

CSC 412/2506:
Probabilistic Learning and Reasoning
Week 10: EM Algorithm & Probabilistic PCA

Murat A. Erdogdu

University of Toronto

Overview of the first hour

- Gaussian mixture models
- EM-algorithm
- Clustering

Mixture of Gaussians

We combine simple models into a complex model by taking a mixture of K multivariate Gaussian densities of the form:

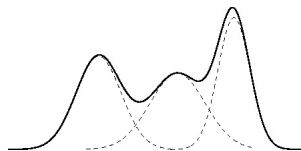
$$p(x) = \sum_{k=1}^K \pi_k N_m(x | \mu_k, \Sigma_k),$$

where $\pi_k \geq 0$ and $\sum_{k=1}^K \pi_k = 1$.

- Each Gaussian component has its own mean vector μ_k and covariance matrix Σ_k .
- The parameters π_k are called the mixing coefficients.

Example:

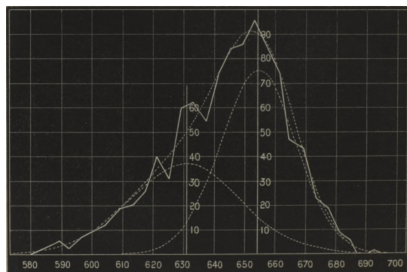
- $K = 3$ (three Gaussian components)
- $m = 1$ (univariate Gaussians)



The crabs from Naples bay

In 1892, scientists collected data on populations of the crab, *Carcinus Moenas*, and observed that the ratio of forehead width to the body length actually showed a highly skewed distribution.

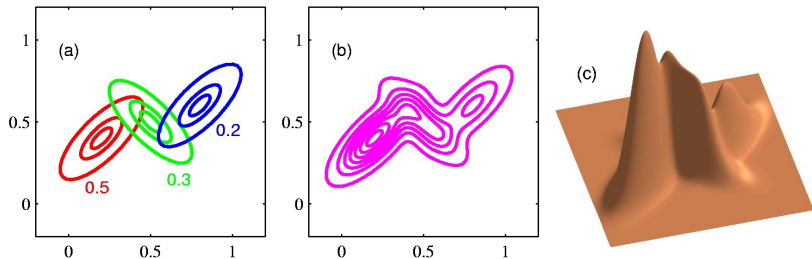
On Certain Correlated Variations in Carcinus maenas (1893) W. F. Weldon



They wondered whether this distribution could be the result of the population being a mix of two different normal distributions (two sub-species).

In **1894**, Karl Pearson proposed a method to fit this model ([read here](#)), whose modern version is the “method of moments”.

- Illustration of a mixture of 3 Gaussians in a 2-dimensional space:



(a) Contours of constant density of each of the mixture components, along with the mixing coefficients

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

(c) A surface plot of the distribution $p(\mathbf{x})$.

Mixture of Gaussians as a latent variable model

Recall: $p(x) = \sum_{k=1}^K \pi_k N_m(x|\mu_k, \Sigma_k)$.

- Consider a latent variable z with K states $z \in \{1, \dots, K\}$.
- The distribution of z given by the mixing coefficients:

$$p(z = k) = \pi_k.$$

- Specify the conditional as $p(x|z = k) = N_m(x|\mu_k, \Sigma_k)$ with joint:

$$p(x, z = k) = p(z = k)p(x|z = k) = \pi_k N_m(x|\mu_k, \Sigma_k).$$

- Then the marginal $p(x)$ satisfies

$$p(x) = \sum_{k=1}^K p(x, z = k) = \sum_{k=1}^K \pi_k N_m(x|\mu_k, \Sigma_k).$$

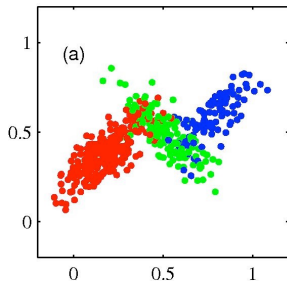
Mixture of Gaussians: inference

- If we have several observations x_1, \dots, x_N , for every observed data point x_n there is a corresponding latent z_n .
- Consider the conditional $p(z|x)$

