Computer vision: models, learning and inference

Chapter 7 Modeling Complex Densities

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Models for machine vision



Table 5.1: Example models in this chapter. These can be categorized into those that arebased on modelling probability density functions, those that are based on linearregression and those that are based on logistic regression.

Face Detection



Type 3: Pr(x|w) - Generative

How to model Pr(**x**|**w**)?

- Choose an appropriate form for Pr(x)
- Make parameters a function of w
- Function takes parameters θ that define its shape

Learning algorithm: learn parameters θ from training data x,w Inference algorithm: Define prior Pr(w) and then compute Pr(w|x) using Bayes' rule

$$Pr(w=1|\mathbf{x}) = \frac{Pr(\mathbf{x}|w=1)Pr(w=1)}{\sum_{k=0}^{1} Pr(\mathbf{x}|w=k)Pr(w=k)}$$

Classification Model

$$Pr(\mathbf{x}|w) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w]$$

Or writing in terms of class conditional density functions $Pr(\mathbf{x}|w=0) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}]$ $Pr(\mathbf{x}|w=1) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_{1}, \boldsymbol{\Sigma}_{1}]$

Parameters μ_0 , Σ_0 learnt just from data S_0 where w=0

$$\hat{\boldsymbol{\mu}}_{0}, \hat{\boldsymbol{\Sigma}}_{0} = \operatorname{argmax}_{\boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}} \left[\prod_{i \in \mathcal{S}_{0}} Pr(\mathbf{x}_{i} | \boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}) \right]$$
$$= \operatorname{argmax}_{\boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}} \left[\prod_{i \in \mathcal{S}_{0}} \operatorname{Norm}_{\mathbf{x}_{i}}[\boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}] \right]$$

Similarly, parameters μ_1 , Σ_1 learnt just from data S_1 where w=1





Inference algorithm: Define prior Pr(w) and then compute Pr(w|x) using Bayes' rule

$$Pr(\mathbf{w}|\mathbf{x}) = \frac{Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})}{\int Pr(\mathbf{x}|\mathbf{w})Pr(\mathbf{w})d\mathbf{w}}$$

Experiment

1000 non-faces 1000 faces

60x60x3 Images =10800 x1 vectors

Equal priors Pr(y=1)=Pr(y=0) = 0.5

75% performance on test set. Not very good!



Results (diagonal covariance)



Figure 7.2 Class conditional density functions for normal model with diagonal covariance. Maximum likelihood fits based on 1000 training examples per class. a) Mean for background data μ_0 (reshaped from 10800×1 vector to 60×60 RGB image). b) Reshaped square root of diagonal covariance for background data Σ_0 . c) Mean for face data μ_1 d) Covariance for face data Σ_1 . The background model has little structure: the mean is uniform and the variance is high everywhere. The mean of the face model clearly captures class-specific information. The covariance of the face is larger at the edges of the image which usually contain hair or background.

The plan... b) Problem 1 Mixture e) models Unimodal Mixture of t-distributions mixture of Gaussians Normal distribution a) c) Mixture of factor analyzers Problem 2 Robust models Sensitive to outliers Mixture of robust subspace models t-distributions d) Robust subspace Problem 3 models Subspace models Too many parameters in high dimensions PPCA, factor analysis

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Hidden (or latent) Variables

Key idea: represent density Pr(x) as marginalization of joint density with another variable h that we do not see

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}, \mathbf{h}) \, d\mathbf{h}$$

Will also depend on some parameters:

$$Pr(\mathbf{x}|\boldsymbol{\theta}) = \int Pr(\mathbf{x}, \mathbf{h}|\boldsymbol{\theta}) \, d\mathbf{h}$$

Hidden (or latent) Variables $Pr(\mathbf{x}|\boldsymbol{\theta}) = \int Pr(\mathbf{x}, \mathbf{h}|\boldsymbol{\theta}) d\mathbf{h}$



Figure 7.4 Using hidden variables to help model complex densities. One way to model the density Pr(x) is to consider the joint probability distribution Pr(x, h) between the observed data x and a hidden variable h. The density Pr(x) can be considered as the marginalization of (integral over) this distribution with respect to the hidden variable h. As we manipulate the parameters $\boldsymbol{\theta}$ of this joint distribution, the marginal changes and the agreements with the observed data $\{x_i\}_{i=1}^{I}$ increases or decreases. Sometimes it is easier to fit the distribution in this indirect way than to directly manipulate Pr(x).

