



CS598:VISUAL INFORMATION RETRIEVAL

Lecture IV: Image Representation:
• Feature Coding and Pooling

RECAP OF LECTURE III

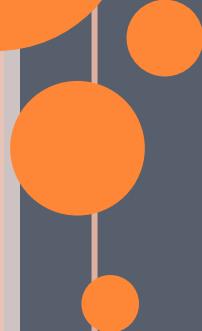
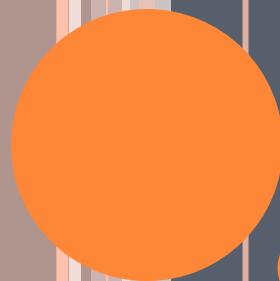
- Blob detection
 - Brief of Gaussian filter
 - Scale selection
 - Lapacian of Gaussian (LoG) detector
 - Difference of Gaussian (DoG) detector
 - Affine co-variant region
- Learning local image descriptors (optional reading)



OUTLINE

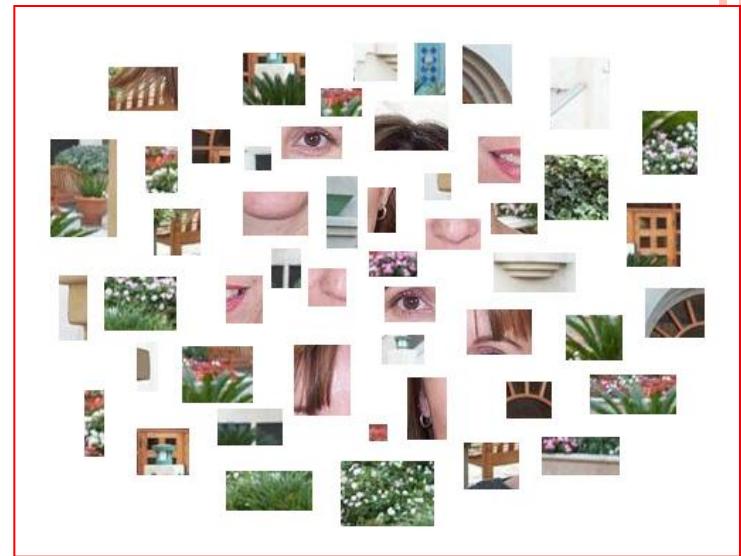
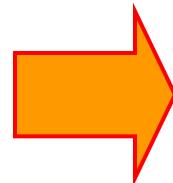
- Histogram of local features
- Bag of words model
- Soft quantization and sparse coding
- Supervector with Gaussian mixture model





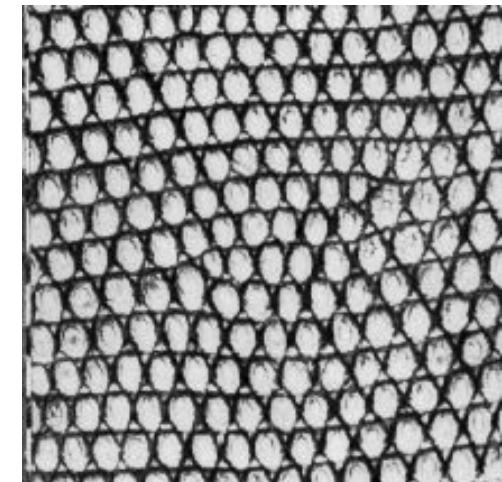
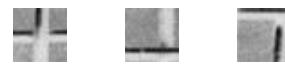
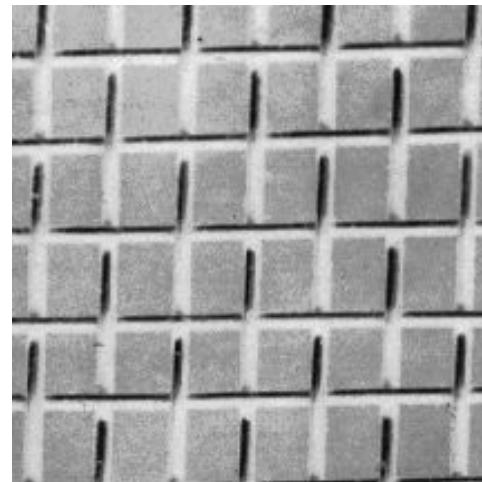
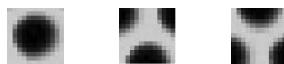
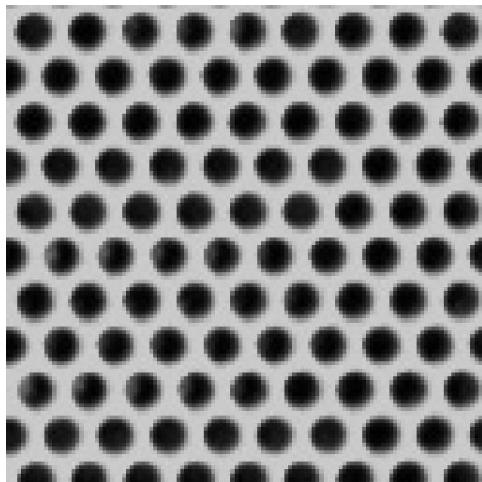
LECTURE IV: PART I

BAG-OF-FEATURES MODELS

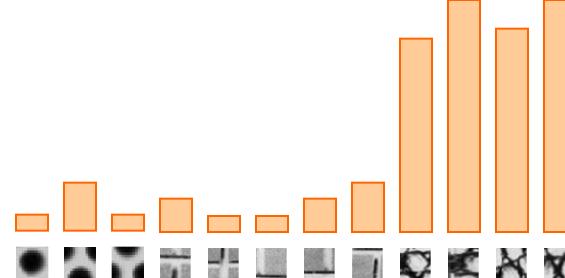
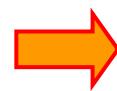
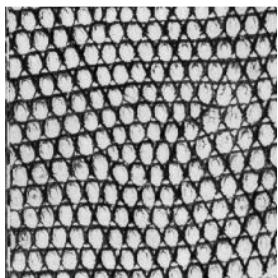
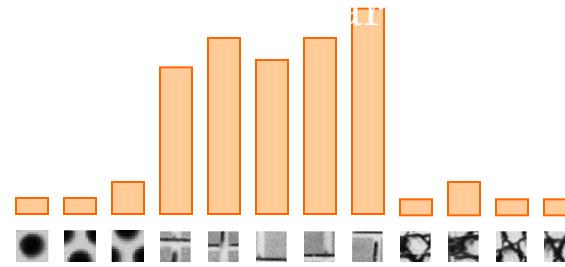
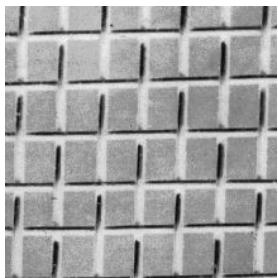
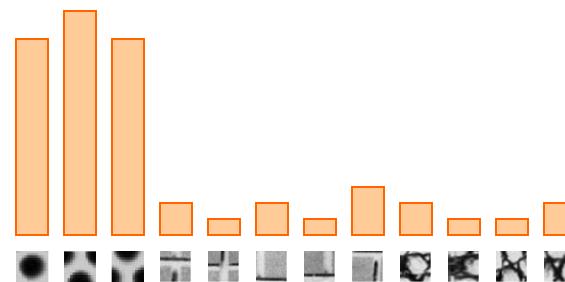
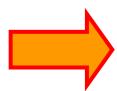
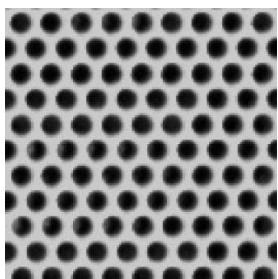


ORIGIN 1: TEXTURE RECOGNITION

- Texture is characterized by the repetition of basic elements or *textons*
- For stochastic textures, it is the identity of the textons, not their spatial arrangement, that matters



ORIGIN 1: TEXTURE RECOGNITION



ORIGIN 2: BAG-OF-WORDS MODELS

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



ORIGIN 2: BAG-OF-WORDS MODELS

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

abandon accountable affordable afghanistan africa aided ally anbar armed army **baghdad** bless challenges chamber chaos choices civilians coalition commanders **commitment** confident confront congressman constitution corps debates deduction deficit deliver **democratic** deploy dikembe diplomacy disruptions earmarks **economy** einstein **elections** eliminates expand **extremists** failing faithful families **freedom** fuel **funding** god haven ideology immigration impose insurgents iran **iraq** islam julie lebanon love madam marine math medicare moderation neighborhoods nuclear offensive palestinian payroll province pursuing **qaeda** radical **regimes** resolve retreat rieman sacrifices science sectarian senate september **shia** stays strength students succeed **sunni** **tax** territories **terrorists** threats uphold victory violence violent **war** washington weapons wesley

ORIGIN 2: BAG-OF-WORDS MODELS

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



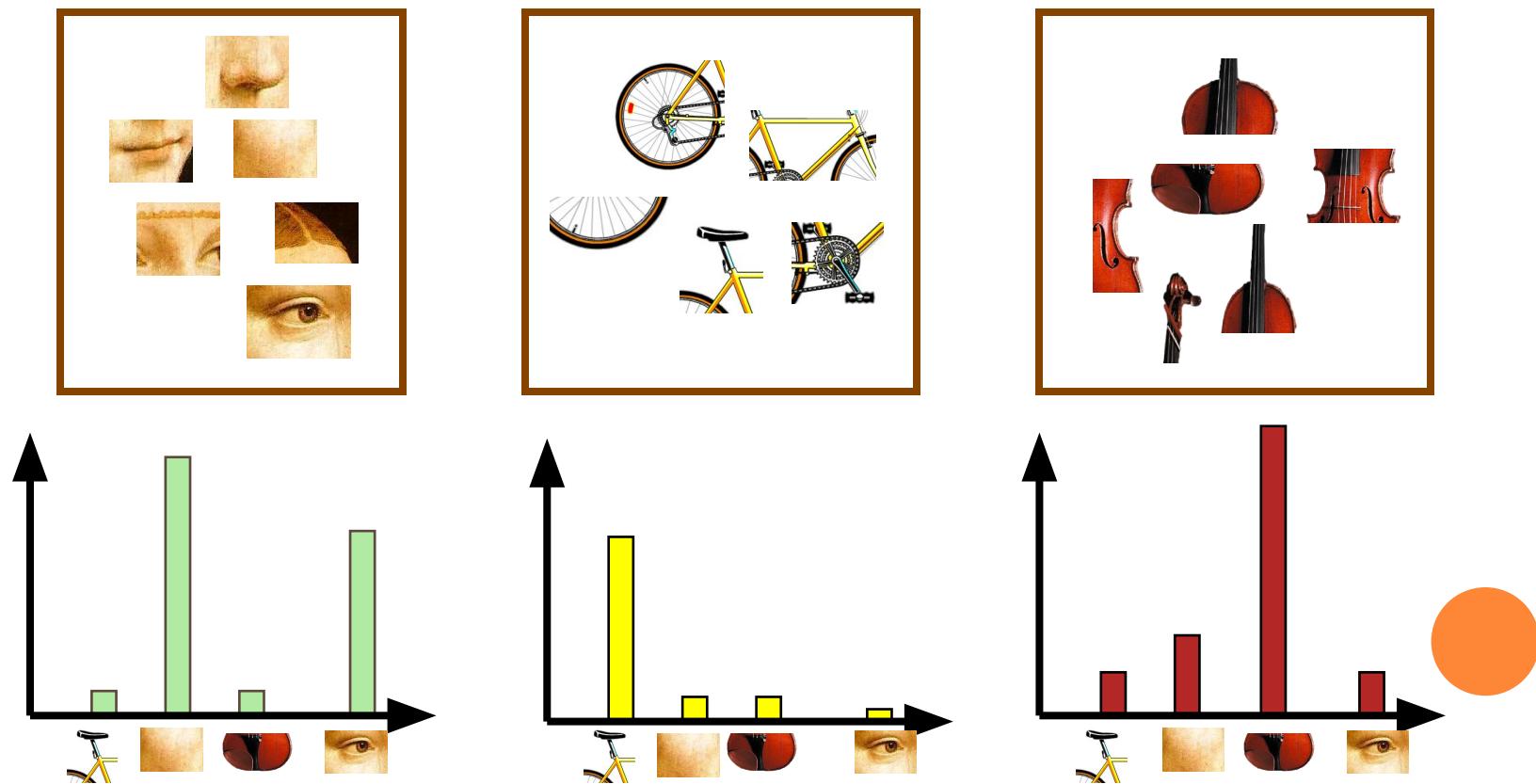
ORIGIN 2: BAG-OF-WORDS MODELS

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



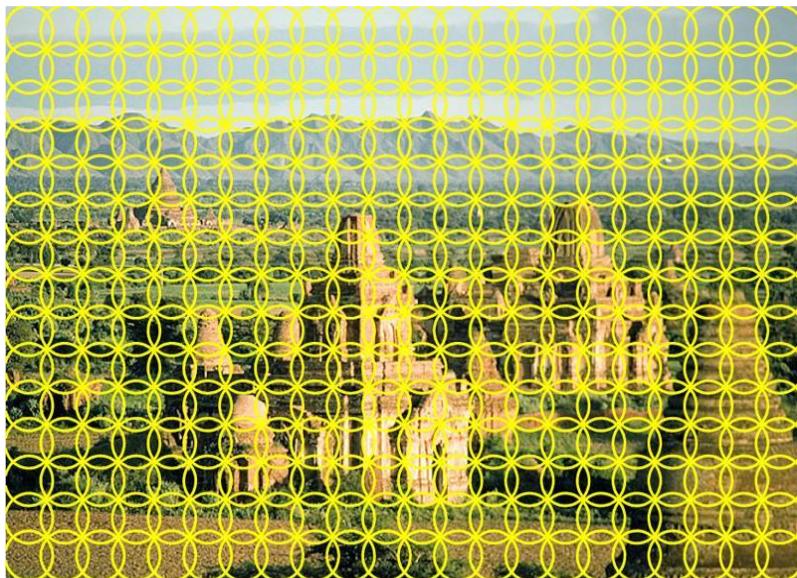
BAG-OF-FEATURES STEPS

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”

