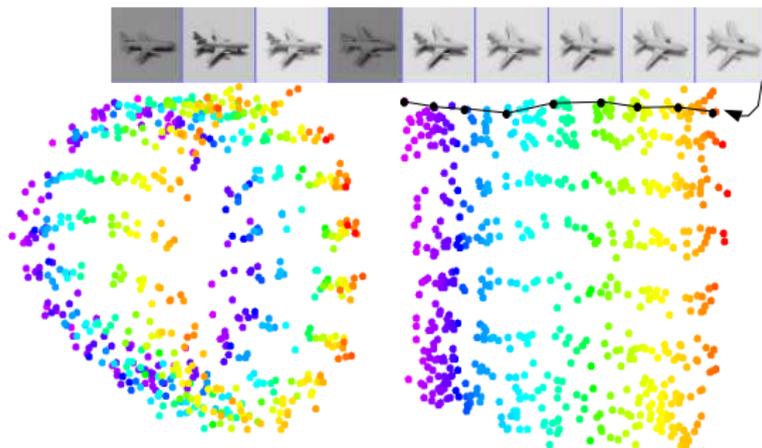


manifold learning: NORB



- 972 images of airplane class: 18 azimuths (every 20°), 9 elevations (in $[30^\circ, 70^\circ]$, every 5°), 6 lighting conditions
- samples similar if taken from contiguous azimuth or elevation, regardless of lighting
- 11k similar pairs, 206M dissimilar pairs
- cylinder in 3d: **azimuth on circumference**, elevation on height

manifold learning: NORB



- 972 images of airplane class: 18 azimuths (every 20°), 9 elevations (in $[30^\circ, 70^\circ]$, every 5°), 6 lighting conditions
- samples similar if taken from contiguous azimuth or elevation, regardless of lighting
- 11k similar pairs, 206M dissimilar pairs
- cylinder in 3d: azimuth on circumference, elevation on height

triplet architecture

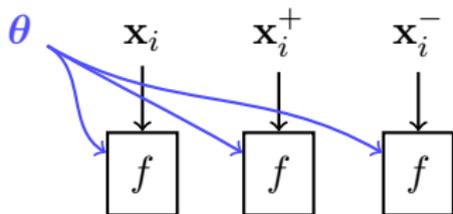
[Wang et al. 2014]

$$\mathbf{x}_i \quad \mathbf{x}_i^+ \quad \mathbf{x}_i^-$$

- an input sample is a **triplet** $(\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-)$
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through the **same** function f with **shared** parameters θ
- loss ℓ_i measured on output triplet $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

triplet architecture

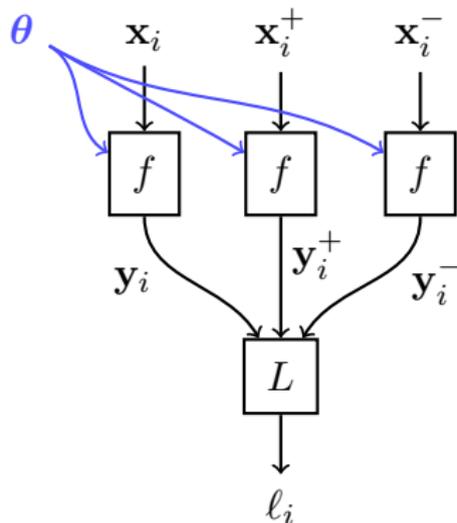
[Wang et al. 2014]



- an input sample is a **triplet** $(\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-)$
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through the **same** function f with **shared** parameters θ
- loss ℓ_i measured on output triplet $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

triplet architecture

[Wang et al. 2014]



- an input sample is a **triplet** (x_i, x_i^+, x_i^-)
- x_i, x_i^+, x_i^- go through the **same** function f with **shared** parameters θ
- loss l_i measured on output triplet (y_i, y_i^+, y_i^-)

triplet loss

- input “anchor” \mathbf{x}_i , output vector $\mathbf{y}_i = f(\mathbf{x}_i; \boldsymbol{\theta})$
- positive $\mathbf{y}_i^+ = f(\mathbf{x}_i^+; \boldsymbol{\theta})$, negative $\mathbf{y}_i^- = f(\mathbf{x}_i^-; \boldsymbol{\theta})$
- triplet loss is a function of distances $\|\mathbf{y}_i - \mathbf{y}_i^+\|$, $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ only

$$\ell_i = L(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-) = \ell(\|\mathbf{y}_i - \mathbf{y}_i^+\|, \|\mathbf{y}_i - \mathbf{y}_i^-\|)$$

$$\ell(x^+, x^-) = [m + (x^+)^2 - (x^-)^2]_+$$

so distance $\|\mathbf{y}_i - \mathbf{y}_i^+\|$ should be less than $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ by margin m

- by taking two pairs $(\mathbf{x}_i, \mathbf{x}_i^+)$ and $(\mathbf{x}_i, \mathbf{x}_i^-)$ at a time with targets 1, 0 respectively, the contrastive loss can be written similarly

$$\ell(x^+, x^-) = (x^+)^2 + [m - x^-]^2_+$$

so distance $\|\mathbf{y}_i - \mathbf{y}_i^+\|$ should small and $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ larger than m

triplet loss

- input “anchor” \mathbf{x}_i , output vector $\mathbf{y}_i = f(\mathbf{x}_i; \theta)$
- positive $\mathbf{y}_i^+ = f(\mathbf{x}_i^+; \theta)$, negative $\mathbf{y}_i^- = f(\mathbf{x}_i^-; \theta)$
- triplet loss is a function of distances $\|\mathbf{y}_i - \mathbf{y}_i^+\|$, $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ only

$$\ell_i = L(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-) = \ell(\|\mathbf{y}_i - \mathbf{y}_i^+\|, \|\mathbf{y}_i - \mathbf{y}_i^-\|)$$

$$\ell(x^+, x^-) = [m + (x^+)^2 - (x^-)^2]_+$$

so distance $\|\mathbf{y}_i - \mathbf{y}_i^+\|$ should be less than $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ by margin m

- by taking two pairs $(\mathbf{x}_i, \mathbf{x}_i^+)$ and $(\mathbf{x}_i, \mathbf{x}_i^-)$ at a time with targets 1, 0 respectively, the contrastive loss can be written similarly

$$\ell(x^+, x^-) = (x^+)^2 + [m - x^-]_+^2$$

so distance $\|\mathbf{y}_i - \mathbf{y}_i^+\|$ should small and $\|\mathbf{y}_i - \mathbf{y}_i^-\|$ larger than m

unsupervised learning by context prediction

[Doersch et al. 2015]

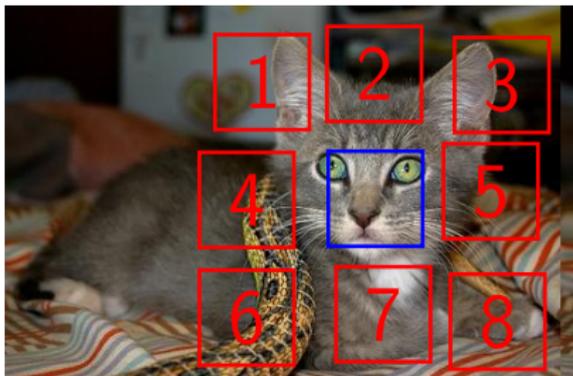


- sample random pairs of patches in one of eight spatial configurations
- patches are randomly jittered and do not overlap
- like solving a puzzle, learn to predict the relative position

$$f\left(\begin{array}{c} \\ \\ \end{array}, \begin{array}{c} \\ \\ \end{array}\right) = 3$$

unsupervised learning by context prediction

[Doersch et al. 2015]

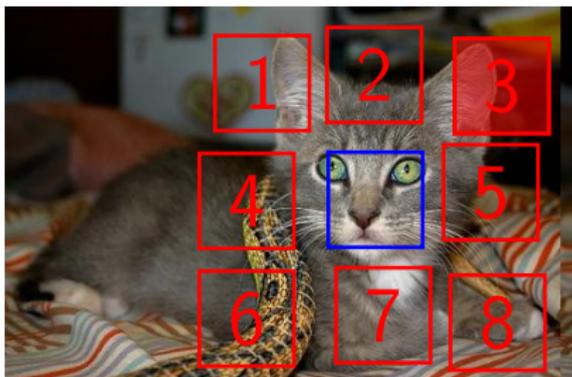


- sample random pairs of patches in one of eight spatial configurations
- patches are randomly jittered and do not overlap
- like solving a puzzle, learn to predict the relative position

$$f\left(\begin{matrix} \text{patch 1} \\ \text{patch 2} \end{matrix}, \begin{matrix} \text{patch 3} \\ \text{patch 4} \end{matrix}\right) = 3$$

unsupervised learning by context prediction

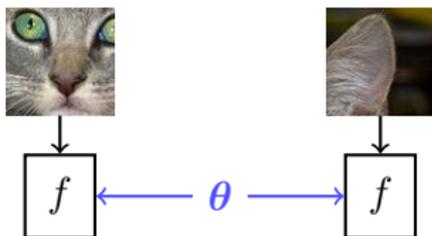
[Doersch et al. 2015]



- sample random pairs of patches in one of eight spatial configurations
- patches are randomly jittered and do not overlap
- like **solving a puzzle**, learn to predict the relative position

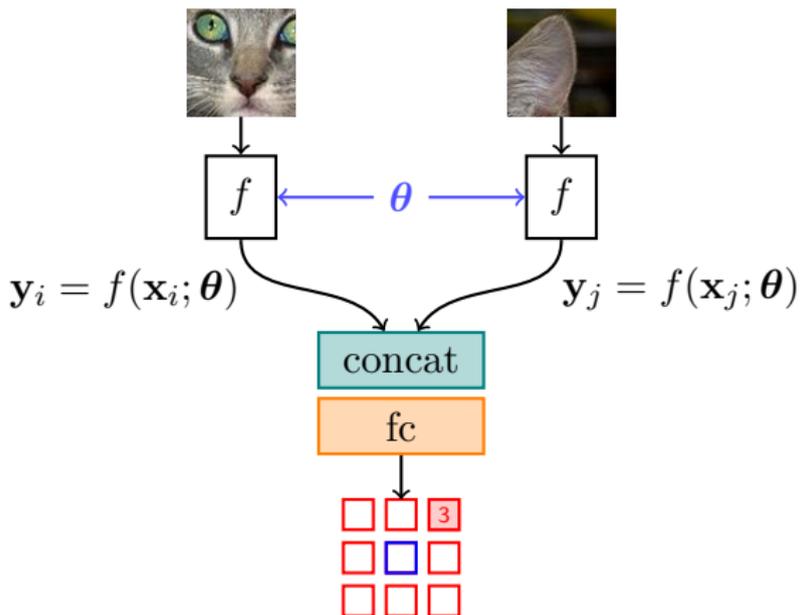
$$f\left(\begin{array}{c} \text{[cat face patch]} \\ \text{[cat ear patch]} \end{array}\right) = 3$$

context prediction: architecture



- network f learned by siamese architecture
- representations are concatenated and followed by softmax classifier, where each spatial configuration is a class

context prediction: architecture



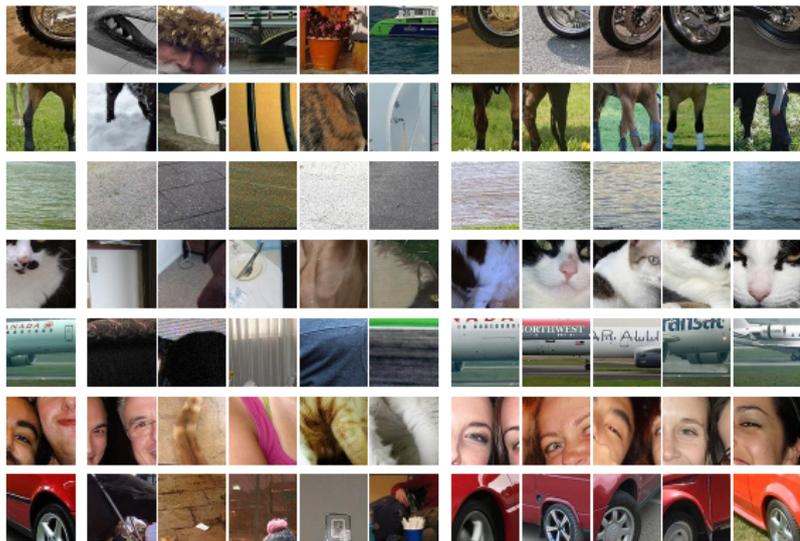
- network f learned by siamese architecture
- representations are concatenated and followed by softmax classifier, where each spatial configuration is a class

context prediction: examples



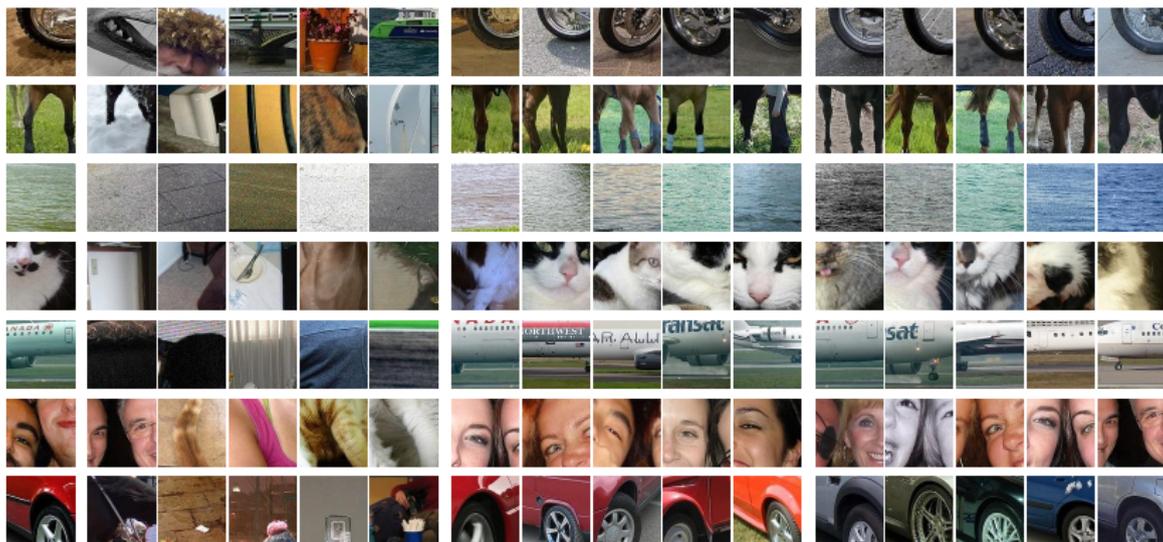
- input image
- nearest neighbors with randomly **initialized** network
- trained by **supervised** classification on ImageNet
- **unsupervised** training from scratch on the context prediction task

context prediction: examples



- input image
- nearest neighbors with randomly **initialized** network
- trained by **supervised** classification on ImageNet
- **unsupervised** training from scratch on the context prediction task

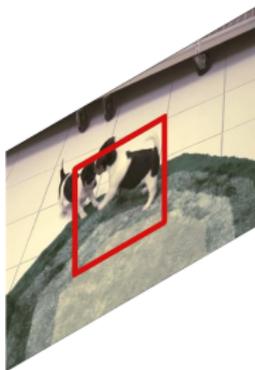
context prediction: examples



- input image
- nearest neighbors with randomly **initialized** network
- trained by **supervised** classification on ImageNet
- **unsupervised** training from scratch on the context prediction task

unsupervised learning on video: tracking

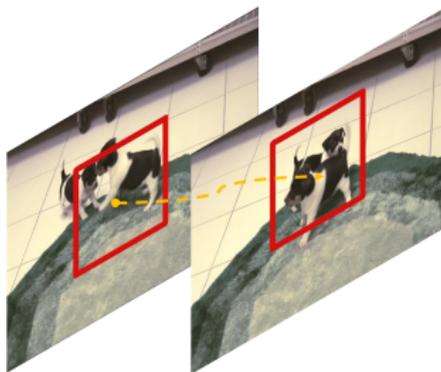
[Wang et al. 2015]



- estimate motion and find the region that contains most motion
- track this region for a number of frames
- generate a pair of matching patches on the first and last frames

unsupervised learning on video: tracking

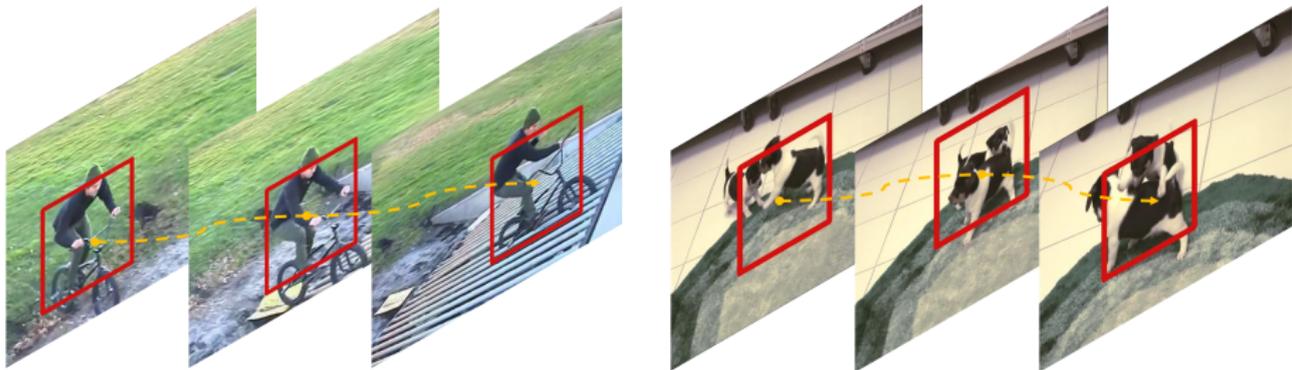
[Wang et al. 2015]