Expectation Maximization

An algorithm specialized to fitting pdfs which are the marginalization of a joint distribution

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left[\sum_{i=1}^{I} \log \left[\int Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) \ d\mathbf{h}_i \right] \right]$$

Defines a lower bound on log likelihood and increases bound iteratively

$$\mathcal{B}[\{q_i(\mathbf{h}_i)\}, \boldsymbol{\theta}] = \sum_{i=1}^{I} \int q_i(\mathbf{h}_i) \log\left[\frac{Pr(\mathbf{x}, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)}\right] d\mathbf{h}_{1...I}$$
$$\leq \sum_{i=1}^{I} \log\left[\int Pr(\mathbf{x}, \mathbf{h} | \boldsymbol{\theta}) d\mathbf{h}_{1...I}\right].$$

Lower bound

$$\mathcal{B}[\{q_i(\mathbf{h}_i)\}, \boldsymbol{\theta}] = \sum_{i=1}^{I} \int q_i(\mathbf{h}_i) \log\left[\frac{Pr(\mathbf{x}, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)}\right] d\mathbf{h}_{1...I}$$

Lower bound is a *function* of parameters θ and a set of probability distributions $q_i(h_i)$

Expectation Maximization (EM) algorithm alternates Esteps and M-Steps

E-Step – Maximize bound w.r.t. distributions $q(\mathbf{h}_i)$ M-Step – Maximize bound w.r.t. parameters $\boldsymbol{\theta}$

Lower bound



E-Step & M-Step

E-Step – Maximize bound w.r.t. distributions $q_i(h_i)$

$$\hat{q}_i(\mathbf{h}_i) = Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta}^{[t]}) = \frac{Pr(\mathbf{x}_i | \mathbf{h}_i, \boldsymbol{\theta}^{[t]}) Pr(\mathbf{h}_i | \boldsymbol{\theta}^{[t]})}{Pr(\mathbf{x}_i)}$$

M-Step – Maximize bound w.r.t. parameters θ

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left[\sum_{i=1}^{I} \int \hat{q}_i(\mathbf{h}_i) \log \left[Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) \right] \, d\mathbf{h}_i \right]$$

E-Step & M-Step



Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Mixture of Gaussians (MoG)



ML Learning

Q. Can we learn this model with maximum likelihood?

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \log \left[Pr(\mathbf{x}_{i} | \boldsymbol{\theta}) \right] \right]$$
$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \log \left[\sum_{k=1}^{K} \lambda_{k} \operatorname{Norm}_{\mathbf{x}_{i}}[\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}] \right] \right]$$

- A. Yes, but using brute force approach is tricky
 - If you take derivative and set to zero, can't solve
 -- the log of the sum causes problems
 - Have to enforce constraints on parameters
 - covariances must be positive definite
 - weights must sum to one

MoG as a marginalization

Define a variable $h \in \{1 \dots K\}$ and then write

$$Pr(\mathbf{x}|h, \boldsymbol{\theta}) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_{h}, \boldsymbol{\Sigma}_{h}]$$
$$Pr(h|\boldsymbol{\theta}) = \operatorname{Cat}_{h}[\boldsymbol{\lambda}]$$

Then we can recover the density by marginalizing Pr(x,h)

$$Pr(\mathbf{x}|\boldsymbol{\theta}) = \sum_{k=1}^{K} Pr(\mathbf{x}, h = k|\boldsymbol{\theta})$$
$$= \sum_{k=1}^{K} Pr(\mathbf{x}|h = k, \boldsymbol{\theta}) Pr(h = k|\boldsymbol{\theta})$$
$$= \sum_{k=1}^{K} \lambda_k \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k].$$



MoG as a marginalization

Define a variable $h \in \{1 \dots K\}$ and then write

$$Pr(\mathbf{x}|h, \boldsymbol{\theta}) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}_{h}, \boldsymbol{\Sigma}_{h}]$$
$$Pr(h|\boldsymbol{\theta}) = \operatorname{Cat}_{h}[\boldsymbol{\lambda}]$$

Note :

- This gives us a method to generate data from MoG
 - First sample Pr(h), then sample Pr(x|h)
- The hidden variable h has a clear interpretation it tells you which Gaussian created data point x

Expectation Maximization for MoG

GOAL: to learn parameters $\theta = \{\lambda_{1...K}, \mu_{1...K}, \Sigma_{1...K}\}$ from training data $\mathbf{x}_{1...I}$

E-Step – Maximize bound w.r.t. distributions q(h_i)

$$\hat{q}_i(\mathbf{h}_i) = Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta}^{[t]}) = \frac{Pr(\mathbf{x}_i | \mathbf{h}_i, \boldsymbol{\theta}^{[t]}) Pr(\mathbf{h}_i | \boldsymbol{\theta}^{[t]})}{Pr(\mathbf{x}_i)}$$

M-Step – Maximize bound w.r.t. parameters θ

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \sum_{k=1}^{K} \hat{q}_i(\mathbf{h}_i = k) \log \left[Pr(\mathbf{x}_i, \mathbf{h}_i = k | \boldsymbol{\theta}) \right] \right]$$

E-Step



Computer vision: models, learning and inference. ©2011 Simon J.D. Prince

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \sum_{k=1}^{K} \hat{q}_i(h_i = k) \log \left[Pr(\mathbf{x}_i, h_i = k | \boldsymbol{\theta}) \right] \right]$$
$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \sum_{k=1}^{K} r_{ik} \log \left[\lambda_k \operatorname{Norm}_{\mathbf{x}_i}[\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k] \right] \right].$$

Take derivative, equate to zero and solve (Lagrange multipliers for λ)

$$\lambda_{k}^{[t+1]} = \frac{\sum_{i=1}^{I} r_{ik}}{\sum_{j=1}^{K} \sum_{i=1}^{I} r_{ij}}$$
$$\mu_{k}^{[t+1]} = \frac{\sum_{i=1}^{I} r_{ik} \mathbf{x}_{i}}{\sum_{i=1}^{I} r_{ik}}$$
$$\boldsymbol{\Sigma}_{k}^{[t+1]} = \frac{\sum_{i=1}^{I} r_{ik} (\mathbf{x}_{i} - \boldsymbol{\mu}_{k}^{[t+1]}) (\mathbf{x}_{i} - \boldsymbol{\mu}_{k}^{[t+1]})^{T}}{\sum_{i=1}^{I} r_{ik}}$$



Update means, covariances and weights according to responsibilities of datapoints

Iterate until no further improvement



E-Step

M-Step



Different flavours...



