

# VISUAL OBJECT RECOGNITION

STATE-OF-THE-ART  
TECHNIQUES AND  
PERFORMANCE EVALUATION

# LECTURE 5: METRICS FOR MATCHING

- ✱ Descriptor distances
  - ✱ Chi<sup>2</sup> distance
  - ✱ Earth Mover's Distance (EMD)
- ✱ Ratio Score
- ✱ Visual Words
- ✱ Learning the metric

# DESCRIPTOR DISTANCES

✱ For a descriptor  $\mathbf{q}$  in a query image. Which prototype in memory  $(\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N)$  is *most likely* to correspond to the same world object.

✱ Assuming additive i.i.d. Gaussian noise on all elements:

$$p(\mathbf{q}|\mathbf{p}_k) \propto \prod_{l=1}^D e^{-.5(p_{kl} - q_l)^2 / \sigma^2}$$

$$\max(p) \Leftrightarrow \min(-\log(p))$$

$$-\log(p(\mathbf{q}|\mathbf{p}_k)) \propto \sum_{l=1}^D (p_{kl} - q_l)^2$$

# DESCRIPTOR DISTANCES

- ✱ So, the match with smallest distance is most likely correct, assuming i.i.d. Gaussian noise.

- ✱ What about the scalar product for normalised vectors/NCC?

$$\|\mathbf{p} - \mathbf{q}\|^2 = \mathbf{p}^T \mathbf{p} + \mathbf{q}^T \mathbf{q} - 2\mathbf{p}^T \mathbf{q} = 2(1 - \mathbf{p}^T \mathbf{q})$$

- ✱ But are all values identically distributed?

- ✱ ...are they independent?

# CHI<sup>2</sup> DISTANCE

- ✻ Many descriptors are histogram-like in their nature.
- ✻ For histograms, the histogram values typically follow the (discrete) *Poisson distribution*:

$$P(k|\mu) = \mu^k e^{-\mu} / k!$$

- ✻ With the statistics:

$$E[P(k)] = \mu \quad E[(P(k) - \mu)^2] = \mu$$

# CHI<sup>2</sup> DISTANCE

- ✱ For large values of  $\mu$ , (e.g. 1000) a (continuous) Gaussian can approximate the Poisson distribution:

$$p(k|\mu) \approx \frac{1}{\mu\sqrt{2\pi}} e^{-.5(k-\mu)^2/\mu}$$

- ✱ Again, assuming independence, this leads to a negative log likelihood proportional to:

$$-\log(p(\mathbf{q}|\mathbf{p}_k)) \propto \sum_{l=1}^D (p_{kl} - q_l)^2 / \mu_l$$

# CHI<sup>2</sup> DISTANCE

- ✱ If we estimate the variance by:

$$\mu_l \approx (p_{kl} + q_l)/2$$

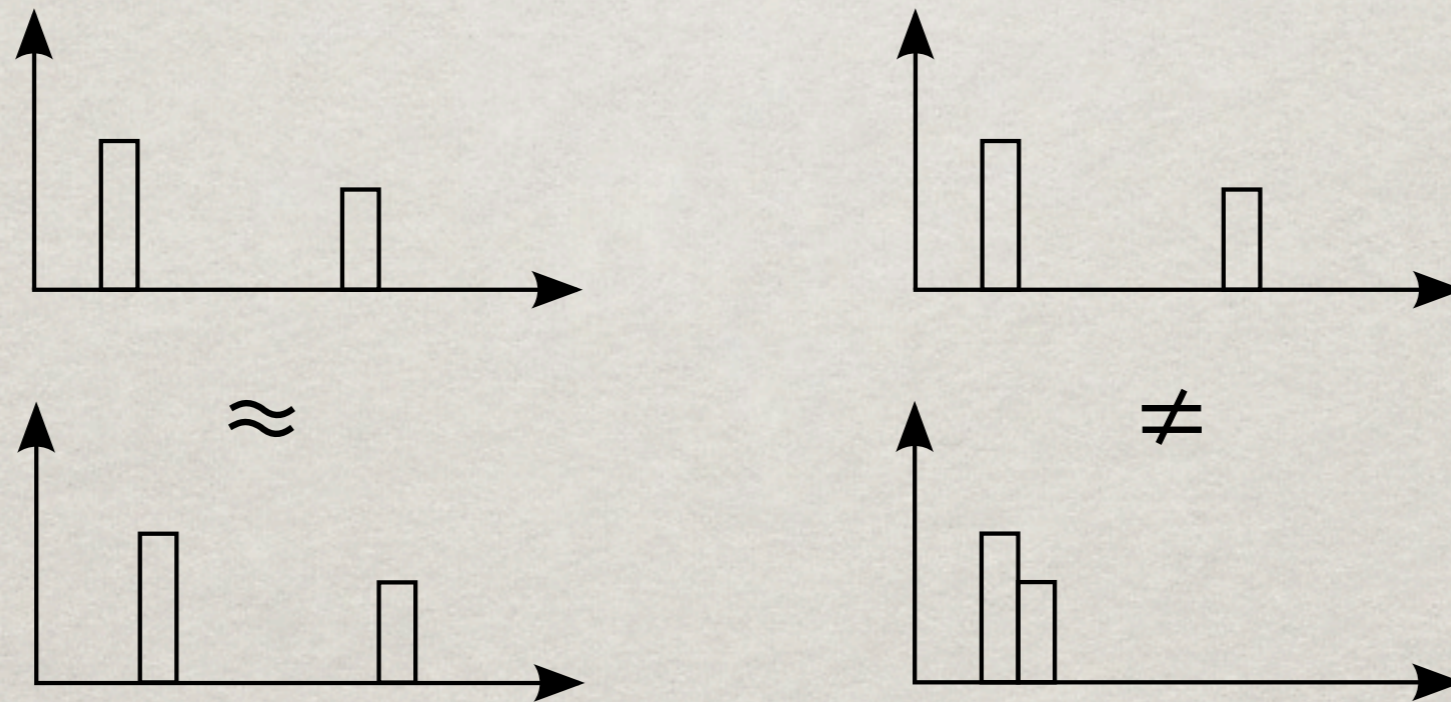
- ✱ We find that the most likely match is the one with the smallest Chi-squared distance:

$$\chi^2(\mathbf{q}, \mathbf{p}_k) = \sum_{l=1}^D \frac{(p_{kl} - q_l)^2}{p_{kl} + q_l}$$

- ✱ This is assuming independence between bins.

# EARTH MOVER'S DISTANCE

- ✻ In histograms, neighbouring bins are typically correlated



- ✻ Instead of falling in bin  $i$ , a sample is likely to fall in bin  $i+1$ .



# EARTH MOVER'S DISTANCE

- ✻ Distance=cost of moving values in  $p$  to  $q$   
cost=amount\*distance
- ✻ First solve a linear programming problem:  
*the transportation problem*, Hitchcock 1941.

$$\min_{f_{ij}} \sum_{i=1}^D \sum_{j=1}^D f_{ij} d_{ij} \quad \text{here} \quad d_{ij} = |i - j|$$

- ✻  $f_{ij}$  amount to move from  $i$  to  $j$ .

# EARTH MOVER'S DISTANCE

✻ Transportation problem, cost function:

$$\min_{f_{ij}} \sum_{i=1}^D \sum_{j=1}^D f_{ij} d_{ij} \quad d_{ij} = |i - j|$$

✻ Constraints:

$$f_{ij} \geq 0 \quad \forall i, j \in [1, D]$$

$$\sum_{i=1}^D f_{ij} = q_j \quad \forall j \in [1, D]$$

$$\sum_{i=1}^D f_{ij} \leq p_j \quad \forall j \in [1, D]$$

# EARTH MOVER'S DISTANCE

- ✻ Now compute EMD as:

$$d(\mathbf{p}, \mathbf{q}) = \min_{f_{ij}} \frac{\sum_{i=1}^D \sum_{j=1}^D f_{ij} |i - j|}{\sum_{i=1}^D \sum_{j=1}^D f_{ij}}$$