- estimate motion and find the region that contains most motion
- track this region for a number of frames
- generate a pair of matching patches on the first and last frames

unsupervised learning on video: tracking

[Wang et al. 2015]



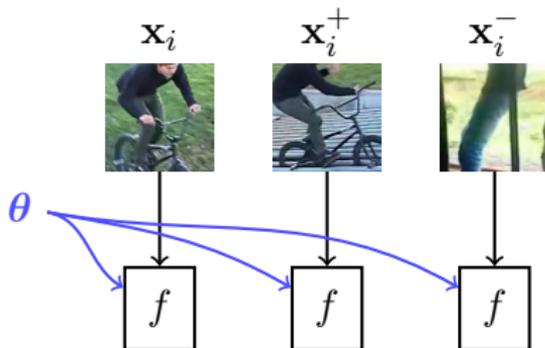
- estimate motion and find the region that contains most motion
- track this region for a number of frames
- generate a pair of matching patches on the first and last frames

unsupervised learning on video: architecture



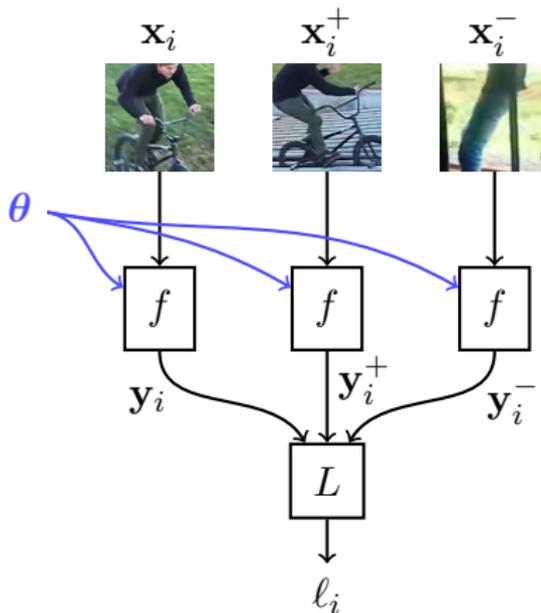
- input query \mathbf{x}_i (first frame), tracked \mathbf{x}_i^+ (last frame), random \mathbf{x}_i^-
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through the same function f with shared parameters θ
- triplet loss ℓ_i measured on output triplet $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

unsupervised learning on video: architecture



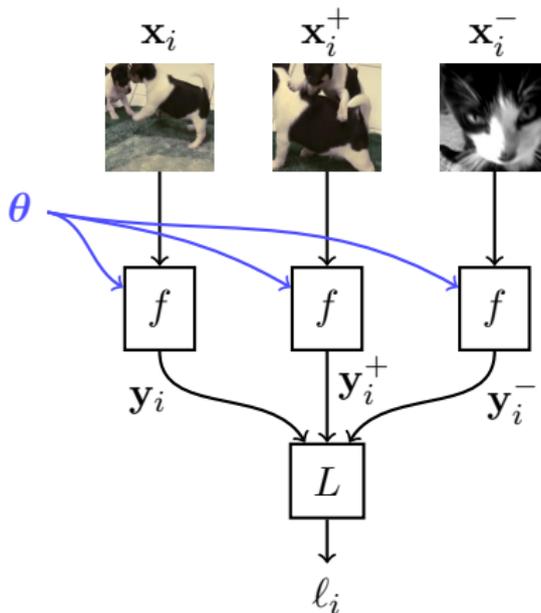
- input query \mathbf{x}_i (first frame), tracked \mathbf{x}_i^+ (last frame), random \mathbf{x}_i^-
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through the **same** function f with **shared** parameters θ
- **triplet loss** ℓ_i measured on output triplet $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

unsupervised learning on video: architecture



- input query x_i (first frame), tracked x_i^+ (last frame), random x_i^-
- x_i, x_i^+, x_i^- go through the **same** function f with **shared** parameters θ
- **triplet loss** l_i measured on output triplet (y_i, y_i^+, y_i^-)

unsupervised learning on video: architecture



- input query x_i (first frame), tracked x_i^+ (last frame), random x_i^-
- x_i, x_i^+, x_i^- go through the **same** function f with **shared** parameters θ
- **triplet loss** l_i measured on output triplet (y_i, y_i^+, y_i^-)

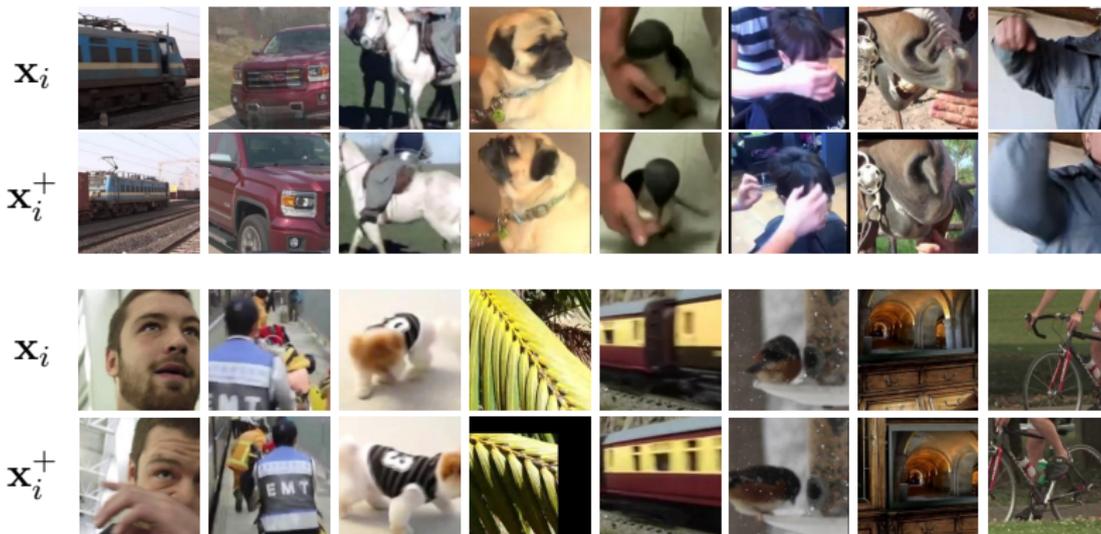
unsupervised learning on video: objective

$$\left\| f\left(\text{img}_1\right) - f\left(\text{img}_2\right) \right\|^2 < \left\| f\left(\text{img}_1\right) - f\left(\text{img}_3\right) \right\|^2 - m$$
$$\left\| f\left(\text{img}_4\right) - f\left(\text{img}_5\right) \right\|^2 < \left\| f\left(\text{img}_4\right) - f\left(\text{img}_6\right) \right\|^2 - m$$

The equations illustrate the objective of unsupervised learning on video. The first equation shows that the squared distance between two frames from the same video (img1 and img2) is less than the squared distance between a frame from the same video (img1) and a frame from a different video (img3), minus a margin m . The second equation shows that the squared distance between two frames from the same video (img4 and img5) is less than the squared distance between a frame from the same video (img4) and a frame from a different video (img6), minus a margin m .

- so, the objective is that squared distance $\|\mathbf{y}_i - \mathbf{y}_i^+\|^2$ is less than $\|\mathbf{y}_i - \mathbf{y}_i^-\|^2$ by margin m

unsupervised learning on video: more examples

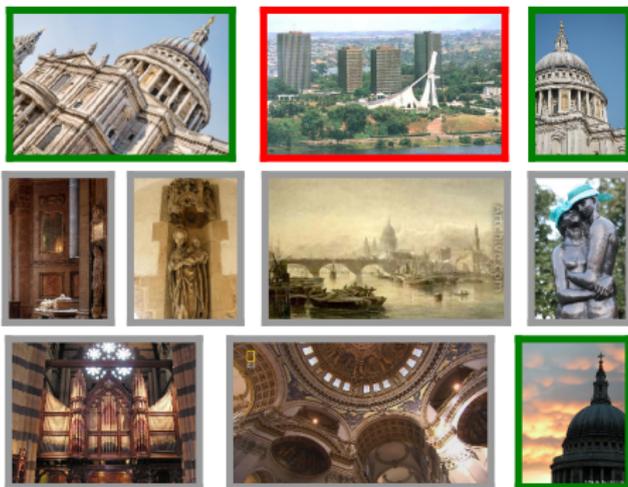


- input query x_i (first frame), tracked x_i^+ (last frame)

fine-tuning

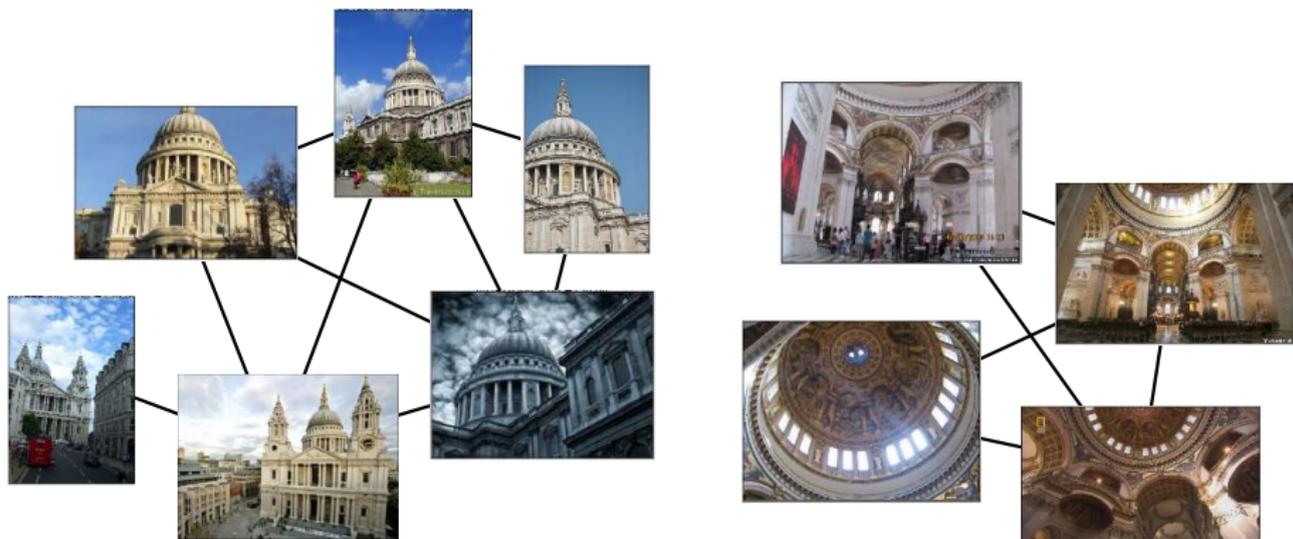
deep image retrieval: dataset cleaning

[Gordo et al. 2016]



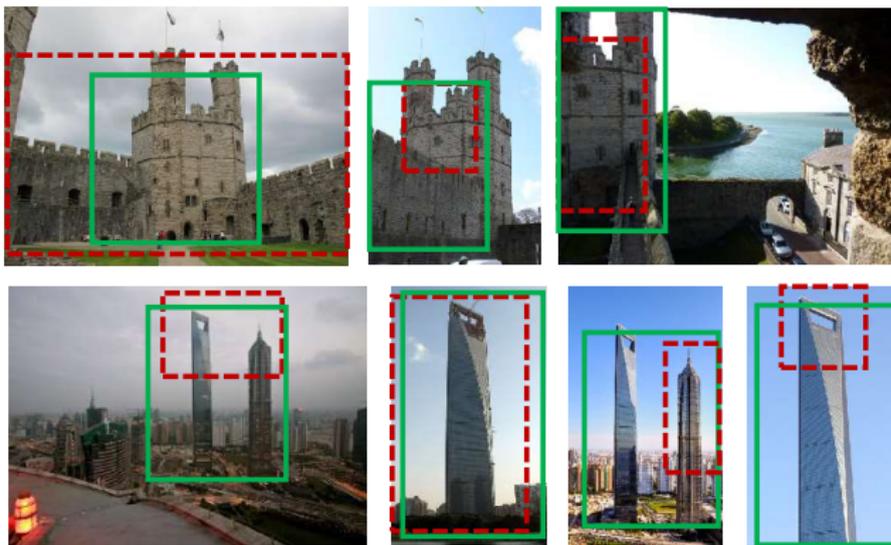
- start from landmark dataset (192k images) and **clean** it (49k images)
- use it to fine-tune a network pre-trained on ImageNet for classification
- **prototypical**, non-prototypical and **incorrect** images per class
- only prototypical are kept to reduce intra-class variability

deep image retrieval: prototypical views



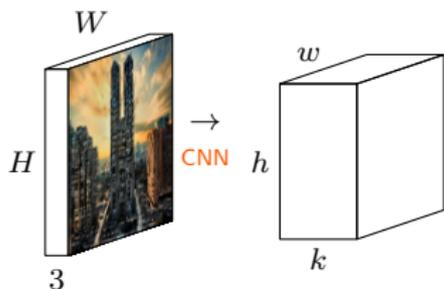
- pairwise match images per class by SIFT descriptors and fast spatial matching
- connect images into a graph and compute the connected components
- keep only the largest component

deep image retrieval: bounding boxes



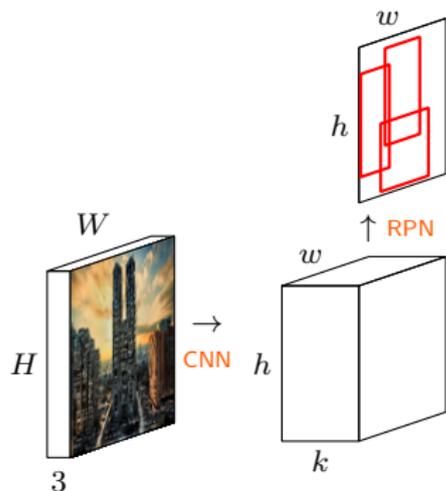
- automatically find object bounding boxes
 - **initialize** with inlier features per image
 - **update** such that boxes are consistent over all matching pairs
- use bounding boxes to train a **region proposal network**

deep image retrieval: network, regions, pooling



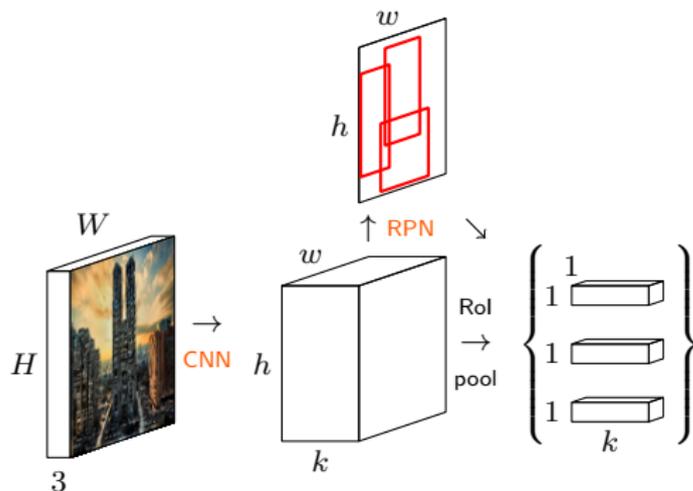
- VGG-16 or ResNet-101 feature maps
- proposals detected on feature maps by RPN and max-pooled
- ℓ_2 -normalization, PCA-whitening (FC layer), ℓ_2 -normalization
- sum-pooling, ℓ_2 -normalization (as in R-MAC)

deep image retrieval: network, regions, pooling



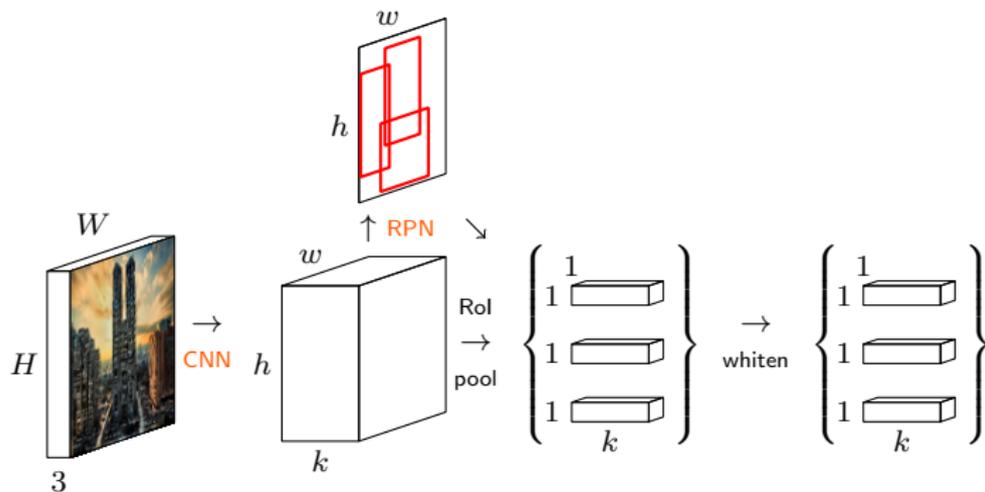
- VGG-16 or ResNet-101 feature maps
- proposals detected on feature maps by RPN and max-pooled
- ℓ_2 -normalization, PCA-whitening (FC layer), ℓ_2 -normalization
- sum-pooling, ℓ_2 -normalization (as in R-MAC)

