

Probabilistic Machine Learning

Week 49, Tuesday

Latent variable models

- Discrete
Mixture models
- Continuous
Probabilistic PCA

Discrete latent variable models: Mixture models

How do we model multimodal distributions?

Complex density functions can be approximated by convex combination of simpler distributions.

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

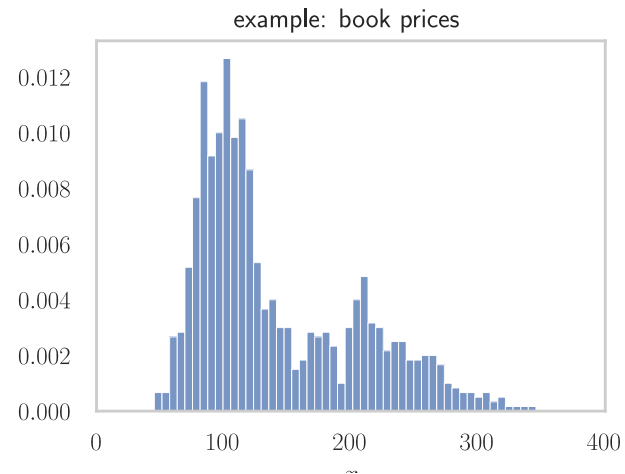
Such models are called mixture models. The simpler distributions in a mixture model are often called *components*. We will focus on Gaussian mixture models.

How do we model multimodal distributions?

Complex density functions can be approximated by convex combination of simpler distributions.

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

Such models are called mixture models. The simpler distributions in a mixture model are often called *components*. We will focus on Gaussian mixture models.

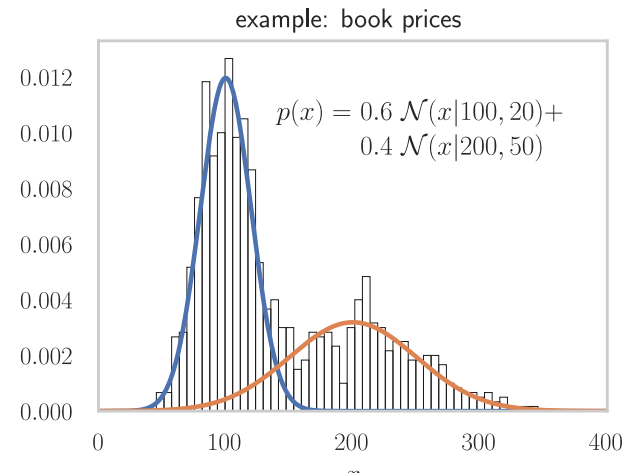


How do we model multimodal distributions?

Complex density functions can be approximated by convex combination of simpler distributions.

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

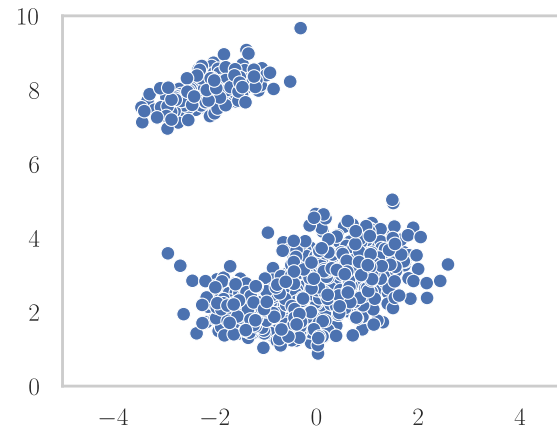
Such models are called mixture models. The simpler distributions in a mixture model are often called *components*. We will focus on Gaussian mixture models.



How do we model multimodal distributions?

Complex density functions can be approximated by convex combination of simpler distributions.

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



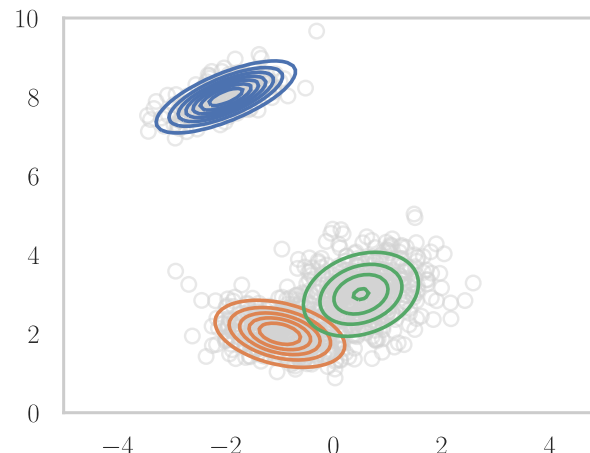
Such models are called mixture models. The simpler distributions in a mixture model are often called *components*. We will focus on Gaussian mixture models.

How do we model multimodal distributions?

Complex density functions can be approximated by convex combination of simpler distributions.

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

Such models are called mixture models. The simpler distributions in a mixture model are often called *components*. We will focus on Gaussian mixture models.



Mixture models are latent variable models

Idea: introduce discrete random variable z which takes K possible outcomes, representing the K components.