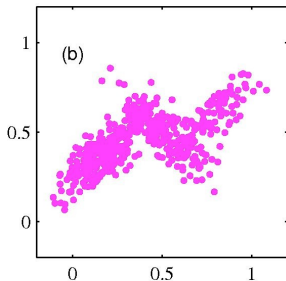
$$\begin{aligned} p(z = k|x) &= \frac{p(z = k)p(x|z = k)}{\sum_{j=1}^K p(z = j)p(x|z = j)} \\ &= \frac{\pi_k N_m(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N_m(x|\mu_j, \Sigma_j)} \end{aligned}$$

- We view π_k as prior probability that $z = k$, and $p(z = k|x)$ is the corresponding posterior once we have observed the data.

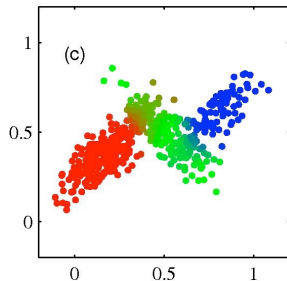
- 500 points drawn from a mixture of 3 Gaussians.



Samples from the **joint** distribution $p(x,z)$.



Samples from the **marginal** distribution $p(x)$.



Same samples where colors represent the value of responsibilities.

The Likelihood function

Parameters: $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$, $\boldsymbol{\Sigma} = (\Sigma_1, \dots, \Sigma_K)$.

Recall: $p(x|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{k=1}^K \pi_k N_m(x|\mu_k, \Sigma_k)$

- Represent the dataset $\{x_1, \dots, x_N\}$ as $\mathbf{X} \in \mathbb{R}^{N \times m}$.
- The latent variable is represented by a vector $\mathbf{z} \in \mathbb{R}^N$.
- The log-likelihood takes the form

$$\log p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N_m(x_n|\mu_k, \Sigma_k) \right)$$

Maximum Likelihood (μ)

Recall: $\log p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N_m(x_n | \mu_k, \Sigma_k) \right)$.

- Differentiating wrt μ_k and setting to zero gives:

$$\begin{aligned} 0 &= \sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)} \Sigma_k^{-1} (x_n - \mu_k) \\ &= \sum_{n=1}^N p(z_n = k | x_n) \Sigma_k^{-1} (x_n - \mu_k). \end{aligned}$$

- Equivalently (as Σ_k is positive definite)

$$\mu_k = \sum_n \frac{p(z = k | x_n)}{N_k} x_n, \quad N_k = \sum_n p(z = k | x_n).$$

- Simple interpretation: the MLE given by the weighted mean of the data weighted by the posterior $p(z = k | x_n)$.

Maximum Likelihood (Σ , π)

Recall: $\log p(\mathbf{X} | \pi, \mu, \Sigma) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k N_m(x_n | \mu_k, \Sigma_k) \right)$.

- Differentiating wrt Σ_k and setting to zero gives:

$$\Sigma_k = \sum_n \frac{p(z = k | x_n)}{N_k} (x_n - \mu_k)(x_n - \mu_k)^\top.$$

- Again data points weighted by posterior probabilities.
- Finally, for the weights π_k the MLE is

$$\pi_k = \frac{N_k}{\sum_{j=1}^K N_j} = \frac{N_k}{N}, \quad N_k = \sum_n p(z = k | x_n).$$

Motivating the EM algorithm

- The MLE **does not have a closed form solution**.
- The estimates depend on the posterior probabilities $p(z = k|x_n)$, which themselves depend on those parameters.
- Indeed, recall that

$$p(z = k|x_n) = \frac{\pi_k N_m(x_n|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N_m(x_n|\mu_j, \Sigma_j)}.$$

- Iterative solution (EM algorithm):
 - ▶ Initialize the parameters to some values.

E-step Update the posteriors $p(z = k|x_n)$.

M-step Update model parameters $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$.

- ▶ Repeat.