deep image retrieval: network, regions, pooling



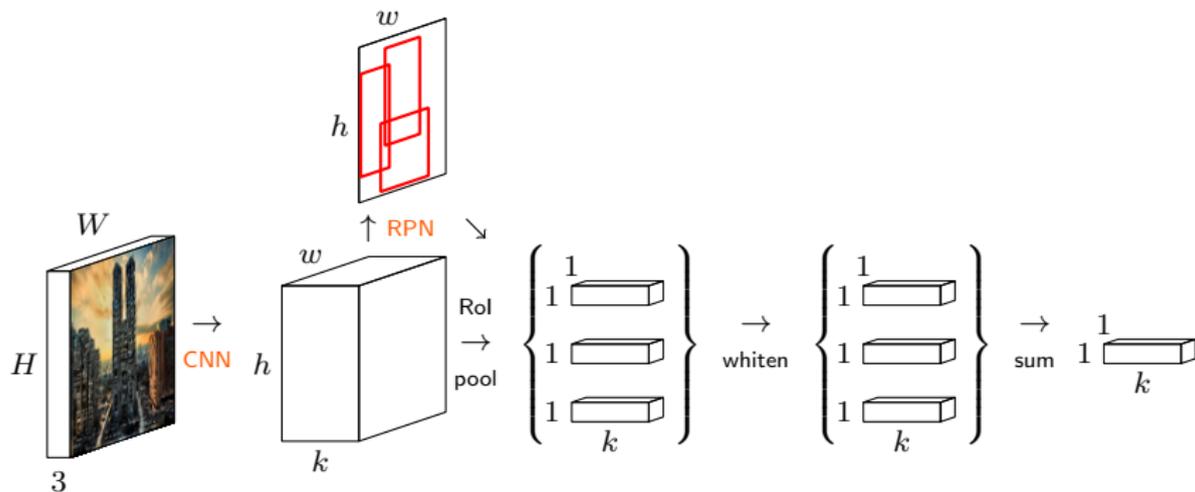
- VGG-16 or ResNet-101 feature maps
- proposals detected on feature maps by RPN and max-pooled
- ℓ_2 -normalization, PCA-whitening (FC layer), ℓ_2 -normalization
- sum-pooling, ℓ_2 -normalization (as in R-MAC)

deep image retrieval: network, regions, pooling



- VGG-16 or ResNet-101 feature maps
- proposals detected on feature maps by RPN and max-pooled
- ℓ_2 -normalization, PCA-whitening (FC layer), ℓ_2 -normalization
- sum-pooling, ℓ_2 -normalization (as in R-MAC)

deep image retrieval: network, regions, pooling



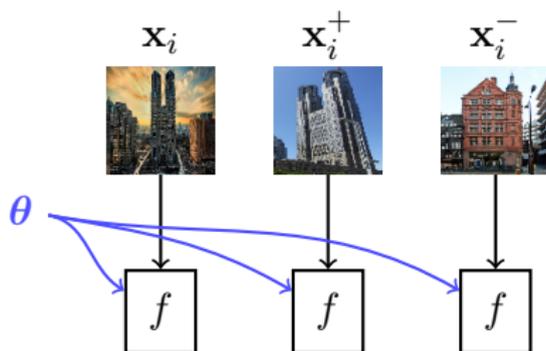
- VGG-16 or ResNet-101 feature maps
- proposals detected on feature maps by RPN and max-pooled
- ℓ_2 -normalization, PCA-whitening (FC layer), ℓ_2 -normalization
- sum-pooling, ℓ_2 -normalization (as in R-MAC)

deep image retrieval: architecture



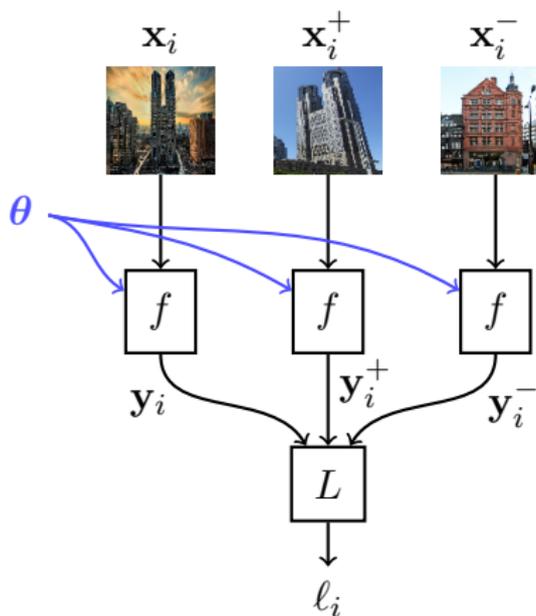
- query \mathbf{x}_i , relevant \mathbf{x}_i^+ (same building), irrelevant \mathbf{x}_i^- (other building)
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through function f including features, RPN, pooling
- triplet loss ℓ_i measured on output $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

deep image retrieval: architecture



- query \mathbf{x}_i , relevant \mathbf{x}_i^+ (same building), irrelevant \mathbf{x}_i^- (other building)
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through function f including features, RPN, pooling
- triplet loss ℓ_i measured on output ($\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-$)

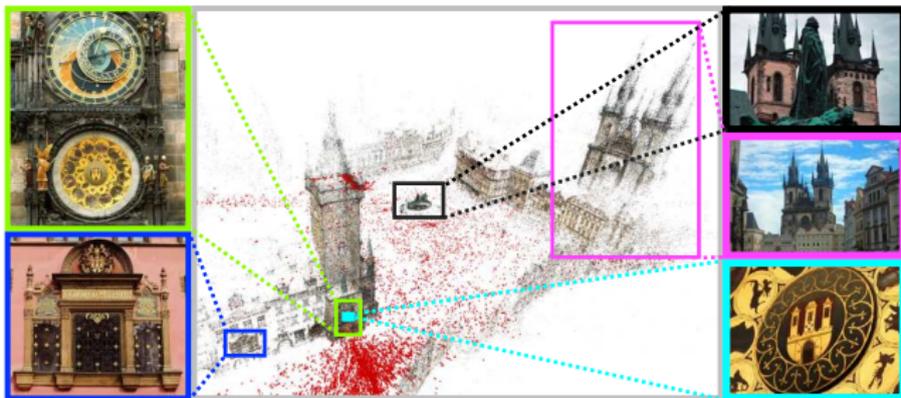
deep image retrieval: architecture



- query \mathbf{x}_i , relevant \mathbf{x}_i^+ (same building), irrelevant \mathbf{x}_i^- (other building)
- $\mathbf{x}_i, \mathbf{x}_i^+, \mathbf{x}_i^-$ go through function f including features, RPN, pooling
- triplet loss ℓ_i measured on output $(\mathbf{y}_i, \mathbf{y}_i^+, \mathbf{y}_i^-)$

learning from bag-of-words: 3d reconstruction

[Radenovic et al. 2016]



- start from an independent dataset of 7.4M images, **no class labels**
- clustering, pairwise matching and **reconstruction** of 713 3d models containing 165k unique images
- **3d models** playing the same role as classes in deep image retrieval
- again, fine-tune a network pre-trained on ImageNet for classification

learning from bag-of-words: positive pairs



- input query
- positive images found in **same model** by minimum MAC distance, maximum inliers, or drawn at random from images having at least a given number of inliers (more challenging)

learning from bag-of-words: positive pairs



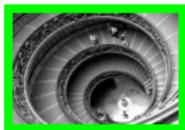
- input query
- positive images found in **same model** by minimum MAC distance, maximum inliers, or drawn at random from images having at least a given number of inliers (more challenging)

learning from bag-of-words: positive pairs



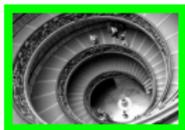
- input query
- positive images found in **same model** by minimum MAC distance, maximum inliers, or drawn at random from images having at least a given number of inliers (more challenging)

learning from bag-of-words: negative pairs



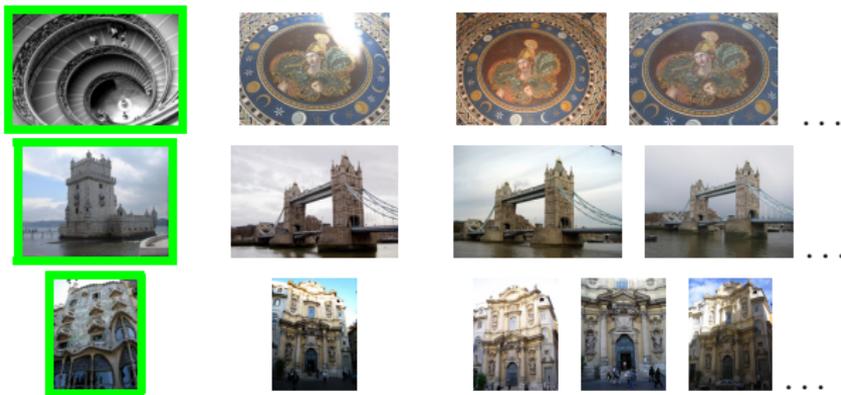
- input query
- negative images found in **different models**
- **hard negatives** are most similar to query, *i.e.* with minimum MAC distance
- hardest negative, nearest neighbors from all other models, or nearest neighbors, one per model (higher variability)

learning from bag-of-words: negative pairs



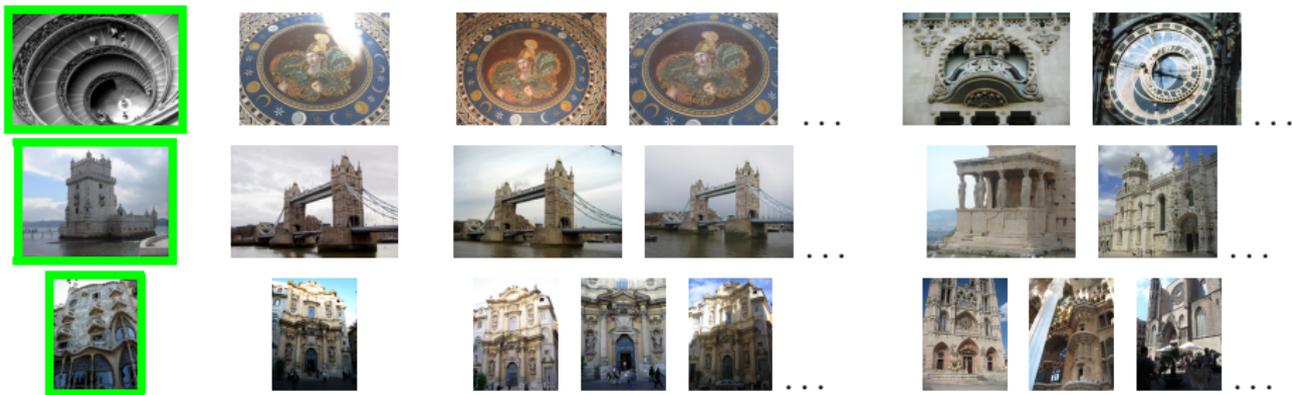
- input query
- negative images found in **different models**
- **hard negatives** are most similar to query, *i.e.* with minimum MAC distance
- **hardest negative**, nearest neighbors from all other models, or nearest neighbors, one per model (higher variability)

learning from bag-of-words: negative pairs



- input query
- negative images found in **different models**
- **hard negatives** are most similar to query, *i.e.* with minimum MAC distance
- hardest negative, nearest neighbors from all other models, or nearest neighbors, one per model (higher variability)

learning from bag-of-words: negative pairs



- input query
- negative images found in **different models**
- **hard negatives** are most similar to query, *i.e.* with minimum MAC distance
- hardest negative, nearest neighbors from all other models, or nearest neighbors, one per model (higher variability)

learning from bag-of-words: architecture

\mathbf{x}_i

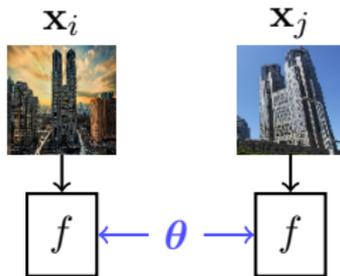


\mathbf{x}_j



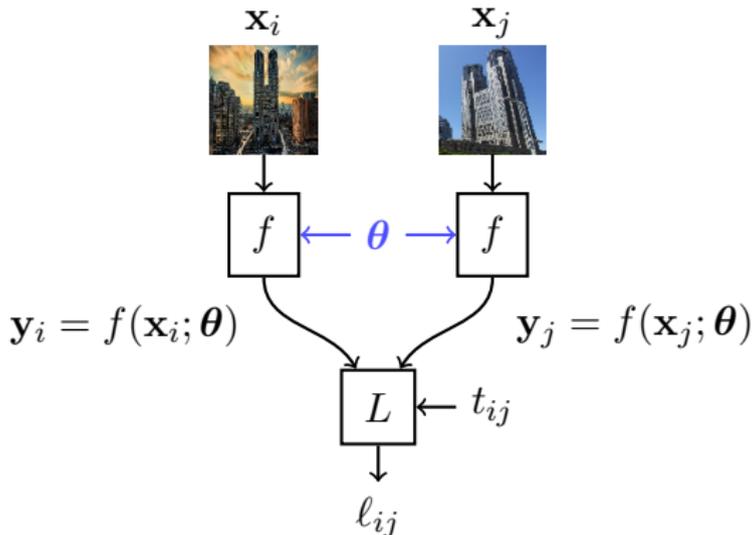
- input $(\mathbf{x}_i, \mathbf{x}_j)$ of relevant or irrelevant images
- both $\mathbf{x}_i, \mathbf{x}_j$ go through function f including features and MAC pooling
- **contrastive loss** ℓ_{ij} measured on output $(\mathbf{y}_i, \mathbf{y}_j)$ and target t_{ij}

learning from bag-of-words: architecture



- input (x_i, x_j) of relevant or irrelevant images
- both x_i, x_j go through function f including features and MAC pooling
- **contrastive loss** l_{ij} measured on output (y_i, y_j) and target t_{ij}

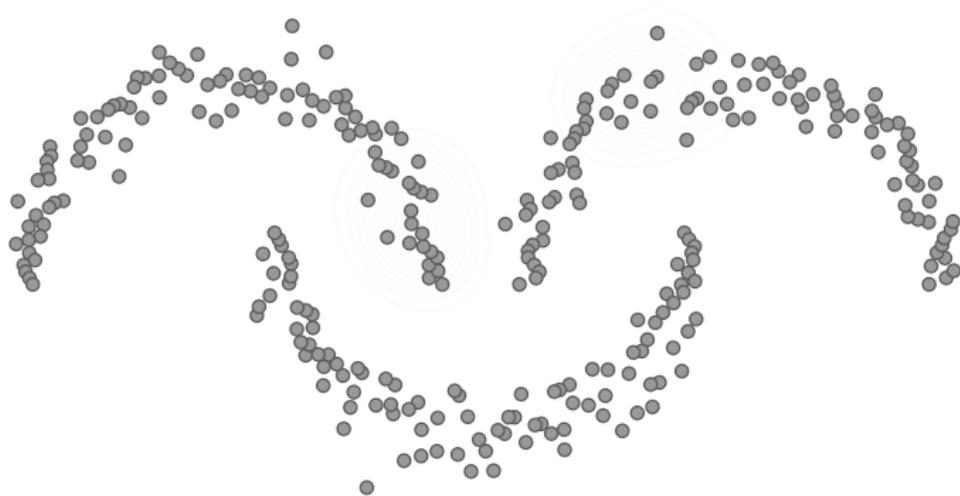
learning from bag-of-words: architecture



- input (x_i, x_j) of relevant or irrelevant images
- both x_i, x_j go through function f including features and MAC pooling
- **contrastive loss** l_{ij} measured on output (y_i, y_j) and target t_{ij}