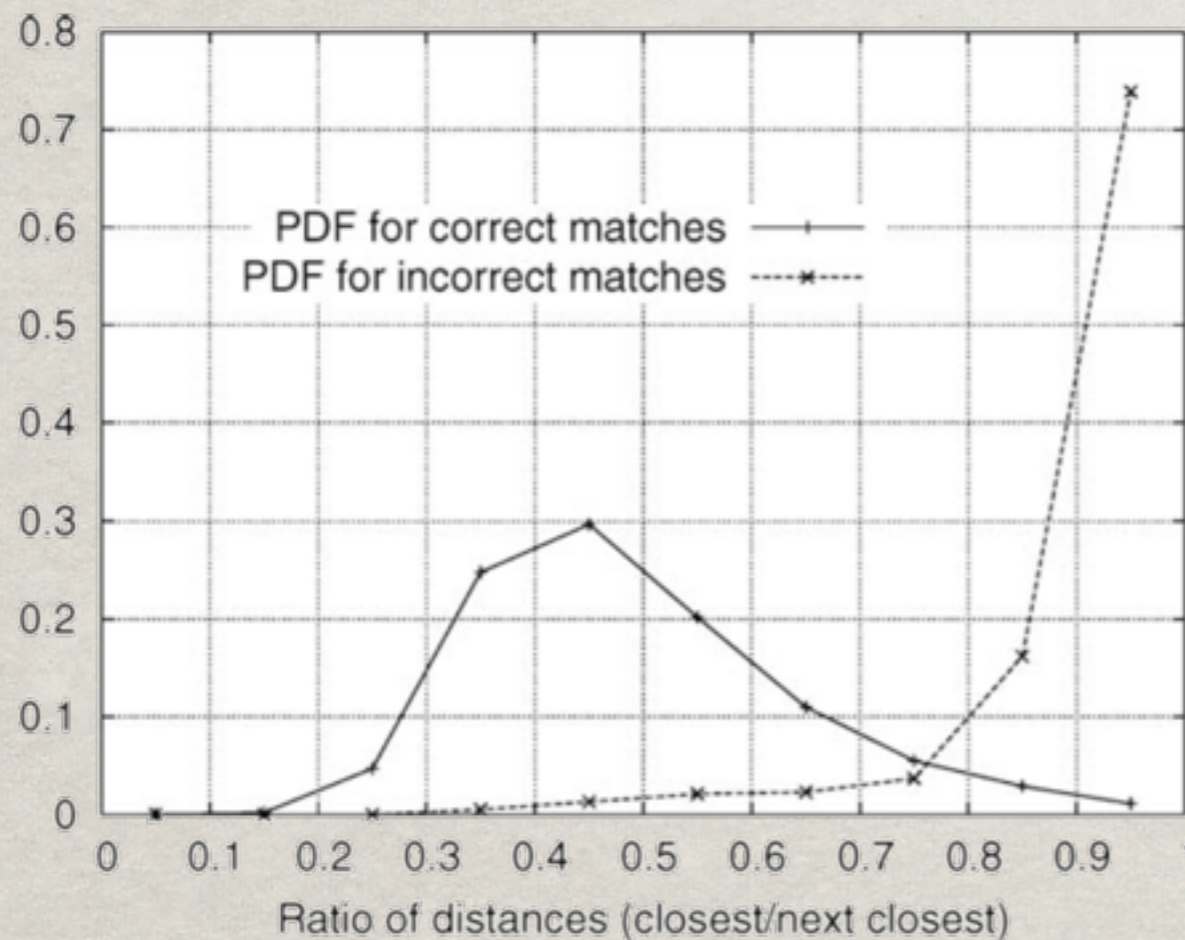
- ✻ The denominator is needed if histograms are computed from different numbers of samples.
- ✻ Introduced in Computer Vision by Rubner&Tomasi&Guibas, at ICCV98
- ✻ Local expert: Thomas Kaijser