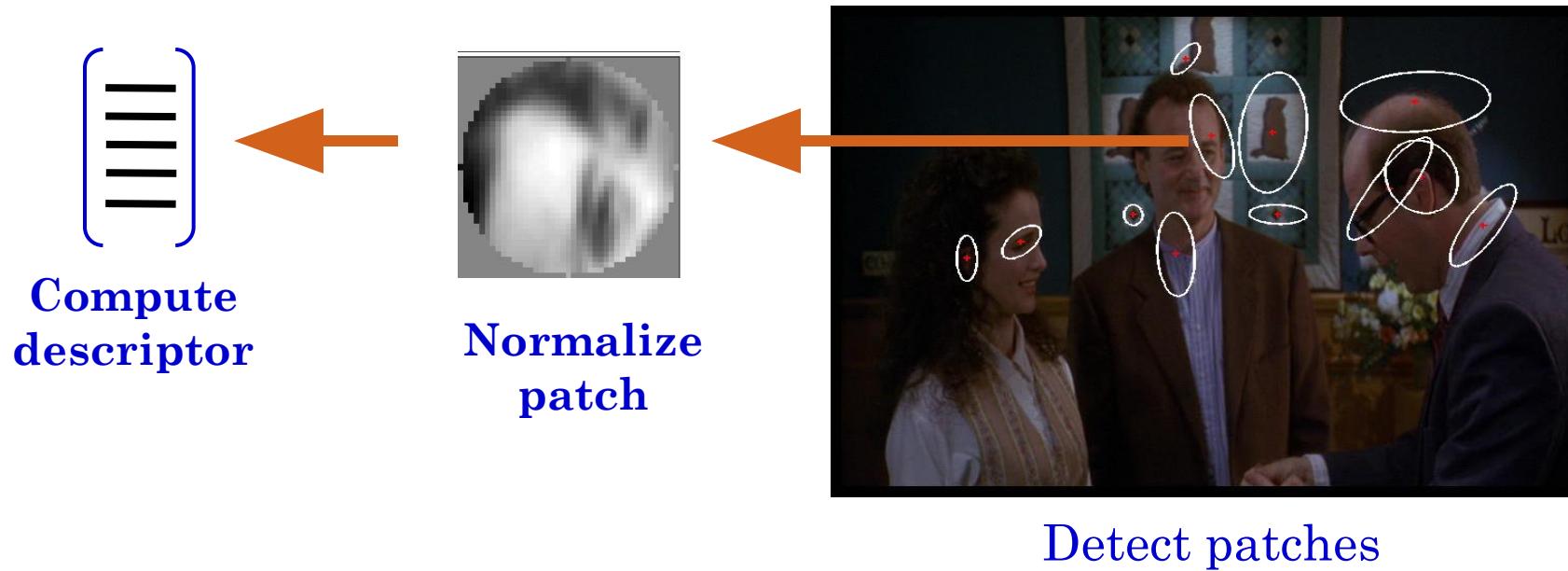


1. FEATURE EXTRACTION

- Regular grid or interest regions

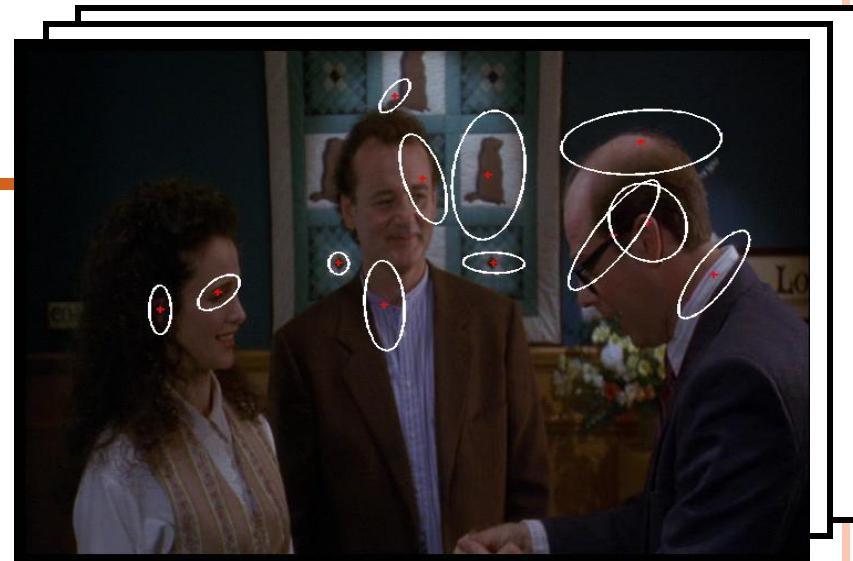
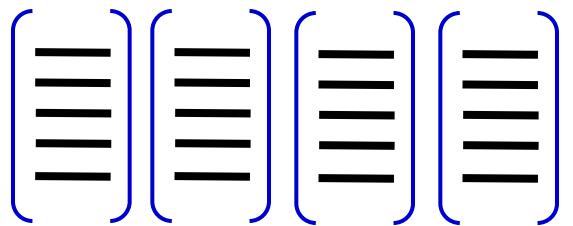


1. FEATURE EXTRACTION



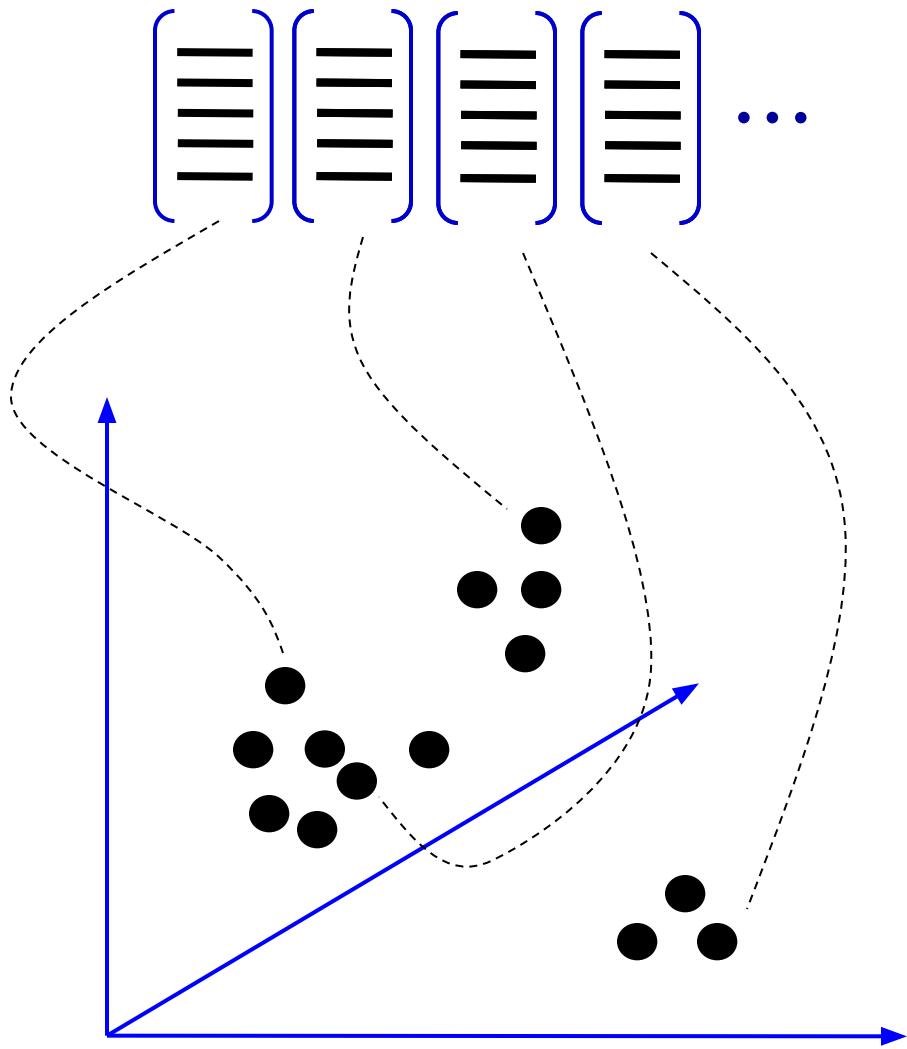
Slide credit: Josef
Sivic

1. FEATURE EXTRACTION

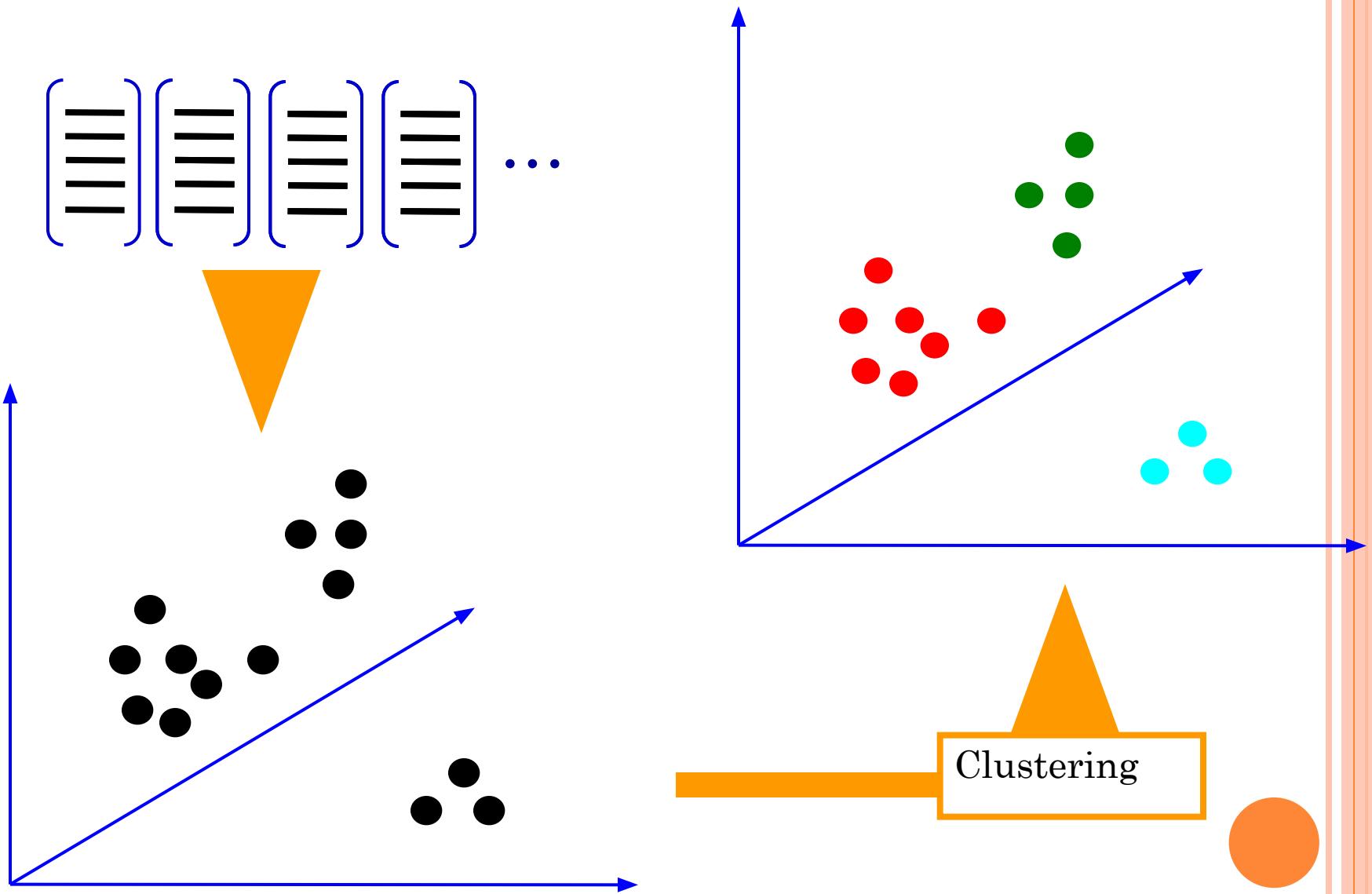


Slide credit: Josef
Gajdala

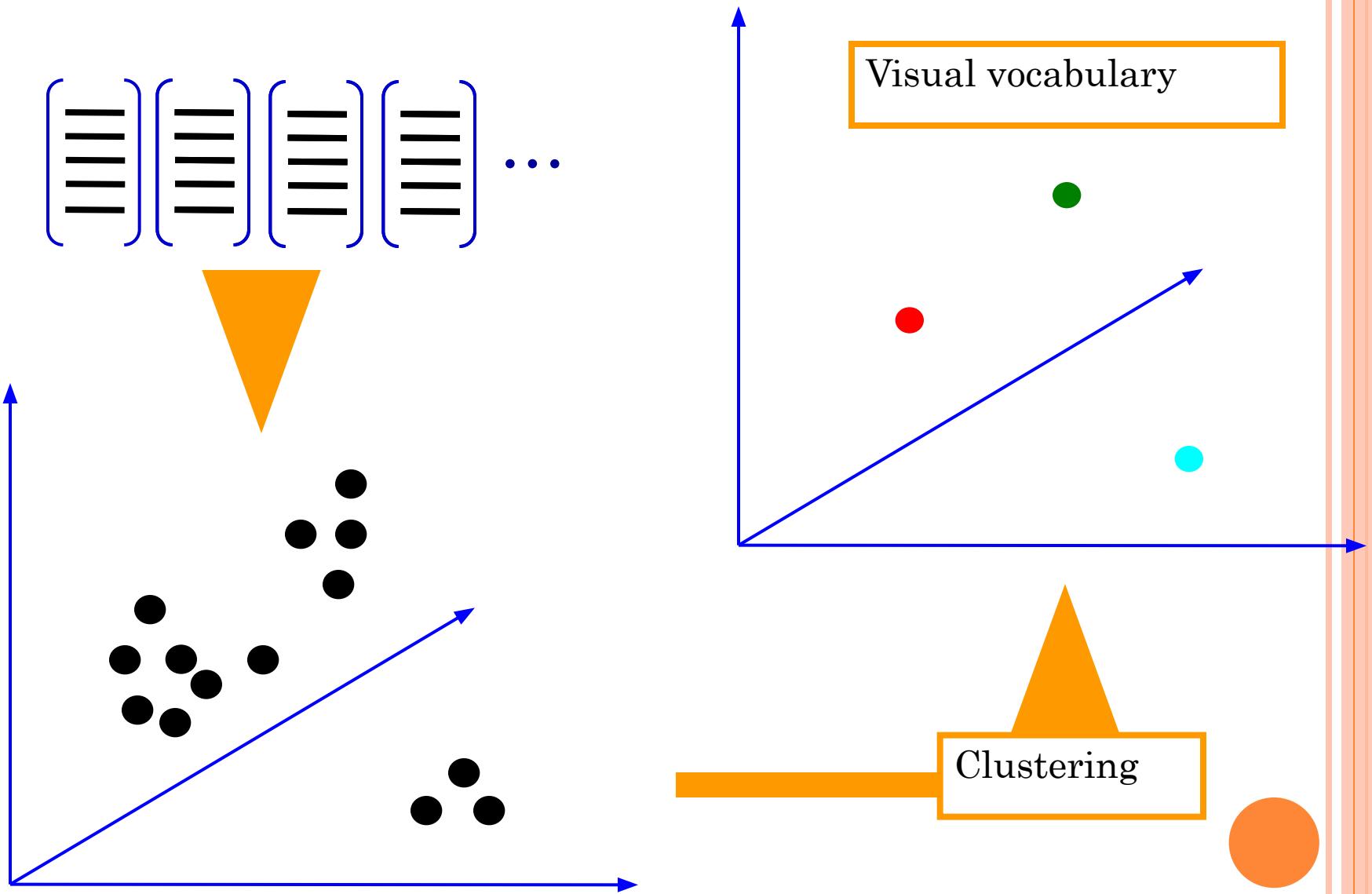
2. LEARNING THE VISUAL VOCABULARY



2. LEARNING THE VISUAL VOCABULARY



2. LEARNING THE VISUAL VOCABULARY



K-MEANS CLUSTERING

- Want to minimize sum of squared Euclidean distances between points x_i and their nearest cluster centers m_k

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\substack{\text{point } i \text{ in} \\ \text{cluster } k}} (x_i - m_k)^2$$

- Algorithm:
 - Randomly initialize K cluster centers
 - Iterate until convergence:
 - Assign each data point to the nearest center
 - Recompute each cluster center as the mean of all points assigned to it

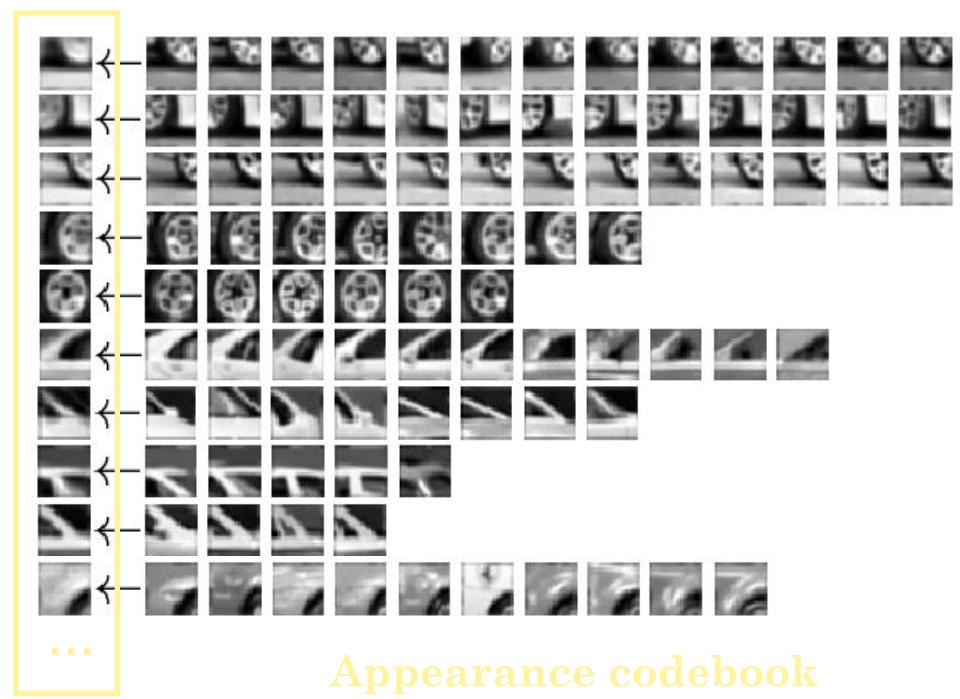


CLUSTERING AND VECTOR QUANTIZATION

- Clustering is a common method for learning a visual vocabulary or codebook
 - Unsupervised learning process
 - Each cluster center produced by k-means becomes a codevector
 - Codebook can be learned on separate training set
 - Provided the training set is sufficiently representative, the codebook will be “universal”
- The codebook is used for quantizing features
 - A *vector quantizer* takes a feature vector and maps it to the index of the nearest codevector in a codebook
 - Codebook = visual vocabulary
 - Codevector = visual word