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

As usual when representing a categorical variable, we will use a *one-hot* (or *one-out-of-K*) representation.

$$z_k \in \{0, 1\} \quad (\text{but only one } z_k \text{ can be nonzero at a time})$$

The marginal probability of z are the mixture components weights

$$p(z_k = 1) = \pi_k$$

Since z has a one-hot representation, we can write

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (\text{only one of the factors will contribute})$$

Mixture models are latent variable models (2)

Model:

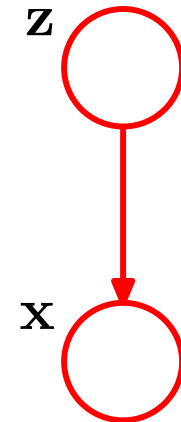
$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= \sum_{\mathbf{z}} \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \pi_k^{z_k} \\ &= \sum_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \underbrace{\pi_k}_{p(z_k=1)} \end{aligned}$$

Mixture models are latent variable models (2)

Model:

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= \sum_{\mathbf{z}} \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \pi_k^{z_k} \\ &= \sum_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \underbrace{\pi_k}_{p(z_k=1)} \end{aligned}$$

Graphical model:



So, we recover the same expression, but with a different interpretation

Why do we need latent variables?

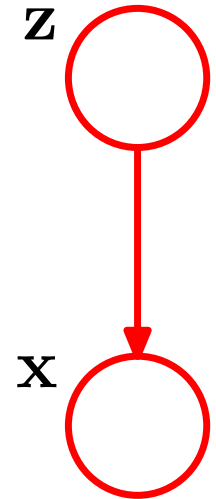
Reasons for introducing latent variables:

Artificial variable introduced to facilitate modelling

Example: Mixture model

Real quantity that is difficult to measure

Example: Underlying causes for disease



Why do we need latent variables?

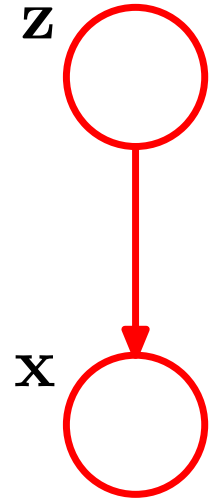
Reasons for introducing latent variables:

Artificial variable introduced to facilitate modelling

Example: Mixture model

Real quantity that is difficult to measure

Example: Underlying causes for disease



Types of latent variables:

Discrete (e.g. mixture models)

Structuring observations/features in groups/clusters

Continuous (PCA, factor analysis)

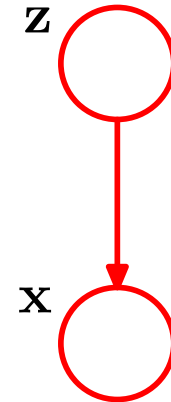
Dimensionality reduction

Sampling from a mixture model

The factorization depicted in the graphical model suggests how we can sample from the model:

1. Draw a sample z from $p(z)$
2. Draw a sample x from $p(x|z)$

(a simple example of *ancestral sampling*)

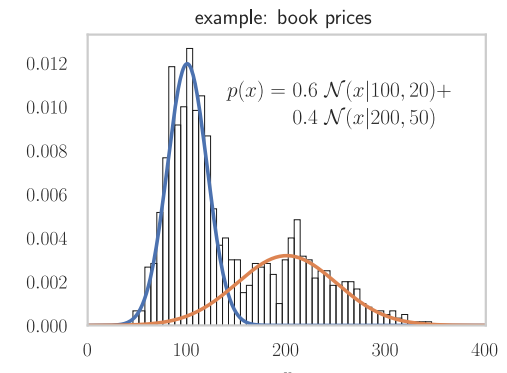
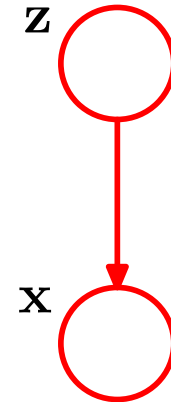


Sampling from a mixture model

The factorization depicted in the graphical model suggests how we can sample from the model:

1. Draw a sample z from $p(z)$
2. Draw a sample x from $p(x|z)$

(a simple example of *ancestral sampling*)



Which component is responsible for a given x ?

We can use Bayes rule to evaluate the probability $p(z_k = 1|\mathbf{x})$

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

These are referred to as the *responsibilities*, $\gamma(z_k)$, of component k to an observation x .

Which component is responsible for a given x ?

We can use Bayes rule to evaluate the probability $p(z_k = 1|\mathbf{x})$

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

These are referred to as the *responsibilities*, $\gamma(z_k)$, of component k to an observation x .

Note that we can evaluate the normalization constant $p(\mathbf{x})$, because z is discrete.

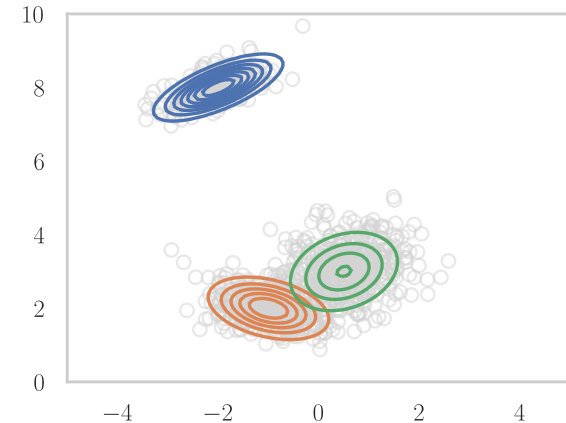
Which component is responsible for a given x ?