FullDiagonalSamecovariancecovariancecovariance

Local Minima

Start from three random positions



Means of face/non-face model



Classification \rightarrow 84% (9% improvement!)

The plan... b) Problem 1 Mixture e) models Unimodal Mixture of t-distributions mixture of Gaussians Normal distribution a) c) Mixture of factor analyzers Problem 2 Robust models Sensitive to outliers Mixture of robust subspace models t-distributions d) Robust subspace Problem 3 models Subspace models Too many parameters in high dimensions PPCA, factor analysis

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Student t-distributions


Student t-distributions motivation

The normal distribution is not very robust – a single outlier can completely throw it off because the tails fall off so fast...



Normal distribution

Normal distribution w/ one extra datapoint!

t-distribution

Student t-distributions

Univariate student t-distribution

$$Pr(x) = \operatorname{Stud}_{\mathbf{x}} \left[\mu, \sigma^{2}, \nu\right]$$
$$= \frac{\Gamma\left[\frac{\nu+1}{2}\right]}{\sqrt{\nu\pi\sigma^{2}}\Gamma\left[\frac{\nu}{2}\right]} \left(1 + \frac{(x-\mu)^{2}}{\nu\sigma^{2}}\right)^{-\frac{\nu+1}{2}}$$

Multivariate student t-distribution

$$Pr(\mathbf{x}) = \operatorname{Stud}_{\mathbf{x}} [\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu]$$

=
$$\frac{\Gamma \left[\frac{\nu+D}{2}\right]}{(\nu\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2} \Gamma \left[\frac{\nu}{2}\right]} \left(1 + \frac{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}{\nu}\right)^{-\frac{\nu+D}{2}}$$

t-distribution as a marginalization

Define hidden variable h

$$Pr(\mathbf{x}|h) = \operatorname{Norm}_{x}[\boldsymbol{\mu}, \boldsymbol{\Sigma}/h]$$
$$Pr(h) = \operatorname{Gam}_{h}[\nu/2, \nu/2]$$

Can be expressed as a marginalization

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}, h) dh = \int Pr(\mathbf{x}|h) Pr(h) dh$$
$$= \int \operatorname{Norm}_{x}[\boldsymbol{\mu}, \boldsymbol{\Sigma}/h] \operatorname{Gam}_{h}[\nu/2, \nu/2] dh$$
$$= \operatorname{Stud}_{\mathbf{x}}[\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu].$$

Gamma distribution



t-distribution as a marginalization

Define hidden variable h

$$Pr(\mathbf{x}|h) = \operatorname{Norm}_{x}[\boldsymbol{\mu}, \boldsymbol{\Sigma}/h]$$
$$Pr(h) = \operatorname{Gam}_{h}[\nu/2, \nu/2]$$

Things to note:

- Again this provides a method to sample from the t-distribution
- Variable h has a clear interpretation:
 - Each datum drawn from a Gaussian, mean μ
 - Covariance depends inversely on h
- Can think of this as an infinite mixture (sum becomes integral) of Gaussians w/ same mean, but different variances

t-distribution



EM for t-distributions

GOAL: to learn parameters $\theta = \{\mu, \sigma^2, \nu\}$ from training data $\mathbf{x}_{1...I}$

E-Step – Maximize bound w.r.t. distributions q(h_i)

$$q_i(h_i) = Pr(h_i | \mathbf{x}_i, \boldsymbol{\theta}^{[t]}) = \frac{Pr(\mathbf{x}_i | h_i, \boldsymbol{\theta}^{[t]}) Pr(h_i)}{Pr(\mathbf{x}_i | \boldsymbol{\theta}^{[t]})}$$

M-Step – Maximize bound w.r.t. parameters θ

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{I} \int \hat{q}_i(h_i) \log \left[Pr(\mathbf{x}, h_i, \boldsymbol{\theta}) \right] dh_i$$

$$\begin{aligned} \mathbf{E}\text{-}\mathbf{Step} \\ q_i(h_i) &= Pr(h_i | \mathbf{x}_i, \boldsymbol{\theta}^{[t]}) &= \frac{Pr(\mathbf{x}_i | h_i, \boldsymbol{\theta}^{[t]}) Pr(h_i)}{Pr(\mathbf{x}_i | \boldsymbol{\theta}^{[t]})} \\ &= \frac{Norm_{\mathbf{x}_i}[\boldsymbol{\mu}, \boldsymbol{\Sigma}/h_i] Gam_{h_i}[\nu/2, \nu/2]}{Pr(\mathbf{x}_i)} \\ &= Gam_{h_i} \left[\frac{\nu + D}{2}, \frac{(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})}{2} + \frac{\nu}{2} \right] \end{aligned}$$