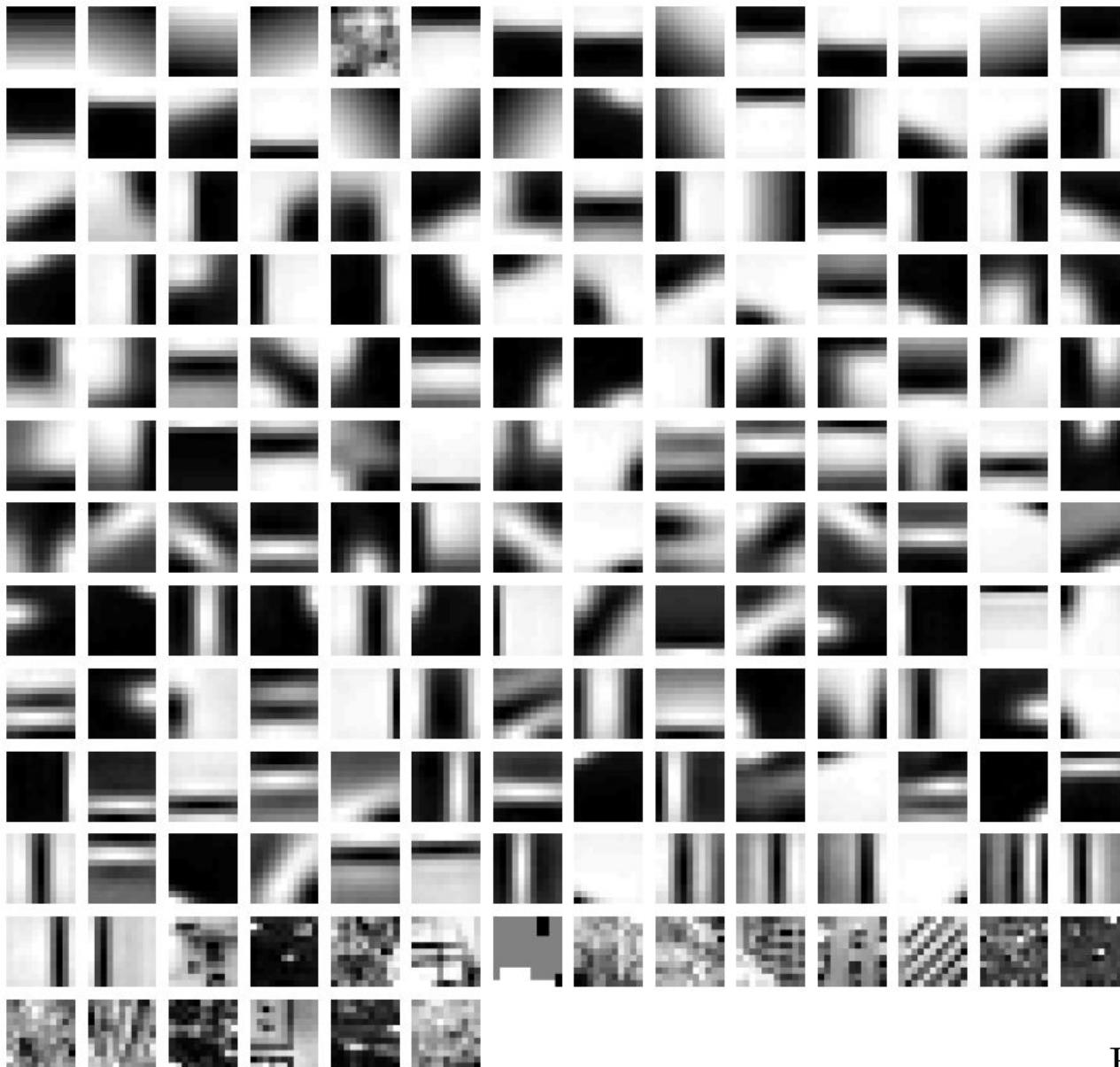
EXAMPLE CODEBOOK



ANOTHER CODEBOOK

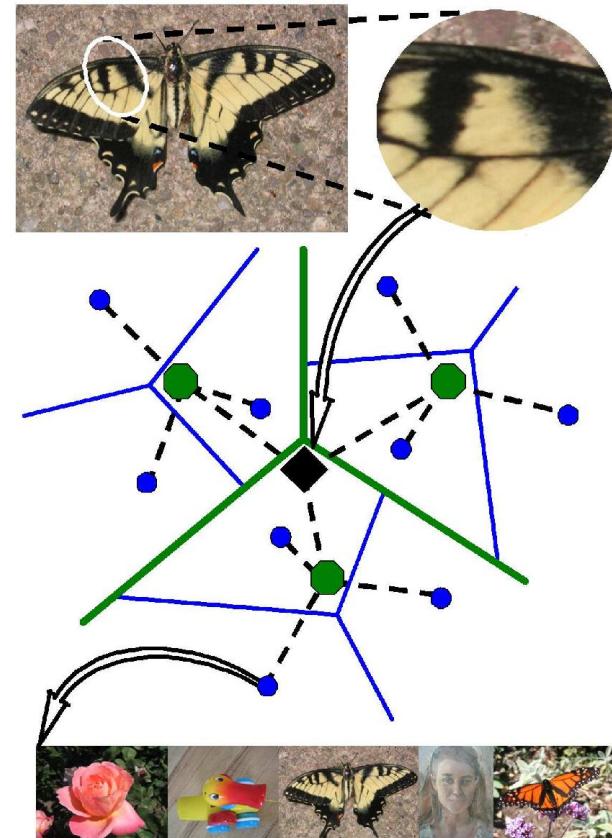


Yet another codebook



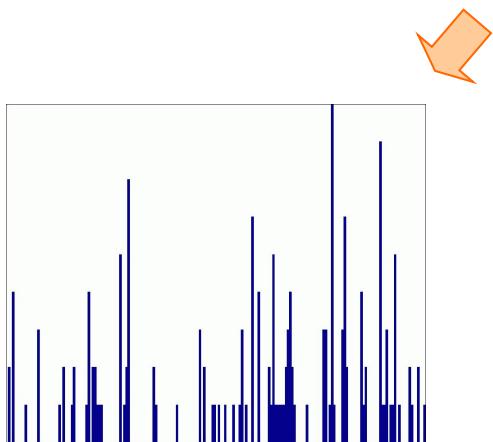
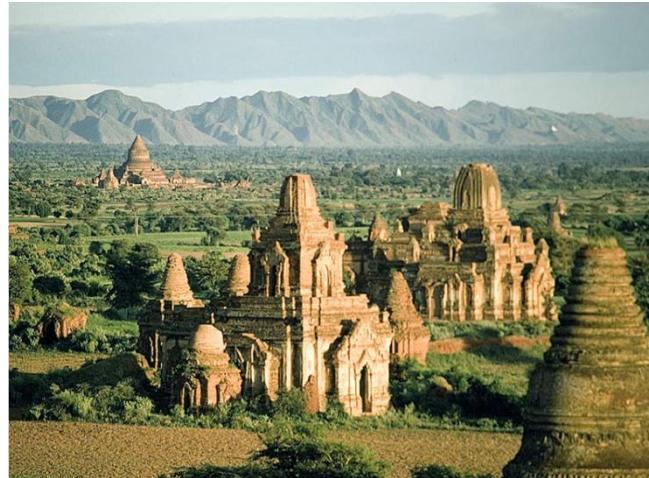
VISUAL VOCABULARIES: ISSUES

- How to choose vocabulary size?
 - Too small: visual words not representative of all patches
 - Too large: quantization artifacts, overfitting
- Computational efficiency
 - Vocabulary trees
(Nister & Stewenius, 2006)



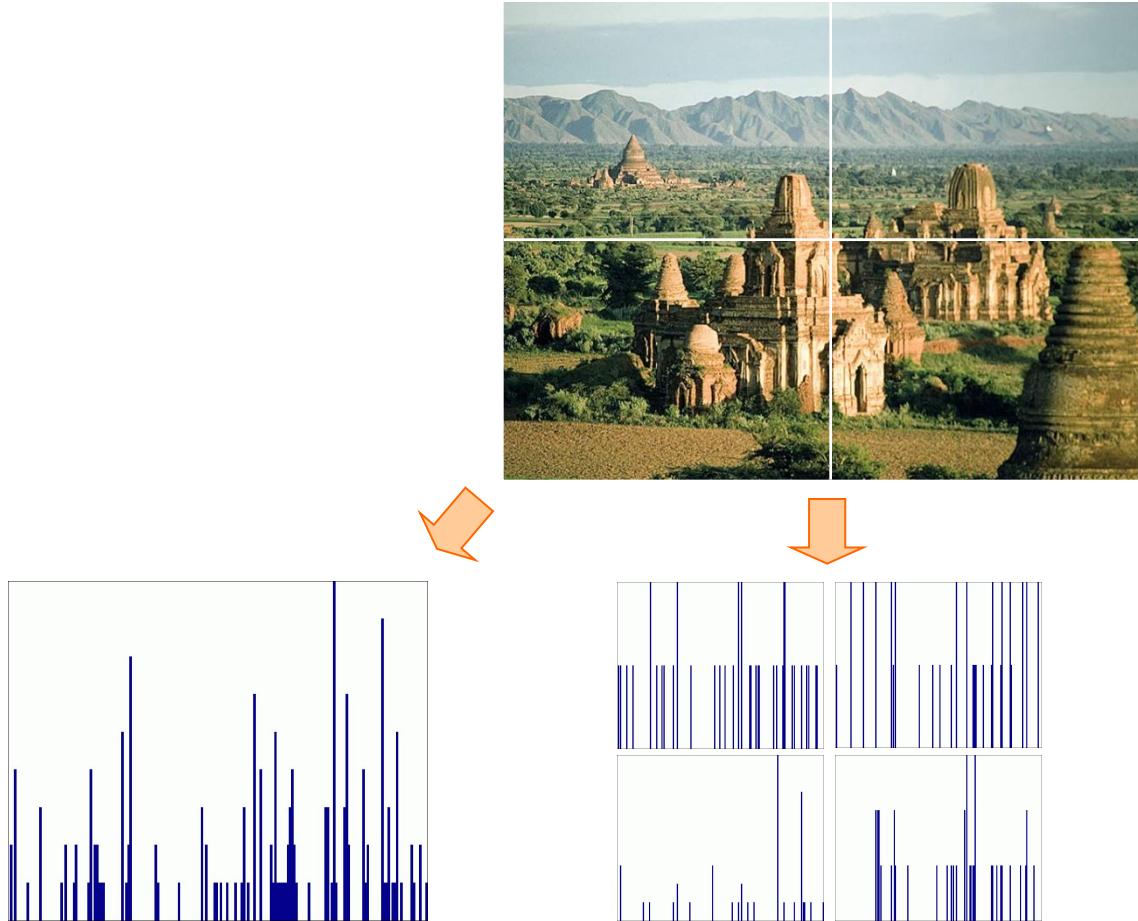
SPATIAL PYRAMID REPRESENTATION

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



SPATIAL PYRAMID REPRESENTATION

- Extension of a bag of features
- Locally orderless representation at several levels of resolution

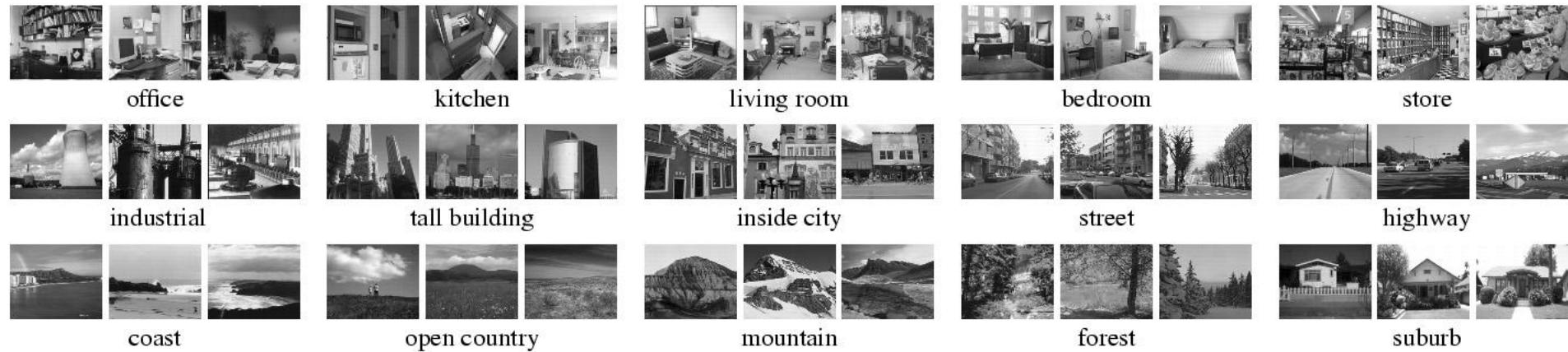


SPATIAL PYRAMID REPRESENTATION

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



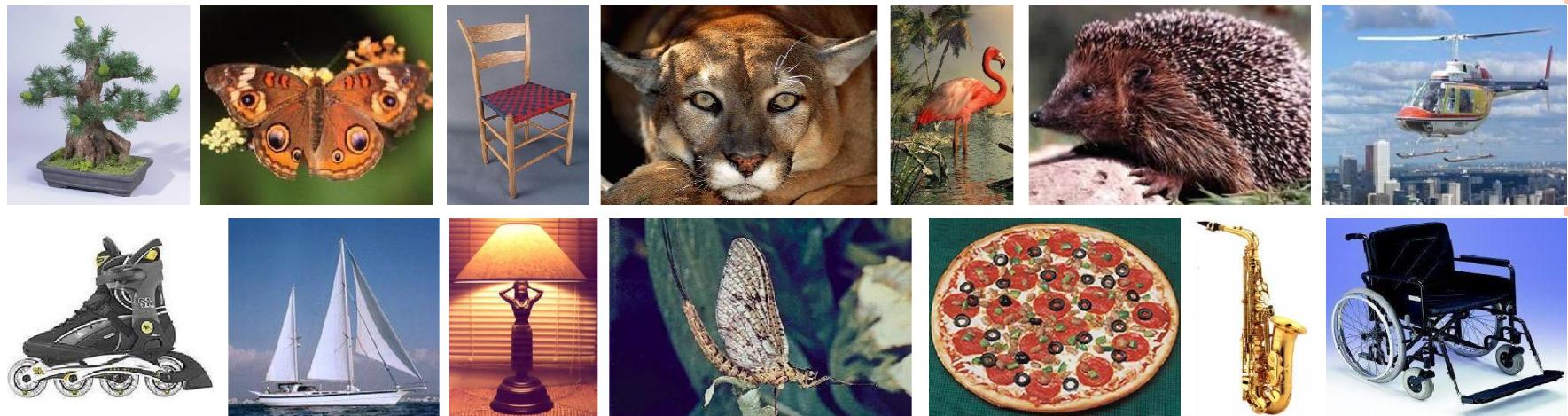
SCENE CATEGORY DATASET



Multi-class classification
results

Level	Weak features (vocabulary size: 16)		Strong features (vocabulary size: 200)	
	Single-level	Pyramid	Single-level	Pyramid
0 (1×1)	45.3 ± 0.5		72.2 ± 0.6	
1 (2×2)	53.6 ± 0.3	56.2 ± 0.6	77.9 ± 0.6	79.0 ± 0.5
2 (4×4)	61.7 ± 0.6	64.7 ± 0.7	79.4 ± 0.3	81.1 ± 0.3
3 (8×8)	63.3 ± 0.8	66.8 ± 0.6	77.2 ± 0.4	80.7 ± 0.3