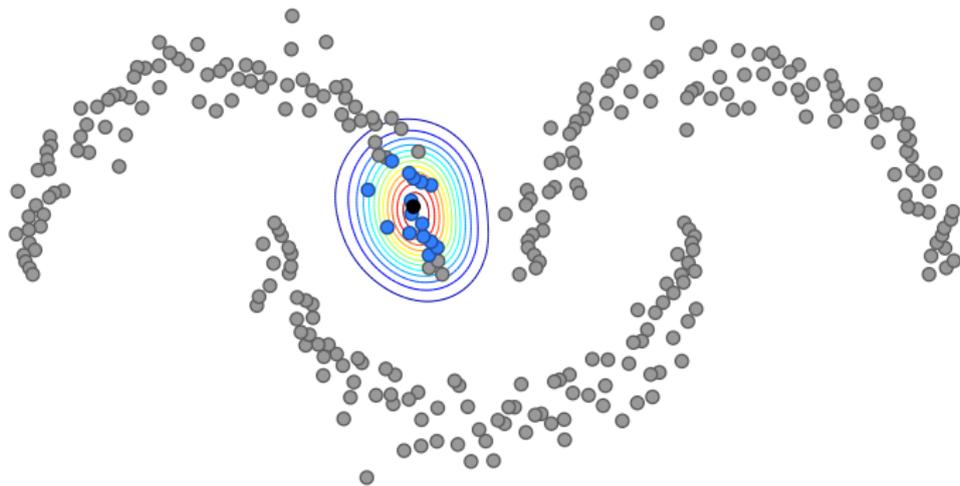
ranking on manifolds

ranking on manifolds: single query



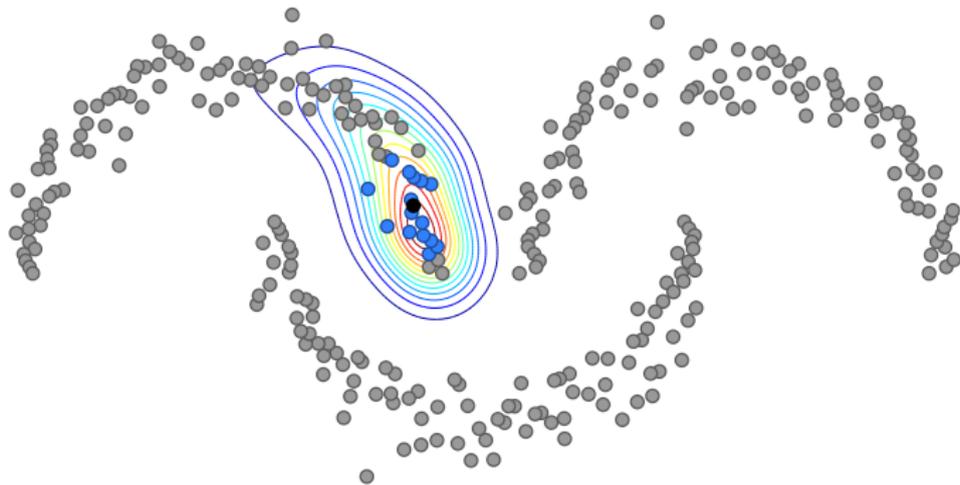
- data points (•), query point (•), nearest neighbors (•)
- iteration $\times 30$

ranking on manifolds: single query



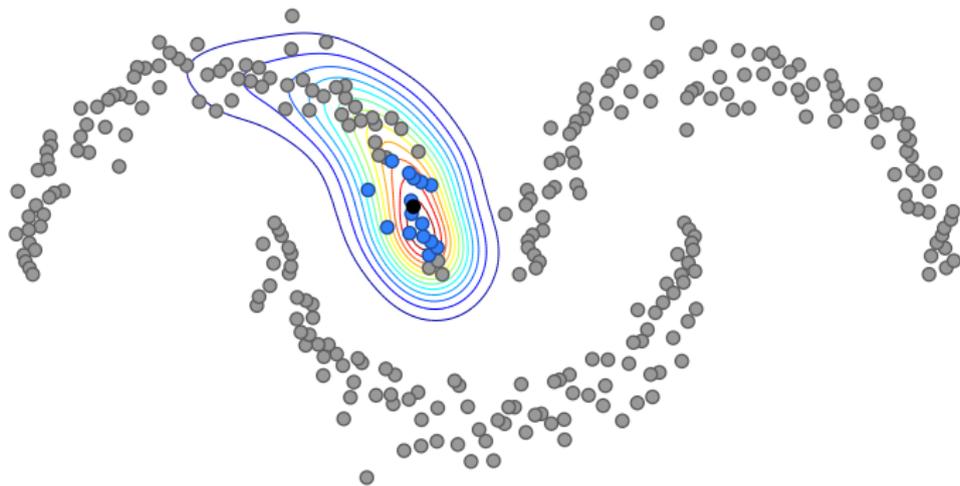
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 0×30

ranking on manifolds: single query



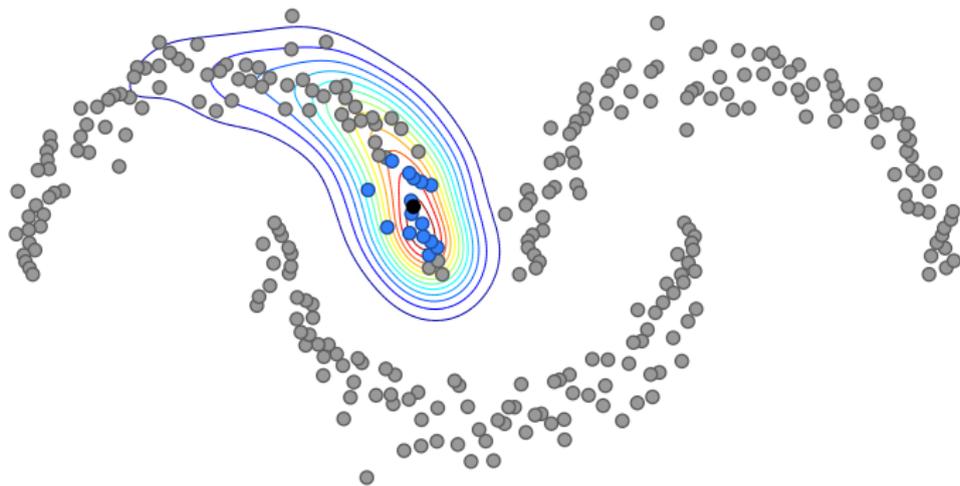
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 1×30

ranking on manifolds: single query



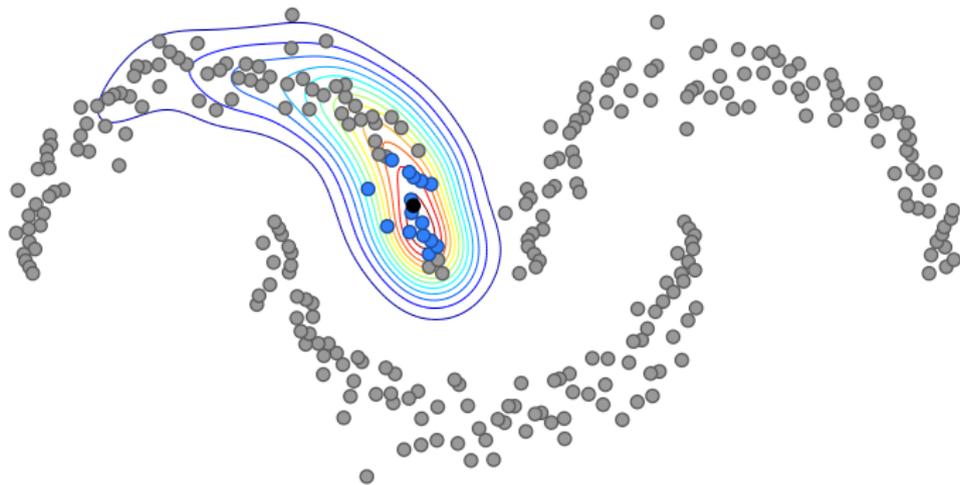
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 2×30

ranking on manifolds: single query



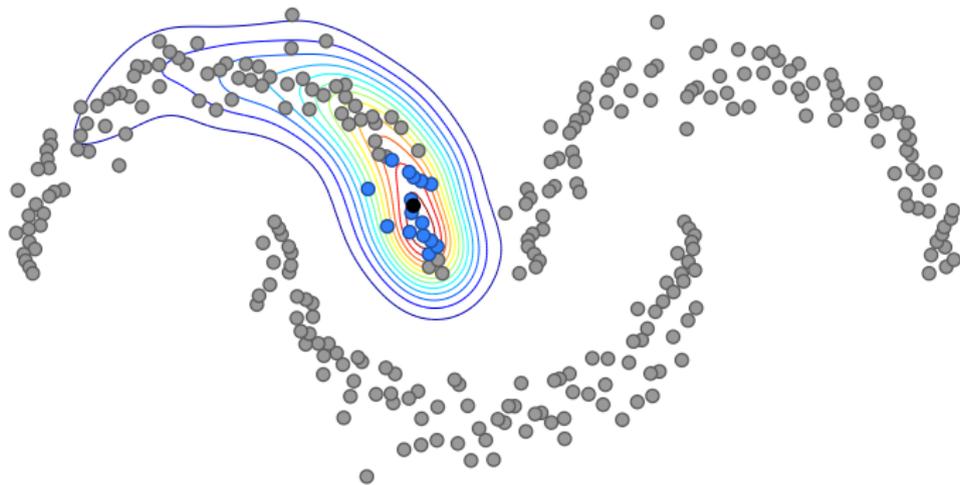
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 3×30

ranking on manifolds: single query



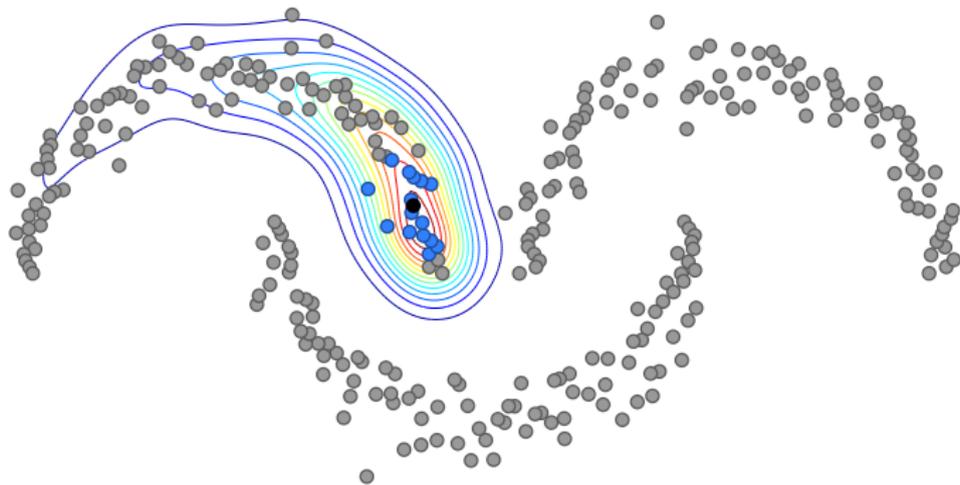
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 4×30

ranking on manifolds: single query



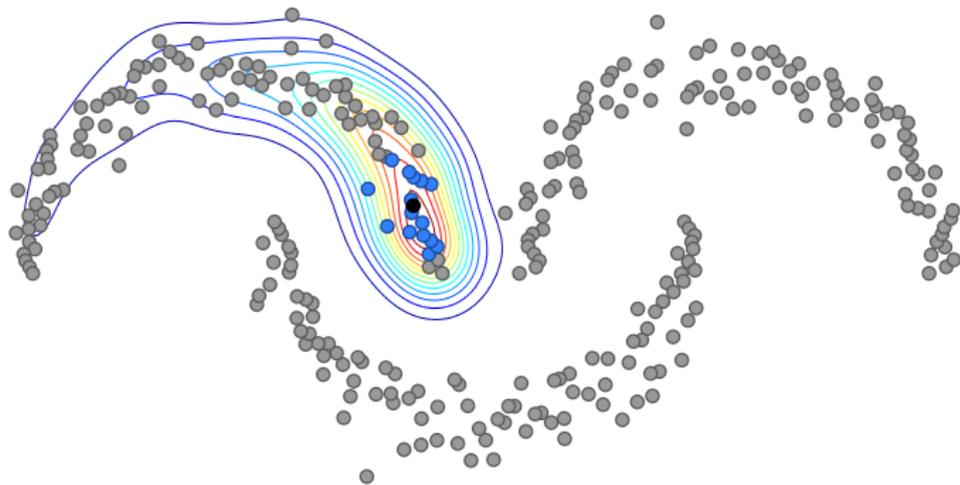
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 5×30

ranking on manifolds: single query



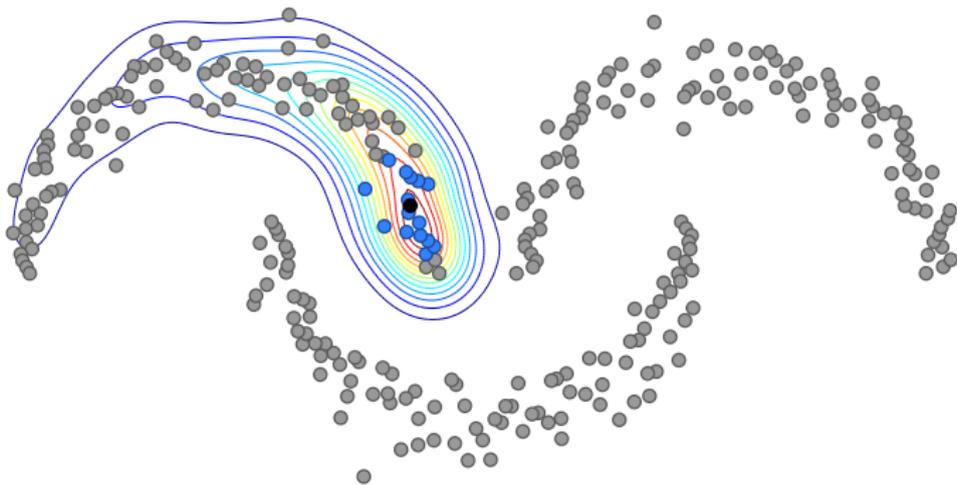
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 6×30

ranking on manifolds: single query



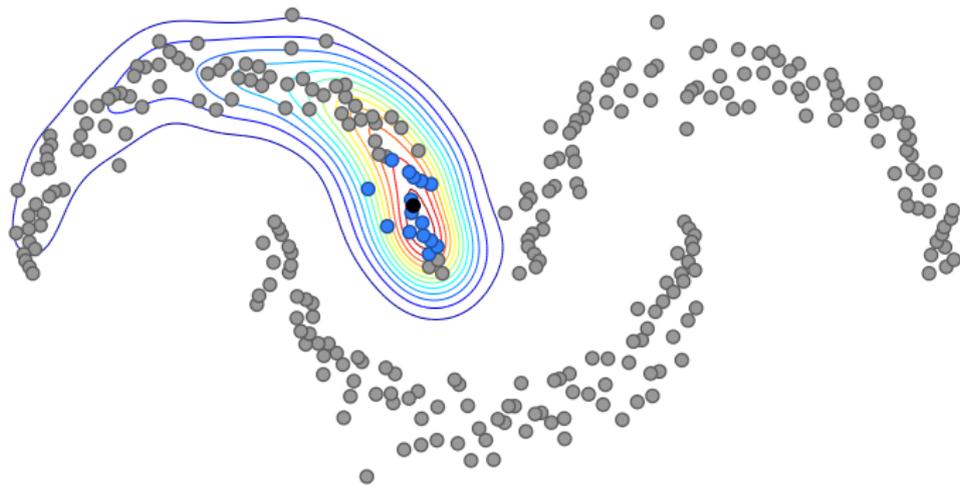
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 7×30

ranking on manifolds: single query



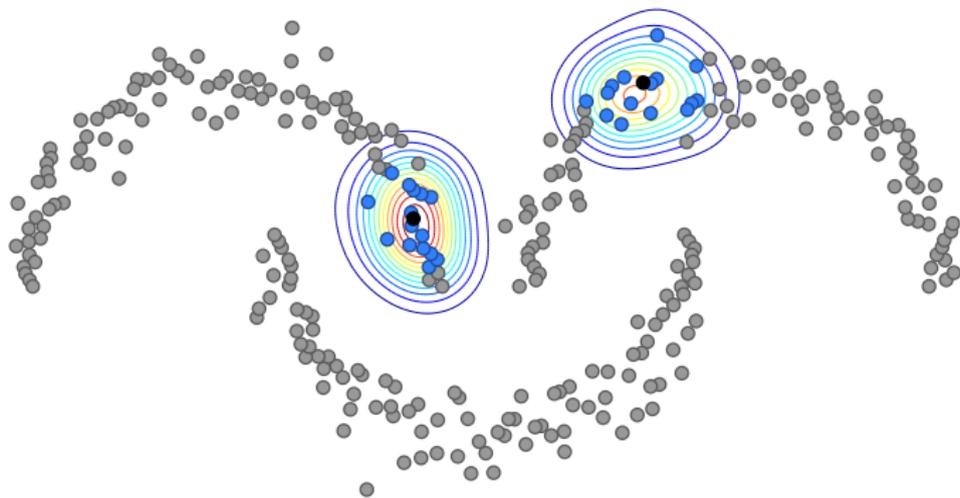
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 8×30

ranking on manifolds: single query



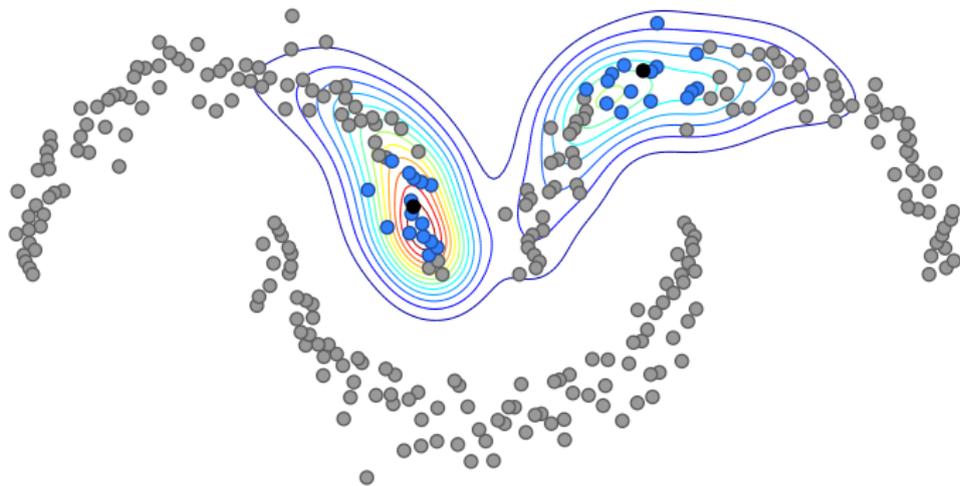
- data points (\bullet), query point (\bullet), nearest neighbors (\bullet)
- iteration 9×30

