

Chapter 6: Generative Probabilistic Models in Multimedia Retrieval

Sami Virpioja
sami.virpioja@tkk.fi

Adaptive Informatics Research Centre
Helsinki University of Technology

February 22, 2008

Outline

- 1 Introduction
- 2 Generative Image Models
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities
- 5 Conclusions

Outline

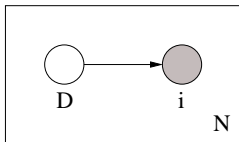
- 1 Introduction
- 2 Generative Image Models
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities
- 5 Conclusions

Classification Tasks in MR

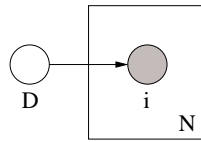
- Many multimedia retrieval tasks can be seen as problems of decision theory:
 - Classification into finite set of classes
(e.g. motorbike or bicycles or people or car)
 - Detection of objects or events (true or false)
 - Retrieval tasks (relevant or not relevant)
- Probabilistic approach is a very natural one:
Given a sample x , return the most probable class c

Graphical Notation for Probabilistic Models

- Nodes are variables
 - Observed variables with solid shading
- Edges are dependencies
- Repetition indicated by boxes



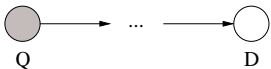
Sampling with repeated choice of dice



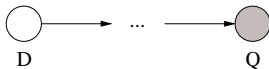
Sampling from a single dice

Discriminative vs. Generative Classification

- Goal: Find the best document $d \in \{d_1, d_2, \dots, d_n\}$ given the query q
- Discriminative classification:
 - Choose the most probable document given the query
 - Build a model for $P(D = d_i | Q = q)$
 - Typical methods: LDA, SVM, perceptrons
- Generative classification:
 - Choose the document that most likely generated the query
 - Build a model for $P(Q = q | D = d_i)$
 - Allows sampling from the joint distribution $P(Q, D)$



Discriminative classification



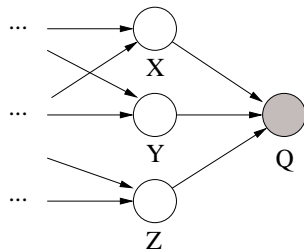
Generative classification

Jigsaw Example

- A set of jigsaw puzzles in their boxes
- One piece is found outside the boxes
- How to select the correct box (without solving all the puzzles)?
- Clue: appearance (color, texture, etc.) of the piece compared to the images of the solved puzzles
- Select the puzzle that maximizes
$$P(\text{piece comes from the puzzle}) = P(\text{piece} | \text{puzzle})$$

Generative Models

- Generative models can be thought as sources that generate samples from a probability distribution
 - Jigsaw example: puzzles = sources, pieces = samples
 - Limited number of sample types → discrete distribution
 - Unlimited sample types → continuous distribution
-
- Samples are from observed variables; a number of hidden variables may affect them
 - Model class \mathcal{M} determines the model structure; parameters θ differentiate the models within a class



Retrieval with Generative Models

- A set of models of class \mathcal{M} with parameters $\theta_1, \theta_2, \dots, \theta_N$
- An observation O is generated from one of the models
- Models can be ranked by sample likelihood $P(O | \theta_i)$
- If observation is a set of samples $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_S\}$ (e.g. words, visual features) and the models are memoryless:

$$P(O | \theta_i) = \prod_{j=1}^S P(\mathbf{v}_j | \theta_i)$$

Estimating Model Parameters - ML

- A set of training samples O from a model of class \mathcal{M}
- How to estimate parameters θ for the model?
- Maximum likelihood (ML) estimate
 - Select such θ that the probability of the training data is maximized
 - Overlearning may be a problem

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} P(O | \theta) \quad (1)$$

Estimating Model Parameters - MAP

- Maximum a posteriori (MAP) estimate
 - Select such θ that the probability of the model given the training data is maximized
 - Model prior $P(\theta)$ is needed

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P(\theta | O) = \arg \max_{\theta} P(\theta)P(O | \theta) \quad (2)$$

Outline

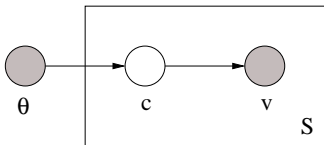
- 1 Introduction
- 2 **Generative Image Models**
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities
- 5 Conclusions