CALTECH101 DATASET



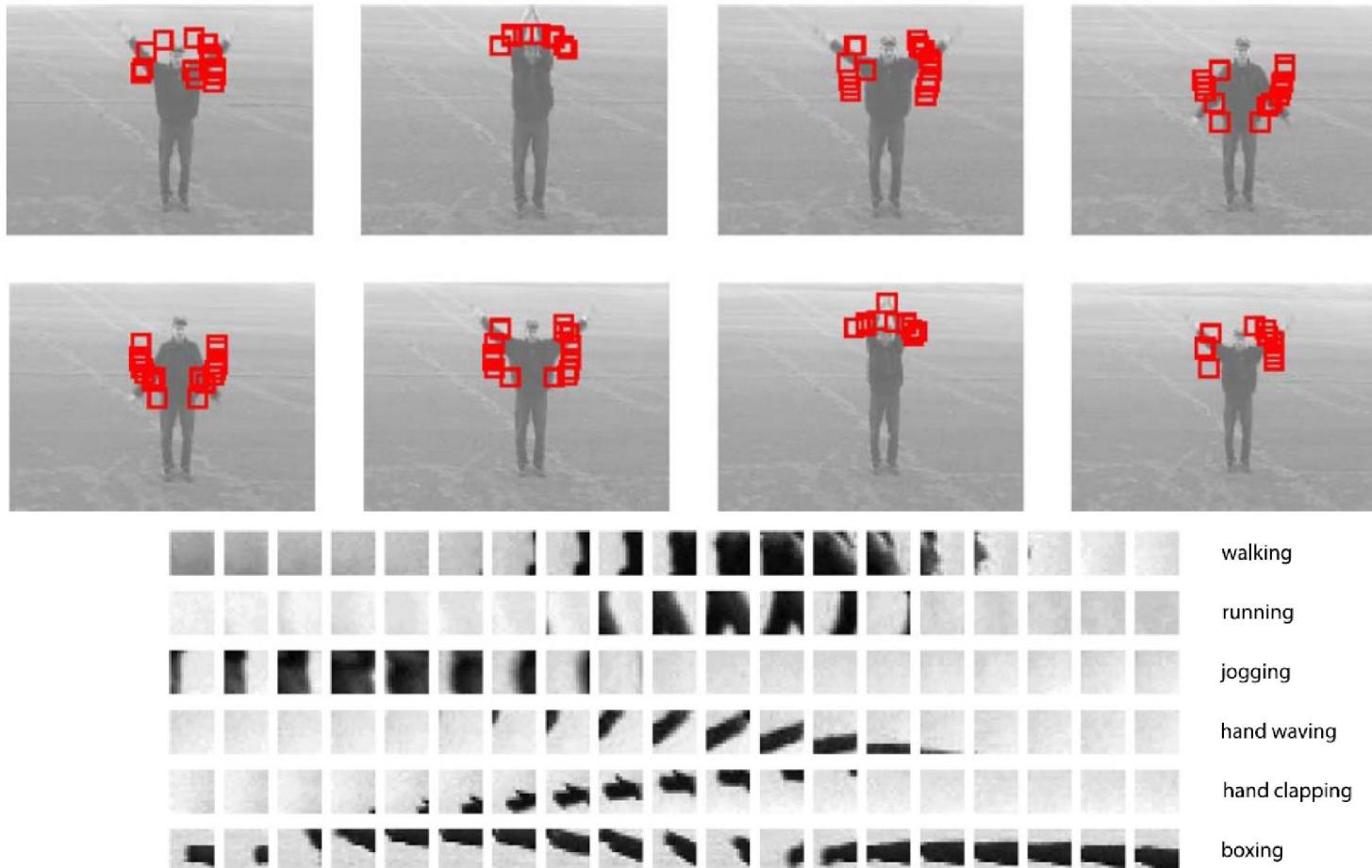
http://www.vision.caltech.edu/Image_Datasets/Caltech101/Caltech101.htm

Multi-class classification results (30 training images per

Level	Weak features (16)		Strong features (200)	
	Single-level	Pyramid	Single-level	Pyramid
0	15.5 ± 0.9		41.2 ± 1.2	
1	31.4 ± 1.2	32.8 ± 1.3	55.9 ± 0.9	57.0 ± 0.8
2	47.2 ± 1.1	49.3 ± 1.4	63.6 ± 0.9	64.6 ± 0.8
3	52.2 ± 0.8	54.0 ± 1.1	60.3 ± 0.9	64.6 ± 0.7

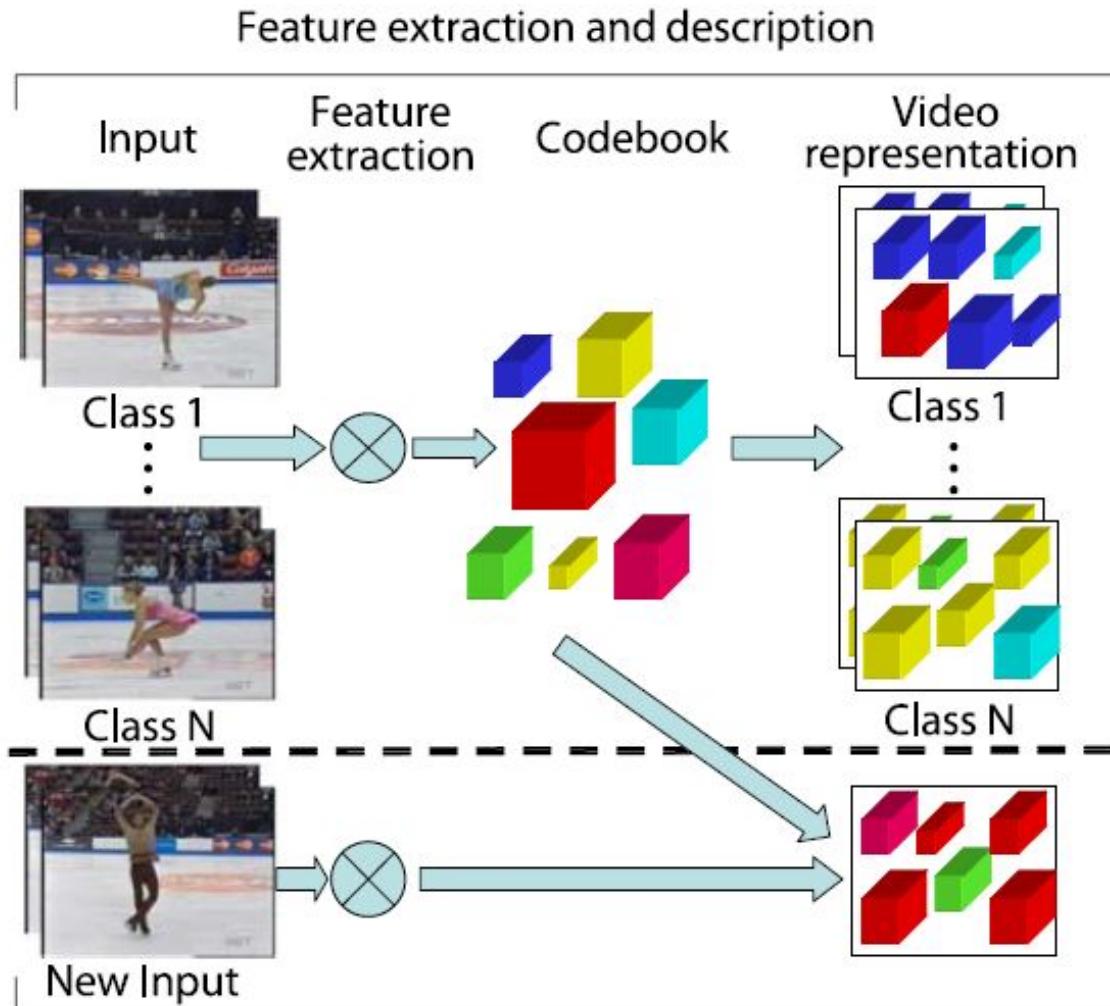
BAGS OF FEATURES FOR ACTION RECOGNITION

Space-time interest points



Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei, Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words, IJCV 2008.

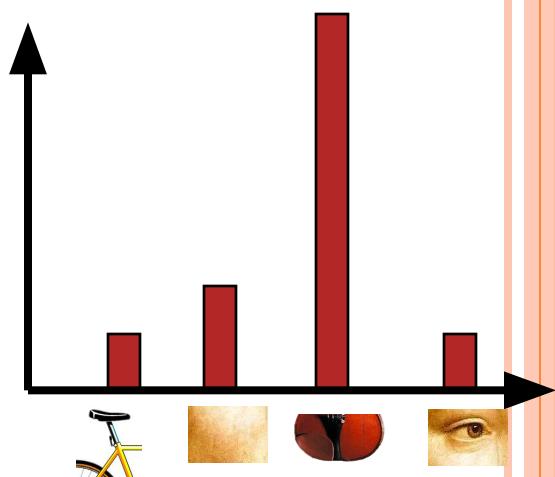
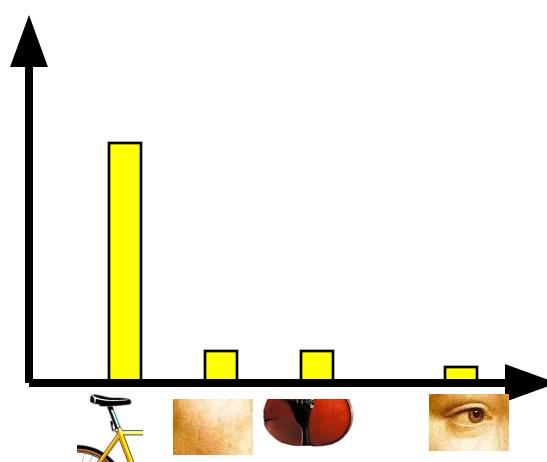
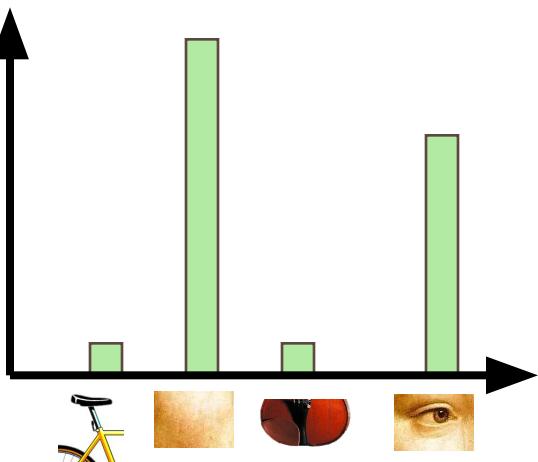
BAGS OF FEATURES FOR ACTION RECOGNITION

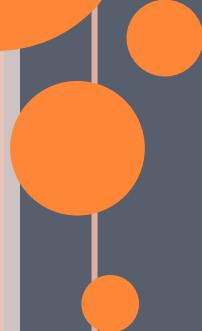
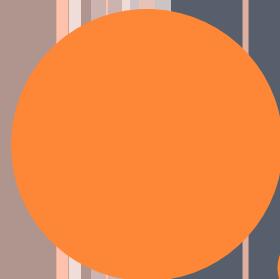


Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei, Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words, IJCV 2008.

IMAGE CLASSIFICATION

- Given the bag-of-features representations of images from different classes, how do we learn a model for distinguishing them?





LECTURE IV: PART II

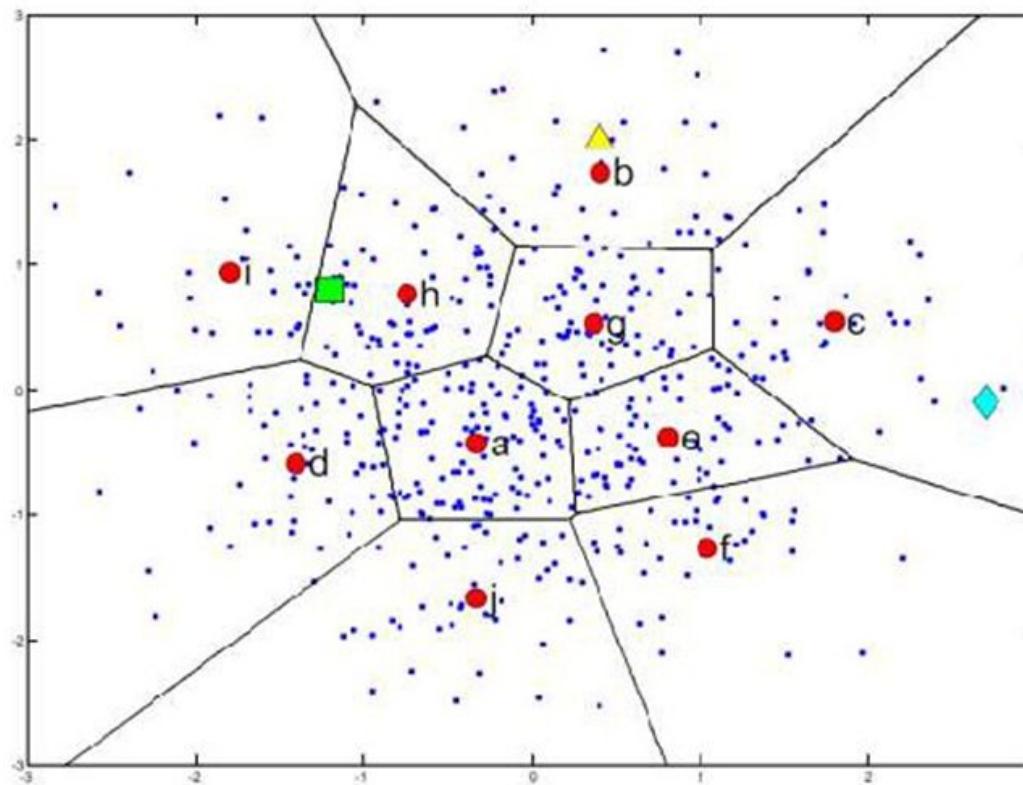
OUTLINE

- Histogram of local features
- Bag of words model
- Soft quantization and sparse coding
- Supervector with Gaussian mixture model



HARD QUANTIZATION

□
$$H(w) = \begin{cases} 1, & \text{if } w = \operatorname{argmin}_{v \in V}(D(v, r_i)) \\ 0, & \text{otherwise} \end{cases}$$



SOFT QUANTIZATION BASED ON UNCERTAINTY

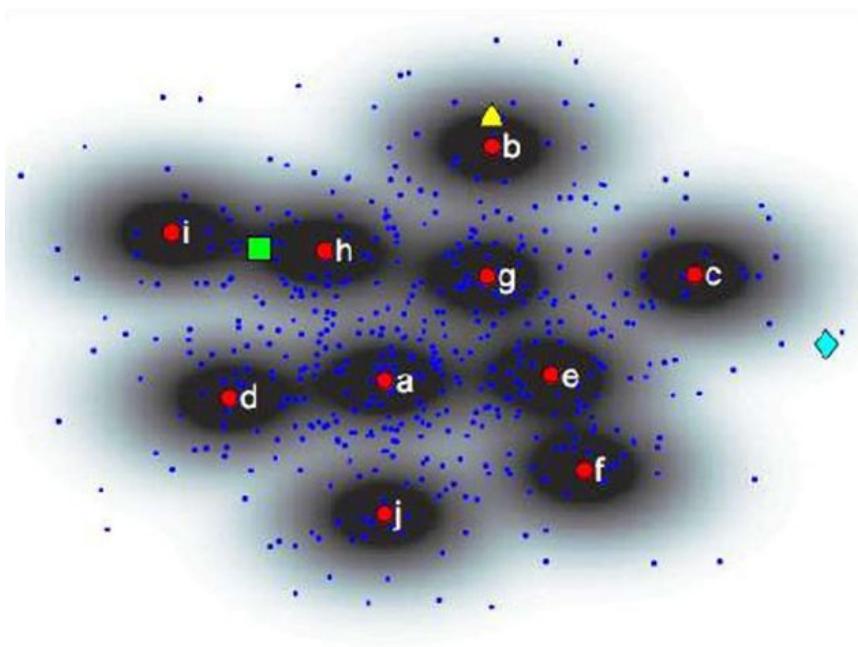
- Quantize each local feature into multiple code-words that are closest to it by appropriately splitting its weight

$$K_\sigma(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

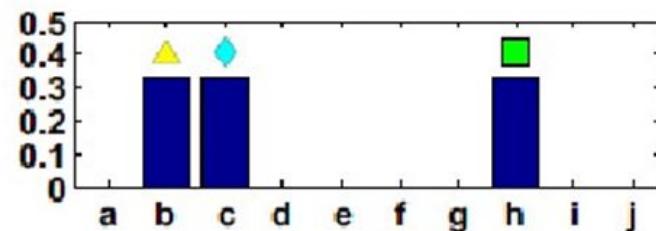
$$UNC(w) = \frac{1}{n} \sum_{i=1}^n \frac{K_\sigma(D(w, f))}{\sum_{j=1}^{|V|} K_\sigma(D(v_j, r_i))}$$