ranking on manifolds: multiple queries



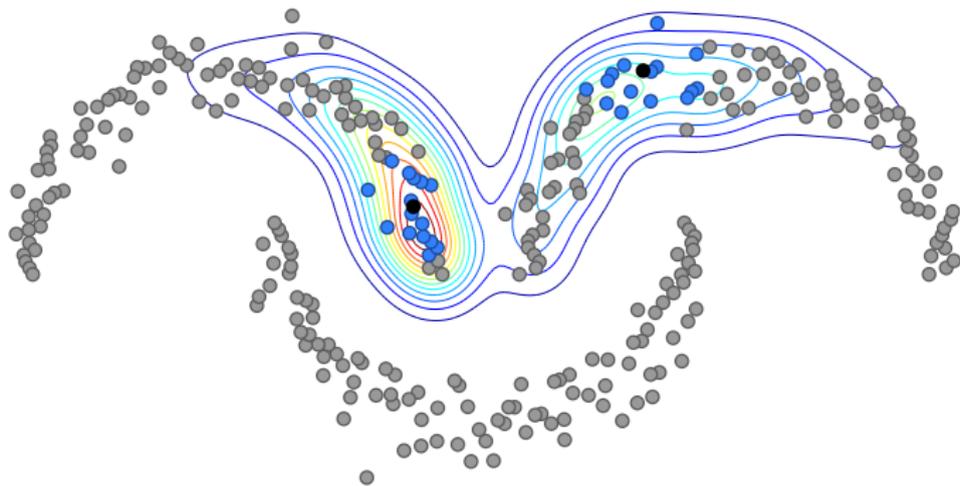
- data points (•), query points (•), nearest neighbors (•)
- iteration 0×30

ranking on manifolds: multiple queries



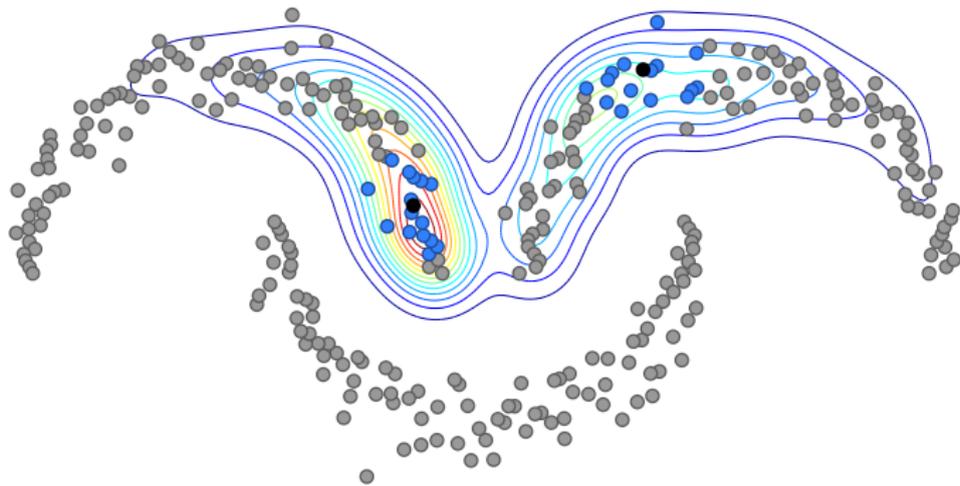
- data points (\circ), query points (\bullet), nearest neighbors (\circ)
- iteration 1×30

ranking on manifolds: multiple queries



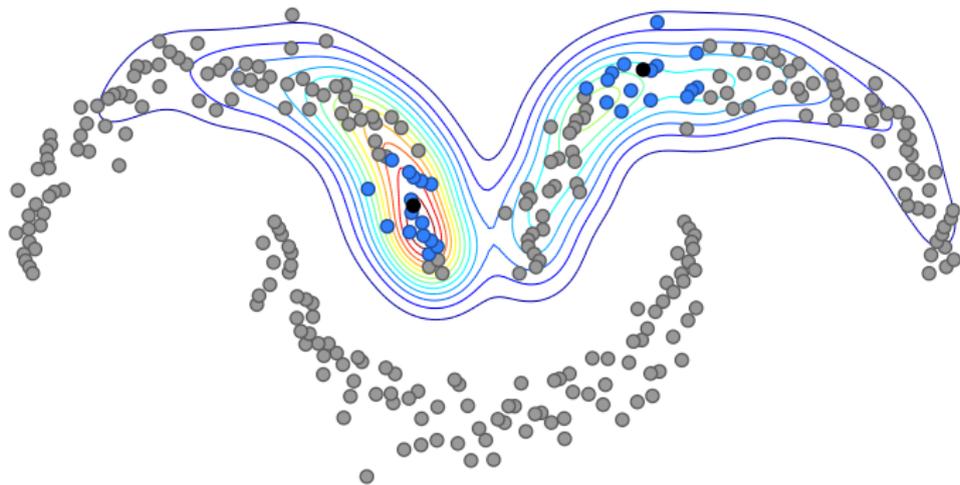
- data points (\circ), query points (\bullet), nearest neighbors (\bullet)
- iteration 2×30

ranking on manifolds: multiple queries



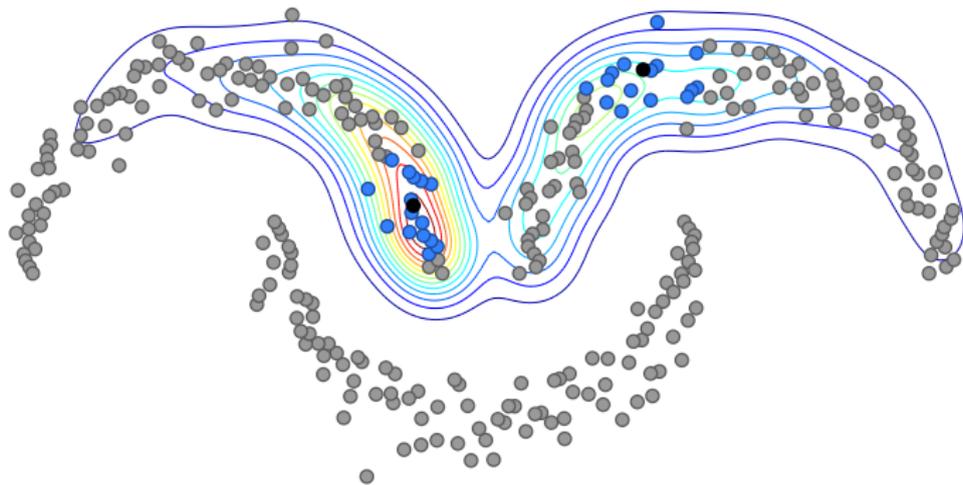
- data points (\bullet), query points (\bullet), nearest neighbors (\bullet)
- iteration 3×30

ranking on manifolds: multiple queries



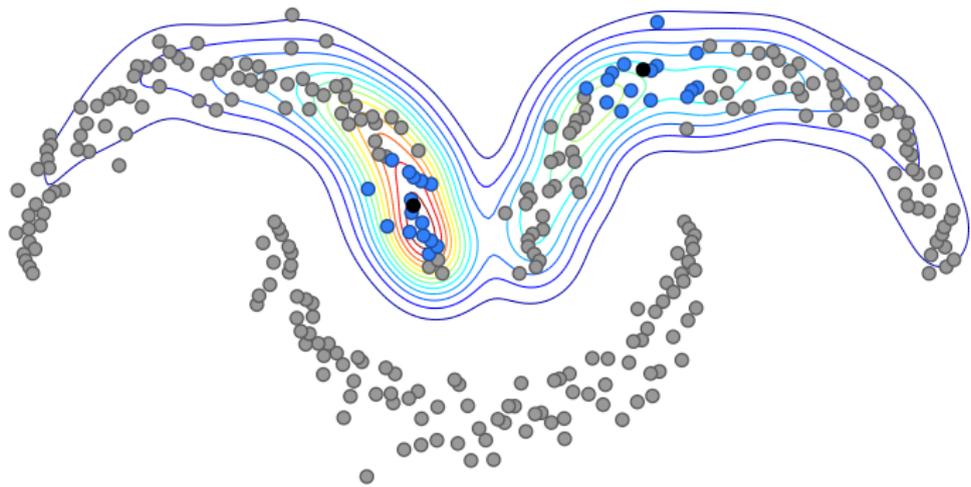
- data points (◦), query points (●), nearest neighbors (◐)
- iteration 4×30

ranking on manifolds: multiple queries



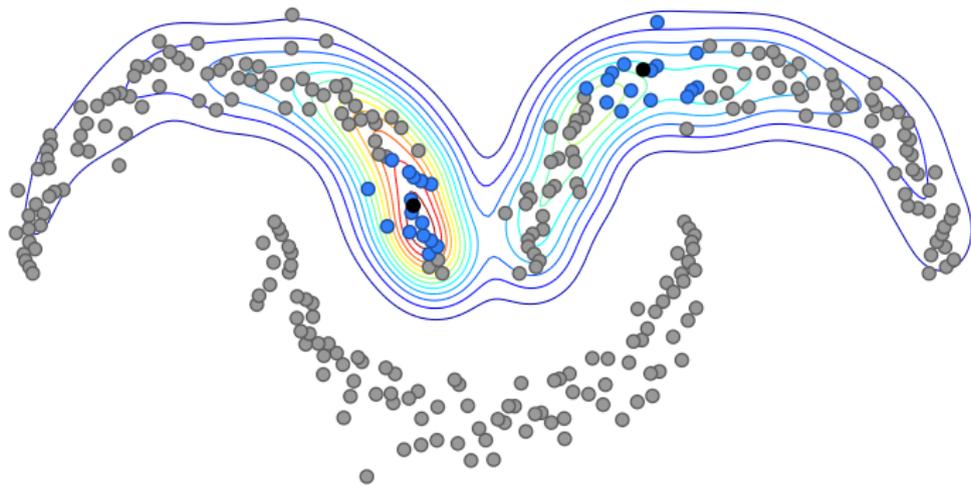
- data points (\bullet), query points (\bullet), nearest neighbors (\bullet)
- iteration 5×30

ranking on manifolds: multiple queries



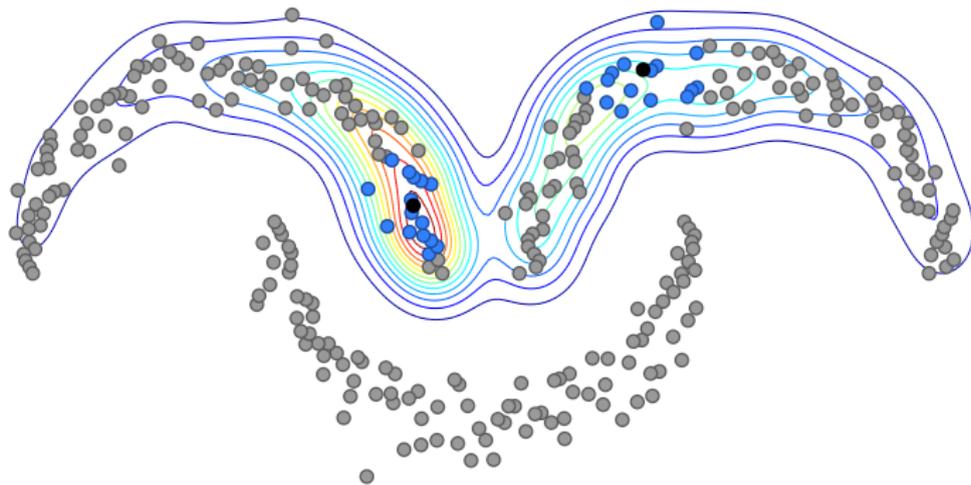
- data points (•), query points (•), nearest neighbors (•)
- iteration 6×30

ranking on manifolds: multiple queries



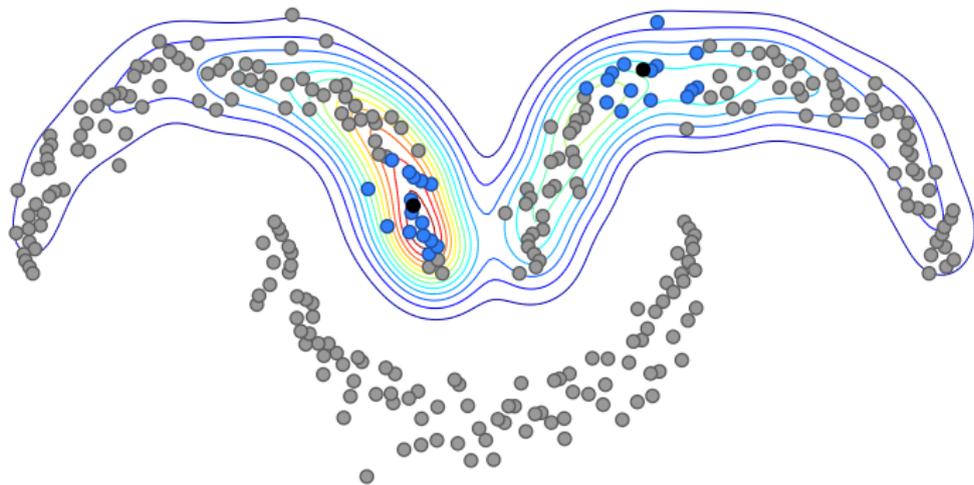
- data points (•), query points (•), nearest neighbors (•)
- iteration 7×30

ranking on manifolds: multiple queries



- data points (•), query points (•), nearest neighbors (•)
- iteration 8×30

ranking on manifolds: multiple queries



- data points (•), query points (•), nearest neighbors (•)
- iteration 9×30

ranking on manifolds: random walk

[Zhou et al. 2003]

- reciprocal k -nearest neighbor graph on n data points
- non-negative, symmetric, sparse adjacency matrix $W \in \mathbb{R}^{n \times n}$, with zero diagonal (no self-loops)
- symmetrically normalized adjacency matrix

$$\mathcal{W} := D^{-1/2} W D^{-1/2}$$

where $D = \text{diag}(W\mathbf{1})$ is the degree matrix

- query: vector $\mathbf{y} \in \mathbb{R}^n$ with $y_i = \mathbb{1}[i \text{ is query}]$
- random walk: starting with any $\mathbf{f}^{(0)} \in \mathbb{R}^n$, iterate

$$\mathbf{f}^{(\tau)} = \alpha \mathcal{W} \mathbf{f}^{(\tau-1)} + (1 - \alpha) \mathbf{y}$$

where $\alpha \in [0, 1)$ (typically close to 1)

- rank data points by descending order of \mathbf{f}

ranking on manifolds: random walk

[Zhou et al. 2003]

- reciprocal k -nearest neighbor graph on n data points
- non-negative, symmetric, sparse adjacency matrix $W \in \mathbb{R}^{n \times n}$, with zero diagonal (no self-loops)
- symmetrically normalized adjacency matrix

$$\mathcal{W} := D^{-1/2} W D^{-1/2}$$

where $D = \text{diag}(W\mathbf{1})$ is the degree matrix

- query: vector $\mathbf{y} \in \mathbb{R}^n$ with $y_i = \mathbb{1}[i \text{ is query}]$
- random walk: starting with any $\mathbf{f}^{(0)} \in \mathbb{R}^n$, iterate

$$\mathbf{f}^{(\tau)} = \alpha \mathcal{W} \mathbf{f}^{(\tau-1)} + (1 - \alpha) \mathbf{y}$$

where $\alpha \in [0, 1)$ (typically close to 1)

- rank data points by descending order of \mathbf{f}

ranking on manifolds: random walk

[Zhou et al. 2003]

- reciprocal k -nearest neighbor graph on n data points
- non-negative, symmetric, sparse adjacency matrix $W \in \mathbb{R}^{n \times n}$, with zero diagonal (no self-loops)
- symmetrically normalized adjacency matrix

$$\mathcal{W} := D^{-1/2} W D^{-1/2}$$

where $D = \text{diag}(W\mathbf{1})$ is the degree matrix

- **query**: vector $\mathbf{y} \in \mathbb{R}^n$ with $y_i = \mathbb{1}[i \text{ is query}]$
- **random walk**: starting with any $\mathbf{f}^{(0)} \in \mathbb{R}^n$, iterate

$$\mathbf{f}^{(\tau)} = \alpha \mathcal{W} \mathbf{f}^{(\tau-1)} + (1 - \alpha) \mathbf{y}$$

where $\alpha \in [0, 1)$ (typically close to 1)

- **rank** data points by descending order of \mathbf{f}

ranking on manifolds: random walk

[Zhou et al. 2003]

- reciprocal k -nearest neighbor graph on n data points
- non-negative, symmetric, sparse adjacency matrix $W \in \mathbb{R}^{n \times n}$, with zero diagonal (no self-loops)
- symmetrically normalized adjacency matrix

$$\mathcal{W} := D^{-1/2} W D^{-1/2}$$

where $D = \text{diag}(W\mathbf{1})$ is the degree matrix

- **query**: vector $\mathbf{y} \in \mathbb{R}^n$ with $y_i = \mathbb{1}[i \text{ is query}]$
- **random walk**: starting with any $\mathbf{f}^{(0)} \in \mathbb{R}^n$, iterate

$$\mathbf{f}^{(\tau)} = \alpha \mathcal{W} \mathbf{f}^{(\tau-1)} + (1 - \alpha) \mathbf{y}$$

where $\alpha \in [0, 1)$ (typically close to 1)

- **rank** data points by descending order of \mathbf{f}

ranking on manifolds: random walk

[Zhou et al. 2003]

- reciprocal k -nearest neighbor graph on n data points
- non-negative, symmetric, sparse adjacency matrix $W \in \mathbb{R}^{n \times n}$, with zero diagonal (no self-loops)
- symmetrically normalized adjacency matrix

$$\mathcal{W} := D^{-1/2} W D^{-1/2}$$

where $D = \text{diag}(W\mathbf{1})$ is the degree matrix

- **query**: vector $\mathbf{y} \in \mathbb{R}^n$ with $y_i = \mathbb{1}[i \text{ is query}]$
- **random walk**: starting with any $\mathbf{f}^{(0)} \in \mathbb{R}^n$, iterate

$$\mathbf{f}^{(\tau)} = \alpha \mathcal{W} \mathbf{f}^{(\tau-1)} + (1 - \alpha) \mathbf{y}$$

where $\alpha \in [0, 1)$ (typically close to 1)