Mixture Models

- Mixture model is a class of generative models
- We assume that samples v come from a mixture of C independent sources:

$$P(v|\theta) = \sum_{i=1}^C P(c_i|\theta)P(v|c_i, \theta) \quad (3)$$

- θ includes the prior probabilities $P(c_i|\theta)$ and parameters of the density functions $P(v|c_i, \theta)$ of the sources c_i



Sampling with a mixture model

Gaussian Mixture Models

- If sources generate samples from normal distributions, we have a Gaussian Mixture Model (GMM)
- Normal distribution is useful approximation especially when the observed variable is affected by many independent variables (Central Limit theorem)
- Observed samples come from a set of different objects/categories → GMM
 - Direct application: C categories really exist
 - Indirect application: just a tool for get a tractable model
 - By increasing the number of mixtures, we can go arbitrarily close to the any distribution

GMM for Images

- Assume that a image are generated as follows:
 - Take a GMM with parameters $\theta = (\pi, \mu, \Sigma)$
 - For each sample v in the image
 - Pick a random component c_i according to the prior distribution $P(c_i|\theta) = \pi_i$
 - Draw a random sample from c_i according to $N(\mu_i, \Sigma_i)$
- Probability of a given sample is thus

$$\begin{aligned} P(v|\theta) &= \sum_{i=1}^C P(c_i|\theta) \times P(v|c_i, \theta) \\ &= \sum_{i=1}^C \pi_i \times \frac{1}{\sqrt{(2\pi)^n |\Sigma_i|}} \exp\left(-\frac{1}{2}(v - \mu_i)^T \Sigma_i (v - \mu_i)\right) \end{aligned}$$

Parameter Estimation for GMMs

- GMM parameters $\theta = (\pi, \mu, \Sigma)$
- Maximum likelihood estimation from a set of samples $O = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_S\}$:

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \arg \max_{\theta} P(O | \theta) \\ &= \arg \max_{\theta} \prod_{j=1}^S P(\mathbf{v}_j | \theta) \\ &= \arg \max_{(\pi, \mu, \Sigma)} \prod_{j=1}^S \sum_{i=1}^C \pi_i |\Sigma_i|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{v}_j - \mu_i)^T \Sigma_i (\mathbf{v}_j - \mu_i)\right)\end{aligned}$$

EM Algorithm

- No analytical solution for θ_{ML}
- Expectation-Maximization (EM) algorithm finds a local optimum iteratively
- A hidden variable \mathbf{H} gives assignments of samples to components: $h_{ij} = P(c_i | \mathbf{v}_j)$
- Initialize \mathbf{H} randomly
- Iterate selecting ML estimates for parameters θ (M-step) and calculating new expectation values for \mathbf{H} (E-step)
- Iterate until likelihood of the training data $P(O | \theta)$ stops increasing

EM Algorithm

- E-step:

$$h_{ij} = \frac{p(\mathbf{v}_j | c_i) \pi_i}{\sum_{c=1}^C p(\mathbf{v}_j | c_c) \pi_c} \quad (4)$$

- M-step:

$$\boldsymbol{\mu}_i = \frac{\sum_j h_{ij} \mathbf{v}_j}{\sum_j h_{ij}} \quad (5)$$

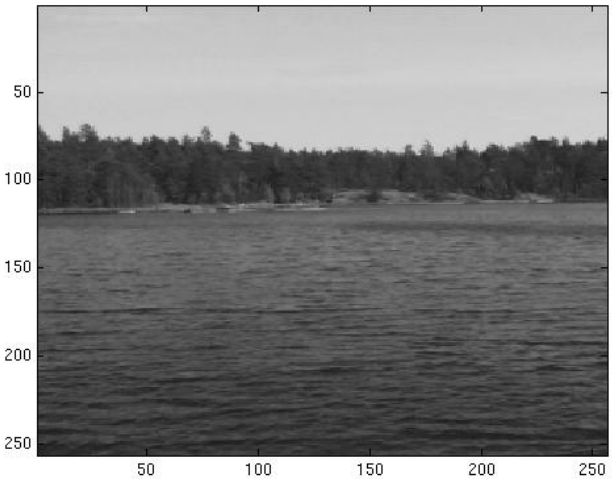
$$\boldsymbol{\Sigma}_i = \frac{\sum_j h_{ij} (\mathbf{v}_j - \boldsymbol{\mu}_i)(\mathbf{v}_j - \boldsymbol{\mu}_i)^T}{\sum_j h_{ij}} \quad (6)$$