We can use Bayes rule to evaluate the probability $p(z_k = 1|\mathbf{x})$

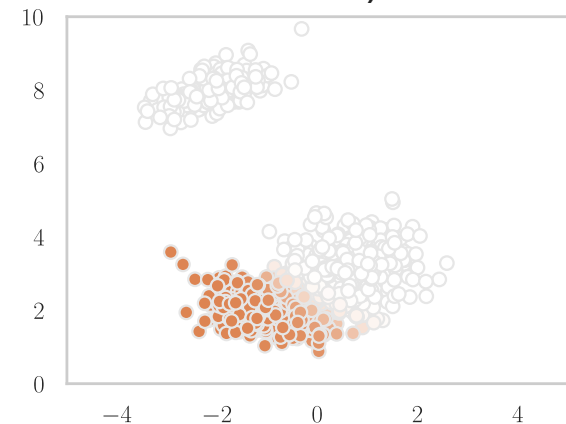
$$\begin{aligned} p(z_k = 1|\mathbf{x}) &= \frac{p(\mathbf{x}|z_k = 1)p(z_k = 1)}{p(\mathbf{x})} \\ &= \frac{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\pi_k}{\sum_{j=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\pi_j} \end{aligned}$$

These are referred to as the *responsibilities*, $\gamma(z_k)$, of component k to an observation x .

Note that we can evaluate the normalization constant $p(\mathbf{x})$, because z is discrete.



$p(z_k = 1|x)$ for $k = 1$ (the orange cluster):



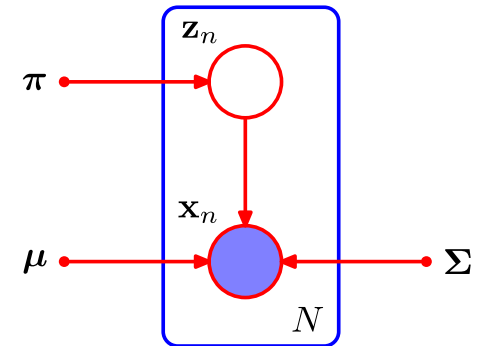
Parameter estimation in mixture models

Log-likelihood for $N \times D$ data matrix \mathbf{X} :

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

We would like to maximize this wrt $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$.

Unfortunately, unlike the Gaussian itself, we cannot optimize this in closed form (due to the sum inside the \ln). ?



Parameter estimation in mixture models

Log-likelihood for $N \times D$ data matrix \mathbf{X} :

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

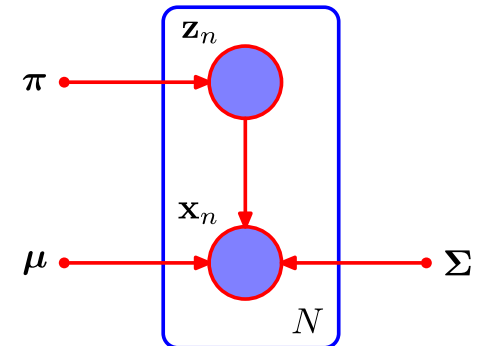
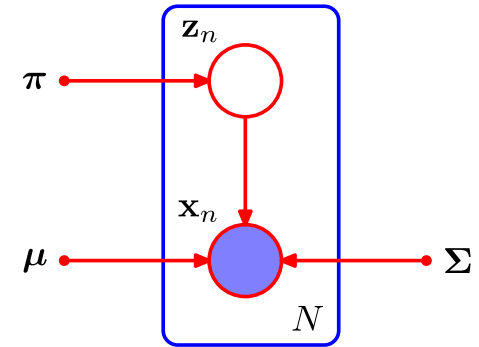
We would like to maximize this wrt $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$.

Unfortunately, unlike the Gaussian itself, we cannot optimize this in closed form (due to the sum inside the \ln). ?

Life would be easier if we had observed all z_n :

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}} \pi_k^{z_{nk}}$$

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



Expectation Maximization (idea)

The complete-data case had a convenient form

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

But we don't know z_{nk} .

Idea:

- For any choice of values for $\boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0$, we can calculate a posterior distribution $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$

Expectation Maximization (idea)

The complete-data case had a convenient form

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

But we don't know z_{nk} .

Idea:

- For any choice of values for $\boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0$, we can calculate a posterior distribution $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$
- We can then calculate the expectation $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})]$ under $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.

Expectation Maximization (idea)

The complete-data case had a convenient form

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

But we don't know z_{nk} .

Idea:

- For any choice of values for $\boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0$, we can calculate a posterior distribution $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$
- We can then calculate the expectation $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})]$ under $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.
- We can now maximize $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})]$ wrt $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$

Expectation Maximization

Procedure:

0. initialize π_0, μ_0, Σ_0 (e.g. randomly)

$t \leftarrow 0$

while not *converged*:

$t \leftarrow t + 1$

1. $\underbrace{\text{calculate } \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \pi, \mu, \Sigma)] \text{ under } p(\mathbf{Z} | \mathbf{X}, \pi_{t-1}, \mu_{t-1})}_{\text{Expectation step}}$

2. $\underbrace{(\pi_t, \mu_t, \Sigma_t) \leftarrow \operatorname{argmax}_{(\pi, \mu, \Sigma)} (\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \pi, \mu, \Sigma)])}_{\text{Maximization step}}$

Let's take a closer look at these steps.

E-step

The expectation goes inside the sum:

$$\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}_{\mathbf{Z}} [z_{nk}] \ln (\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k)$$

The expectation of the binary indicator variable z_{nk} is its probability:

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}} [z_{nk}] &= p(z_{nk} = 1 | \mathbf{x}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \frac{p(\mathbf{x}_n | z_{nk}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(z_{nk})}{p(\mathbf{x}_n)} = \frac{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k}{\sum_j^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \pi_j} = \gamma(z_{nk}) \end{aligned}$$

So, we see that $\mathbb{E}_{\mathbf{Z}} [z_{nk}]$ is simply the responsibility of component k to datapoint \mathbf{x}_n .