- **rank** data points by descending order of \mathbf{f}

ranking as solving a linear system

[Iscen et al. 2017]

- **query**: sparse vector $\mathbf{y} \in \mathbb{R}^n$ with nearest neighbor similarities

$$y_i = \sum_{\mathbf{q} \in Q} s(\mathbf{q}, \mathbf{x}_i) \times \mathbb{1}[\mathbf{x}_i \in \text{NN}_X^k(\mathbf{q})]$$

where $Q, X \subset \mathbb{R}^d$ query/data points, $\mathbf{x}_i \in X$, s similarity function

- regularized **Laplacian**

$$\mathcal{L}_\alpha = \frac{I - \alpha \mathcal{W}}{1 - \alpha}$$

- solve **linear system**

$$\mathcal{L}_\alpha \mathbf{f} = \mathbf{y}$$

by conjugate gradient method

ranking as solving a linear system

[Iscen et al. 2017]

- **query**: sparse vector $\mathbf{y} \in \mathbb{R}^n$ with nearest neighbor similarities

$$y_i = \sum_{\mathbf{q} \in Q} s(\mathbf{q}, \mathbf{x}_i) \times \mathbb{1}[\mathbf{x}_i \in \text{NN}_X^k(\mathbf{q})]$$

where $Q, X \subset \mathbb{R}^d$ query/data points, $\mathbf{x}_i \in X$, s similarity function

- regularized **Laplacian**

$$\mathcal{L}_\alpha = \frac{I - \alpha W}{1 - \alpha}$$

- solve **linear system**

$$\mathcal{L}_\alpha \mathbf{f} = \mathbf{y}$$

by conjugate gradient method

ranking as solving a linear system

[Iscen et al. 2017]

- **query**: sparse vector $\mathbf{y} \in \mathbb{R}^n$ with nearest neighbor similarities

$$y_i = \sum_{\mathbf{q} \in Q} s(\mathbf{q}, \mathbf{x}_i) \times \mathbb{1}[\mathbf{x}_i \in \text{NN}_X^k(\mathbf{q})]$$

where $Q, X \subset \mathbb{R}^d$ query/data points, $\mathbf{x}_i \in X$, s similarity function

- regularized **Laplacian**

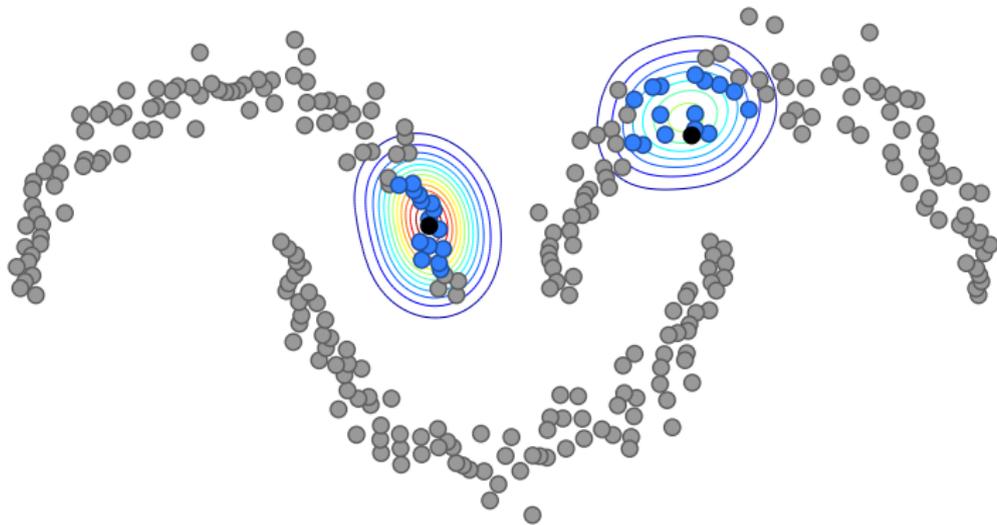
$$\mathcal{L}_\alpha = \frac{I - \alpha W}{1 - \alpha}$$

- solve **linear system**

$$\mathcal{L}_\alpha \mathbf{f} = \mathbf{y}$$

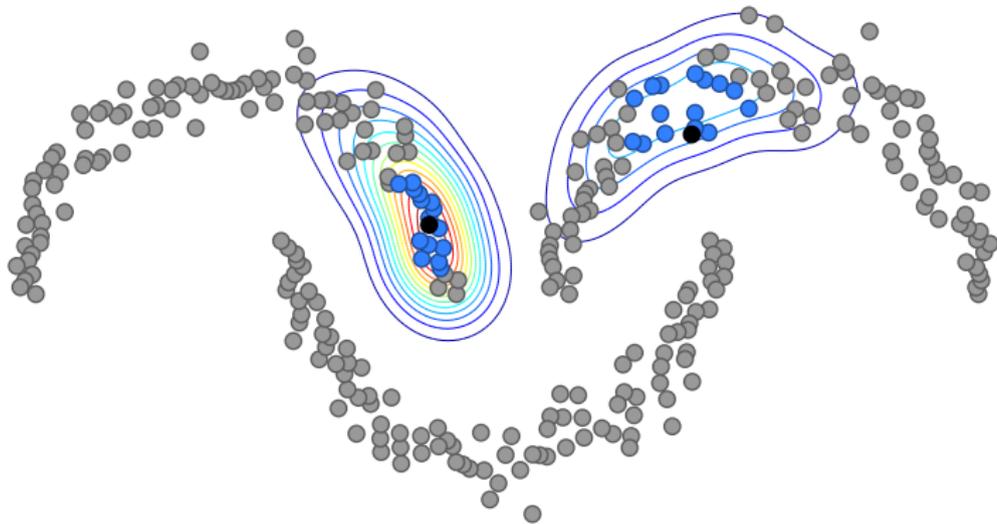
by conjugate gradient method

ranking by conjugate gradient



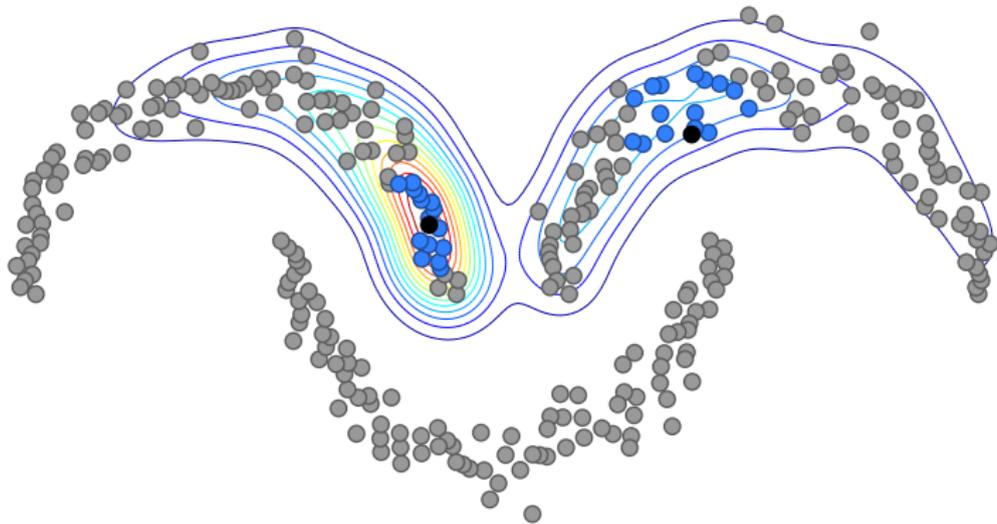
- data points (◦), query points (•), nearest neighbors (◐)
- iteration 0×2

ranking by conjugate gradient



- data points (◦), query points (•), nearest neighbors (◐)
- iteration 2×2

ranking by conjugate gradient

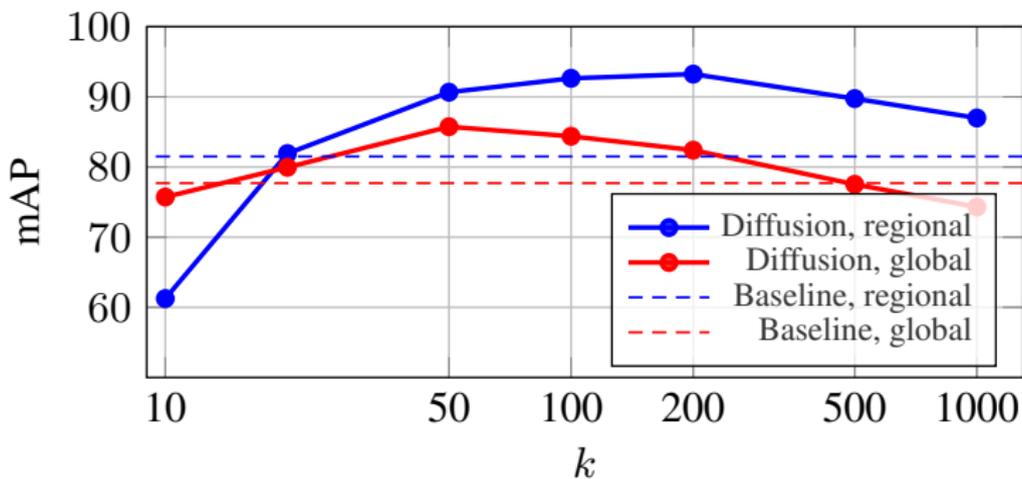


- data points (◦), query points (●), nearest neighbors (◦)
- iteration 6×2

ranking as solving a linear system

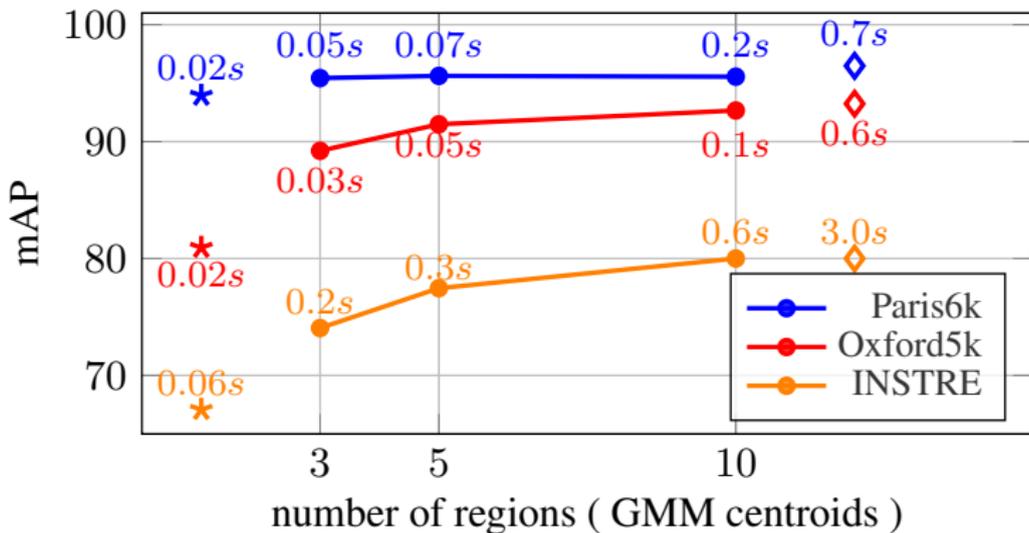
- represent image by global descriptor or multiple regional descriptors
- perform initial query based on Euclidean nearest neighbors
- re-rank by solving linear system
- ResNet-101 fine-tuned by BoW + R-MAC + re-ranking:
 - mAP 87.1 (95.8) on Oxford5k, 96.5 (96.9) on Paris6k
 - 1 (21) descriptors/image \times 2048 dimensions

dependence on neighbors, k (Oxford5k)

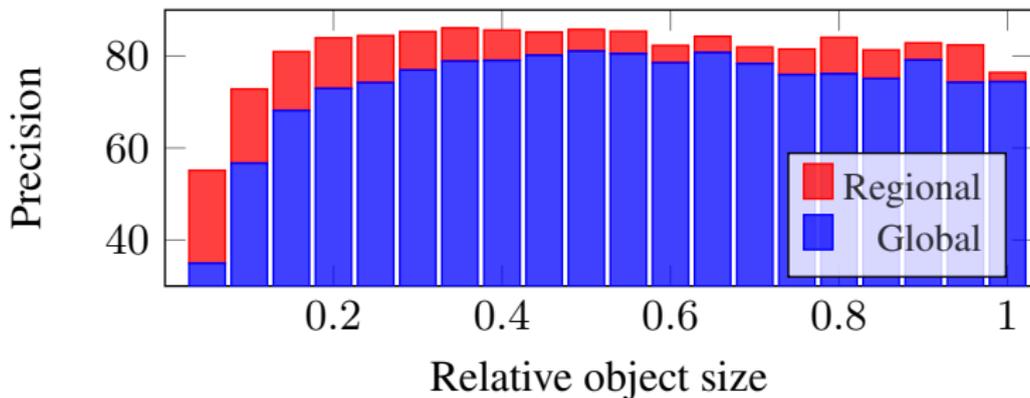


“small patterns appear more frequently than entire images”

global \rightarrow regional



small objects (INSTRE)



fast spectral ranking

faster than CG?

- want to solve $\mathcal{L}_\alpha \mathbf{x} = \mathbf{y}$
- could invert \mathcal{L}_α offline, but it wouldn't be sparse
- could approximate \mathcal{L}_α^{-1} by $\Phi^\top \Phi$ where Φ is a (sparse) $r \times n$ matrix with $r \ll n$; then

$$\mathbf{x} \approx \Phi^\top \Phi \mathbf{y}$$

- but how to compute Φ without ever inverting \mathcal{L}_α ?
- still, there is no generalization; even α is given in advance

faster than CG?

- want to solve $\mathcal{L}_\alpha \mathbf{x} = \mathbf{y}$
- could invert \mathcal{L}_α offline, but it wouldn't be sparse
- could approximate \mathcal{L}_α^{-1} by $\Phi^\top \Phi$ where Φ is a (sparse) $r \times n$ matrix with $r \ll n$; then

$$\mathbf{x} \approx \Phi^\top \Phi \mathbf{y}$$

- but how to compute Φ without ever inverting \mathcal{L}_α ?
- still, there is no generalization; even α is given in advance

faster than CG?

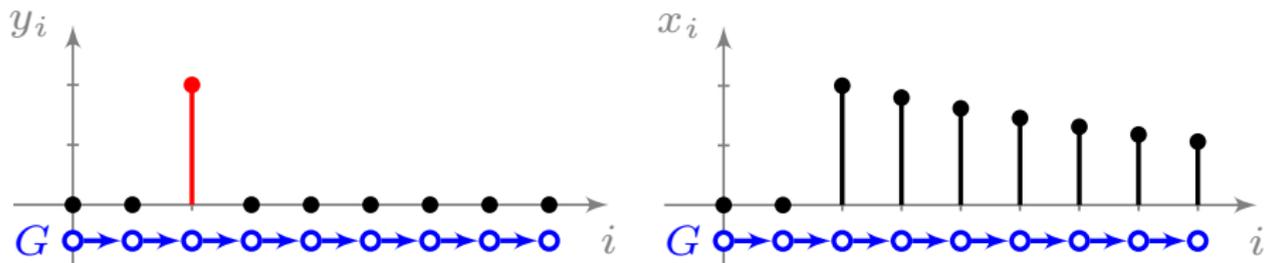
- want to solve $\mathcal{L}_\alpha \mathbf{x} = \mathbf{y}$
- could invert \mathcal{L}_α offline, but it wouldn't be sparse
- could approximate \mathcal{L}_α^{-1} by $\Phi^\top \Phi$ where Φ is a (sparse) $r \times n$ matrix with $r \ll n$; then

$$\mathbf{x} \approx \Phi^\top \Phi \mathbf{y}$$

- but how to compute Φ without ever inverting \mathcal{L}_α ?
- still, there is no generalization; even α is given in advance

searching on manifolds as smoothing

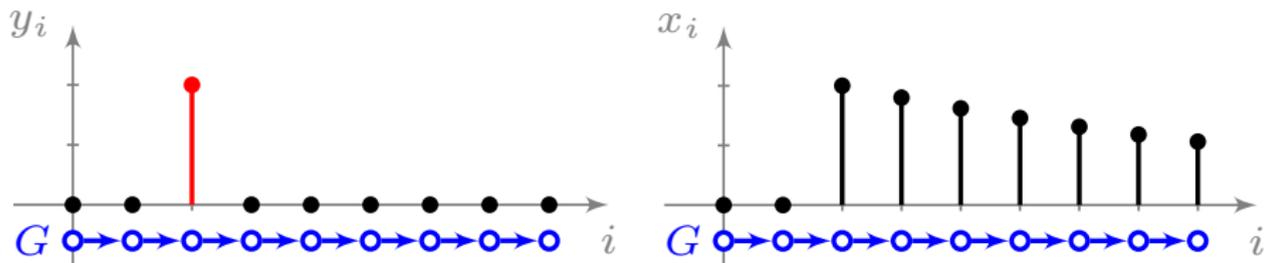
[Isken et al. 2018]