EM algorithm for Gaussian mixtures

- Initialize π, μ, Σ .
- **E-step**: for each k, n compute the posterior probabilities

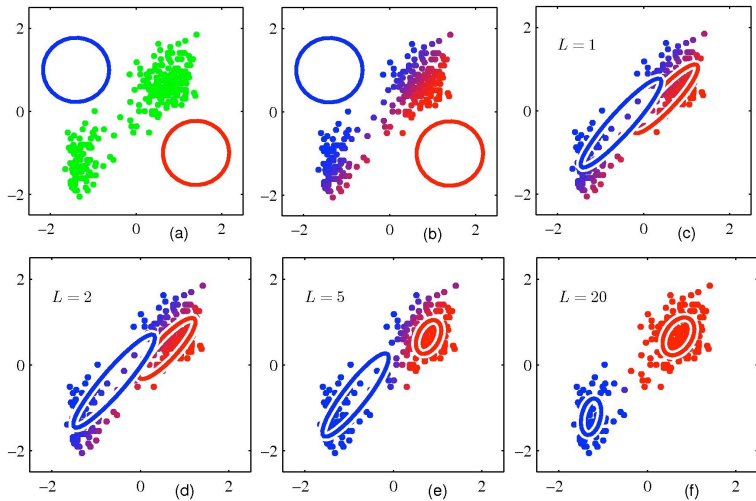
$$p(z = k | x_n) = \frac{\pi_k N_m(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N_m(x_n | \mu_j, \Sigma_j)}.$$

- **M-step**: Re-estimate model parameters

$$\begin{aligned}\mu_k^{\text{new}} &= \sum_{n=1}^N \frac{p(z = k | x_n)}{N_k} x_n, & N_k &= \sum_{n=1}^N p(z = k | x_n), \\ \Sigma_k^{\text{new}} &= \sum_{n=1}^N \frac{p(z = k | x_n)}{N_k} (x_n - \mu_k^{\text{new}})(x_n - \mu_k^{\text{new}})^\top, \\ \pi_k^{\text{new}} &= \frac{N_k}{N}.\end{aligned}$$

- Evaluate the log-likelihood and check for convergence.

Illustration of the EM algorithm:



The General EM algorithm

Consider a general setting with latent variables.

- Observed dataset $\mathbf{X} \in \mathbb{R}^{N \times D}$, latent variables $\mathbf{Z} \in \mathbb{R}^{N \times K}$.

Maximize the log-likelihood $\log p(\mathbf{X}|\theta) = \log(\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta))$.

- Initialize parameters θ^{old} .
- **E-step:** use θ^{old} to compute the posterior $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$.
- **M-step:** $\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$, where

$$\begin{aligned} Q(\theta, \theta^{\text{old}}) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\theta) \\ &= \mathbb{E}\left(\log p(\mathbf{X}, \mathbf{Z}|\theta) \middle| \mathbf{X}, \theta^{\text{old}}\right) \end{aligned}$$

which is tractable in many applications.

- Replace $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$. Repeat until convergence.

Example: Gaussian mixture

- If z was observed, the MLE would be trivial

$$\log p(\mathbf{X}, \mathbf{Z}|\theta) = \sum_{n=1}^N \log p(x_n, z_n|\theta) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}(z_n = k) \log(\pi_k N(x_n|\mu_k, \Sigma_k)).$$

For the E-step: $p(\mathbf{Z}|\mathbf{X}, \theta) = \prod_{n=1}^N p(z_n|\mathbf{X}, \theta)$ we have

$$p(z_n = k|\mathbf{X}, \theta) = p(z_n = k|x_n, \theta) = \frac{\pi_k N_m(x_n|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N_m(x_n|\mu_j, \Sigma_j)}.$$

For the M-step: $\mathbb{E}(\mathbb{1}(z_n = k)|\mathbf{X}, \theta^{\text{old}}) = p(z_n = k|\mathbf{X}, \theta^{\text{old}})$ and so

$$\mathbb{E}\left(\log p(\mathbf{X}, \mathbf{Z}|\theta)\middle|\mathbf{X}, \theta^{\text{old}}\right) = \sum_{n=1}^N \sum_{k=1}^K p(z_n = k|\mathbf{X}, \theta^{\text{old}}) \log(\pi_k N(x_n|\mu_k, \Sigma_k)).$$

Maximizing gives the formulas on Slide 13.