$$\pi_i = \frac{1}{S} \sum_j h_{ij} \quad (7)$$

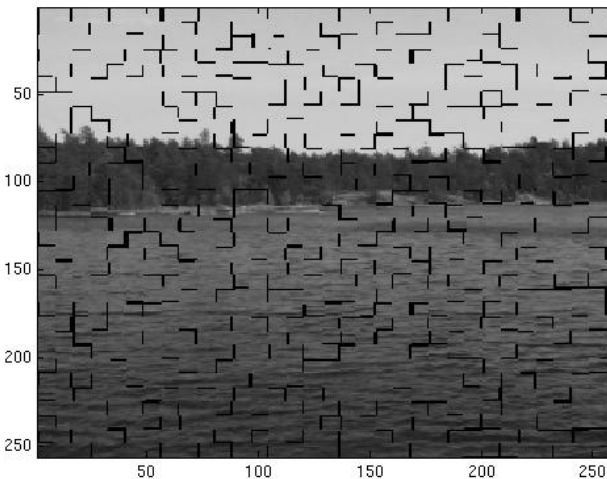
Example: GMM for a Single Image

- 256×256 grayscale image
- Samples: 8×8 blocks uniformly from the image
- Features: (x, y) coordinates of the middle point, DCT transform of the pixel intensities

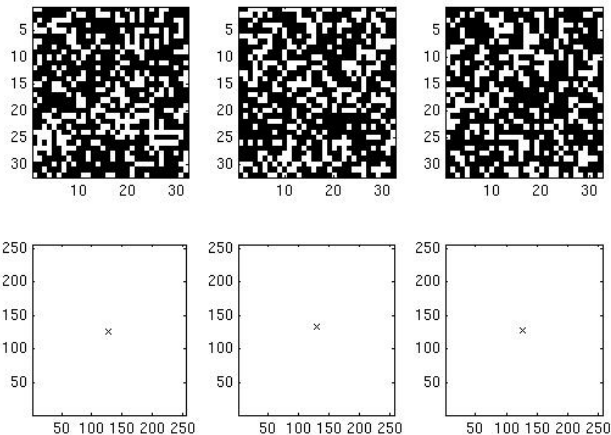
Example: GMM for a Single Image



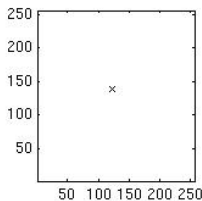
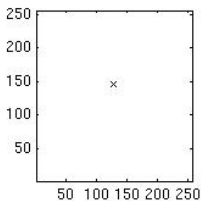
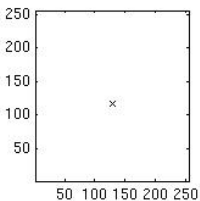
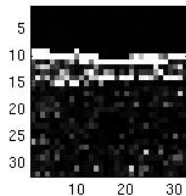
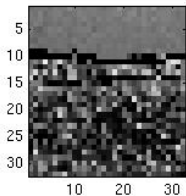
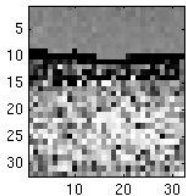
Example: GMM for a Single Image



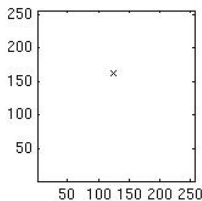
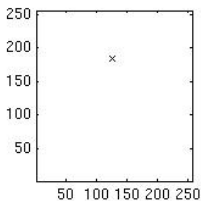
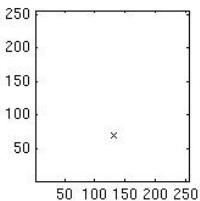
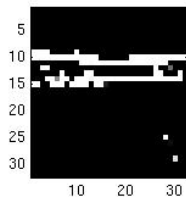
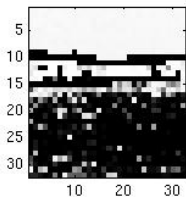
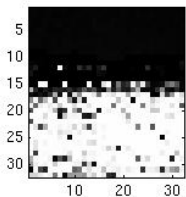
Example: GMM for a Single Image