- exponential moving average filter
- output given by $x_i := (1 - \alpha) \sum_{t=0}^{\infty} \alpha^t y_{i-t}$
- or by recurrence $x_i = \alpha x_{i-1} + (1 - \alpha) y_i$
- impulse response $h_i = (1 - \alpha) \alpha^i u_i$
- transfer function $H(z) := (1 - \alpha) \sum_{t=0}^{\infty} (\alpha z^{-1})^t = (1 - \alpha) / (1 - \alpha z^{-1})$

searching on manifolds as smoothing

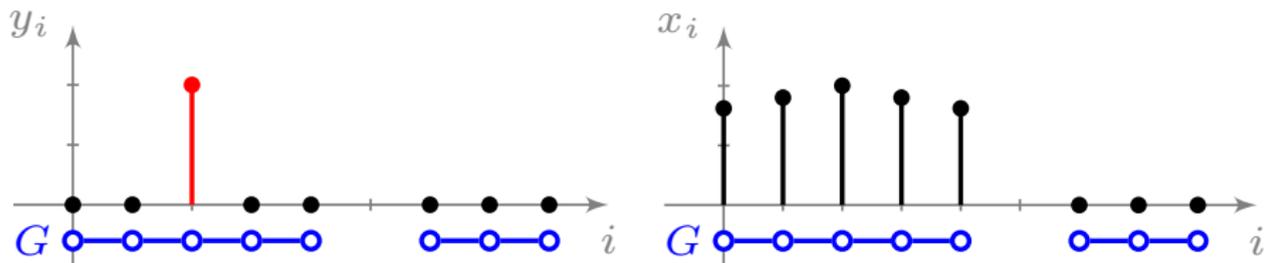
[Isken et al. 2018]



- exponential moving average filter
- output given by $x_i := (1 - \alpha) \sum_{t=0}^{\infty} \alpha^t y_{i-t}$
- or by recurrence $x_i = \alpha x_{i-1} + (1 - \alpha) y_i$
- impulse response $h_i = (1 - \alpha) \alpha^i u_i$
- transfer function $H(z) := (1 - \alpha) \sum_{t=0}^{\infty} (\alpha z^{-1})^t = (1 - \alpha) / (1 - \alpha z^{-1})$

searching on manifolds as smoothing

[Iscen et al. 2018]



- using a weighted undirected graph G instead
- information “flows” in all directions, controlled by edge weights

searching on manifolds as smoothing

- express \mathcal{L}_α^{-1} using a transfer function

$$\mathcal{L}_\alpha^{-1} = h_\alpha(\mathcal{W}) = (1 - \alpha)(I - \alpha\mathcal{W})^{-1}$$

- given any matrix function h , we want to compute

$$\mathbf{x} = h(\mathcal{W})\mathbf{y}$$

without computing $h(\mathcal{W})$

searching on manifolds as smoothing

- express \mathcal{L}_α^{-1} using a transfer function

$$\mathcal{L}_\alpha^{-1} = h_\alpha(\mathcal{W}) = (1 - \alpha)(I - \alpha\mathcal{W})^{-1}$$

- given any matrix function h , we want to compute

$$\mathbf{x} = h(\mathcal{W})\mathbf{y}$$

without computing $h(\mathcal{W})$

searching on manifolds as smoothing

$$\mathbf{x} = h \left(\mathcal{W} \right) \mathbf{y}$$

- eigenvalue decomposition of \mathcal{W}
- low-rank approximation
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

searching on manifolds as smoothing

$$\mathbf{x} = h \left(\begin{array}{|c|} \hline U \\ \hline \end{array} \begin{array}{|c|} \hline \Lambda \\ \hline \end{array} \begin{array}{|c|} \hline U^\top \\ \hline \end{array} \right) \mathbf{y}$$

- eigenvalue decomposition of \mathcal{W}
- low-rank approximation
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

searching on manifolds as smoothing

$$\mathbf{x} \approx h \left(\begin{array}{c} U \\ \Lambda \\ U^T \end{array} \right) \mathbf{y}$$

- eigenvalue decomposition of \mathcal{W}
- **low-rank approximation**
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

searching on manifolds as smoothing

$$\mathbf{x} \approx U h \left(\Lambda \right) U^T \mathbf{y}$$

- eigenvalue decomposition of \mathcal{W}
- low-rank approximation
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

searching on manifolds as smoothing

$$\mathbf{x} \approx U h \left(\Lambda \right) U^T \mathbf{y}$$

diagonal sparse

- eigenvalue decomposition of \mathcal{W}
- low-rank approximation
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

searching on manifolds as smoothing

$$\mathbf{x} \approx \mathcal{F}^{-1} h \left(\begin{array}{c} \Lambda \end{array} \right) \mathcal{F} \mathbf{y}$$

- eigenvalue decomposition of \mathcal{W}
- low-rank approximation
- (under conditions on h and Λ)
- dot-product search
- linear graph filter in frequency domain

interpretation: graph signal processing

- in discrete signal processing, a **signal** of period n is a vector $\mathbf{s} \in \mathbb{R}^n$ where $s_{\bar{i}} := s_{(i \bmod n)+1}$ for $i \in 1, \dots, n$
- a **shift** of \mathbf{s} is the mapping $s_{\bar{i}} \mapsto s_{\overline{i-1}}$; also represented by $\mathbf{s} \mapsto C_n \mathbf{s}$ where C_n is an $n \times n$ circulant zero-one matrix, e.g. for $n = 5$

$$C_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

interpretation: graph signal processing

- in discrete signal processing, a **signal** of period n is a vector $\mathbf{s} \in \mathbb{R}^n$ where $s_{\bar{i}} := s_{(i \bmod n)+1}$ for $i \in 1, \dots, n$
- a **shift** of \mathbf{s} is the mapping $s_{\bar{i}} \mapsto s_{\overline{i-1}}$; also represented by $\mathbf{s} \mapsto C_n \mathbf{s}$ where C_n is an $n \times n$ circulant zero-one matrix, e.g. for $n = 5$

$$C_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

interpretation: graph signal processing

- a linear, shift invariant, causal **filter** is the mapping $s \mapsto Hs$ where

$$H := h(C_n) = \sum_{t=0}^{\infty} h_t C_n^t$$

- matrix C_n has the eigenvalue decomposition $U\Lambda U^\top$ where U^\top is the $n \times n$ **discrete Fourier transform** matrix \mathcal{F}
- if the series $h(C_n)$ converges, filtering $s \mapsto Hs$ is written as

$$s \mapsto \mathcal{F}^{-1}h(\Lambda)\mathcal{F}s$$

- **graph signal processing** generalizes the above by replacing C_n with a matrix determined by a graph

interpretation: graph signal processing

- a linear, shift invariant, causal **filter** is the mapping $\mathbf{s} \mapsto H\mathbf{s}$ where

$$H := h(C_n) = \sum_{t=0}^{\infty} h_t C_n^t$$

- matrix C_n has the eigenvalue decomposition $U\Lambda U^\top$ where U^\top is the $n \times n$ **discrete Fourier transform** matrix \mathcal{F}
- if the series $h(C_n)$ converges, filtering $\mathbf{s} \mapsto H\mathbf{s}$ is written as

$$\mathbf{s} \mapsto \mathcal{F}^{-1}h(\Lambda)\mathcal{F}\mathbf{s}$$

- **graph signal processing** generalizes the above by replacing C_n with a matrix determined by a graph

interpretation: graph signal processing

- a linear, shift invariant, causal **filter** is the mapping $\mathbf{s} \mapsto H\mathbf{s}$ where

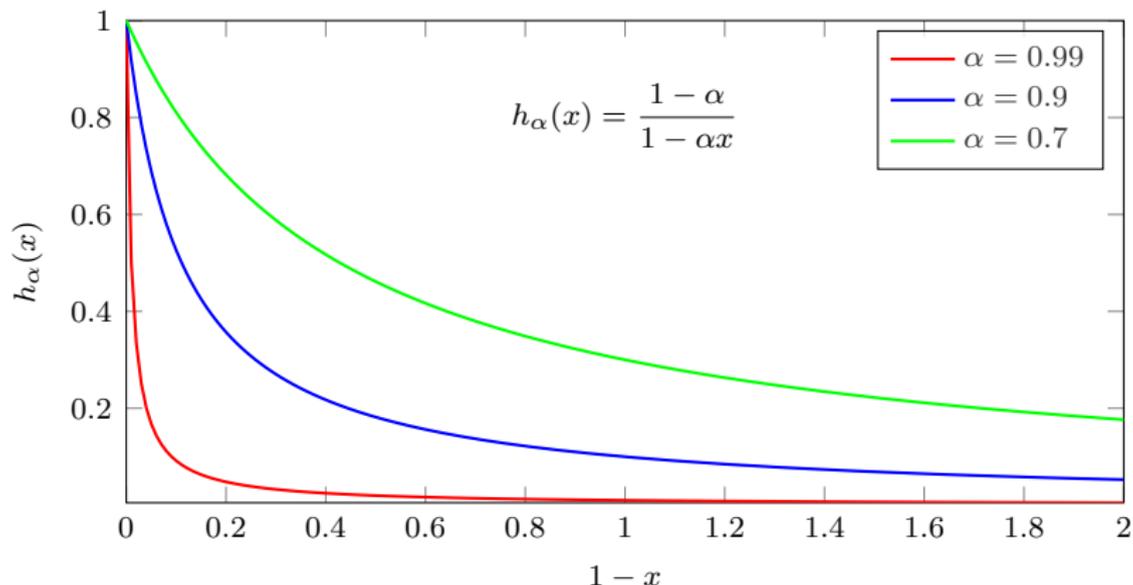
$$H := h(C_n) = \sum_{t=0}^{\infty} h_t C_n^t$$

- matrix C_n has the eigenvalue decomposition $U\Lambda U^\top$ where U^\top is the $n \times n$ **discrete Fourier transform** matrix \mathcal{F}
- if the series $h(C_n)$ converges, filtering $\mathbf{s} \mapsto H\mathbf{s}$ is written as

$$\mathbf{s} \mapsto \mathcal{F}^{-1}h(\Lambda)\mathcal{F}\mathbf{s}$$

- **graph signal processing** generalizes the above by replacing C_n with a matrix determined by a graph

interpretation: graph signal processing



- low-pass filtering in the frequency domain
- or, “soft” dimensionality reduction

interpretation: random fields

- a **Gaussian Markov random field** (GMRF) with precision A and mean $\boldsymbol{\mu}$ can be parametrized as

$$p(\mathbf{x}) := \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, A^{-1}) \propto e^{-E(\mathbf{x}|\mathbf{b}, A)}$$

where $E(\mathbf{x}|\mathbf{b}, A) := \frac{1}{2}\mathbf{x}^\top A\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ is a quadratic **energy**

- its expectation $\boldsymbol{\mu} = A^{-1}\mathbf{b}$ is the minimizer of this energy
- our solution $\mathbf{x}^* = \mathcal{L}_\alpha^{-1}\mathbf{y}$ is the expectation of a GMRF with energy

$$f_\alpha(\mathbf{x}) := E(\mathbf{x}|\mathbf{y}, \mathcal{L}_\alpha) = \frac{1}{2}\mathbf{x}^\top \mathcal{L}_\alpha \mathbf{x} - \mathbf{y}^\top \mathbf{x}$$

- if $\hat{\mathbf{x}} := D^{-1/2}\mathbf{x}$, this energy has the same minimizer as

$$\alpha \sum_{i,j} w_{ij} \|\hat{x}_i - \hat{x}_j\|^2 + (1 - \alpha) \|\mathbf{x} - \mathbf{y}\|^2$$

interpretation: random fields

- a **Gaussian Markov random field** (GMRF) with precision A and mean $\boldsymbol{\mu}$ can be parametrized as

$$p(\mathbf{x}) := \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, A^{-1}) \propto e^{-E(\mathbf{x}|\mathbf{b}, A)}$$

where $E(\mathbf{x}|\mathbf{b}, A) := \frac{1}{2}\mathbf{x}^\top A\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ is a quadratic **energy**

- its expectation $\boldsymbol{\mu} = A^{-1}\mathbf{b}$ is the minimizer of this energy
- our solution $\mathbf{x}^* = \mathcal{L}_\alpha^{-1}\mathbf{y}$ is the expectation of a GMRF with energy

$$f_\alpha(\mathbf{x}) := E(\mathbf{x}|\mathbf{y}, \mathcal{L}_\alpha) = \frac{1}{2}\mathbf{x}^\top \mathcal{L}_\alpha \mathbf{x} - \mathbf{y}^\top \mathbf{x}$$

- if $\hat{\mathbf{x}} := D^{-1/2}\mathbf{x}$, this energy has the same minimizer as

$$\alpha \sum_{i,j} w_{ij} \|\hat{x}_i - \hat{x}_j\|^2 + (1 - \alpha) \|\mathbf{x} - \mathbf{y}\|^2$$

interpretation: random fields

- a **Gaussian Markov random field** (GMRF) with precision A and mean $\boldsymbol{\mu}$ can be parametrized as

$$p(\mathbf{x}) := \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, A^{-1}) \propto e^{-E(\mathbf{x}|\mathbf{b}, A)}$$

where $E(\mathbf{x}|\mathbf{b}, A) := \frac{1}{2}\mathbf{x}^\top A\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ is a quadratic **energy**

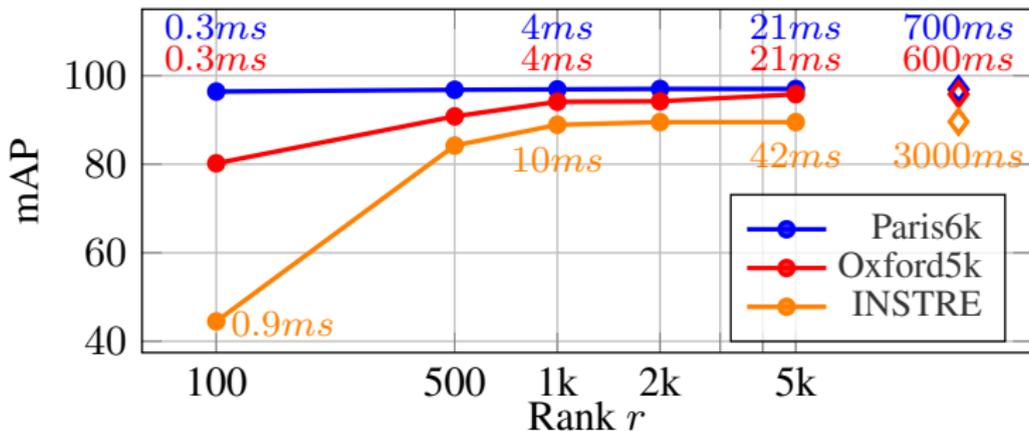
- its expectation $\boldsymbol{\mu} = A^{-1}\mathbf{b}$ is the minimizer of this energy
- our solution $\mathbf{x}^* = \mathcal{L}_\alpha^{-1}\mathbf{y}$ is the expectation of a GMRF with energy

$$f_\alpha(\mathbf{x}) := E(\mathbf{x}|\mathbf{y}, \mathcal{L}_\alpha) = \frac{1}{2}\mathbf{x}^\top \mathcal{L}_\alpha \mathbf{x} - \mathbf{y}^\top \mathbf{x}$$

- if $\hat{\mathbf{x}} := D^{-1/2}\mathbf{x}$, this energy has the same minimizer as

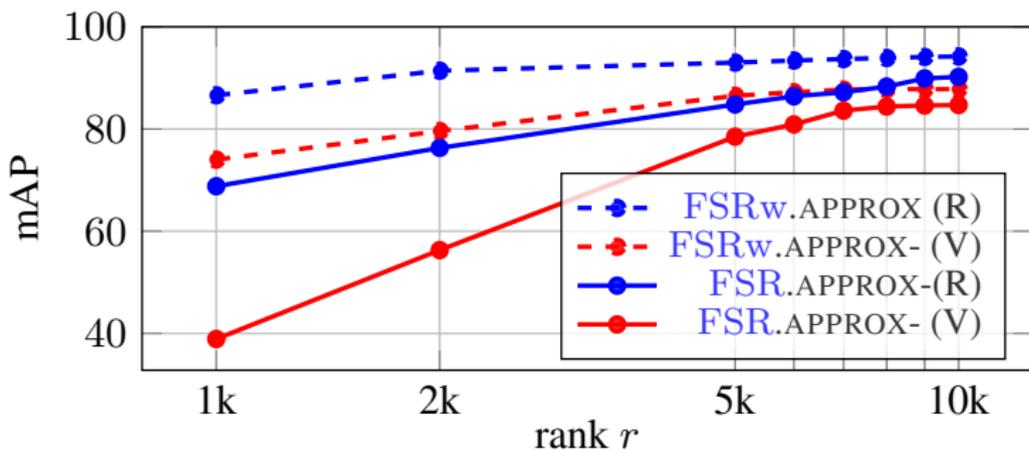
$$\alpha \sum_{i,j} w_{ij} \|\hat{x}_i - \hat{x}_j\|^2 + (1 - \alpha) \|\mathbf{x} - \mathbf{y}\|^2$$

small scale (21 descriptors/image)



- in summary: same performance as CG, two orders of magnitude faster, but $3\times$ space needed

Oxford105k (5 descriptors/image)



- in summary: same performance as CG, two orders of magnitude faster, but $3\times$ space needed

hard examples?