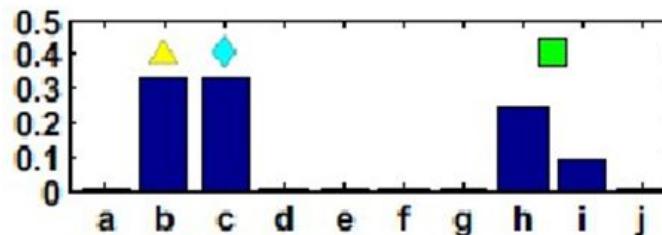
SOFT QUANTIZATION



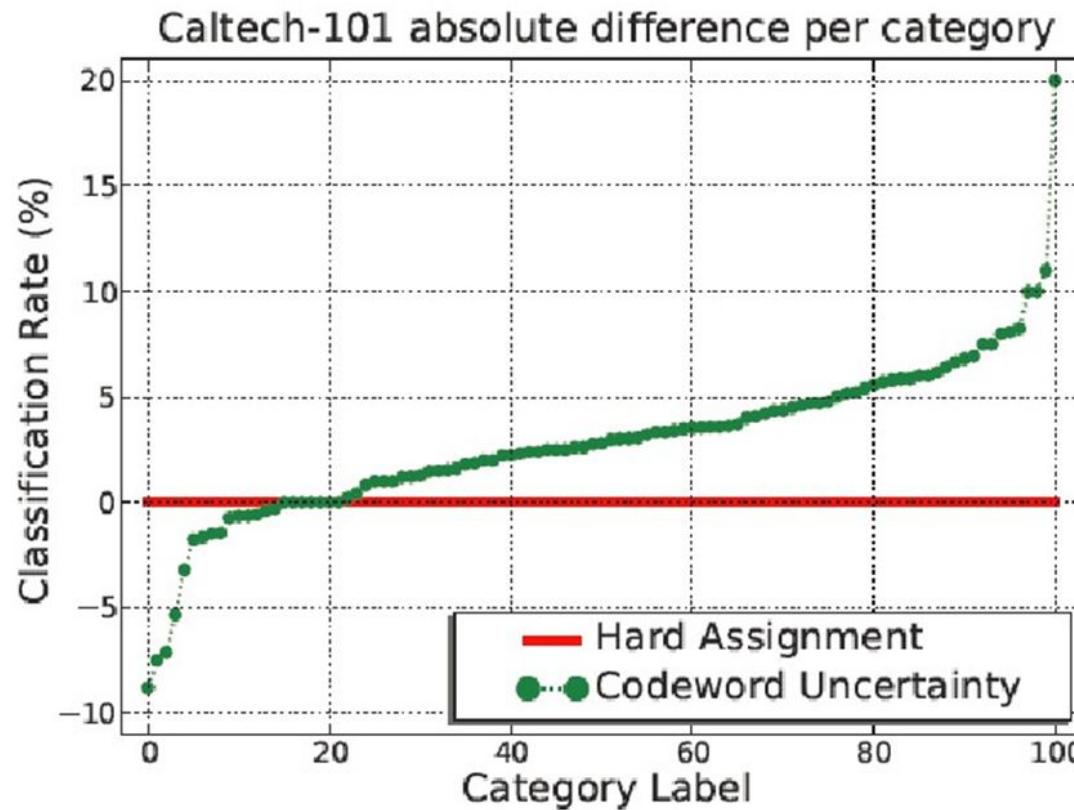
Hard quantization



Soft quantization



SOME EXPERIMENTS ON SOFT QUANTIZATION



- Improvement on classification rate from soft quantization

SPARSE CODING

- Hard quantization is an “extremely sparse representation

$$\min \|\mathbf{x} - \mathbf{Dz}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{z}\|_0 = 1$$

- Generalizing from it, we may consider to solve

$$\min \|\mathbf{x} - \mathbf{Dz}\|_2^2 + \lambda \|\mathbf{z}\|_0$$

for soft quantization, but it is hard to solve

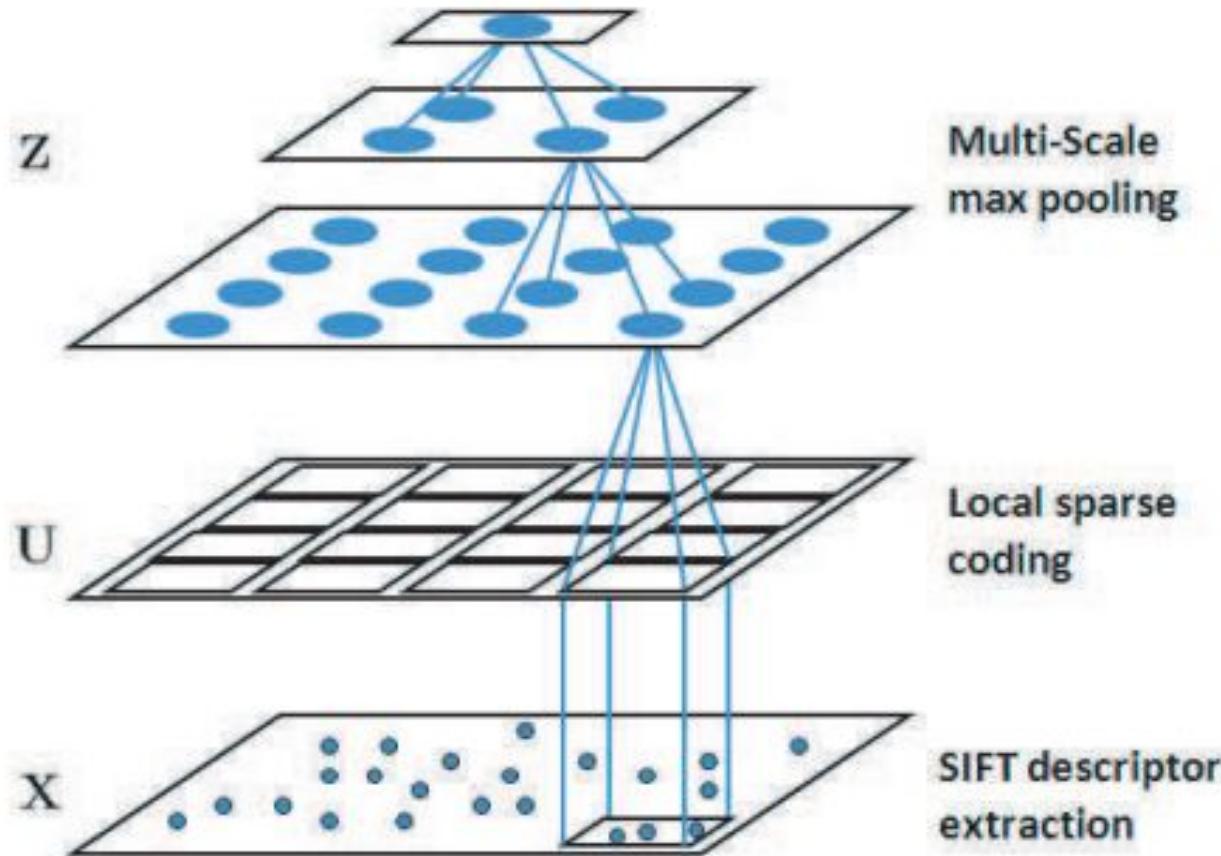
- In practice, we consider solving the sparse coding as

$$\min \|\mathbf{x} - \mathbf{Dz}\|_2^2 + \lambda \|\mathbf{z}\|_1$$

for quantization



SOFT QUANTIZATION AND POOLING



[Yang et al, 2009]:
Linear Spatial Pyramid Matching using Sparse Coding for Image Classification

OUTLINE

- Histogram of local features
- Bag of words model
- Soft quantization and sparse coding
- Supervector with Gaussian mixture model



QUIZ

- What is the potential shortcoming of bag-of-feature representation based on quantization and pooling?



MODELING THE FEATURE DISTRIBUTION

- The bag-of-feature histogram represents the distribution of a set of features
 - The quantization steps introduces information loss!
- How about modeling the distribution without quantization?
 - Gaussian mixture model

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



QUIZ

- Given a set of features $\{x_1, x_2, \dots, x_n\}$, how can we fit the GMM distribution?



THE GAUSSIAN DISTRIBUTION

- ## □ Multivariate Gaussian

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi|\boldsymbol{\Sigma}|)^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$


- Define precision to be the inverse of the covariance

$$\Lambda \equiv \Sigma^{-1}$$

- ## □ In 1-dimension

$$\tau = \frac{1}{\sigma^2}$$



LIKELIHOOD FUNCTION

- Data set

$$D = \{\mathbf{x}_n\} \quad n = 1, \dots, N$$

- Assume observed data points generated independently

$$p(D|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Viewed as a function of the parameters, this is known as the *likelihood function*



MAXIMUM LIKELIHOOD

- Set the parameters by maximizing the likelihood function
- Equivalently maximize the log likelihood

$$\begin{aligned}\ln p(D|\mu, \Sigma) &= -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \ln(2\pi) \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$



MAXIMUM LIKELIHOOD SOLUTION

- Maximizing w.r.t. the mean gives the *sample mean*

$$\boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Maximizing w.r.t covariance gives the *sample covariance*

$$\boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T$$



BIAS OF MAXIMUM LIKELIHOOD

- Consider the expectations of the maximum likelihood estimates under the Gaussian distribution

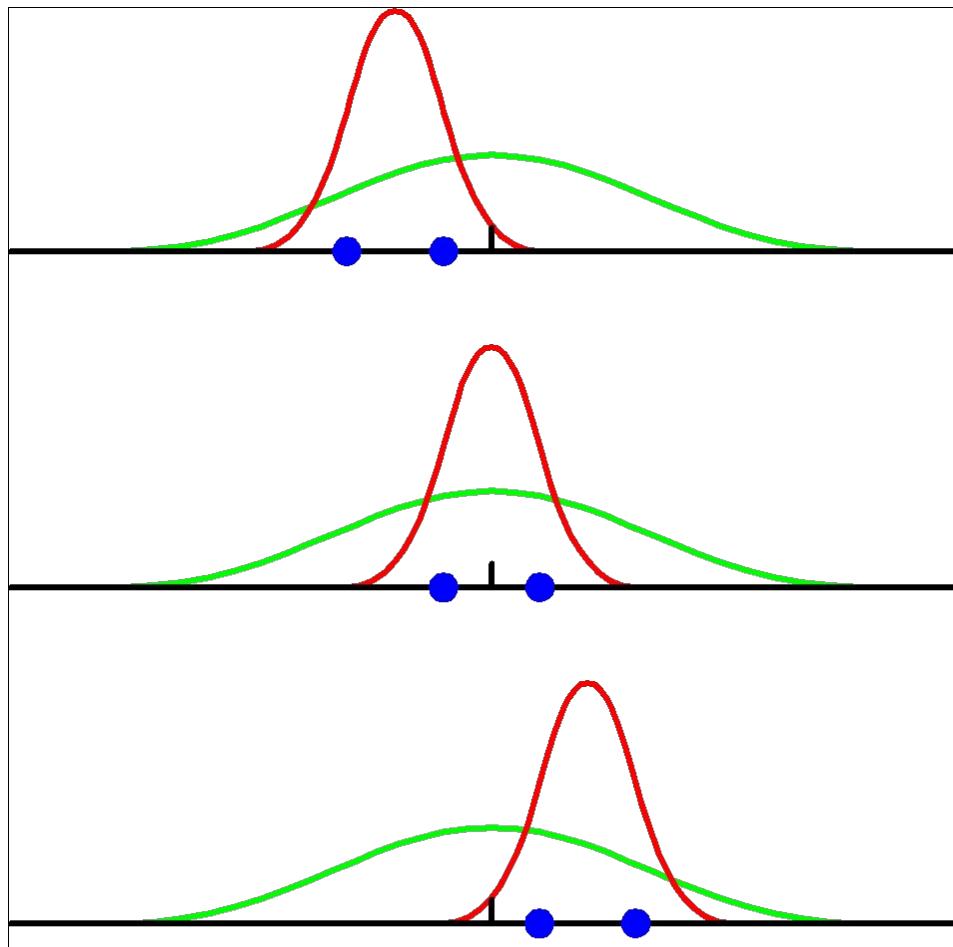
$$\mathbb{E}[\mu_{\text{ML}}] = \mu$$

$$\mathbb{E}[\Sigma_{\text{ML}}] = \left(\frac{N-1}{N}\right) \Sigma$$

- The maximum likelihood solution systematically under-estimates the covariance
- This is an example of *over-fitting*



INTUITIVE EXPLANATION OF OVER-FITTING



UNBIASED VARIANCE ESTIMATE

- Clearly we can remove the bias by using

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T$$

since this gives

$$\mathbb{E} [\tilde{\Sigma}] = \Sigma$$

- For an infinite data set the two expressions are equal



GAUSSIAN MIXTURES

- Linear super-position of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Normalization and positivity require

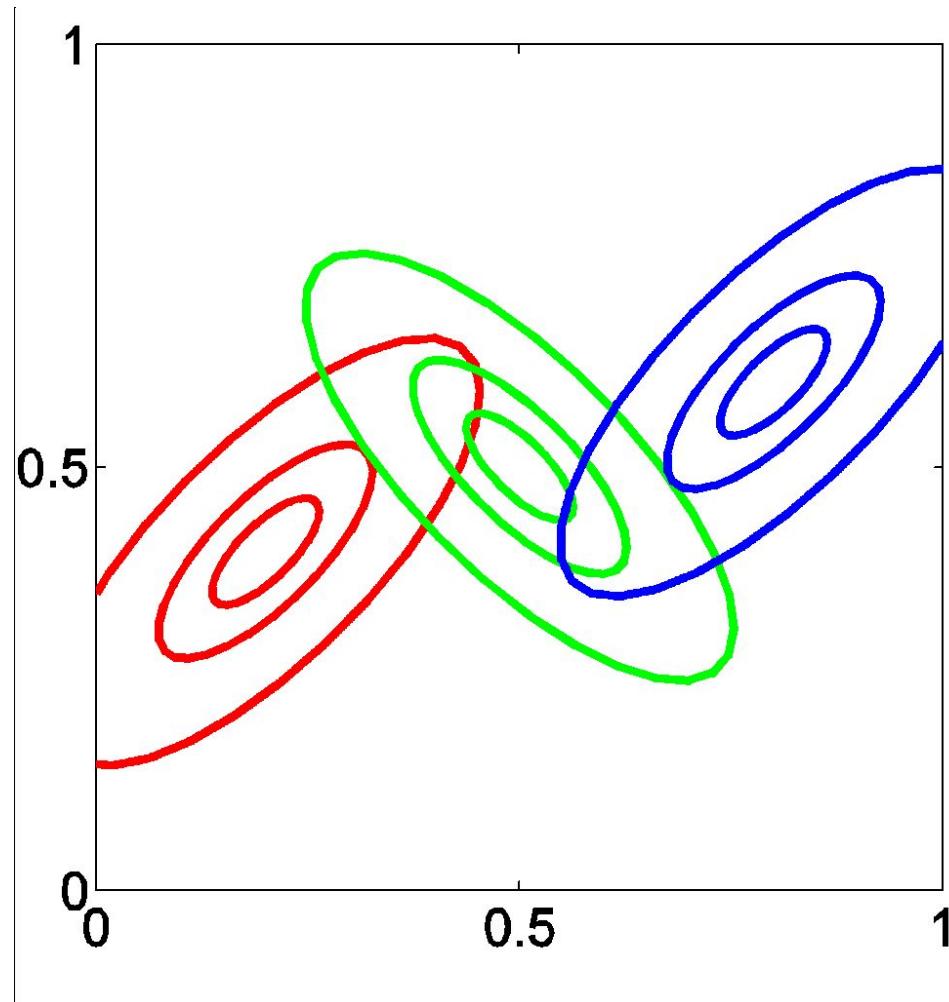
$$\sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$$