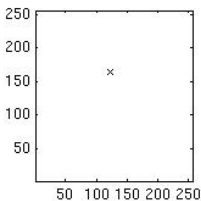
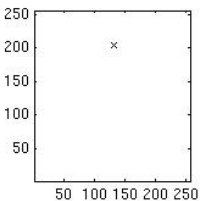
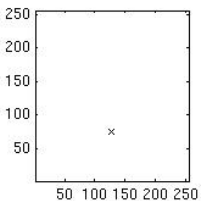
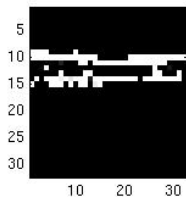
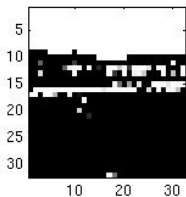
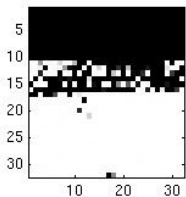
Example: GMM for a Single Image



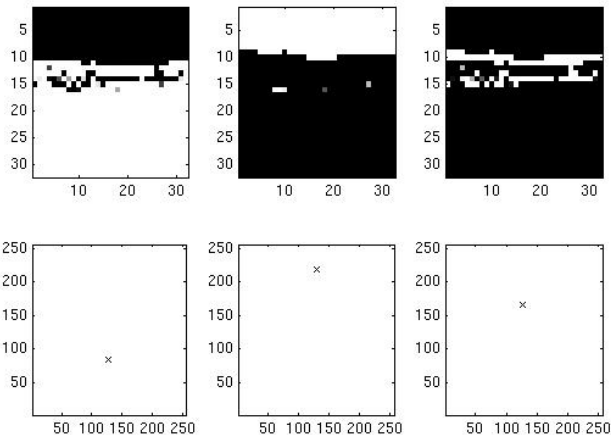
Example: GMM for a Single Image



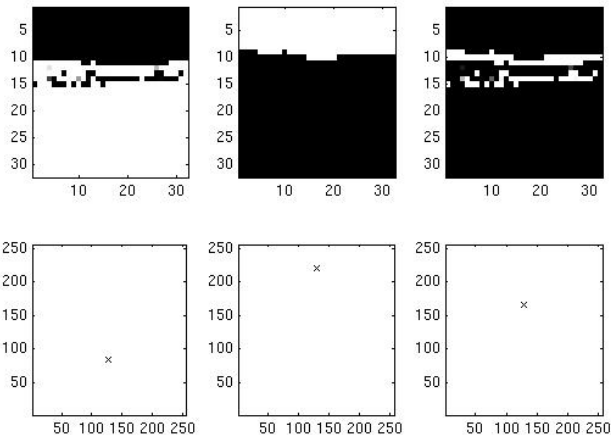
Example: GMM for a Single Image



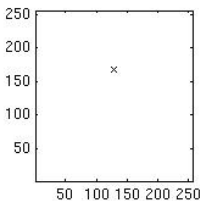
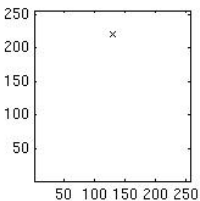
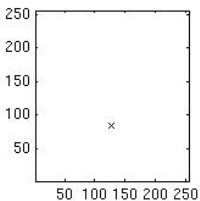
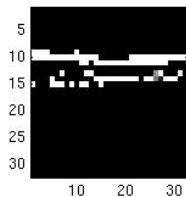
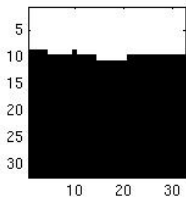
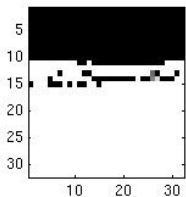
Example: GMM for a Single Image