(AP: 92.1)



#5



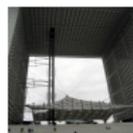
#32



#51



#70



#71



#76



#79



#126



(AP: 92.7)



#2



#4



#8



#61



#68



#72



#75



#108

- red: drift
- blue: incorrect annotations

mining on manifolds

embeddings for manifold similarity

- the $n \times n$ kernel matrix

$$K = h(\mathcal{W}) = Uh(\Lambda)U^\top$$

expresses the pairwise manifold similarity of all data points

- if $h(x) > 0$ for $x \in \mathbb{R}$, which holds for h_α , then K is positive-definite and there is an $n \times n$ matrix Φ such that $K = \Phi^\top \Phi$
- a particular choice is

$$\Phi = h(\Lambda)^{1/2}U^\top$$

- if we choose a rank- r approximation instead, then Φ is $r \times n$ and defines a low-dimensional embedding onto \mathbb{R}^r
- so why not learn an embedding on a training set such that it generalizes manifold similarity to other data sets?

embeddings for manifold similarity

- the $n \times n$ kernel matrix

$$K = h(\mathcal{W}) = Uh(\Lambda)U^\top$$

expresses the pairwise manifold similarity of all data points

- if $h(x) > 0$ for $x \in \mathbb{R}$, which holds for h_α , then K is positive-definite and there is an $n \times n$ matrix Φ such that $K = \Phi^\top \Phi$
- a particular choice is

$$\Phi = h(\Lambda)^{1/2}U^\top$$

- if we choose a rank- r approximation instead, then Φ is $r \times n$ and defines a low-dimensional embedding onto \mathbb{R}^r
- so why not learn an embedding on a training set such that it generalizes manifold similarity to other data sets?

embeddings for manifold similarity

- the $n \times n$ kernel matrix

$$K = h(\mathcal{W}) = Uh(\Lambda)U^\top$$

expresses the pairwise manifold similarity of all data points

- if $h(x) > 0$ for $x \in \mathbb{R}$, which holds for h_α , then K is positive-definite and there is an $n \times n$ matrix Φ such that $K = \Phi^\top \Phi$
- a particular choice is

$$\Phi = h(\Lambda)^{1/2}U^\top$$

- if we choose a rank- r approximation instead, then Φ is $r \times n$ and defines a low-dimensional embedding onto \mathbb{R}^r
- so why not learn an embedding on a training set such that it generalizes manifold similarity to other data sets?

embeddings for manifold similarity

- the $n \times n$ kernel matrix

$$K = h(\mathcal{W}) = Uh(\Lambda)U^\top$$

expresses the pairwise manifold similarity of all data points

- if $h(x) > 0$ for $x \in \mathbb{R}$, which holds for h_α , then K is positive-definite and there is an $n \times n$ matrix Φ such that $K = \Phi^\top \Phi$
- a particular choice is

$$\Phi = h(\Lambda)^{1/2}U^\top$$

- if we choose a rank- r approximation instead, then Φ is $r \times n$ and defines a low-dimensional embedding onto \mathbb{R}^r
- so why not learn an embedding on a training set such that it generalizes manifold similarity to other data sets?

mining on manifolds

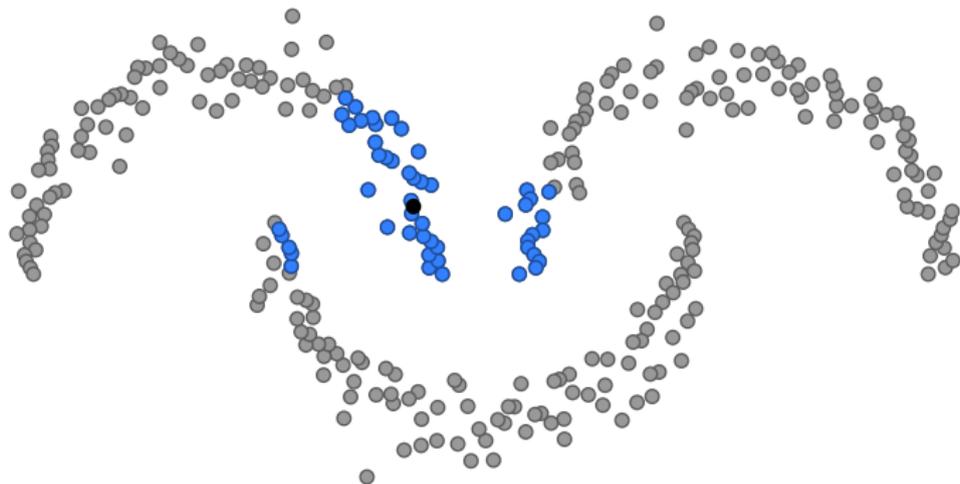
[Iscen et al. 2018]



- data points (\circ), query point x (\bullet)

mining on manifolds

[Iscen et al. 2018]



- data points (\circ), query point x (\bullet)
- Euclidean nearest neighbors $E(x)$ (\bullet)

mining on manifolds

[Iscen et al. 2018]



- data points (\bullet), query point x (\bullet)
- manifold nearest neighbors $M(x)$ (\bullet)

mining on manifolds

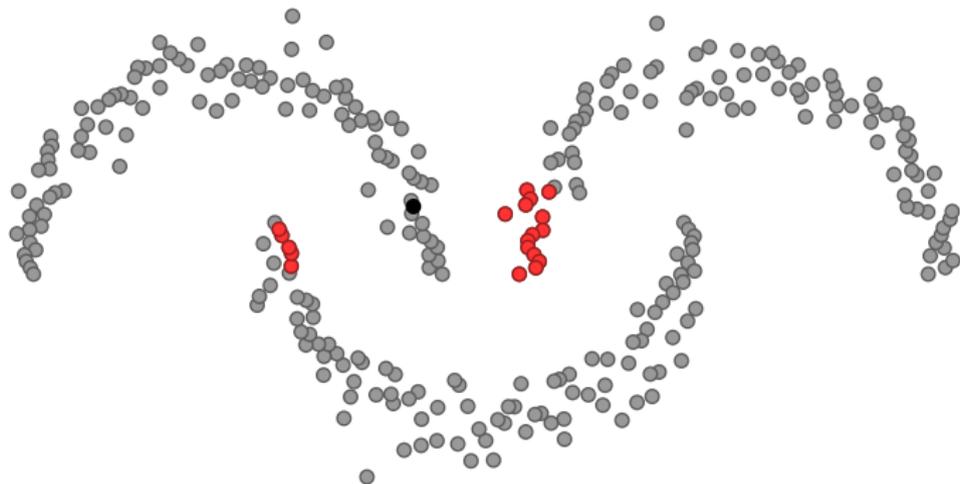
[Iscen et al. 2018]



- data points (\circ), query point \mathbf{x} (\bullet)
- **hard positives** $S^+ = M(\mathbf{x}) \setminus E(\mathbf{x})$ (\bullet)

mining on manifolds

[Iscen et al. 2018]



- data points (\circ), query point \mathbf{x} (\bullet)
- **hard negatives** $S^- = E(\mathbf{x}) \setminus M(\mathbf{x})$ (\circ)

hard positive/negative examples



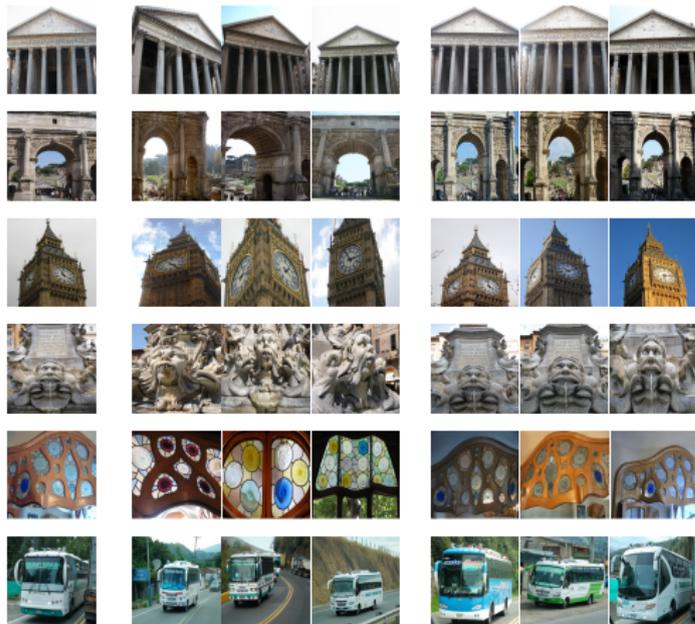
- query (anchor) (\mathbf{x})
- positives $S^+(\mathbf{x})$ vs. Euclidean neighbors $E(\mathbf{x})$
- negatives $S^-(\mathbf{x})$ vs. Euclidean non-neighbors $X \setminus E(\mathbf{x})$

hard positive/negative examples



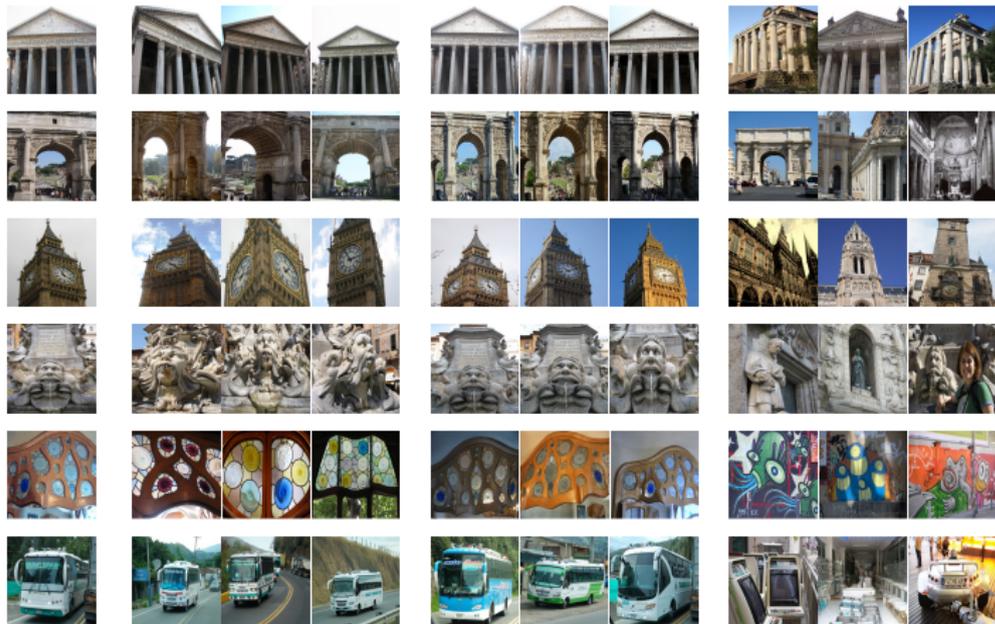
- query (anchor) (\mathbf{x})
- positives $S^+(\mathbf{x})$ vs. Euclidean neighbors $E(\mathbf{x})$
- negatives $S^-(\mathbf{x})$ vs. Euclidean non-neighbors $X \setminus E(\mathbf{x})$

hard positive/negative examples



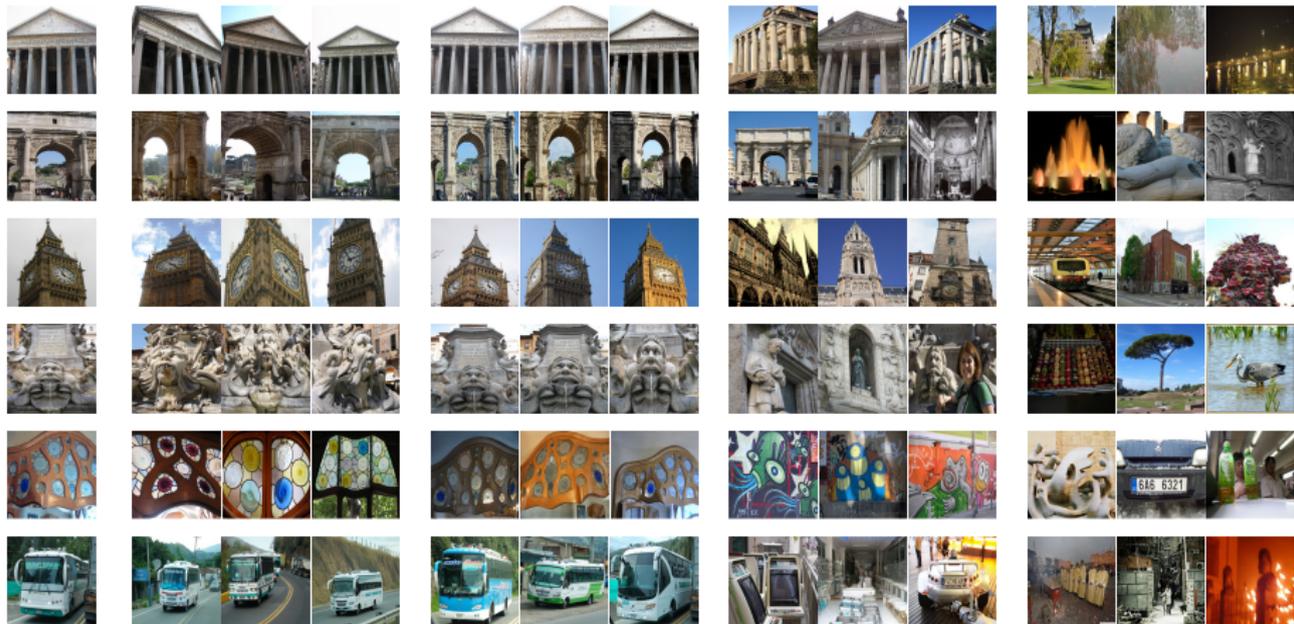
- query (anchor) (\mathbf{x})
- positives $S^+(\mathbf{x})$ vs. Euclidean neighbors $E(\mathbf{x})$
- negatives $S^-(\mathbf{x})$ vs. Euclidean non-neighbors $X \setminus E(\mathbf{x})$

hard positive/negative examples



- query (anchor) (\mathbf{x})
- positives $S^+(\mathbf{x})$ vs. Euclidean neighbors $E(\mathbf{x})$
- negatives $S^-(\mathbf{x})$ vs. Euclidean non-neighbors $X \setminus E(\mathbf{x})$

hard positive/negative examples



- query (anchor) (\mathbf{x})
- positives $S^+(\mathbf{x})$ vs. Euclidean neighbors $E(\mathbf{x})$
- negatives $S^-(\mathbf{x})$ vs. Euclidean non-neighbors $X \setminus E(\mathbf{x})$

fine-tuning with hard example mining

- pre-train network
- extract descriptors on **unlabeled** dataset
- construct nearest neighbor graph
- sample **anchors**, measure Euclidean and manifold distances
- sample **positives** and **negatives**
- fine-tune using **contrastive** or **triplet** loss

fine-tuning with hard example mining

- pre-train network
- extract descriptors on **unlabeled** dataset
- construct nearest neighbor graph
- sample **anchors**, measure Euclidean and manifold distances
- sample **positives** and **negatives**
- fine-tune using **contrastive** or **triplet** loss

fine-grained categorization results

Method	Labels	R@1	R@2	R@4	R@8	NMI
Initial	No	35.0	46.8	59.3	72.0	48.1
Triplet+semi-hard	Yes	42.3	55.0	66.4	77.2	55.4
Lifted-Structure	Yes	43.6	56.6	68.6	79.6	56.5
Triplet+	Yes	45.9	57.7	69.6	79.8	58.1
Clustering	Yes	48.2	61.4	71.8	81.9	59.2
Triplet+++	Yes	49.8	62.3	74.1	83.3	59.9
Cyclic match	No	40.8	52.8	65.1	76.0	52.6
Ours	No	45.3	57.8	68.6	78.4	55.0

- CUB200-2011 dataset, 200 bird species, 100 training / 100 testing
- GoogLeNet pre-trained on ImageNet, then fine-tuned with triplet loss

particular object retrieval results

Model	Pooling	Labels	Oxf5k	Oxf105k	Par6k	Par106k	Hol	Instre
ImageNet	MAC	Human	58.5	50.3	73.0	59.0	79.4	48.5
From BoW		SfM	79.7	73.9	82.4	74.6	81.4	48.5
Ours		—	78.7	74.3	83.1	75.6	82.6	55.5
ImageNet	R-MAC	Human	68.0	61.0	76.6	72.1	87.0	55.6
From BoW		SfM	77.8	70.1	84.1	76.8	84.4	47.7
Ours		—	78.2	72.6	85.1	78.0	87.5	57.7

- VGG-16 pre-trained on ImageNet, then fine-tuned with constrastive loss on a 1M unlabeled dataset with MAC representation
- at test time, either MAC or R-MAC used

summary

- **pooling** CNN representations is best at last convolutional layers
- **fine-tuning** with **contrastive** or **triplet** loss allows **transferring** to a new domain and learning to **rank** as opposed to classify
- now that images are represented by a global descriptor or just a few regional descriptors, **graph methods** are more applicable than ever
- it turns out that query expansion is not just “post processing” but at the core of **ranking on manifolds**
- there is at least one low-dimensional **embedding** of manifold similarity, but is dataset-specific
- modeling the manifold explicitly allows **unsupervised fine-tuning** without labels, auxiliary systems (e.g. SIFT pipeline), or other information (e.g. temporal neighborhood in video)

summary

- **pooling** CNN representations is best at last convolutional layers
- **fine-tuning** with **contrastive** or **triplet** loss allows **transferring** to a new domain and learning to **rank** as opposed to classify
- now that images are represented by a global descriptor or just a few regional descriptors, **graph methods** are more applicable than ever
- it turns out that query expansion is not just “post processing” but at the core of **ranking on manifolds**
- there is at least one low-dimensional **embedding** of manifold similarity, but is dataset-specific
- modeling the manifold explicitly allows **unsupervised fine-tuning** without labels, auxiliary systems (e.g. SIFT pipeline), or other information (e.g. temporal neighborhood in video)

summary

- **pooling** CNN representations is best at last convolutional layers
- **fine-tuning** with **contrastive** or **triplet** loss allows **transferring** to a new domain and learning to **rank** as opposed to classify
- now that images are represented by a global descriptor or just a few regional descriptors, **graph methods** are more applicable than ever
- it turns out that query expansion is not just “post processing” but at the core of **ranking on manifolds**
- there is at least one low-dimensional **embedding** of manifold similarity, but is dataset-specific
- modeling the manifold explicitly allows **unsupervised fine-tuning** without labels, auxiliary systems (e.g. SIFT pipeline), or other information (e.g. temporal neighborhood in video)