M-step

Now we just need to optimize:

$$\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \underbrace{\ln(\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k)}_{\ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \ln(\pi_k)} =: f(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

M-step

Now we just need to optimize:

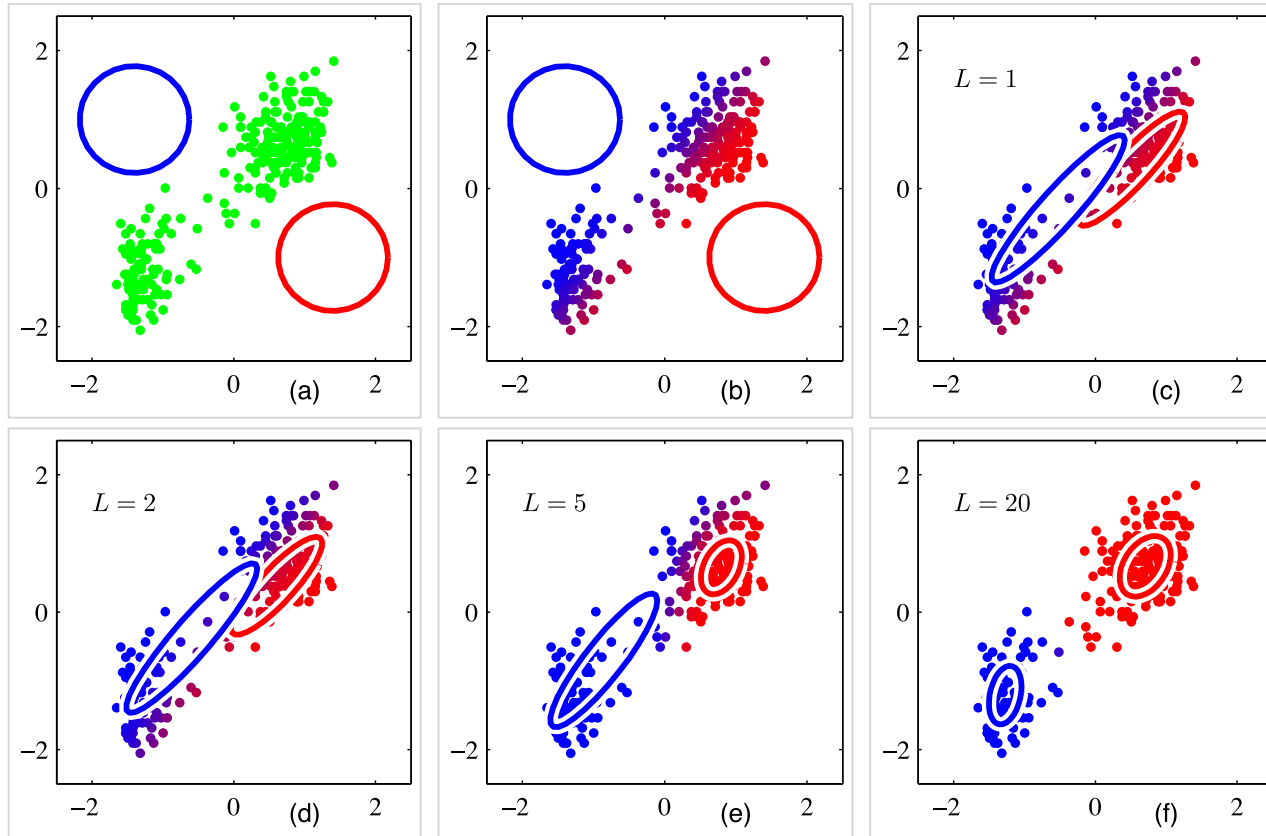
$$\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \underbrace{\ln(\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k)}_{\ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \ln(\pi_k)} =: f(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Considering single component $\boldsymbol{\mu}_k$. Set derivative to zero:

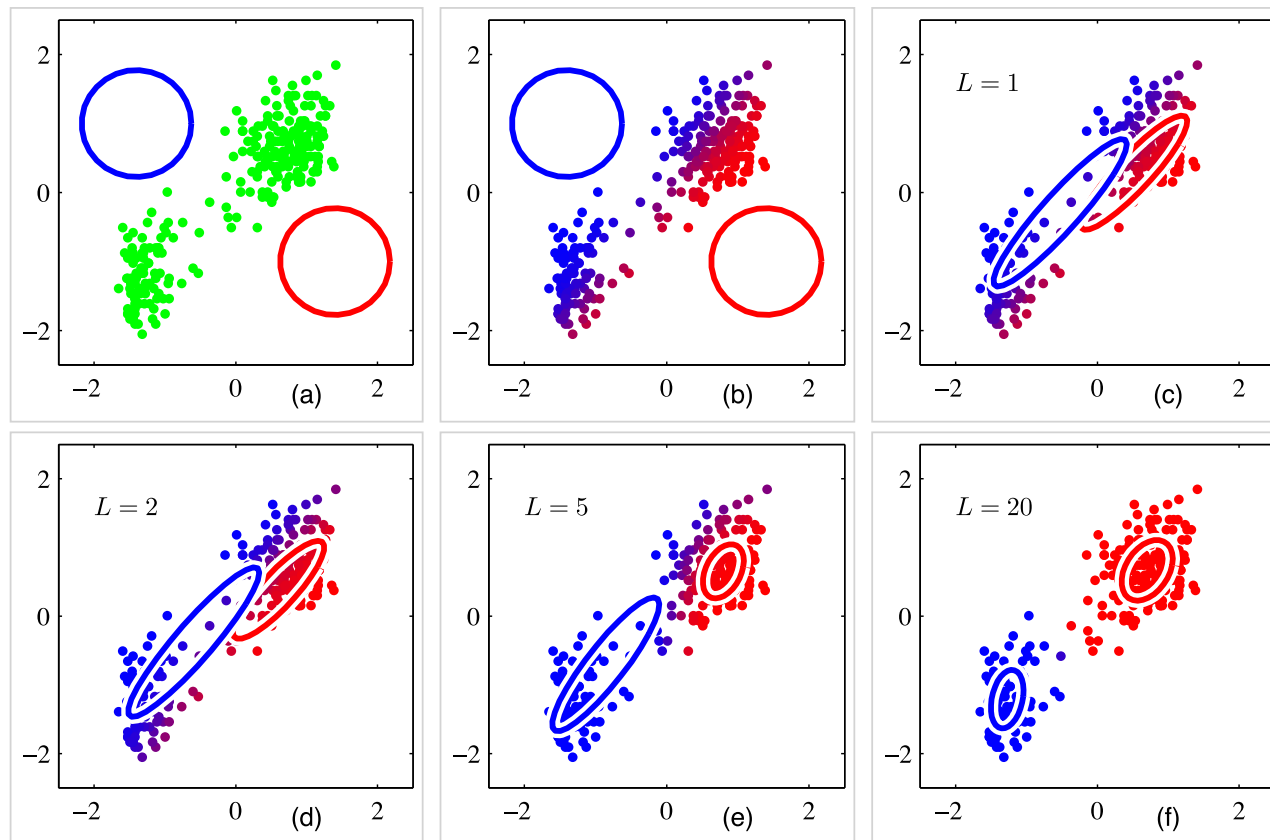
$$\begin{aligned} \frac{\partial f}{\partial \boldsymbol{\mu}_k} &= \sum_{n=1}^N \gamma(z_{nk}) \frac{\partial}{\partial \boldsymbol{\mu}_k} (\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \\ &= \sum_{n=1}^N \gamma(z_{nk}) \underbrace{\frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right)}_{\text{Look up in Matrix cookbook: eq 86}} \\ &= \sum_{n=1}^N \gamma(z_{nk}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad \Rightarrow \quad \boldsymbol{\mu}_k = \frac{1}{\sum_{n=1}^N \gamma(z_{nk})} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \end{aligned}$$

similar for the other parameters

EM illustration



EM illustration



Note: you can get faster convergence using Kmeans as an initializer.

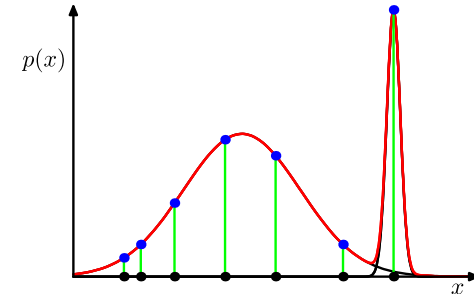
Is maximum likelihood a good idea?

Does it makes sense to optimize the likelihood?

Is maximum likelihood a good idea?

Does it makes sense to optimize the likelihood?

Not really. It's ill-posed.



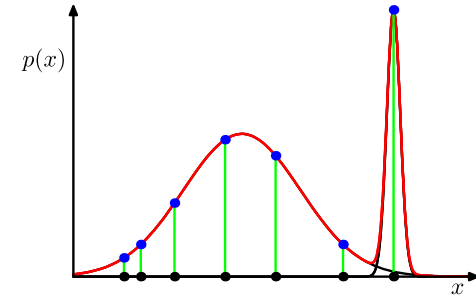
Is maximum likelihood a good idea?

Does it makes sense to optimize the likelihood?

Not really. It's ill-posed.

Why does EM work then?

Reasonable initialization + local optimization.



Is maximum likelihood a good idea?

Does it make sense to optimize the likelihood?

Not really. It's ill-posed.

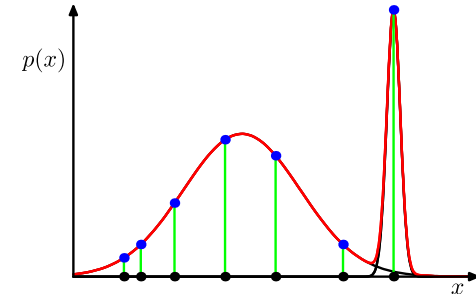
Why does EM work then?

Reasonable initialization + local optimization.

What is the proper way to fix this?

We can apply regularization to avoid the covariances in Σ becoming too small.

In the Bayesian setting: add a prior on Σ



EM in terms of log-likelihood bound

The product rule states $p(\mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})}$, so

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$$

General setting:

1. $p(\mathbf{X}|\boldsymbol{\theta})$ difficult to optimize
2. $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ easy to optimize

EM in terms of log-likelihood bound

The product rule states $p(\mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})}$, so

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$$

Add and subtract $\ln q(\mathbf{Z})$ (valid for any choice of q):

General setting:

1. $p(\mathbf{X}|\boldsymbol{\theta})$ difficult to optimize
2. $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ easy to optimize

EM in terms of log-likelihood bound

The product rule states $p(\mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})}$, so

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$$

Add and subtract $\ln q(\mathbf{Z})$ (valid for any choice of q):

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\theta}) &= \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln q(\mathbf{Z}) - (\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\ &= \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} - \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})}\end{aligned}$$

General setting:

1. $p(\mathbf{X}|\boldsymbol{\theta})$ difficult to optimize
2. $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ easy to optimize

EM in terms of log-likelihood bound

The product rule states $p(\mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})}$, so