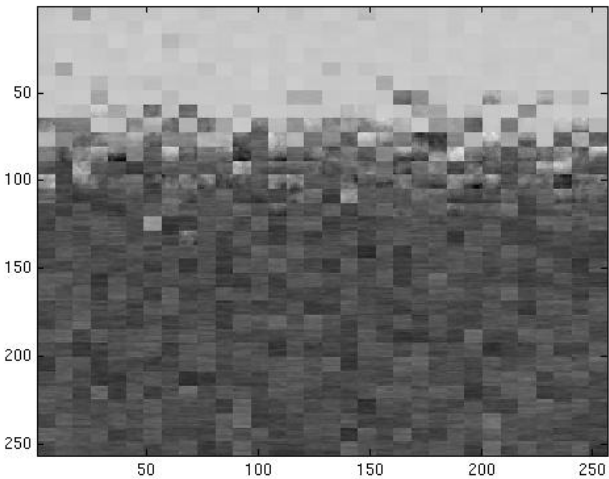
Example: GMM for a Single Image



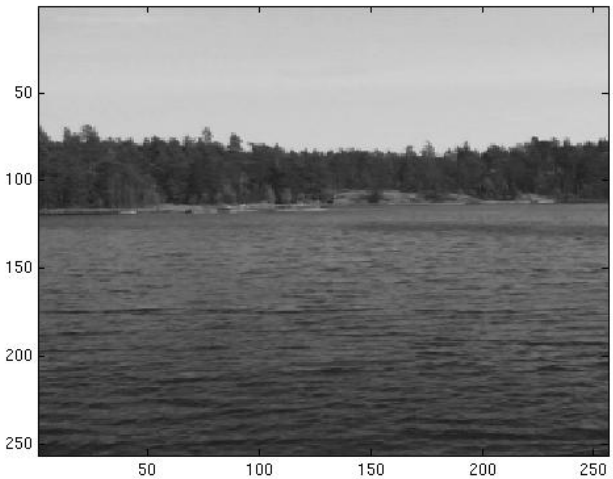
Example: GMM for a Single Image



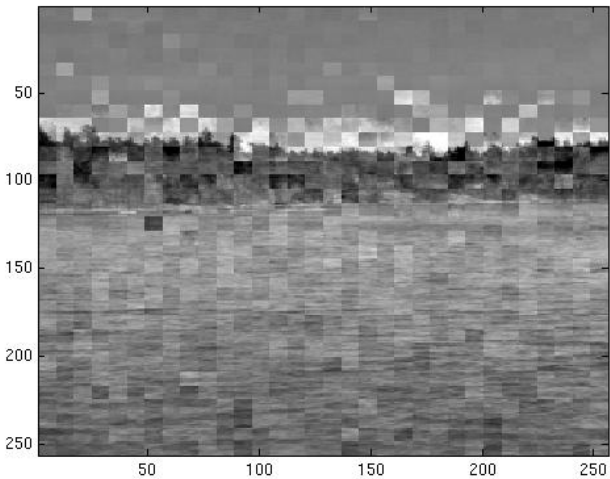
Example: GMM for a Single Image



Example: GMM for a Single Image



Example: GMM for a Single Image



Outline

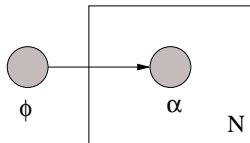
- 1 Introduction
- 2 Generative Image Models
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities
- 5 Conclusions

Statistical Language Models

- Statistical language model gives a probability distribution $P(S)$ over texts strings S of a language
- Usually based on words: $P(w_1, w_2, \dots, w_N)$
- N-gram assumption:
$$P(w_i | w_1, \dots, w_{i-2}, w_{i-1}) = P(w_i | w_{i-n+1}, \dots, w_{i-2}, w_{i-1})$$
- n -gram model is a $(n - 1)$:th order Markov model
- High-order n -gram models are needed in speech recognition and statistical machine translation
- Unigram (1-gram) model seems to be enough for information retrieval

Unigram Model for Text Retrieval

- A lexicon of T words, α_i is the i :th word in lexicon
- Model parameters $\phi = (\phi_1, \phi_2, \dots, \phi_T)$
- $P(w = \alpha_i | \phi) = \phi_i$
- Generation process: For each term of the document, draw a random term from ϕ according to multinomial distribution $\text{mult}(\phi)$



Sampling with a unigram model

Sample Likelihood for Unigram Model

- Sampled document is denoted as $\mathbf{t} = (t_1, t_2, \dots, t_T)$, where t_i is the number of occurrences of α_i
- Probability of \mathbf{t} given the model ϕ is

$$P(\mathbf{t} | \phi) = \frac{(\sum_{i=1}^T t_i)!}{\prod_{i=1}^T t_i!} \prod_{i=1}^T \phi_i^{t_i} \quad (8)$$

- $\prod_{i=1}^T \phi_i^{t_i}$ gives the probability of words in one order
- The messy first term gives the number of permutations for the set of $N = \sum_{i=1}^T t_i$ words (bag-of-words model)