# RATIO SCORE

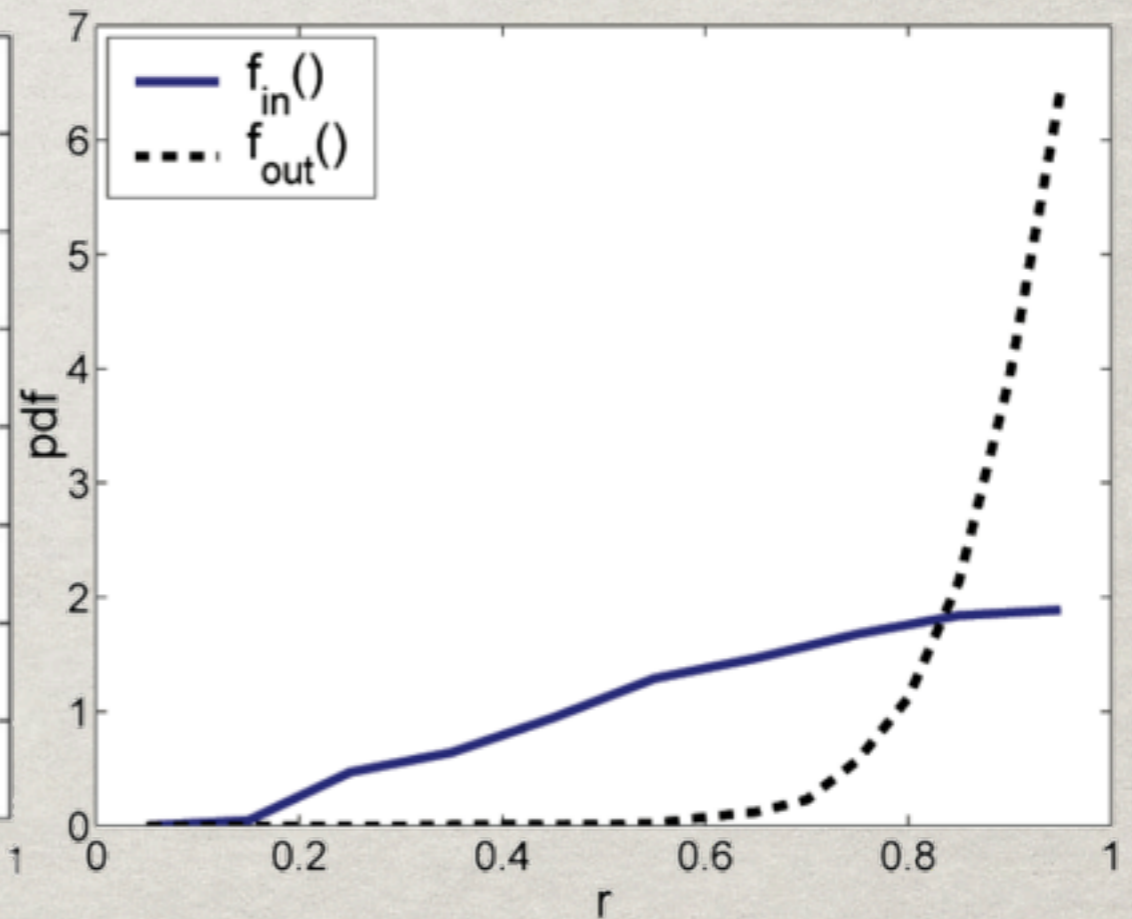
- ✱ If we have best matches for descriptors  $q_1$  and  $q_2$  in the image. Which one is better?
- ✱ Some features are *more common than others*, and by scoring the match for  $q_1$ , according to the ratio between the best, and the second best match we can compensate for this:

$$r = d_1 / d_2$$

# RATIO SCORE



Lowe IJCV04



Goshen&Shimshoni PAMI08

# BAGS-OF-FEATURES AND VISUAL WORDS

- ✻ A common technique for matching in large datasets. Sivic&Zisserman “Video Google” ICCV03.

**Object**

**Bag of ‘words’**



Illustration by Li Fei-Fei, <http://people.csail.mit.edu/torralba/shortCourseRLOC/>

- ✻ Completely disregards spatial relationships among features.

# VISUAL WORDS

- ✱ Vector quantize feature space to into  $K$  parts using e.g.  $K$ -means clustering (e.g. function `kmeans` in *Matlab*) on large training set.

- ✱ Clustering is done in whitened space:

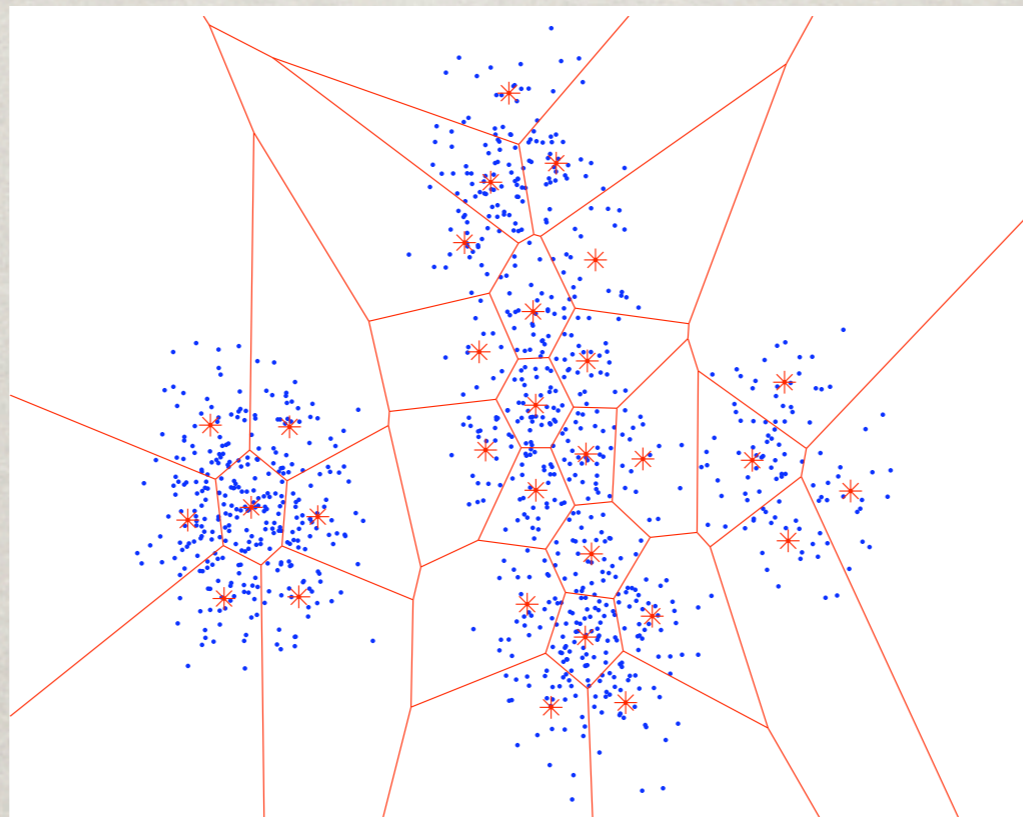
$$\hat{\mathbf{x}} = \mathbf{C}^{-1/2}(\mathbf{x} - \mu)$$

- ✱ Each descriptor is then approximated by the most similar prototype/visual word.

# VISUAL WORDS

- ✱ K-means finds a local min of the following objective function ( $\mathbf{x}$ -samples,  $\mathbf{p}$ -prototypes):