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$$

Add and subtract $\ln q(\mathbf{Z})$ (valid for any choice of q):

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\theta}) &= \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) - \ln q(\mathbf{Z}) - (\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\ &= \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} - \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})}\end{aligned}$$

General setting:

1. $p(\mathbf{X}|\boldsymbol{\theta})$ difficult to optimize
2. $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ easy to optimize

Since this holds for any \mathbf{Z} , it also holds for $\mathbb{E}_{\mathbf{Z}}$ under any distribution over \mathbf{Z} . We choose $q(\mathbf{Z})$:

$$\underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \ln p(\mathbf{X}|\boldsymbol{\theta})}_{=\ln p(\mathbf{X}|\boldsymbol{\theta})} = \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(\ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(-\ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\text{KL}(q||p)}$$

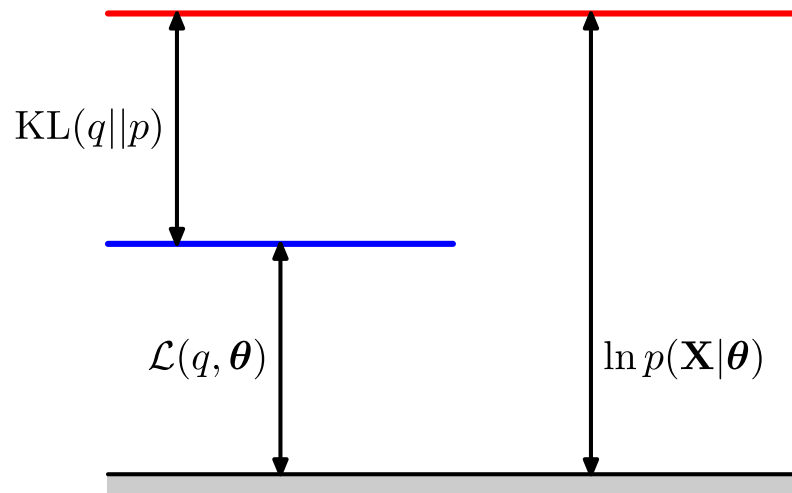
EM in terms of log-likelihood bound (2)

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(\ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(-\ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\text{KL}(q||p)}$$

The Kullback Leibler divergence, $\text{KL}(q||p)$ specifies how "close" $q(\mathbf{Z})$ is to $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$.

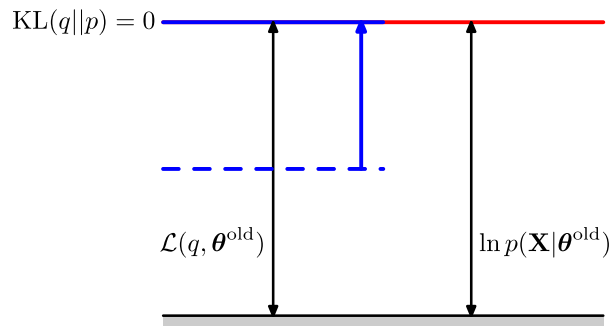
Since $\text{KL}(q||p) \geq 0$, $\mathcal{L}(q, \boldsymbol{\theta})$ is a lower bound on $\ln p(\mathbf{X}|\boldsymbol{\theta})$.

\mathcal{L} is sometimes referred to as an ELBO (evidence lower bound)



EM in terms of log-likelihood bound (3)

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(\ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(- \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\text{KL}(q||p)}$$



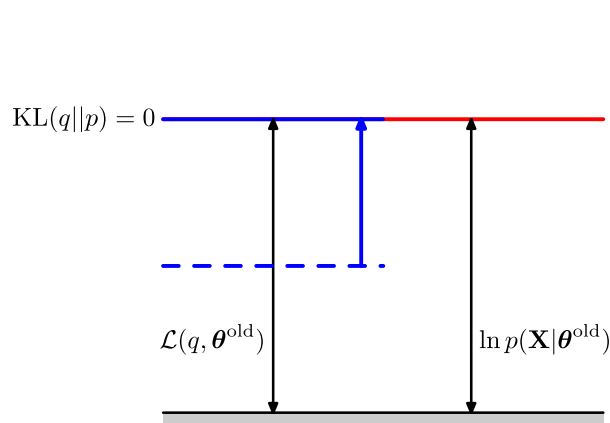
E-step (iteration t)

Optimize $\mathcal{L}(q, \boldsymbol{\theta}_{t-1})$ with respect to q , keeping $\boldsymbol{\theta}_{t-1}$ fixed. ?

$$\Rightarrow q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{t-1})$$

EM in terms of log-likelihood bound (3)

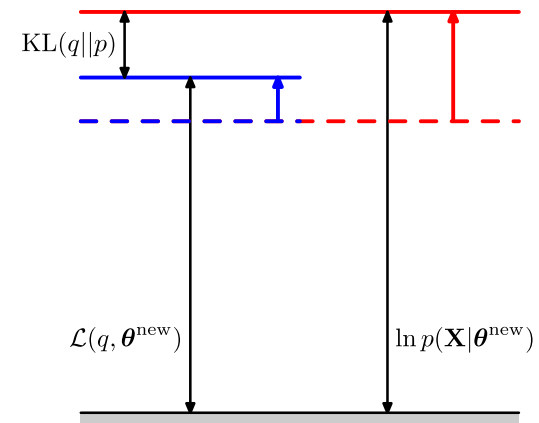
$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(\ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} \left(-\ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right)}_{\text{KL}(q||p)}$$



E-step (iteration t)