Extract expectations

$$E[h_i] = \frac{(\nu + D)}{\nu + (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}$$
$$E[\log[h_i]] = \Psi\left[\frac{\nu + D}{2}\right] - \log\left[\frac{\nu + (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}{2}\right]$$

$$\begin{split} \mathbf{M}-\mathbf{Step} \\ \hat{\boldsymbol{\theta}}^{[t+1]} &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int \hat{q}_{i}(h_{i}) \log \left[Pr(\mathbf{x}_{i}, h_{i} | \boldsymbol{\theta}) \right] dh_{i} \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int \hat{q}_{i}(h_{i}) \left(\log \left[Pr(\mathbf{x}_{i} | h_{i}, \boldsymbol{\theta}) \right] + \log \left[Pr(h_{i}) \right] \right) dh_{i} \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int Pr(h_{i} | \mathbf{x}_{i}, \boldsymbol{\theta}^{[t]}) \left(\log \left[Pr(\mathbf{x}_{i} | h_{i}, \boldsymbol{\theta}) \right] + \log \left[Pr(h_{i}) \right] \right) dh_{i} \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int Pr(h_{i} | \mathbf{x}_{i}, \boldsymbol{\theta}^{[t]}) \left(\log \left[Pr(\mathbf{x}_{i} | h_{i}, \boldsymbol{\theta}) \right] + \log \left[Pr(h_{i}) \right] \right) dh_{i} \right] \end{split}$$

Where...

$$E\left[\log\left[Pr(\mathbf{x}_{i}|h_{i},\boldsymbol{\theta})\right]\right] = \frac{D\mathrm{E}[\log h_{i}] - D\log 2\pi - \log|\boldsymbol{\Sigma}| - (\mathbf{x}_{i} - \boldsymbol{\mu})^{T}\boldsymbol{\Sigma}^{-1}(\mathbf{x}_{i} - \boldsymbol{\mu})\mathrm{E}[h_{i}]}{2}$$
$$E\left[\log\left[Pr(h_{i})\right]\right] = \frac{\nu}{2}\log\left[\frac{\nu}{2}\right] - \log\Gamma\left[\frac{\nu}{2}\right] + \left(\frac{\nu}{2} - 1\right)\mathrm{E}[\log h_{i}] - \frac{\nu}{2}\mathrm{E}[h_{i}].$$

Updates

$$\boldsymbol{\mu}^{[t+1]} = \frac{\sum_{i=1}^{I} \mathrm{E}[h_i] \mathbf{x}_i}{\sum_{i=1}^{I} \mathrm{E}[h_i]}$$
$$\boldsymbol{\Sigma}^{[t+1]} = \frac{\sum_{i=1}^{I} \mathrm{E}[h_i] (\mathbf{x}_i - \boldsymbol{\mu}^{[t+1]}) (\mathbf{x}_i - \boldsymbol{\mu}^{[t+1]})^T}{\sum_{i=1}^{I} \mathrm{E}[h_i]}$$

No closed form solution for v. Must optimize bound – since it is only one number we can optimize by just evaluating with a fine grid of values and picking the best

EM algorithm for t-distributions



Figure 7.18 Expectation maximization for fitting t-distributions. a) Estimate of distribution before update. b) In the E-Step we calculate the posterior distribution $Pr(h_i|x_i)$ over the hidden variable h_i for each data point x_i . The color of each curve corresponds to that of the original data point in (a). c) In the M-Step we use these distributions over h to update the estimate of the parameters $\theta = \{\mu, \sigma^2, \nu\}$.

The plan... b) Problem 1 Mixture e) models Unimodal Mixture of t-distributions mixture of Gaussians Normal distribution a) c) Mixture of factor analyzers Problem 2 Robust models Sensitive to outliers Mixture of robust subspace models t-distributions d) Robust subspace Problem 3 models Subspace models Too many parameters in high dimensions PPCA, factor analysis

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Factor Analysis

Compromise between

- Full covariance matrix (desirable but many parameters)
- Diagonal covariance matrix (no modelling of covariance)

Models full covariance in a subspace Φ Mops everything else up with diagonal component Σ

$$Pr(\mathbf{x}) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}, \boldsymbol{\Phi}\boldsymbol{\Phi}^T + \boldsymbol{\Sigma}]$$

Data Density



Consider modelling distribution of 100x100x3 RGB Images

Gaussian w/ spherical covariance:	1	covariance matrix parameters
Gaussian w/ diagonal covariance:	D _x	covariance matrix parameters
Full Gaussian:	$\sim D_x^2$	covariance matrix parameters
PPCA:	$\sim D_x D_h$	covariance matrix parameters
Factor Analysis:	$\sim D_x(D_h+1)$	covariance matrix parameters
Computer vision: models, learning and inference. ©2011 Simon J.D. Prince (full covariance in subspace + diagonal)		

Subspaces



Factor Analysis

Compromise between

- Full covariance matrix (desirable but many parameters)
- Diagonal covariance matrix (no modelling of covariance)