$$J = \sum_{n=1}^N \min_{k \in [1 \dots K]} \|\mathbf{x}_n - \mathbf{p}_k\|^2$$

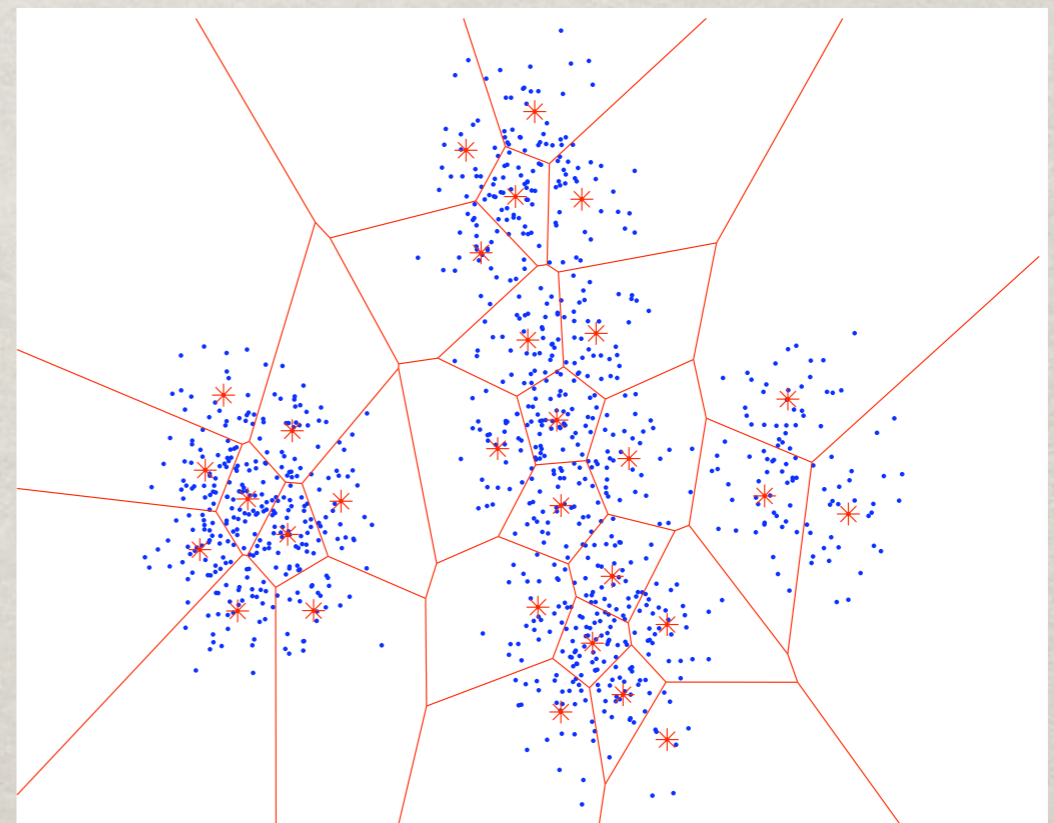
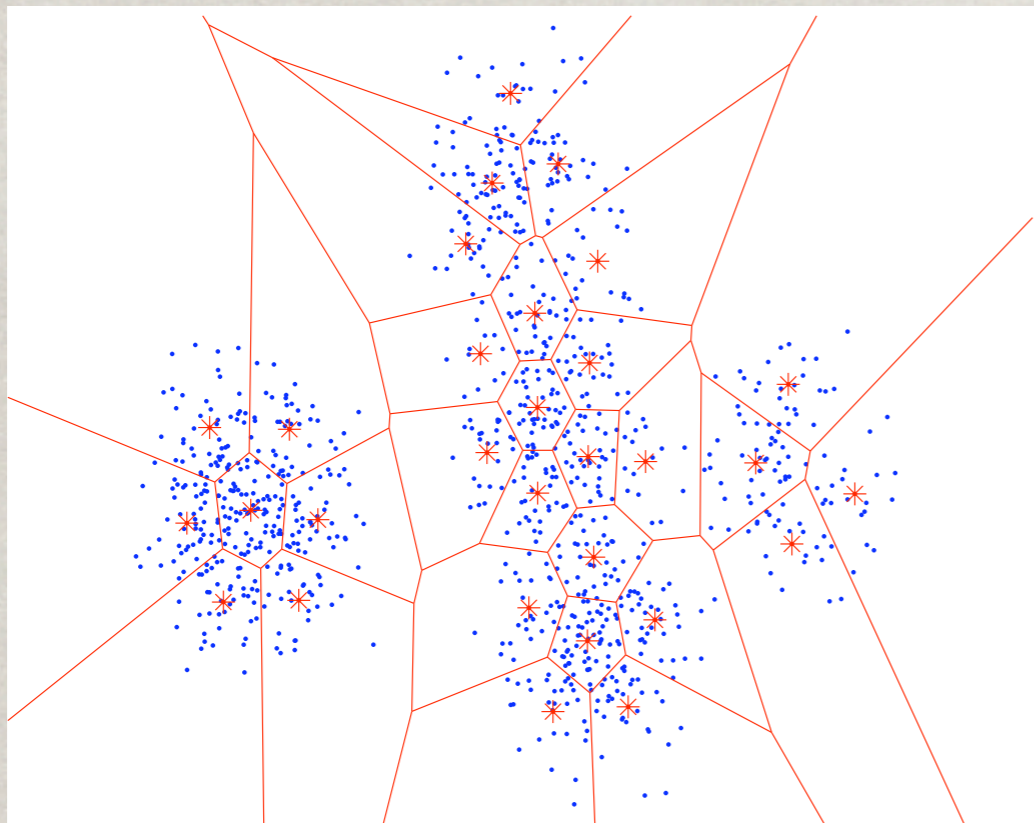


(C) 2008 PER-ERIK FORSSÉN



# VISUAL WORDS

- ✱ Probability of visual words is somewhat equalized. cf. ratio score and histogram eq.
- ✱ K-means is not perfectly repeatable.  
(Try several times and pick highest J.)



(C) 2008 PER-ERIK FORSSÉN

# VISUAL WORDS

- ☼ Analogy with text document matching.

- ☼ Each document (i.e. image) is represented as a vector of (TF-IDF) word frequencies

$$\mathbf{v}_d = (v_1 \quad \dots \quad v_K)^T \quad v_k = \frac{N_{kd}}{N_d} \log \frac{N}{N_k}$$

- ☼ term frequency:  $N_{kd}/N_d$  (word  $k$ , document  $d$ )  
Nistér&Stewenius CVPR06: skip  $N_d$ .

- ☼ inverse document frequency:  $N/N_k$  - inverse frequency of word  $k$  in whole database.