Optimize $\mathcal{L}(q, \boldsymbol{\theta}_{t-1})$ with respect to q , keeping $\boldsymbol{\theta}_{t-1}$ fixed. ?

$$\Rightarrow q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{t-1})$$

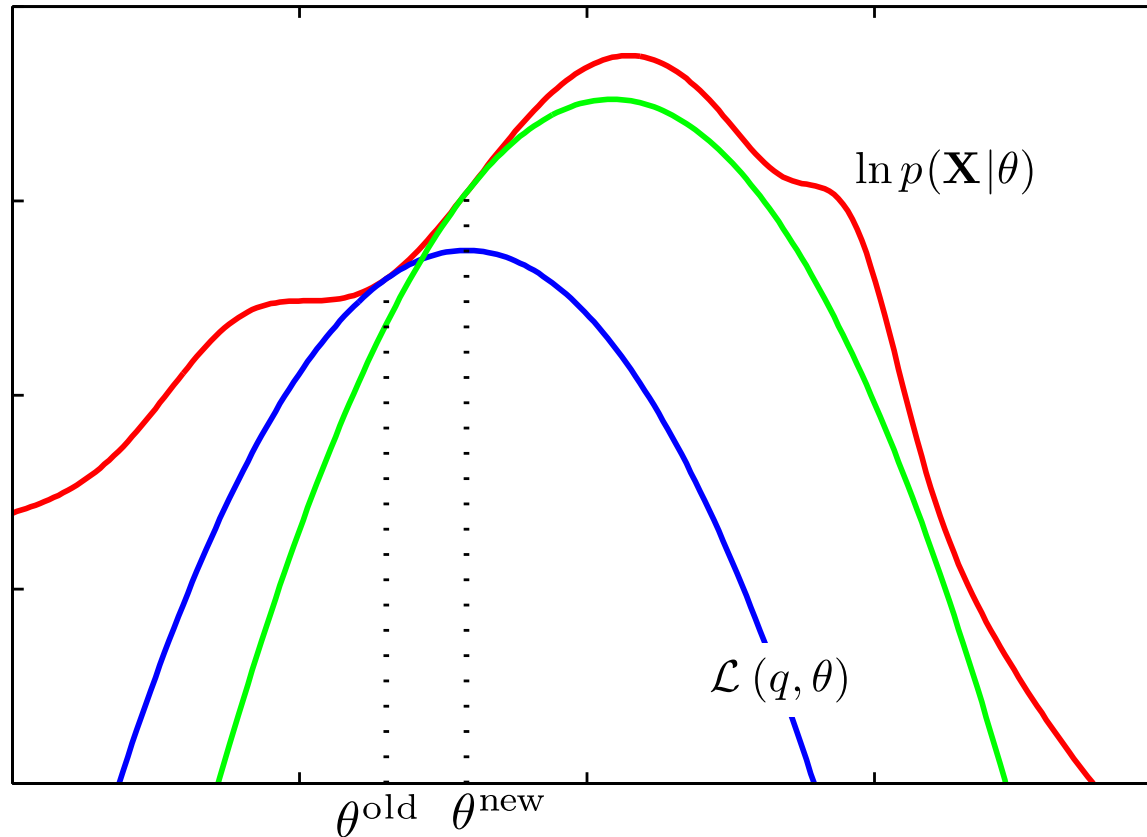


M-step (iteration t)

Optimize $\mathcal{L}(q, \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$, keeping q fixed. ?

$$\mathcal{L}(q, \boldsymbol{\theta}_t) \geq \mathcal{L}(q, \boldsymbol{\theta}_{t-1}) \Rightarrow p(\mathbf{X}|\boldsymbol{\theta}_t) \geq p(\mathbf{X}|\boldsymbol{\theta}_{t-1})$$

EM in terms of log-likelihood bound (4)



Mixture models - Concluding remarks

Mixture models allow modelling of complex distributions.

Mixture models - Concluding remarks

Mixture models allow modelling of complex distributions.

The latent variable provides structure/groups in the data

Mixture models - Concluding remarks

Mixture models allow modelling of complex distributions.

The latent variable provides structure/groups in the data

Gaussian mixtures are somewhat similar to Kmeans, but estimating (co)variances in addition to means.

Mixture models - Concluding remarks

Mixture models allow modelling of complex distributions.

The latent variable provides structure/groups in the data

Gaussian mixtures are somewhat similar to Kmeans, but estimating (co)variances in addition to means.

Mixture models can be trained using gradient descent, but the EM algorithm is often more efficient.

Continuous latent variable models

Remember PCA?

Core idea: Reduce the number of dimensions by finding linear combinations of features that capture most of the variance in the data.

Remember PCA?

Core idea: Reduce the number of dimensions by finding linear combinations of features that capture most of the variance in the data.

Two formal definitions:

- Maximize the variance of the projected data.
- Minimize the sum-of-squared of the projection errors.

Remember PCA?

Core idea: Reduce the number of dimensions by finding linear combinations of features that capture most of the variance in the data.

Two formal definitions:

- Maximize the variance of the projected data.
- Minimize the sum-of-squared of the projection errors.

These are equivalent.

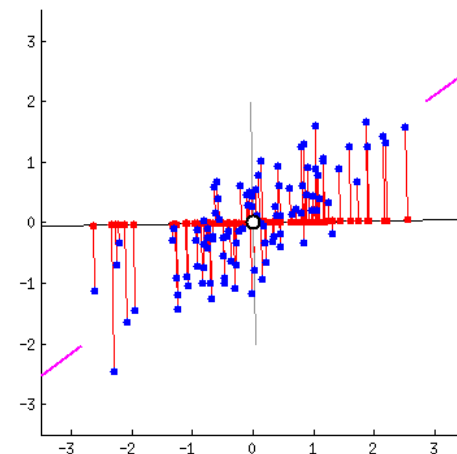


Figure from
<https://stats.stackexchange.com/a/140579>.