Models full covariance in a subspace ${f \Phi}$ Mops everything else up with diagonal component Σ

$$Pr(\mathbf{x}) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}, \boldsymbol{\Phi}\boldsymbol{\Phi}^T + \boldsymbol{\Sigma}]$$

Factor Analysis as a Marginalization

Let's define

$$Pr(\mathbf{x}|\mathbf{h}) = \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu} + \boldsymbol{\Phi}\mathbf{h}, \boldsymbol{\Sigma}]$$
$$Pr(\mathbf{h}) = \operatorname{Norm}_{\mathbf{h}}[\mathbf{0}, \mathbf{I}]$$

Then it can be shown (not obvious) that

$$Pr(\mathbf{x}) = \int Pr(\mathbf{x}, \mathbf{h}) dh = \int Pr(\mathbf{x}|\mathbf{h}) Pr(\mathbf{h}) d\mathbf{h}$$
$$= \int \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu} + \boldsymbol{\Phi}\mathbf{h}, \boldsymbol{\Sigma}] \operatorname{Norm}_{\mathbf{h}}[\mathbf{0}, \mathbf{I}] d\mathbf{h}$$
$$= \operatorname{Norm}_{\mathbf{x}}[\boldsymbol{\mu}, \boldsymbol{\Phi}\boldsymbol{\Phi}^{T} + \boldsymbol{\Sigma}]$$

Sampling from factor analyzer



Factor analysis vs. MoG



E-Step

• Compute posterior over hidden variables

$$\hat{q}(\mathbf{h}_{i}) = Pr(\mathbf{h}_{i}|\mathbf{x}_{i}, \boldsymbol{\theta}^{[t]})$$

$$= \frac{Pr(\mathbf{x}_{i}|\mathbf{h}_{i}, \boldsymbol{\theta}^{[t]})Pr(\mathbf{h}_{i})}{Pr(\mathbf{x}_{i}|\boldsymbol{\theta}^{[t]})}$$

$$= \frac{Norm_{\mathbf{x}_{i}}[\boldsymbol{\mu} + \boldsymbol{\Phi}\mathbf{h}_{i}, \boldsymbol{\Sigma}]Norm_{\mathbf{h}_{i}}[\mathbf{0}, \mathbf{I}]}{Pr(\mathbf{x}_{i}|\boldsymbol{\theta}^{[t]})}$$

$$= Norm_{\mathbf{h}_{i}}[(\boldsymbol{\Phi}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\Phi} + \mathbf{I})^{-1}\boldsymbol{\Phi}^{T}\boldsymbol{\Sigma}^{-1}(\mathbf{x}_{i} - \boldsymbol{\mu}), (\boldsymbol{\Phi}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\Phi} + \mathbf{I})]$$

• Compute expectations needed for M-Step

$$\begin{aligned} \mathbf{E}[\mathbf{h}_i] &= (\mathbf{\Phi}^T \mathbf{\Sigma}^{-1} \mathbf{\Phi} + \mathbf{I})^{-1} \mathbf{\Phi}^T \mathbf{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \\ \mathbf{E}[\mathbf{h}_i \mathbf{h}_i^T] &= E\left[(\mathbf{h}_i - \mathbf{E}[\mathbf{h}_i]) (\mathbf{h}_i - \mathbf{E}[\mathbf{h}_i])^T \right] + \mathbf{E}[\mathbf{h}_i] \mathbf{E}[\mathbf{h}_i]^T \\ &= (\mathbf{\Phi}^T \mathbf{\Sigma}^{-1} \mathbf{\Phi} + \mathbf{I})^{-1} + \mathbf{E}[\mathbf{h}_i] \mathbf{E}[\mathbf{h}_i]^T. \end{aligned}$$

E-Step



M-Step

• Optimize parameters w.r.t. EM bound

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int \hat{q}_{i}(\mathbf{h}_{i}) \log \left[Pr(\mathbf{x}, \mathbf{h}_{i}, \boldsymbol{\theta}) \right] d\mathbf{h}_{i} \right]$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int \hat{q}_{i}(\mathbf{h}_{i}) \left[\log \left[Pr(\mathbf{x} | \mathbf{h}_{i}, \boldsymbol{\theta}) \right] + \log \left[Pr(\mathbf{h}_{i}) \right] \right] d\mathbf{h}_{i} \right]$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int \hat{q}_{i}(\mathbf{h}_{i}) \log \left[Pr(\mathbf{x} | \mathbf{h}_{i}, \boldsymbol{\theta}) \right] d\mathbf{h}_{i} \right]$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} E \left[\log Pr(\mathbf{x} | \mathbf{h}_{i}, \boldsymbol{\theta}) \right] \right],$$

where....

$$\log Pr(\mathbf{x}_i|\mathbf{h}_i) = -\frac{D\log(2\pi) + \log|\mathbf{\Sigma}| + (\mathbf{x}_i - \boldsymbol{\mu} - \boldsymbol{\Phi}\mathbf{h}_i)^T \mathbf{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu} - \boldsymbol{\Phi}\mathbf{h}_i)}{2}$$