Relationship to K-Means (CSC 311)

- Consider a Gaussian mixture, s.t. $\Sigma_k = \epsilon I$ for all $k = 1, \dots, K$.
- We have

$$p(x|\mu_k, \Sigma_k) = \frac{1}{(2\pi\epsilon)^{m/2}} \exp\left(-\frac{1}{2\epsilon}\|x - \mu_k\|^2\right).$$

- Consider the EM algorithm in this special case, $\theta = (\boldsymbol{\pi}, \boldsymbol{\mu})$.
- The posterior probabilities take the form

$$p(z_n = k|\mathbf{X}, \theta) = \frac{\pi_k \exp(-\|x_n - \mu_k\|^2/2\epsilon)}{\sum_{j=1}^K \pi_j \exp(-\|x_n - \mu_j\|^2/2\epsilon)}.$$

- If $\epsilon \rightarrow 0$, the term with smallest $\|x_n - \mu_j\|$ tends to zero most slowly.

- Thus $p(z_n = k|\mathbf{X}, \theta) \rightarrow r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\| \\ 0 & \text{otherwise} \end{cases}$

Relationship to K-Means

Recall: $\mathbb{E}(\log p(\mathbf{X}, \mathbf{Z}|\theta) | \mathbf{X}, \theta^{\text{old}}) = \sum_{n=1}^N \sum_{k=1}^K p(z_n = k | \mathbf{X}, \theta^{\text{old}}) \log(\pi_k N(x_n | \mu_k, \Sigma_k))$.

As $\epsilon \rightarrow 0$, we have

$$p(z_n = k | \mathbf{X}, \theta) \rightarrow r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\| \\ 0 & \text{otherwise} \end{cases}$$

which gives

$$\mathbb{E}(\log p(\mathbf{X}, \mathbf{Z}|\theta) | \mathbf{X}, \theta^{\text{old}}) \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2 + \text{const.}$$

- In the limit, maximizing the expected log-likelihood is equivalent to minimizing the distortion measure in the K-means algorithm.
- The EM-algorithm is slower but more flexible and accurate.

- The ELBO is given as

$$\mathcal{L}(x; \theta, \phi) = E_{z \sim q_\phi} \left[\log p_\theta(x, z) \right] + H(q_\phi)$$

- ▶ This maximizes expected complete data log-likelihood while penalizing low entropy distributions.
- ▶ We perform alternating gradient descent (ascent).
- Expectation in EM algorithm maximizes

$$\mathcal{Q}(\phi, \phi^{\text{old}}) = E_{z \sim q_\phi^{\text{old}}} \left[\log p_\phi(x, z) \right]$$

- ▶ This maximizes expected complete data log-likelihood while the expectation is over the posterior.
- ▶ We perform maximization at each iteration.

Summary

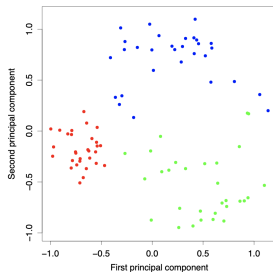
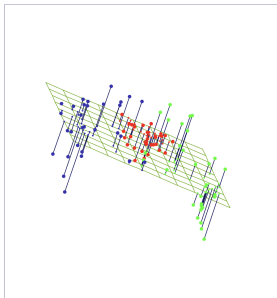
- EM algorithm is a classical method in statistics.
- It can be used in the presence of latent variables.
- When applied to Gaussian mixtures, compared to k-means, it captures the covariance structure of the data.

Overview

- A probabilistic model for continuous latent variables.
 - ▶ Probabilistic interpretation of the PCA
- Earlier formulation of PCA was motivated geometrically.
- We will show that it can be expressed as the maximum likelihood estimate of a certain probabilistic model.