# VISUAL WORDS

- ✻ Image matching is done by a normalised scalar product:

$$\hat{\mathbf{v}}_q^T \hat{\mathbf{v}}_p = \cos \phi$$

- ✻ An *inverted file* makes real-time matching possible on very large datasets:

word1: frame 3, frame 17, frame 243...

word2: frame 2, frame 23, frame 33...

...

# BAG-OF-FEATURES

- ✱ Instead of the TF-IDF vector in the visual words method, one could simply compute a histogram of visual word occurrences.
- ✱ This is called a *bag-of-features*, or *bag-of-keypoints*
- ✱ With IDF as metric this is equivalent to the TF-IDF vector (when  $TF = N_{kd}$ ).

# BAG-OF-FEATURES

- ✱ The bag-of-features vector is often fed into a machine learning algorithm. (E.g. in today's article).
- ✱ Typically  $K$  is large and most values are zero.

Csurka et al.  $K=1000$

Sivic&Zisserman  $K=6000$  and  $10,000$

Nistér&Stewenius  $K=16e6$

# LEARNING THE METRIC

- ✱ What we ultimately want is to distinguish good feature matches from bad.
- ✱ Collect known corresponding descriptors:  
 $\{(\mathbf{p}_k, \mathbf{q}_k)\}_1^K$  and set  $\mathbf{d}_k = \mathbf{p}_k - \mathbf{q}_k$
- ✱ We now want to find a linear transformation that makes the noise equal in magnitude in all directions:

$$\mathbf{y}_k = \mathbf{T}\mathbf{p}_k \quad \text{assuming} \quad \mathbf{d}_k \sim \mathcal{N}(0, \mathbf{C})$$

# LEARNING THE METRIC

- ✱ Find a whitening transform  $\mathbf{T}$  from the covariance matrix:

$$\mathbf{C} = \frac{1}{K} \sum_{k=1}^K \mathbf{d}_k \mathbf{d}_k^T \quad \text{with} \quad \mathbf{T} \mathbf{C} \mathbf{T}^T = \mathbf{I}$$

- ✱ Valid solutions:

$$\mathbf{T} = \mathbf{R} \mathbf{C}^{-1/2} \quad \text{where} \quad \mathbf{R} \mathbf{R}^T = \mathbf{I}$$

- ✱ If we only use the first few dimensions we should choose  $\mathbf{R}$  such that it selects dimensions where we “see things happen”.

# LEARNING THE METRIC

- ✱ Find  $\mathbf{R}$  from PCA of transformed SIFT feature space:

$$\mathbf{C}_b = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^T - \mathbf{m} \mathbf{m}^T \quad \mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$$

$$\mathbf{R} \mathbf{D} \mathbf{R}^T = \mathbf{C}$$

- ✱ Final contraction operator:

$$\mathbf{P} = \bar{\mathbf{I}}_k \mathbf{R} \mathbf{C}^{-1/2}$$

- ✱ Where  $\mathbf{I}_k$  is a  $k * 128$  truncated identity matrix.



# LEARNING THE METRIC

- ✻ This Mahalanobis metric for features was published at ICCV07 by Mikolajczyk & Matas, SIFT 128  $\rightarrow$  40 dim
- ✻ A similar method that only finds a rotation called linear discriminant embedding (LDE) also at ICCV07 by Hua & Brown & Winder, SIFT 128  $\rightarrow$  14/18 dim
- ✻ Besides reducing dimensionality, these techniques also improve matching results.

# LEARNING THE METRIC

✱ Linear Discriminant Embedding(LDE)

✱ Maximise

$$\mathbf{J}(\mathbf{w}) = \frac{\sum_{\text{outlier}(i,j)} \mathbf{w}^T (\mathbf{p}_i - \mathbf{q}_j)^2}{\sum_{\text{inlier}(i,j)} \mathbf{w}^T (\mathbf{p}_i - \mathbf{q}_j)^2}$$

$$\mathbf{J}(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{A} \mathbf{w}}{\mathbf{w}^T \mathbf{B} \mathbf{w}}, \quad \|\mathbf{w}\| = 1$$

✱ Where  $\mathbf{A}$  covariance for outliers and  $\mathbf{B}$  inliers.