- Can interpret the mixing coefficients as prior probabilities

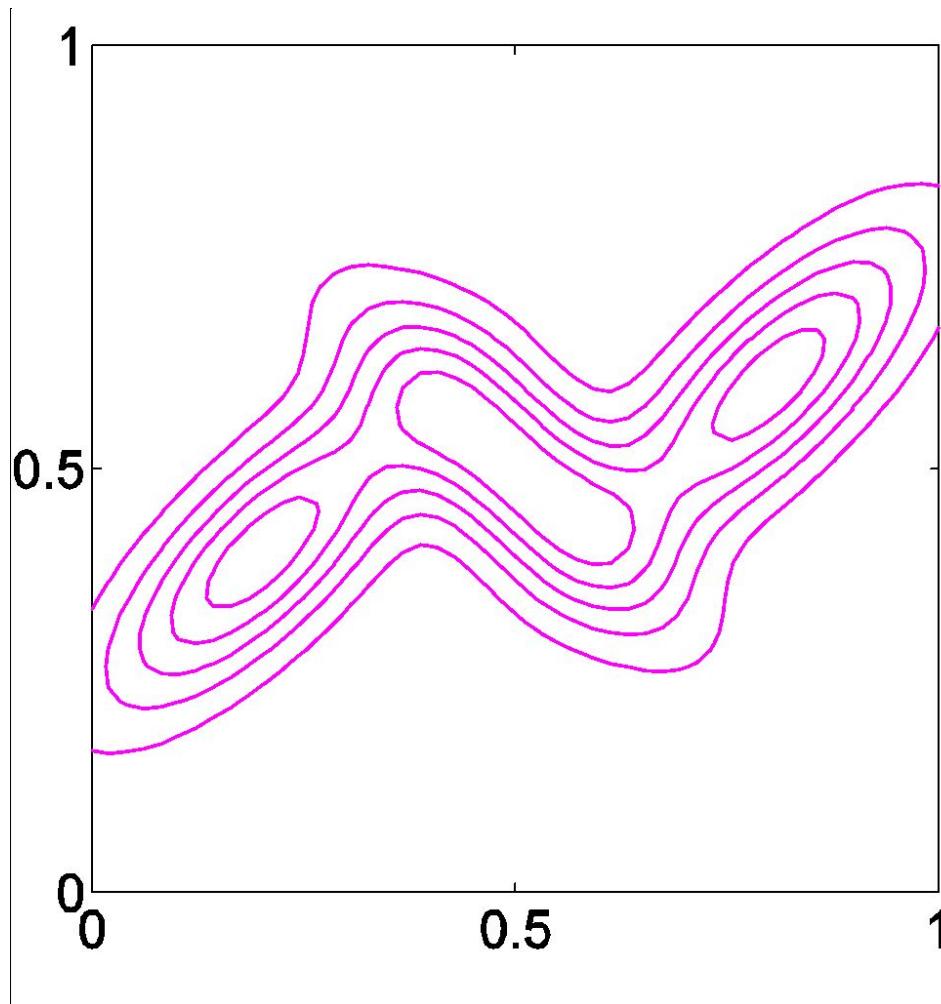
$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x} | k)$$

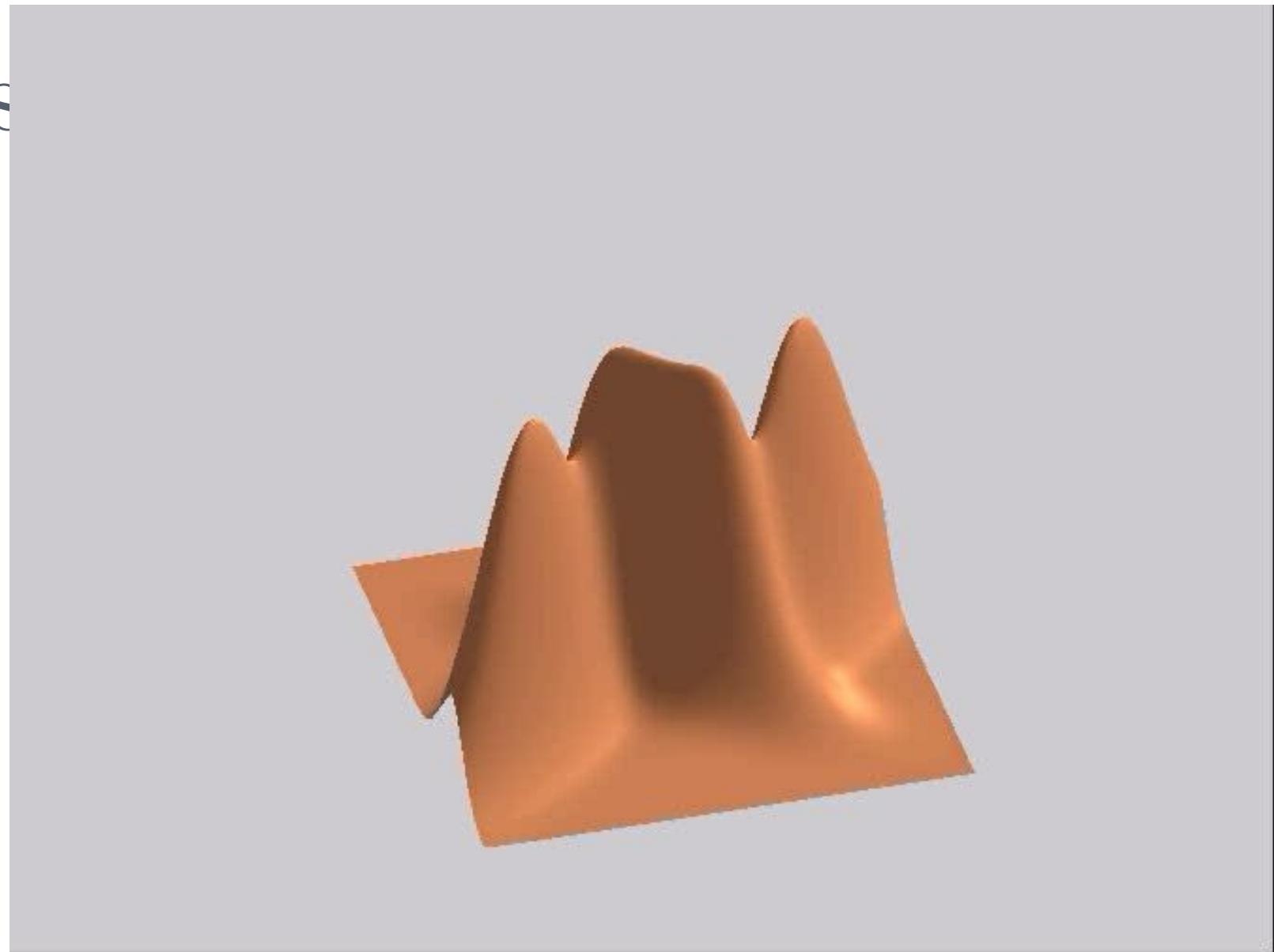


EXAMPLE: MIXTURE OF 3 GAUSSIANS



CONTOURS OF PROBABILITY DISTRIBUTION



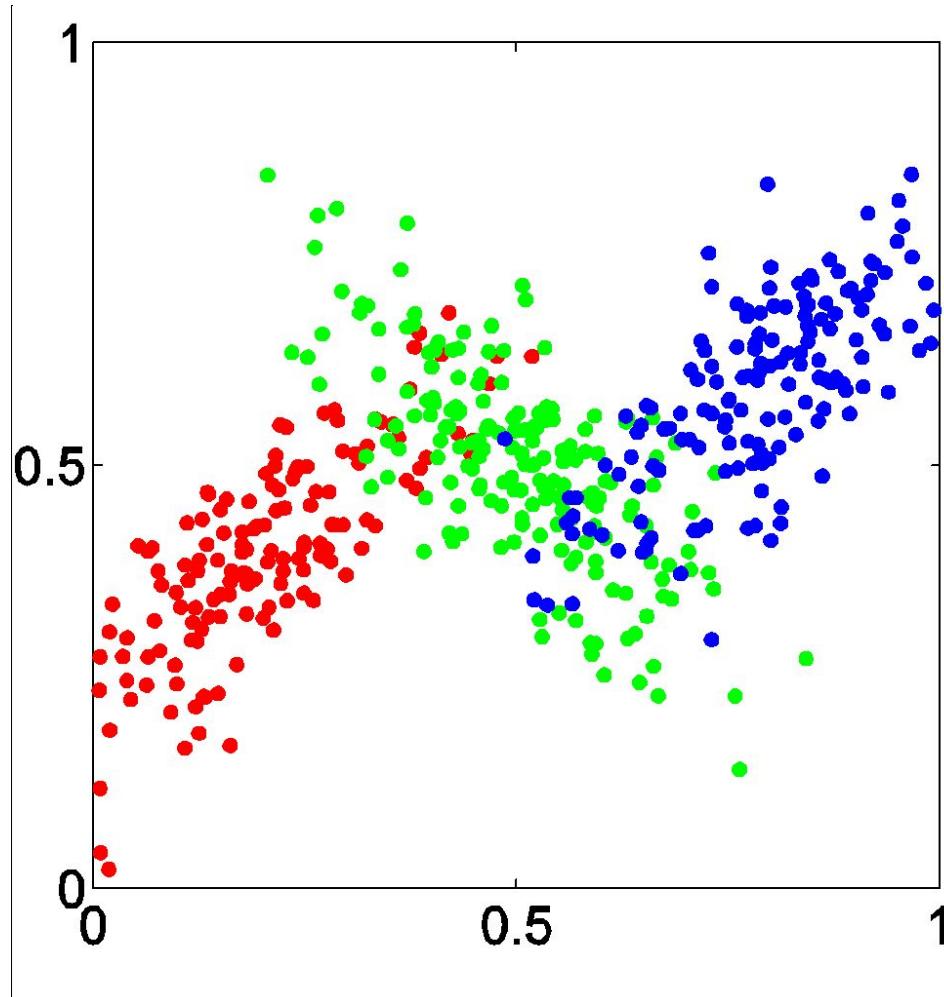


SAMPLING FROM THE GAUSSIAN

- To generate a data point:
 - first pick one of the components with probability π_k
 - then draw a sample \mathbf{x}_n from that component
- Repeat these two steps for each new data point



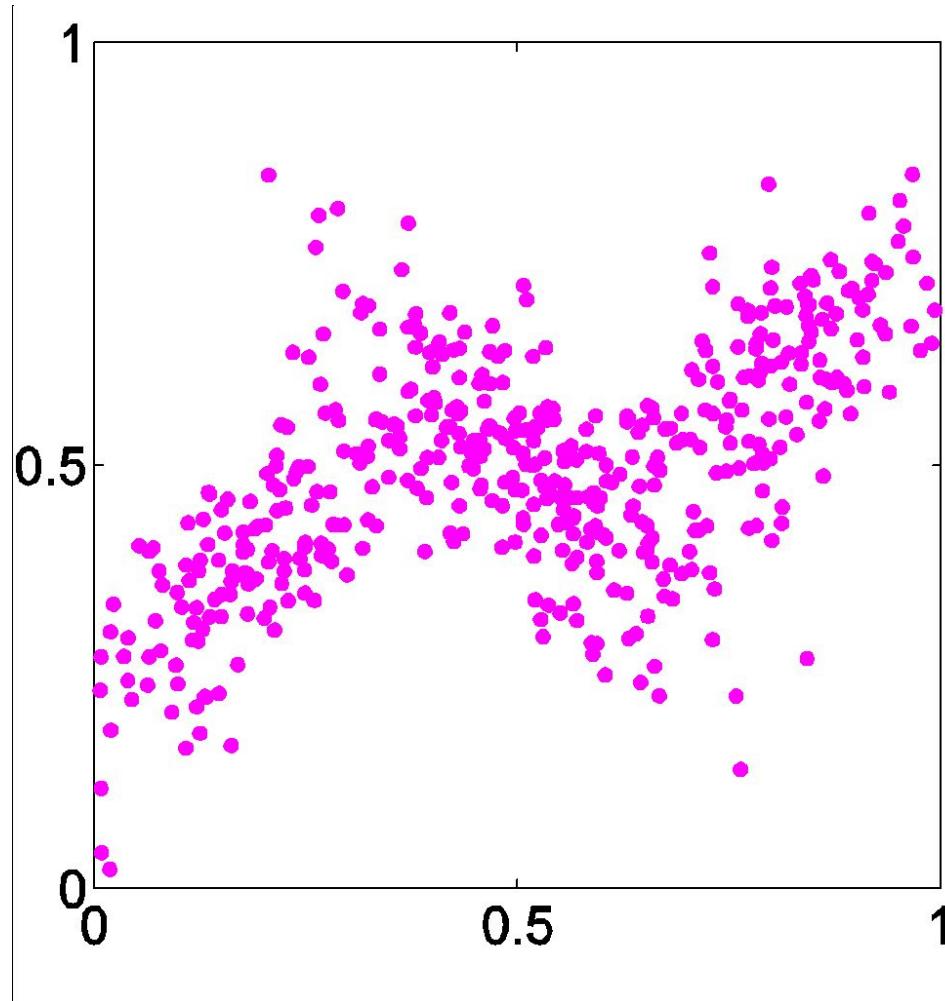
SYNTHETIC DATA SET



FITTING THE GAUSSIAN MIXTURE

- We wish to invert this process – given the data set, find the corresponding parameters:
 - mixing coefficients, means, and covariances
- If we knew which component generated each data point, the maximum likelihood solution would involve fitting each component to the corresponding cluster
- Problem: the data set is unlabelled
- We shall refer to the labels as *latent* (= hidden) variables

SYNTHETIC DATA SET WITHOUT LABELS



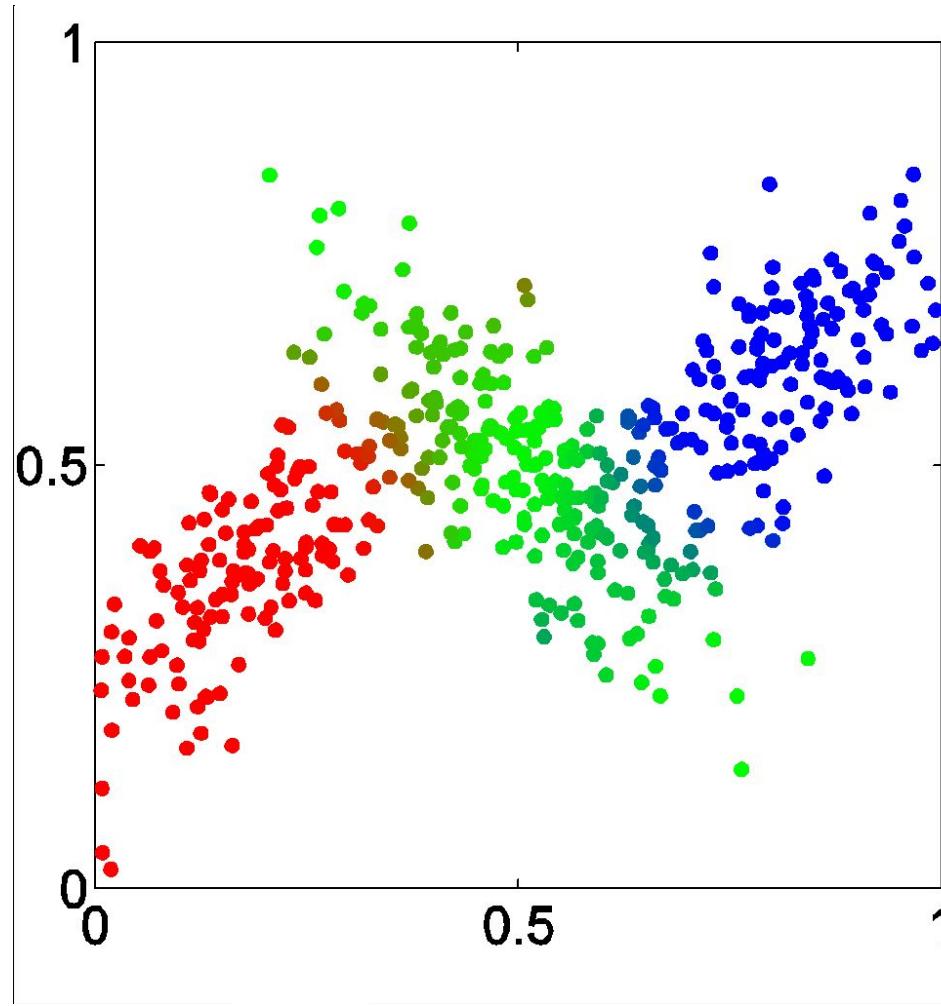
POSTERIOR PROBABILITIES

- We can think of the mixing coefficients as prior probabilities for the components
- For a given value of \mathbf{x} we can evaluate the corresponding posterior probabilities, called *responsibilities*
- These are given from Bayes' theorem by