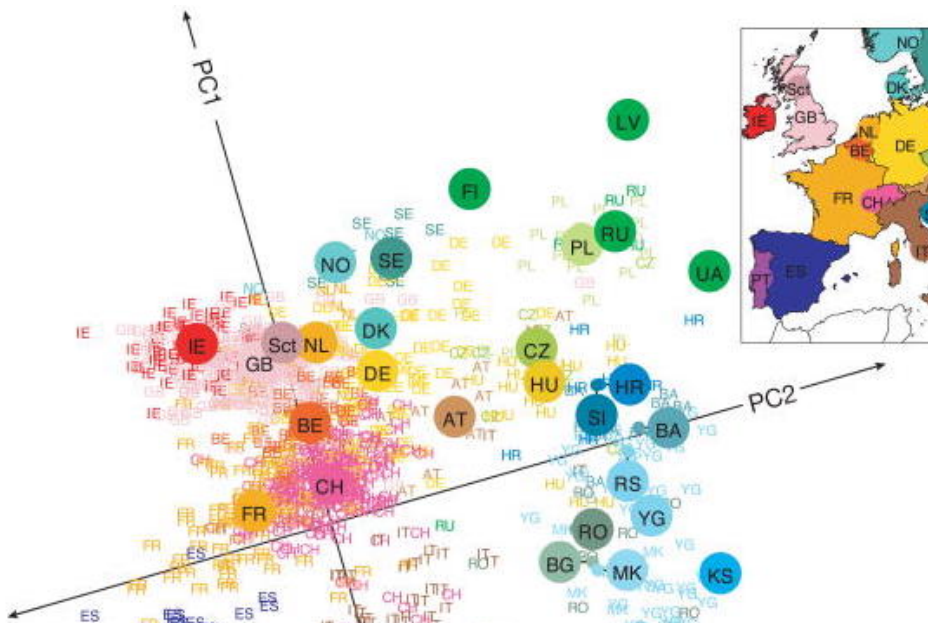
Low dimensional representation

- In practice, even though data is very high dimensional, its important features can be accurately captured in a low dimensional subspace.



- Find a low dimensional representation of your data.
 - ▶ Computational benefits
 - ▶ Interpretability, visualization
 - ▶ Generalization

Nice example



Recall: Principal Component Analysis (PCA)

- Data set $\{\mathbf{x}^{(i)}\}_{i=1}^N$
- Each input vector $\mathbf{x}^{(i)} \in \mathbb{R}^D$ is approximated as $\bar{\mathbf{x}} + \mathbf{U}\mathbf{z}^{(i)}$,

$$\mathbf{x}^{(i)} \approx \tilde{\mathbf{x}}^{(i)} = \bar{\mathbf{x}} + \mathbf{U}\mathbf{z}^{(i)}$$

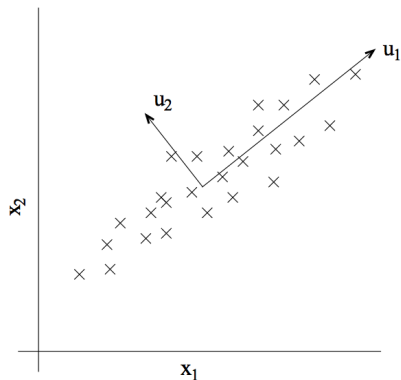
where $\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}^{(i)}$ is the data mean, $\mathbf{U} \in \mathbb{R}^{D \times K}$ is the orthogonal basis for the principal subspace, and $\mathbf{z}^{(i)} \in \mathbb{R}^K$ is the code vector

$$\mathbf{z}^{(i)} = \mathbf{U}^\top (\mathbf{x}^{(i)} - \bar{\mathbf{x}})$$

- \mathbf{U} is chosen to minimize the reconstruction error

$$\mathbf{U}^* = \arg \min_{\mathbf{U}} \sum_i \|\mathbf{x}^{(i)} - \bar{\mathbf{x}} - \mathbf{U}\mathbf{U}^\top (\mathbf{x}^{(i)} - \bar{\mathbf{x}})\|^2$$

We are looking for directions



- For example, in a 2-dimensional problem, we are looking for the direction u_1 along which the data is **well represented**: (?)
 - ▶ e.g. direction of higher variance
 - ▶ e.g. direction of minimum reconstruction error
 - ▶ Recall: they are the same!

Probabilistic PCA

Consider the following latent variable model.