M-Step

• Optimize parameters w.r.t. EM bound

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} E\left[\log Pr(\mathbf{x}|\mathbf{h}_{i}, \boldsymbol{\theta}) \right] \right]$$

where...

$$\log Pr(\mathbf{x}_i|\mathbf{h}_i) = -\frac{D\log(2\pi) + \log|\mathbf{\Sigma}| + (\mathbf{x}_i - \boldsymbol{\mu} - \boldsymbol{\Phi}\mathbf{h}_i)^T \mathbf{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu} - \boldsymbol{\Phi}\mathbf{h}_i)}{2}$$

giving expressions:

$$\hat{\boldsymbol{\mu}} = \frac{\sum_{i=1}^{I} \mathbf{x}_{i}}{I}$$

$$\hat{\boldsymbol{\Phi}} = \left(\sum_{i=1}^{I} (\mathbf{x}_{i} - \hat{\boldsymbol{\mu}}) \mathbf{E}[\mathbf{h}_{i}]^{T}\right) \left(\sum_{i=1}^{I} \mathbf{E}[\mathbf{h}_{i}\mathbf{h}_{i}^{T}]\right)^{-1}$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{I} \sum_{i=1}^{I} \operatorname{diag} \left[(\mathbf{x}_{i} - \hat{\boldsymbol{\mu}})^{T} (\mathbf{x}_{i} - \hat{\boldsymbol{\mu}}) - \hat{\boldsymbol{\Phi}} \mathbf{E}[\mathbf{h}_{i}] \mathbf{x}^{T} \right]$$

Face model



Sampling from 10 parameter model

To generate:

- Choose factor loadings, **h**_i from standard normal distribution
- Multiply by factors, Φ
- Add mean, μ
- (should add random noise component ϵ_i w/ diagonal cov Σ)



Computer vision: models, learning and inference. ©2011 Simon J.D. Prince

Combining Models

Can combine three models in various combinations

- Mixture of t-distributions = mixture of Gaussians + t-distributions
- Robust subspace models = t-distributions + factor analysis
- Mixture of factor analyzers = mixture of Gaussians + factor analysis

Or combine all models to create mixture of robust subspace models

$$Pr(\mathbf{x}) = \sum_{k=1}^{K} \lambda_k \operatorname{Stud}_{\mathbf{x}} \left[\boldsymbol{\mu}_k, \boldsymbol{\Phi}_k \boldsymbol{\Phi}_k^T + \boldsymbol{\Sigma}_k, \nu_k \right]$$

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Expectation Maximization

Problem: Optimize cost functions of the form

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{I} \log \left[\sum_{h} Pr(\mathbf{x}_{i}, h_{i}) \right] \qquad \qquad \text{Discrete case}$$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{I} \log \left[\int Pr(\mathbf{x}_{i}, \mathbf{h}_{i}) d\mathbf{h}_{i} \right] \qquad \qquad \qquad \text{Continuous case}$$

Solution: Expectation Maximization (EM) algorithm (Dempster, Laird and Rubin 1977)

Key idea: Define lower bound on log-likelihood and increase at each iteration

Expectation Maximization

Defines a lower bound on log likelihood and increases bound iteratively

$$\mathcal{B}[\{q_i(\mathbf{h}_i)\}, \boldsymbol{\theta}] = \sum_{i=1}^{I} \int q_i(\mathbf{h}_i) \log\left[\frac{Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)}\right] d\mathbf{h}_i$$

$$\leq \sum_{i=1}^{I} \log\left[\int q_i(\mathbf{h}_i) \frac{Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} d\mathbf{h}_i\right]$$

$$= \sum_{i=1}^{I} \log\left[\int Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) d\mathbf{h}_i\right],$$

Lower bound



Lower bound is a *function* of parameters θ and a *set* of probability distributions {q_i(**h**_i)}

E-Step & M-Step

E-Step – Maximize bound w.r.t. distributions $\{q_i(\mathbf{h}_i)\}$

$$q_i^{[t]}[\mathbf{h}_i] = \operatorname*{argmax}_{q_i[\mathbf{h}_i]} \left[\mathcal{B}[\{q_i(\mathbf{h}_i)\}, \theta^{[t-1]}] \right]$$

M-Step – Maximize bound w.r.t. parameters θ

$$\boldsymbol{\theta}^{[t]} = \operatorname*{argmax}_{\boldsymbol{\theta}} \left[\mathcal{B}[\{q_i^{[t]}(\mathbf{h}_i)\}, \boldsymbol{\theta}] \right]$$



E-Step: Update $\{q_i[h_i]\}$ so that bound equals log likelihood for this θ



M-Step: Update θ to maximum

Missing Parts of Argument



Computer vision: models, learning and inference. ©2011 Simon J.D. Prince

Missing Parts of Argument

- 1. Show that this is a lower bound
- 2. Show that the E-Step update is correct
- 3. Show that the M-Step update is correct
Jensen's Inequality