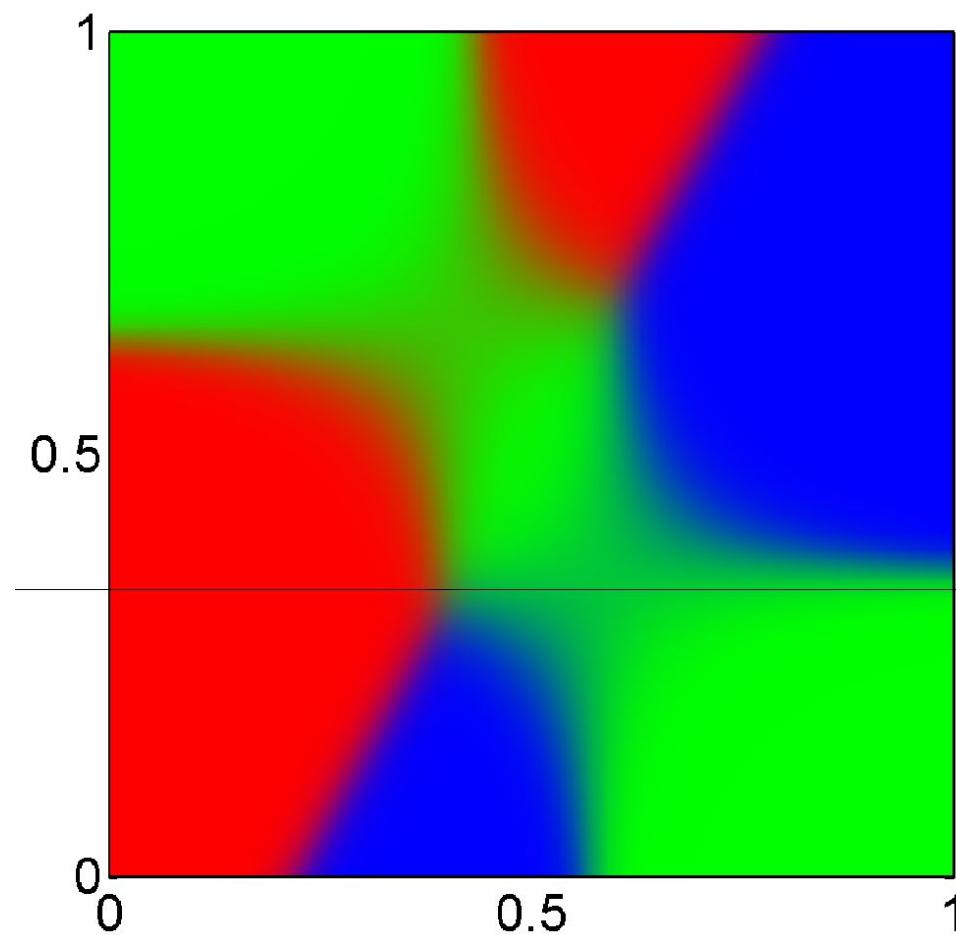
$$\begin{aligned}\gamma_k(\mathbf{x}) \equiv p(k|\mathbf{x}) &= \frac{p(k)p(\mathbf{x}|k)}{p(\mathbf{x})} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}\end{aligned}$$



POSTERIOR PROBABILITIES (COLOUR CODED)



POSTERIOR PROBABILITY MAP



MAXIMUM LIKELIHOOD FOR THE GMM

- The log likelihood function takes the form

$$\ln p(D|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- Note: sum over components appears *inside* the log
- There is no closed form solution for maximum likelihood!



OVER-FITTING IN GAUSSIAN MIXTURE MODELS

- Singularities in likelihood function when a component ‘collapses’ onto a data point:

$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}$$

then consider $\sigma_j \rightarrow 0$

- Likelihood function gets larger as we add more components (and hence parameters) to the model
 - not clear how to choose the number K of components



PROBLEMS AND SOLUTIONS

- How to maximize the log likelihood
 - solved by expectation-maximization (EM) algorithm
- How to avoid singularities in the likelihood function
 - solved by a Bayesian treatment
- How to choose number K of components
 - also solved by a Bayesian treatment



EM ALGORITHM – INFORMAL DERIVATION

- Let us proceed by simply differentiating the log likelihood
- Setting derivative with respect to μ_j equal to zero gives

$$-\sum_{n=1}^N \frac{\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\underbrace{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\gamma_j(\mathbf{x}_n)}} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j) = 0$$

giving

$$\boldsymbol{\mu}_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)}$$

- which is simply the weighted mean of the data

EM ALGORITHM – INFORMAL DERIVATION

- Similarly for the co-variances

$$\Sigma_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)(\mathbf{x}_n - \boldsymbol{\mu}_j)(\mathbf{x}_n - \boldsymbol{\mu}_j)^\top}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)}$$

- For mixing coefficients use a Lagrange multiplier to give

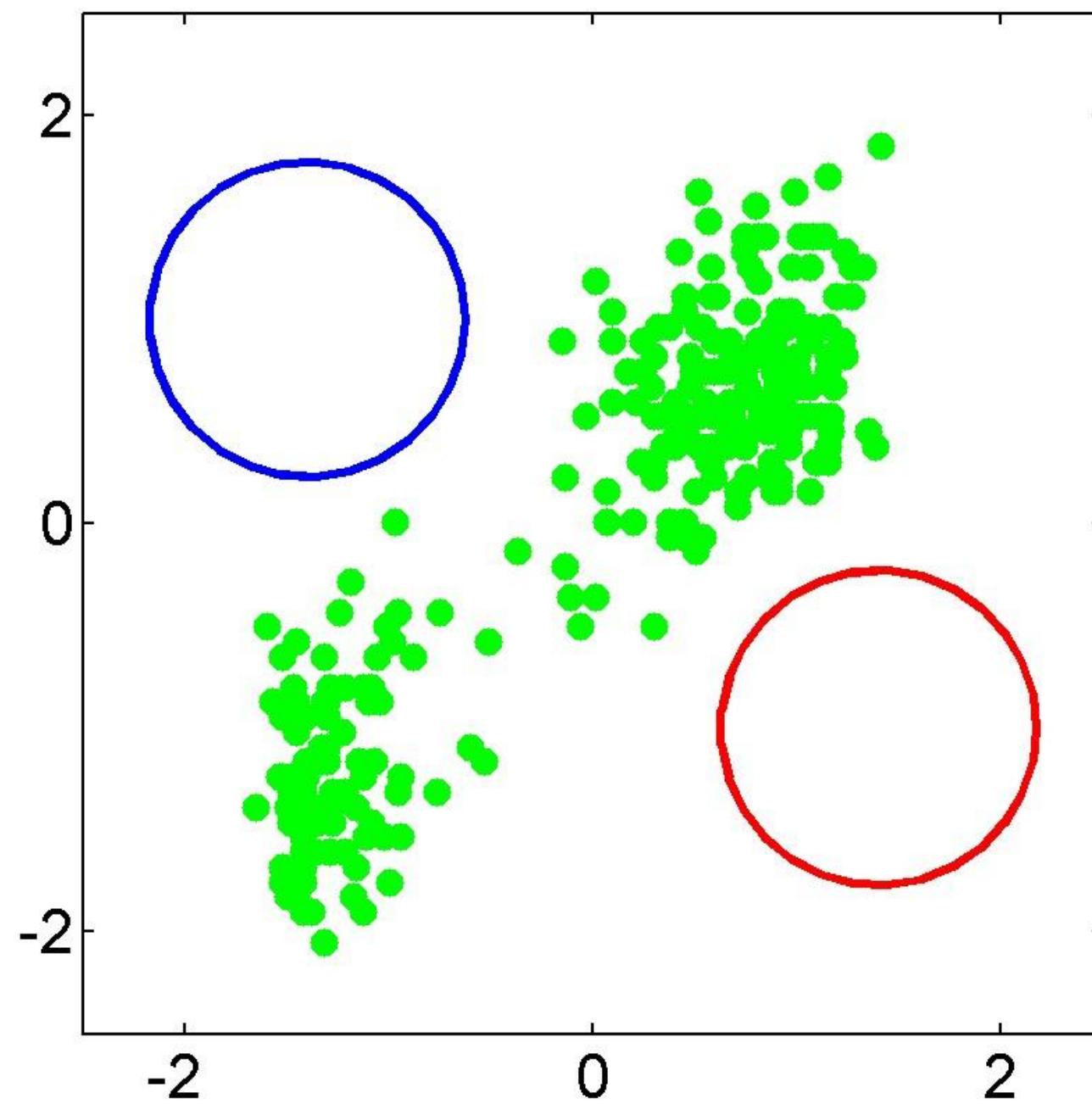
$$\pi_j = \frac{1}{N} \sum_{n=1}^N \gamma_j(\mathbf{x}_n)$$