ML Estimate for Unigram Parameters

- Let's find the ML estimate for the parameters:

$$\begin{aligned}\hat{\phi}_{\text{ML}} &= \arg \max_{\phi} P(\mathbf{t} | \phi) \\ &= \arg \max_{\phi} \frac{(\sum_{i=1}^T t_i)!}{\prod_{i=1}^T t_i!} \prod_{i=1}^T \phi_i^{t_i} \\ &= \arg \max_{\phi} \prod_{i=1}^T \phi_i^{t_i} \\ &= \arg \min_{\phi} \left[- \sum_{i=1}^T t_i \log \phi_i \right]\end{aligned}$$

Unigram Parameters

- $-\sum_{i=1}^T t_i \log \phi_i$ is (empirical) cross-entropy of the source that outputs words according to distribution ϕ
- Minimized with $\phi_i = \frac{t_i}{\sum_i t_i} = \frac{t_i}{N}$
- Thus ML estimates are just relative frequencies of words in the document

Smoothing and Interpolation

- However, the ML estimate has two fundamental problems
 - Words that did not occur in the document will have zero probabilities
 - Words that occurred only few times will have overestimated probabilities
- Traditional solutions:
 - **Smoothing** explicitly moves probability mass from the low frequencies to the zero frequencies (more or less heuristically)
 - **Interpolation** or **back-off** with/to more general (e.g. lower order) models
- Here we discuss only interpolation

Linear Interpolation with a Background Model

- Make a background model ϕ_{BG} using all the documents as training data
- Then interpolate linearly with a coefficient λ :

$$\phi_i = \lambda\phi_{i\text{ML}} + (1 - \lambda)\phi_{i\text{BG}} \quad (9)$$

- If $t_{d,i}$ is the number of times α_i occurred in document d ,
 $\phi_{i\text{BG}} = \sum_d t_{d,i} / \sum_d \sum_j t_{d,j}$
- Or with document frequencies: $\phi_{i\text{BG}} = df_i / \sum_j df_j$

Relation to Weighting with *tf.idf*

- Let's try this in text retrieval for a query t
- With some of manipulation (see p. 192 in book), we get

$$P(\mathbf{t} | \phi) \propto \prod_{t_i > 0} \left[\frac{\lambda \phi_{iML}}{(1 - \lambda) \phi_{iBG}} + 1 \right]^{t_i} \quad (10)$$

- Background probabilities ϕ_{BG} weight down words that are common in the collection, as the *idf* component does in term weighting

Hierarchical Interpolation

- Even more models can be combined with linear interpolation
- E.g. for a video subdivided into scenes which are subdivided into shots,

$$\phi_i = \lambda_{\text{Shot}}\phi_{i\text{Shot}} + \lambda_{\text{Scene}}\phi_{i\text{Scene}} + \lambda_{\text{Coll}}\phi_{i\text{Coll}} \quad (11)$$

where $\lambda_{\text{Shot}} + \lambda_{\text{Scene}} + \lambda_{\text{Coll}} = 1$

- λ :s can be estimated from a separate development set
- Interpolation can be used as well for generative image or video models

Outline

- 1 Introduction
- 2 Generative Image Models
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities**
- 5 Conclusions

Combining Multiple Modalities

- Generative models for different modalities can naturally be combined
- E.g. multimedia query $\mathbf{q} = (\mathbf{v}, \mathbf{t})$ with visual part \mathbf{v} and textual part \mathbf{t}
- Rank documents by joint probability $P(\mathbf{q} | d) = P(\mathbf{v}, \mathbf{t} | d)$
- Assuming that \mathbf{v} and \mathbf{t} are independent
$$P(\mathbf{q} | d) = P(\mathbf{v} | \theta_d)P(\mathbf{t} | \phi_d)$$
- The independence assumption is of course totally unrealistic
- However, may work well enough in practice

Outline

- 1 Introduction
- 2 Generative Image Models
 - Gaussian Mixture Model for Images
 - Parameter Estimation
 - Example
- 3 Generative Language Models
 - Language Models for Text Retrieval
 - Parameter Estimation
 - Interpolation
- 4 Combining Multiple Modalities
- 5 Conclusions

Conclusions

- In IR, generative models describe the process of generating samples from documents
- Given the query, documents can be ranked with the sample likelihood of their models
- Two examples were given:
 - Gaussian Mixture Model for images
 - EM algorithm for training
 - Unigram model for text documents
 - Direct ML estimates, but smoothing/interpolation required