- Similar to the Gaussian mixture model but with Gaussian latents:

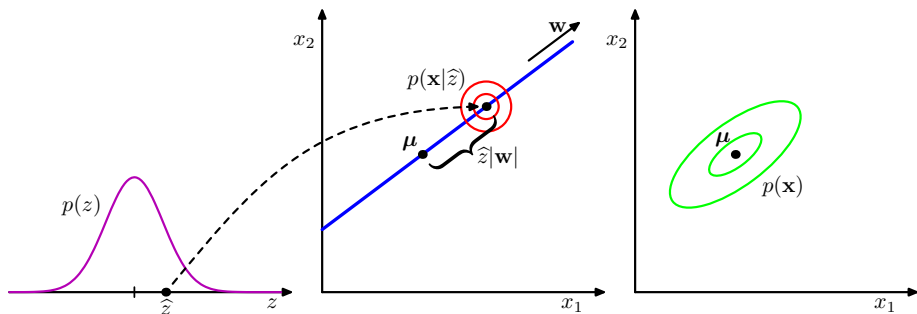
$$\begin{aligned}\mathbf{z} &\sim \mathcal{N}_K(\mathbf{0}, \mathbf{I}_K) \\ \mathbf{x} \mid \mathbf{z} &\sim \mathcal{N}_D(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}_D)\end{aligned}$$

- This is similar to naive Bayes graphical model, because $p(\mathbf{x} \mid \mathbf{z})$ factorizes with respect to the dimensions of \mathbf{x} .
- What sort of data does this model produce?

Matrix-vector multiplication: $\mathbf{W}\mathbf{z}$ is a linear combination of the columns of \mathbf{W} with coefficients $\mathbf{z} = (z_1, \dots, z_K)$.

Probabilistic PCA

- \mathbf{Wz} is a random linear combination of the columns of \mathbf{W}
- To get the random variable \mathbf{x} , we sample a standard normal \mathbf{z} and then add a small amount of isotropic noise to $\mathbf{Wz} + \boldsymbol{\mu}$.



The column span of \mathbf{W} refers to the principal subspace in PCA.

Probabilistic PCA : The Likelihood function

- To perform maximum likelihood in this model, we need to maximize the following:

$$\max_{\mathbf{W}, \boldsymbol{\mu}, \sigma^2} \log p(\mathbf{x} | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \max_{\mathbf{W}, \boldsymbol{\mu}, \sigma^2} \log \int p(\mathbf{x} | \mathbf{z}, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}) d\mathbf{z}$$

- This is easier than the Gaussian mixture model.
- $\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \epsilon$ (\mathbf{x} is an affine transformations of Gaussian vars)
- $p(\mathbf{x} | \mathbf{W}, \boldsymbol{\mu}, \sigma^2)$ is Gaussian
 - ▶ Only need to compute $\mathbb{E}[\mathbf{x}]$ and $\text{Cov}[\mathbf{x}]$.

Probabilistic PCA : Maximum Likelihood

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\mu}$$

$$\begin{aligned}\text{Cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})^\top] \\ &= \mathbb{E}[(\mathbf{W}\mathbf{z}\mathbf{z}^\top \mathbf{W}^\top] + \text{Cov}[\boldsymbol{\epsilon}] \\ &= \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D\end{aligned}$$

Recall: \mathbf{R} orthogonal if $\mathbf{R}\mathbf{R}^\top = \mathbf{I}$.

This model is not identifiable because $\mathbf{W}\mathbf{W}^\top = (\mathbf{W}\mathbf{R})(\mathbf{W}\mathbf{R})^\top$.

Probabilistic PCA : Maximum Likelihood

Thus, the log-likelihood of the data under this model is given by

$$-\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu})$$

where $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D$.

Here the MLE $(\hat{\boldsymbol{\mu}}, \hat{\mathbf{W}}, \hat{\sigma}^2)$ is given in a closed-form!

Check Tipping and Bishop (Probabilistic PCA, 1999) for details.

The maximum likelihood estimates

The maximum likelihood estimator is:

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$$

$$\widehat{\mathbf{W}} = \widehat{\mathbf{U}}(\widehat{\mathbf{L}} - \hat{\sigma}^2 \mathbf{I}_K)^{\frac{1}{2}} \mathbf{R}$$

$$\hat{\sigma}^2 = \frac{1}{D - K} \sum_{i=K+1}^D \lambda_i$$

- The columns of $\widehat{\mathbf{U}} \in \mathbb{R}^{D \times K}$ are the K unit eigenvectors of the empirical covariance matrix $\widehat{\boldsymbol{\Sigma}}$ that have the largest eigenvalues,
- $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$ are the eigenvalues of $\widehat{\boldsymbol{\Sigma}}$.
- $\widehat{\mathbf{L}} = \text{diag}(\lambda_1, \dots, \lambda_K)$ is the diagonal matrix whose elements are the corresponding eigenvalues, and \mathbf{R} is any orthogonal matrix.