Lower Bound for EM Algorithm

$$\mathcal{B}[\{q_i(\mathbf{h}_i)\}, \boldsymbol{\theta}] = \sum_{i=1}^{I} \int q_i(\mathbf{h}_i) \log\left[\frac{Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)}\right] d\mathbf{h}_i$$

$$\leq \sum_{i=1}^{I} \log\left[\int q_i(\mathbf{h}_i) \frac{Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} d\mathbf{h}_i\right]$$

$$= \sum_{i=1}^{I} \log\left[\int Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) d\mathbf{h}_i\right],$$

Where we have used Jensen's inequality

$$\int q(y) \log[y] dy \le \int \log[q(y)] dy$$

E-Step – Optimize bound w.r.t {q_i(**h**_i)}

$$\mathcal{B}[\{q_{i}(\mathbf{h}_{i})\},\boldsymbol{\theta}] = \sum_{i=1}^{I} \int q_{i}(\mathbf{h}_{i}) \log \left[\frac{Pr(\mathbf{x}_{i},\mathbf{h}_{i}|\boldsymbol{\theta})}{q_{i}(\mathbf{h}_{i})}\right] d\mathbf{h}_{i}$$

$$= \sum_{i=1}^{I} \int q_{i}(\mathbf{h}_{i}) \log \left[\frac{Pr(\mathbf{h}_{i}|\mathbf{x}_{i},\boldsymbol{\theta})Pr(\mathbf{x}_{i}|\boldsymbol{\theta})}{q_{i}(\mathbf{h}_{i})}\right] d\mathbf{h}_{i}$$

$$= \sum_{i=1}^{I} \int q_{i}(\mathbf{h}_{i}) \log \left[Pr(\mathbf{x}_{i}|\boldsymbol{\theta})\right] d\mathbf{h}_{i} - \sum_{i=1}^{I} \int q_{i}(\mathbf{h}_{i}) \log \left[\frac{q_{i}(\mathbf{h}_{i})}{Pr(\mathbf{h}_{i}|\mathbf{x}_{i},\boldsymbol{\theta})}\right] d\mathbf{h}_{i}$$

$$= \sum_{i=1}^{I} \log \left[Pr(\mathbf{x}_{i}|\boldsymbol{\theta})\right] - \sum_{i=1}^{I} \int q_{i}(\mathbf{h}_{i}) \log \left[\frac{q_{i}(\mathbf{h}_{i})}{Pr(\mathbf{h}_{i}|\mathbf{x}_{i},\boldsymbol{\theta})}\right] d\mathbf{h}_{i} \qquad (7.46)$$
Constant w.r.t. q(h) Only this term matters

$$\begin{aligned} \mathbf{E}\text{-}\mathbf{Step}\\ \hat{q}_{i}(\mathbf{h}_{i}) &= \operatorname{argmax}_{q_{i}(\mathbf{h}_{i})} \left[-\int q_{i}(\mathbf{h}_{i}) \log \left[\frac{q_{i}(\mathbf{h}_{i})}{Pr(\mathbf{h}_{i} | \mathbf{x}_{i}, \boldsymbol{\theta})} \right] d\mathbf{h}_{i} \right] \\ &= \operatorname{argmax}_{q_{i}(\mathbf{h}_{i})} \left[\int q_{i}(\mathbf{h}_{i}) \log \left[\frac{Pr(\mathbf{h}_{i} | \mathbf{x}_{i}, \boldsymbol{\theta})}{q_{i}(\mathbf{h}_{i})} \right] d\mathbf{h}_{i} \right] \\ &= \operatorname{argmin}_{q_{i}(\mathbf{h}_{i})} \left[-\int q_{i}(\mathbf{h}_{i}) \log \left[\frac{Pr(\mathbf{h}_{i} | \mathbf{x}_{i}, \boldsymbol{\theta})}{q_{i}(\mathbf{h}_{i})} \right] d\mathbf{h}_{i} \right] \end{aligned}$$

Kullback Leibler divergence – distance between probability distributions . We are maximizing the negative distance (i.e. Minimizing distance)

Use this relation



Kullback Leibler Divergence

$$\int q_i(\mathbf{h}_i) \log \left[\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_i \leq \int q_i(\mathbf{h}_i) \left(\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} - 1 \right) d\mathbf{h}_i$$
$$= \int Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta}) - q_i(\mathbf{h}_i) d\mathbf{h}_i$$
$$= 1 - 1 = 0,$$

So the cost function must be positive

$$\hat{q}_i(\mathbf{h}_i) = \operatorname{argmin}_{q_i(\mathbf{h}_i)} \left[-\int q_i(\mathbf{h}_i) \log \left[\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_i \right]$$

In other words, the best we can do is choose $q_i(\boldsymbol{h}_i)$ so that this is zero

E-Step

So the cost function must be positive

$$\hat{q}_i(\mathbf{h}_i) = \operatorname{argmin}_{q_i(\mathbf{h}_i)} \left[-\int q_i(\mathbf{h}_i) \log \left[\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_i \right]$$

The best we can do is choose $q_i(\mathbf{h}_i)$ so that this is zero.