# LEARNING THE METRIC

✱  $J(\mathbf{w})$  is maximised by eigenvectors with large eigenvalues in  $\mathbf{B}^{-1}\mathbf{A}$

✱ Eigenvalues of  $\mathbf{B}$  are set to  $\tilde{\lambda}_i = \max(\lambda_i, \lambda_r)$

$$r = \arg \min_n \frac{\sum_{i=n}^N \lambda_i}{\sum_{i=1}^N \lambda_i} \geq \alpha$$

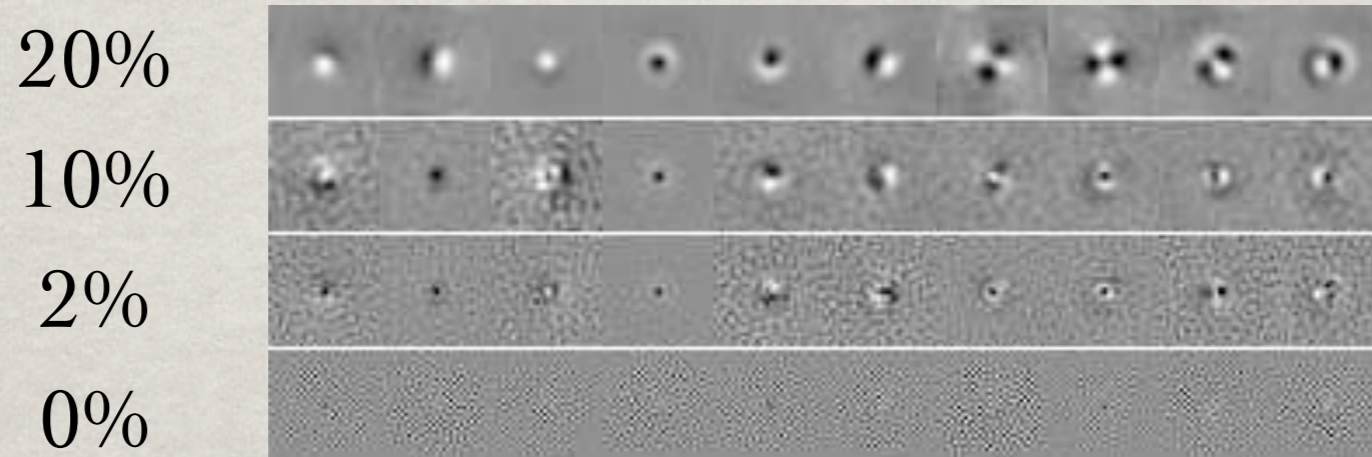
✱  $\alpha$  can be interpreted as a threshold on SNR.  
This is called *Power Regularisation*

✱ Many variations of the algorithm in the paper.

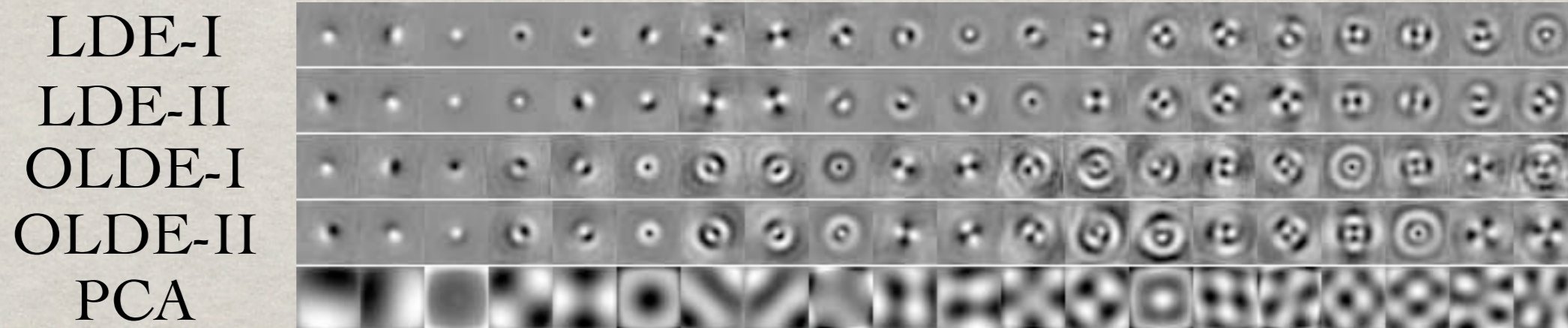
# LEARNING THE METRIC

✱ Some LDE results on grey-scale patches:

✱ Reducing the amount of power reg:



✱ Linear filters found on grey-scale patches:



(C) 2008 PER-ERIK FORSSÉN

# DISCUSSION

☼ Questions/comments on paper and lecture.

(C) 2008 PER-ERIK FORSSÉN