Probabilistic PCA : Maximum Likelihood

- That seems complex, to get an intuition about how this model behaves when it is fit to data, lets consider the MLE density.
- Recall that the marginal distribution on \mathbf{x} in our fitted model is a Gaussian with mean

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$$

and covariance

$$\widehat{\mathbf{C}} = \widehat{\mathbf{W}}\widehat{\mathbf{W}}^{\top} + \hat{\sigma}^2\mathbf{I} = \widehat{\mathbf{U}}(\widehat{\mathbf{L}} - \hat{\sigma}^2\mathbf{I})\widehat{\mathbf{U}}^{\top} + \hat{\sigma}^2\mathbf{I}$$

- The covariance gives us a nice intuition about the model.

Probabilistic PCA : Maximum Likelihood

- Center the data and check the variance along one of the unit eigenvectors \mathbf{u}_i , which are the vectors forming the columns of $\widehat{\mathbf{U}}$:

$$\begin{aligned}\text{Var}(\mathbf{u}_i^\top (\mathbf{x} - \bar{\mathbf{x}})) &= \mathbf{u}_i^\top \text{Cov}[\mathbf{x}] \mathbf{u}_i \\ &= \mathbf{u}_i^\top \widehat{\mathbf{U}} (\widehat{\mathbf{L}} - \hat{\sigma}^2 \mathbf{I}) \widehat{\mathbf{U}}^\top \mathbf{u}_i + \hat{\sigma}^2 \\ &= \lambda_i - \hat{\sigma}^2 + \hat{\sigma}^2 = \lambda_i\end{aligned}$$

- Now, center the data and check the variance along any unit vector orthogonal to the subspace spanned by $\widehat{\mathbf{U}}$ ($i > K$):

$$\begin{aligned}\text{Var}(\mathbf{u}_i^\top (\mathbf{x} - \bar{\mathbf{x}})) &= \mathbf{u}_i^\top \widehat{\mathbf{U}} (\widehat{\mathbf{L}} - \hat{\sigma}^2 \mathbf{I}) \widehat{\mathbf{U}}^\top \mathbf{u}_i + \hat{\sigma}^2 \\ &= \hat{\sigma}^2\end{aligned}$$

- The model captures the variance along the principle axes and approximates it in all remaining directions with a single variance.

How does it relate to PCA?

- The posterior mean is given by

$$\mathbb{E}[\mathbf{z} | \mathbf{x}] = (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu})$$

- Posterior variance:

$$\text{Cov}[\mathbf{z} | \mathbf{x}] = \sigma^{-2} (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})$$

- In the limit $\sigma^2 \rightarrow 0$, we get

$$\mathbb{E}[\mathbf{z} | \mathbf{x}] \xrightarrow{\sigma^2 \rightarrow 0} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu})$$

- Plugging in the MLEs, this limit recovers the standard PCA.

Why Probabilistic PCA (PPCA)?

- Fitting a full-covariance Gaussian model of data requires $D(D + 1)/2 + D$ parameters. With PPCA we model only the K most significant correlations and this only requires $\mathcal{O}(KD)$ parameters as long as K is small.
- Bayesian PCA gives us a Bayesian method for determining the low dimensional principal subspace.
- Existence of likelihood functions allows direct comparison with other probabilistic models.
- Instead of solving directly, we can also use EM. The EM can be scaled to very large high- dimensional datasets.

Summary: Some Gaussian models

- Gaussian mixture model.
 - ▶ Gaussian latent variable model $p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$ used for clustering.
- Probabilistic PCA.
 - ▶ Gaussian latent variable model $p(\mathbf{x}) = \int_z p(\mathbf{x}, z)$ used for dimensionality reduction.
- Bayesian linear regression (next lecture).
 - ▶ Gaussian discriminative model $p(y | \mathbf{x})$ used for regression with a Bayesian analysis for the weights.