How can we do this? Easy – choose posterior Pr(h|x)

$$\int q_i(\mathbf{h}_i) \log \left[\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_i = \int Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta}) \log \left[\frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})}{Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta})} \right] d\mathbf{h}_i$$
$$= \int Pr(\mathbf{h}_i | \mathbf{x}_i, \boldsymbol{\theta}) \log [1] d\mathbf{h}_i = 0.$$

M-Step – Optimize bound w.r.t. θ

$$\begin{aligned} \boldsymbol{\theta}^{[t]} &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\mathcal{B}[\{q_i^{[t]}(\mathbf{h}_i)\}, \boldsymbol{\theta}] \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int q_i^{[t]}(\mathbf{h}_i) \log \left[\frac{Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta})}{q_i^{[t]}(\mathbf{h}_i)} \right] d\mathbf{h}_i \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int q_i^{[t]}(\mathbf{h}_i) \log \left[Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) \right] - q_i^{[t]}(\mathbf{h}_i) \log \left[q_i^{[t]}(\mathbf{h}_i) \right] d\mathbf{h}_i \right] \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \left[\sum_{i=1}^{I} \int q_i^{[t]}(\mathbf{h}_i) \log \left[Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) \right] d\mathbf{h}_i \right] \end{aligned}$$

In the M-Step we optimize expected joint log likelihood with respect to parameters θ (Expectation w.r.t distribution from E-Step)

E-Step & M-Step

E-Step – Maximize bound w.r.t. distributions $q_i(h_i)$

$$\hat{q}_i(\mathbf{h}_i) = Pr(\mathbf{h}_i|\mathbf{x}_i, \boldsymbol{\theta}^{[t]}) = \frac{Pr(\mathbf{x}_i|\mathbf{h}_i, \boldsymbol{\theta}^{[t]})Pr(\mathbf{h}_i)}{Pr(\mathbf{x}_i)}$$

M-Step – Maximize bound w.r.t. parameters θ

$$\hat{\boldsymbol{\theta}}^{[t+1]} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left[\sum_{i=1}^{I} \int \hat{q}_i(\mathbf{h}_i) \log \left[Pr(\mathbf{x}_i, \mathbf{h}_i | \boldsymbol{\theta}) \right] \, d\mathbf{h}_i \right]$$

Structure

- Densities for classification
- Models with hidden variables
 - Mixture of Gaussians
 - t-distributions
 - Factor analysis
- EM algorithm in detail
- Applications

Face Detection

-86		To			
a)			b)	25	

	Color	Grayscale	Equalized
Single Gaussian Mixture of 10 Gaussians	$76\% \\ 81\%$	$79\% \\ 85\%$	$rac{80\%}{89\%}$

(not a very realistic model – face detection really done using discriminative methods)

Object recognition

Aeschliman et al. (2010)



Segmentation



















 $Pr(w_n) = \operatorname{Cat}_{w_n}[\boldsymbol{\lambda}] \qquad \text{Sfikas et al. (2007) @IEEE 2007} \\ Pr(\mathbf{x}_i | w_i = k) = \operatorname{Norm}_{\mathbf{x}_i}[\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k].$

Face Recognition

Figure 6.28 Face recognition. Our goal is to take the RGB values of a facial image \mathbf{x} and assign a label $y \in \{1 \dots K\}$. Since the data is high dimensional, we model the class conditional density function $Pr(\mathbf{x}|y=k)$ for each individual in the database as a factor analyzer. To classify a new face, we apply Bayes' rule with suitable priors Pr(y) to compute the posterior distribution $Pr(y|\mathbf{x})$. Finally, we choose the label \hat{y} = $\arg \max_k Pr(y=k|\mathbf{x})$ that maximizes the posterior. This is not suitable if there are not plentiful training examples of each individual or where there are significant pose or lighting changes.



Pose Regression





Conditional Distributions



$$Pr(\mathbf{x}) = Pr\left(\begin{bmatrix}\mathbf{x}_1\\\mathbf{x}_2\end{bmatrix}\right) = \operatorname{Norm}_{\mathbf{x}}\left(\begin{bmatrix}\boldsymbol{\mu}_1\\\boldsymbol{\mu}_2\end{bmatrix}, \begin{bmatrix}\boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12}^T\\\boldsymbol{\Sigma}_{12} & \boldsymbol{\Sigma}_{22}\end{bmatrix}\right)$$

then

$$Pr(\mathbf{x}_{1}|\mathbf{x}_{2}) = \operatorname{Norm}_{\mathbf{x}_{1}} \left(\boldsymbol{\mu}_{1} + \boldsymbol{\Sigma}_{12}^{T} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_{2} - \boldsymbol{\mu}_{2}), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}^{T} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12} \right)$$

$$Pr(\mathbf{x}_{2}|\mathbf{x}_{1}) = \operatorname{Norm}_{\mathbf{x}_{2}} \left(\boldsymbol{\mu}_{2} + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_{1} - \boldsymbol{\mu}_{1}), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12}^{T} \right)$$



Figure 7.30 Modeling transformations with hidden variables. a) Original set of digit images are only weakly aligned. b) Mean and standard deviation images are consequently blurred out. The probability density model does not fit well. c) Each possible value of a discrete hidden variable represents a different transformation (here inverse transformations are shown). Red square highlights most likely choice of hidden variable after 10 iterations. d) Transformed digits (based on most likely hidden variable). d) New mean and standard deviation images are more focussed: the probability density function fits better.