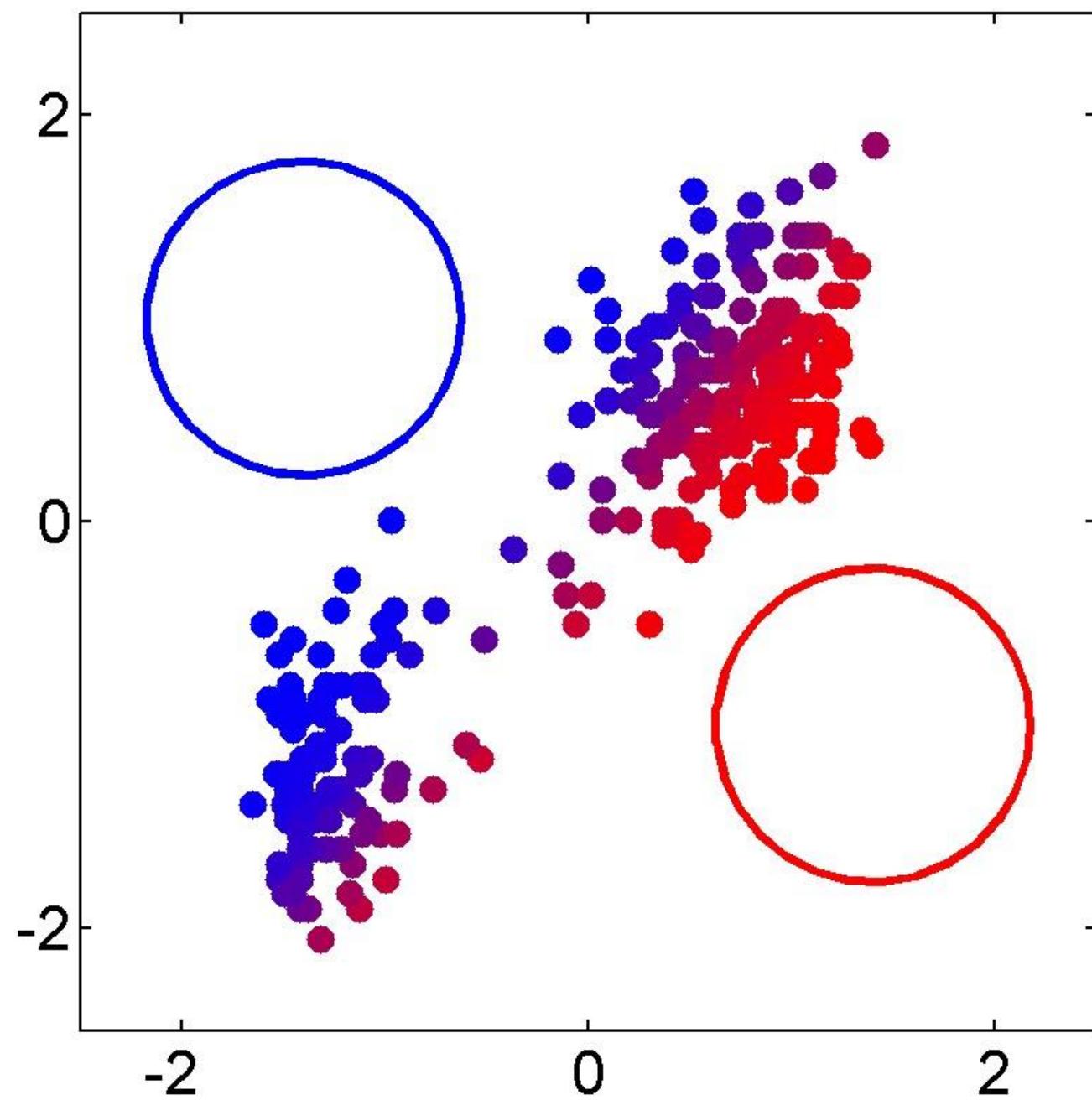


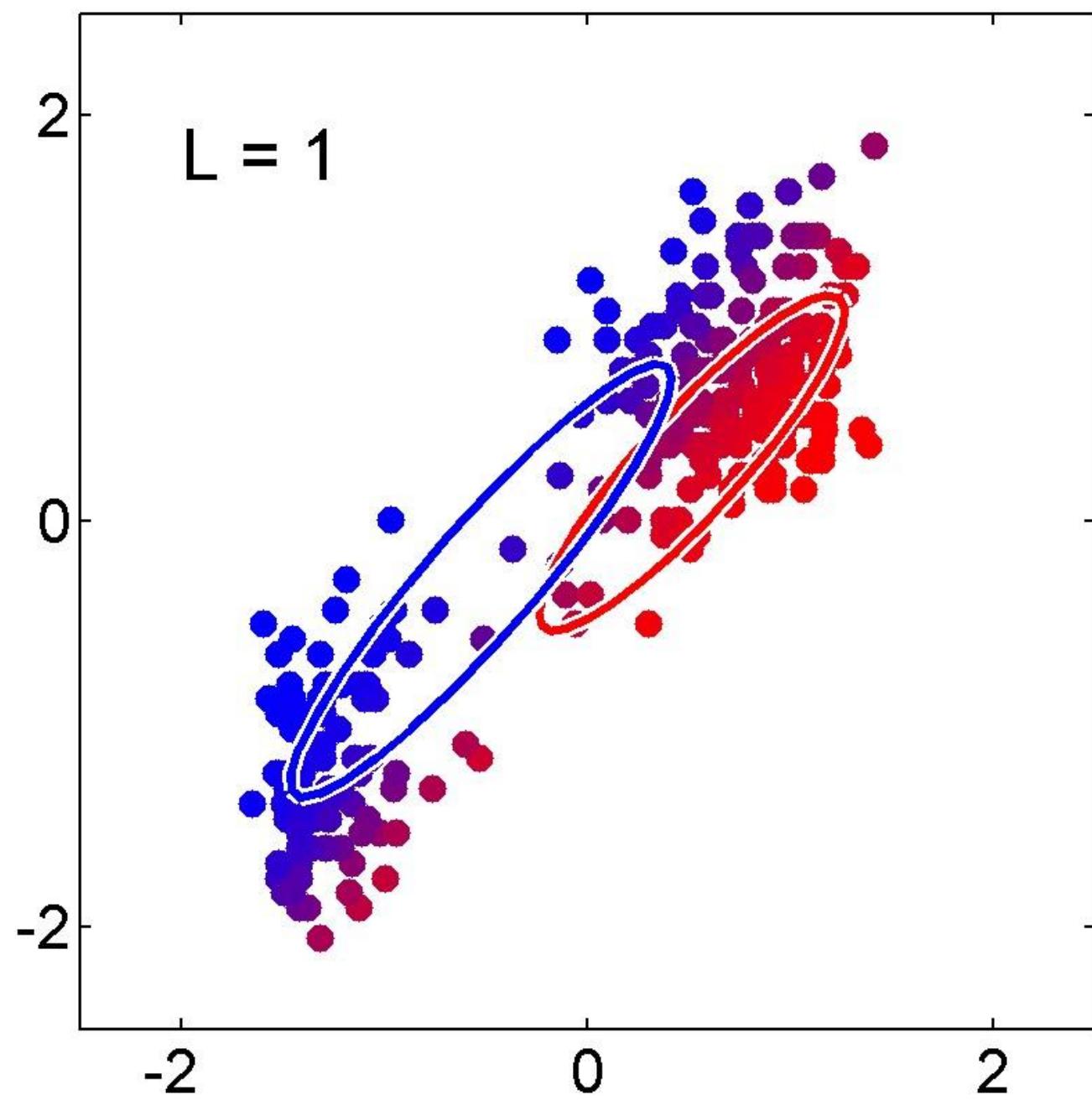
EM ALGORITHM – INFORMAL DERIVATION

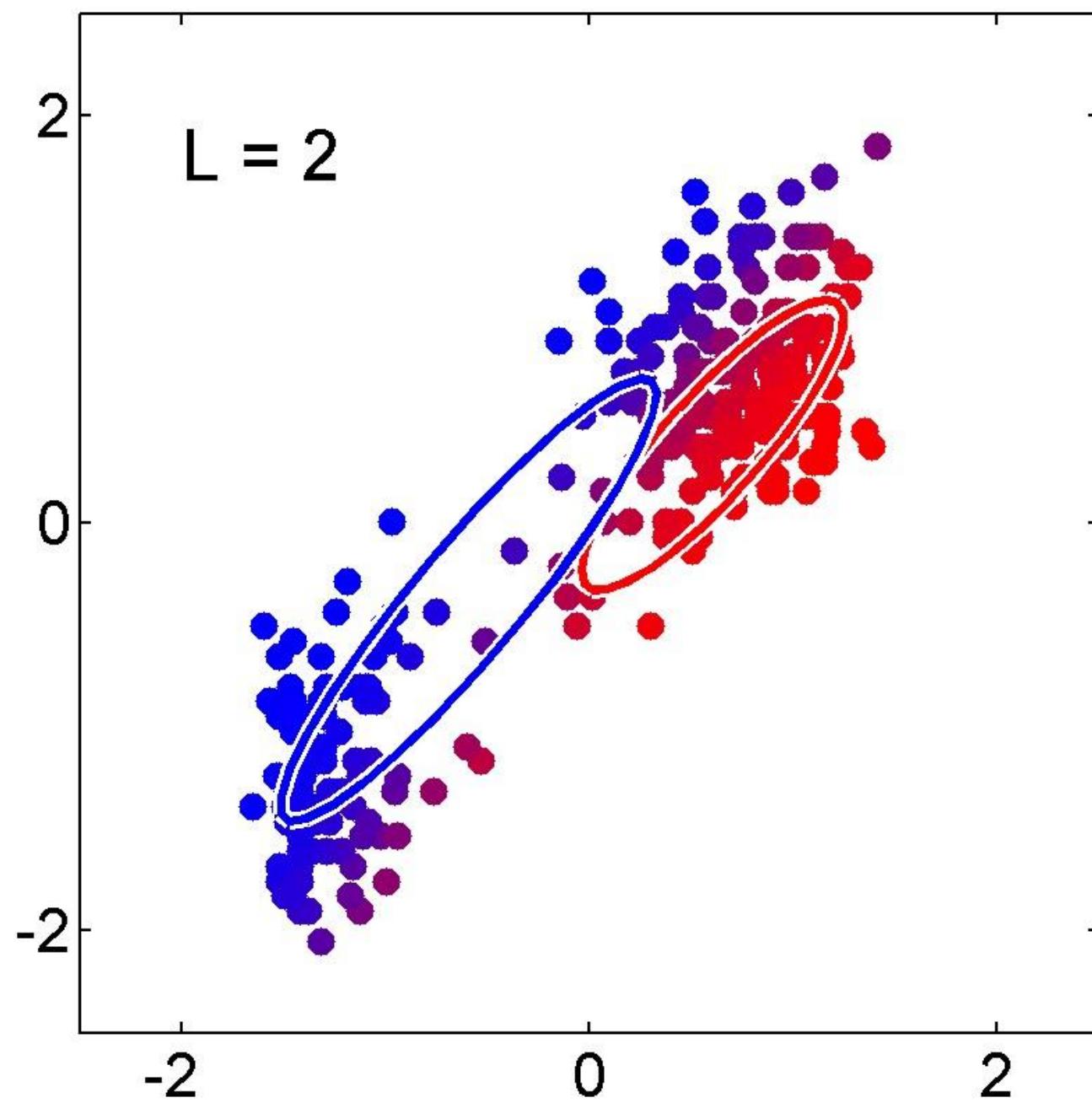
- The solutions are not closed form since they are coupled
- Suggests an iterative scheme for solving them:
 - Make initial guesses for the parameters
 - Alternate between the following two stages:
 1. E-step: evaluate responsibilities
 2. M-step: update parameters using ML results

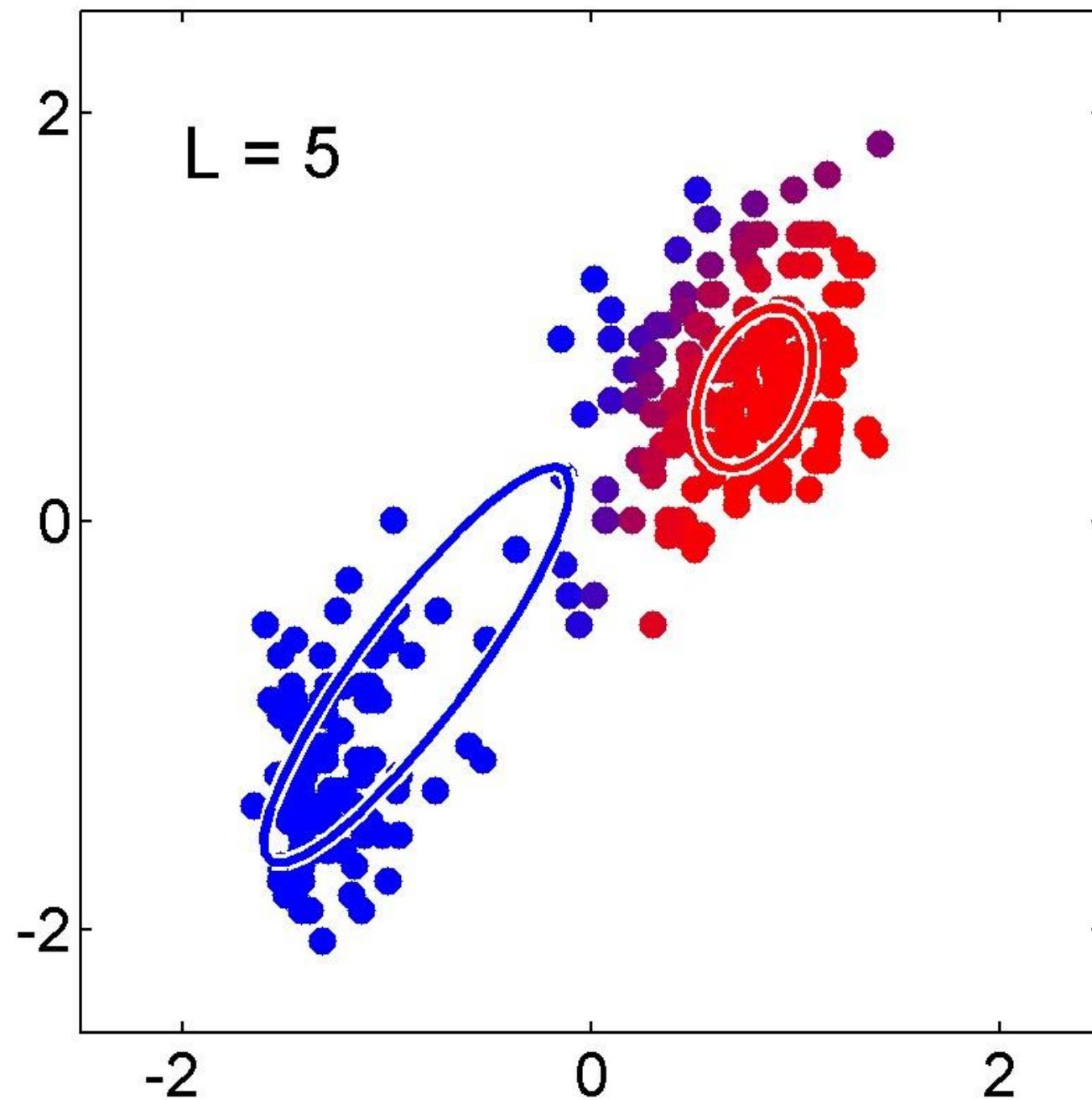


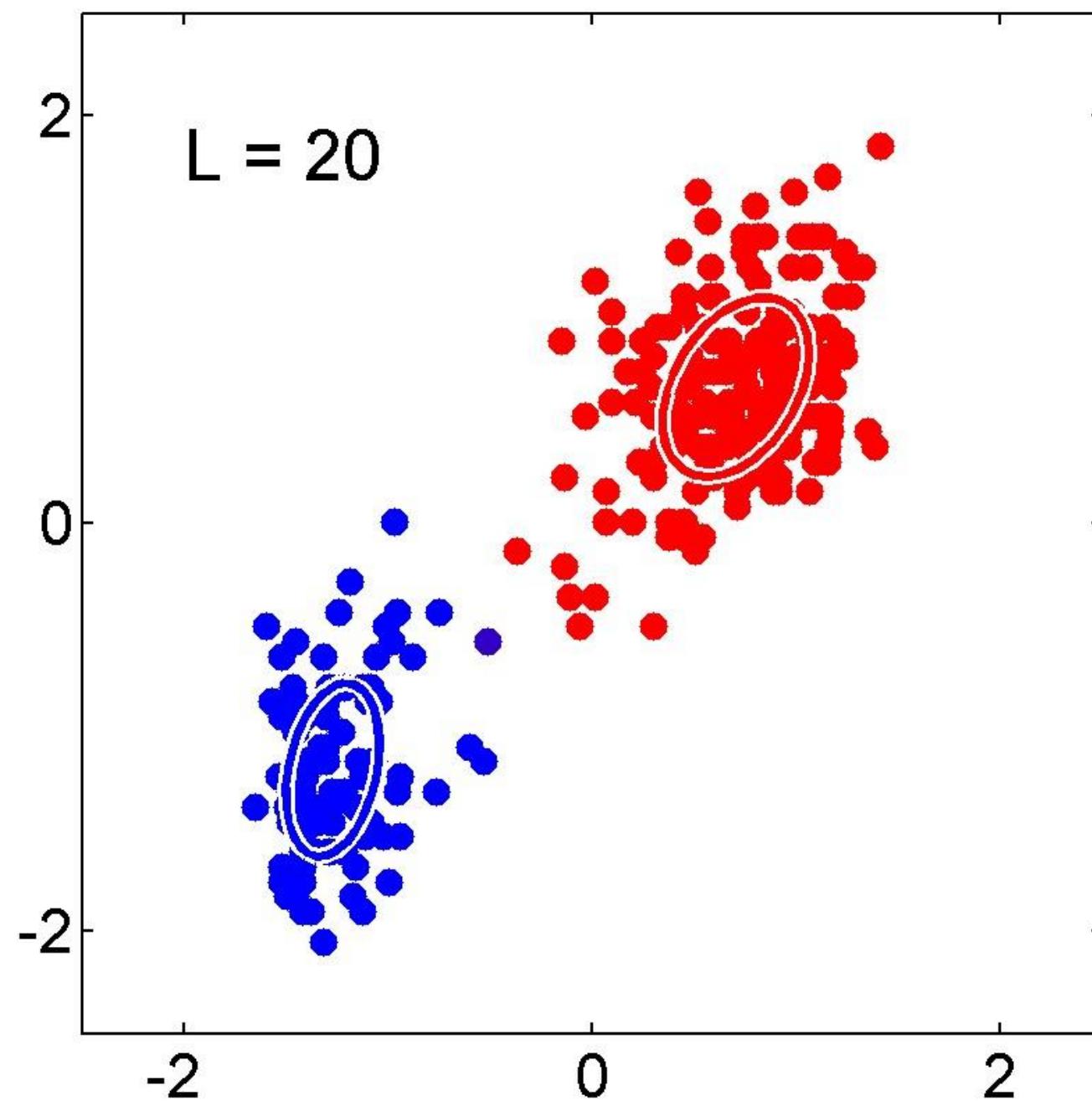






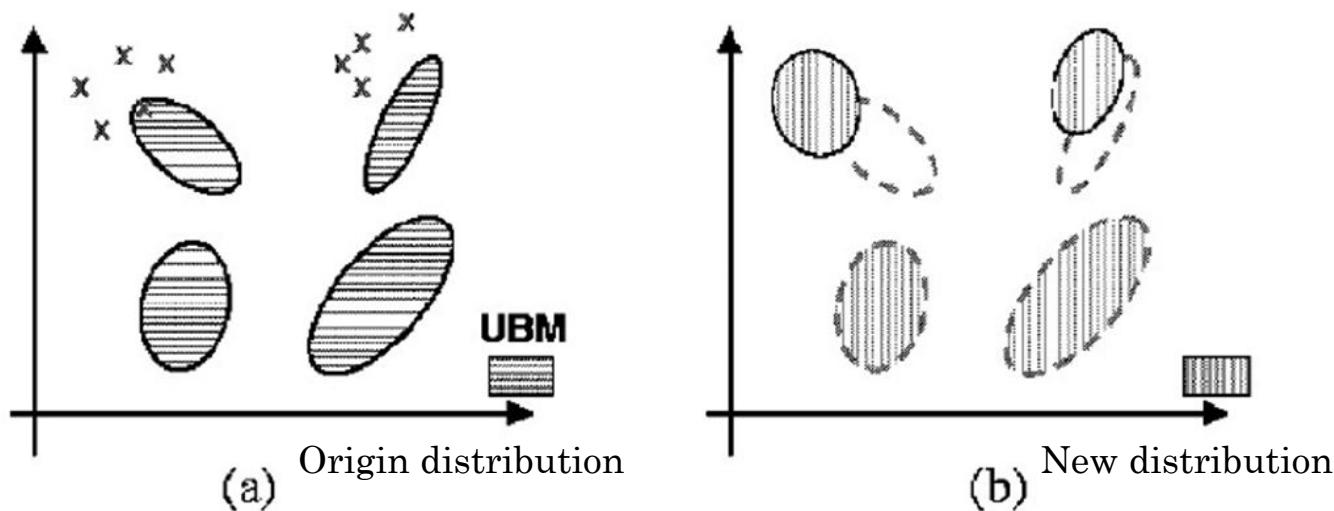






THE SUPERVECTOR REPRESENTATION (1)

- Given a set of features from a set of images, train a Gaussian mixture model. This is called a Universal Background Model.
- Given the UBM and a set of features from a single image, adapt the UBM to the image feature set by Bayesian EM (check equations in paper below).



THE SUPERVECTOR REPRESENTATION (2)

- Supervector:

$$[\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_K]$$

- Normalized supervector:

$$[\sqrt{\omega_1} \Sigma_1^{-1/2} \mathbf{m}_1, \sqrt{\omega_2} \Sigma_2^{-1/2} \mathbf{m}_2, \dots, \sqrt{\omega_K} \Sigma_K^{-1/2} \mathbf{m}_K]$$

- Centralized and normalized supervector

$$\left[\sqrt{\omega_1} \Sigma_1^{-\frac{1}{2}} (\mathbf{m}_1 - \boldsymbol{\mu}_1), \sqrt{\omega_2} \Sigma_2^{-1/2} (\mathbf{m}_2 - \boldsymbol{\mu}_2), \dots, \sqrt{\omega_K} \Sigma_K^{-\frac{1}{2}} (\mathbf{m}_K - \boldsymbol{\mu}_K) \right]$$