Remember PCA?

Core idea: Reduce the number of dimensions by finding linear combinations of features that capture most of the variance in the data.

Two formal definitions:

- Maximize the variance of the projected data.
- Minimize the sum-of-squared of the projection errors.

These are equivalent.

PCA is calculated as the eigenvectors of the data covariance matrix.

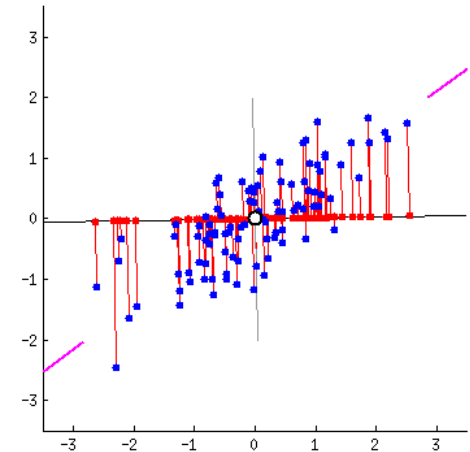


Figure from
<https://stats.stackexchange.com/a/140579>.

Why eigenvectors of the covariance matrix?

Consider a multivariate Gaussian with covariance Σ and mean 0

Since Σ is square, symmetric, we can diagonalize it

$$\begin{aligned}\Sigma &= R\tilde{\Sigma}R^T \\ &= \begin{bmatrix} | & | & \cdots & | \\ v_1 & v_2 & \cdots & v_D \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{bmatrix} \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ & \vdots & \\ - & v_D^T & - \end{bmatrix}\end{aligned}$$

where v_1, \dots, v_D are the eigenvectors and $\lambda_1, \dots, \lambda_D$ are the eigenvalues of Σ .

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) \propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right)$$

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) \quad \text{inverse of matrix product} \end{aligned}$$

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) && \text{inverse of matrix product} \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}^{-1}\mathbf{R}^T)\mathbf{x}\right) && \text{since } \mathbf{R} \text{ is orthogonal} \end{aligned}$$

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) && \text{inverse of matrix product} \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}^{-1}\mathbf{R}^T)\mathbf{x}\right) && \text{since } \mathbf{R} \text{ is orthogonal} \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{R}^T\mathbf{x})^T\tilde{\Sigma}^{-1}(\mathbf{R}^T\mathbf{x})\right) && \text{since } (AB)^T = B^T A^T \end{aligned}$$

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) && \text{inverse of matrix product} \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}^{-1}\mathbf{R}^T)\mathbf{x}\right) && \text{since } \mathbf{R} \text{ is orthogonal} \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{R}^T\mathbf{x})^T\tilde{\Sigma}^{-1}(\mathbf{R}^T\mathbf{x})\right) && \text{since } (AB)^T = B^T A^T \end{aligned}$$

\mathbf{R} can thus be interpreted as the rotation matrix which we can apply to our data such that our new feature dimensions have no covariance.

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) && \text{inverse of matrix product} \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}^{-1}\mathbf{R}^T)\mathbf{x}\right) && \text{since } \mathbf{R} \text{ is orthogonal} \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{R}^T\mathbf{x})^T\tilde{\Sigma}^{-1}(\mathbf{R}^T\mathbf{x})\right) && \text{since } (AB)^T = B^T A^T \end{aligned}$$

\mathbf{R} can thus be interpreted as the rotation matrix which we can apply to our data such that our new feature dimensions have no covariance.

Note that the *eigenvalues* are the variances along these dimensions.

Why eigenvectors of the covariance matrix? (2)

Plugging $\Sigma = \mathbf{R}\tilde{\Sigma}\mathbf{R}^T$ into the Gaussian PDF, we have

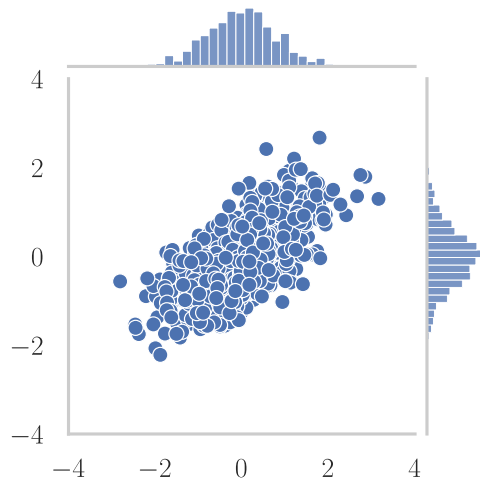
$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}\mathbf{R}^T)^{-1}\mathbf{x}\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T((\mathbf{R}^T)^{-1}\tilde{\Sigma}^{-1}\mathbf{R}^{-1})\mathbf{x}\right) && \text{inverse of matrix product} \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^T(\mathbf{R}\tilde{\Sigma}^{-1}\mathbf{R}^T)\mathbf{x}\right) && \text{since } \mathbf{R} \text{ is orthogonal} \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{R}^T\mathbf{x})^T\tilde{\Sigma}^{-1}(\mathbf{R}^T\mathbf{x})\right) && \text{since } (AB)^T = B^T A^T \end{aligned}$$

\mathbf{R} can thus be interpreted as the rotation matrix which we can apply to our data such that our new feature dimensions have no covariance.

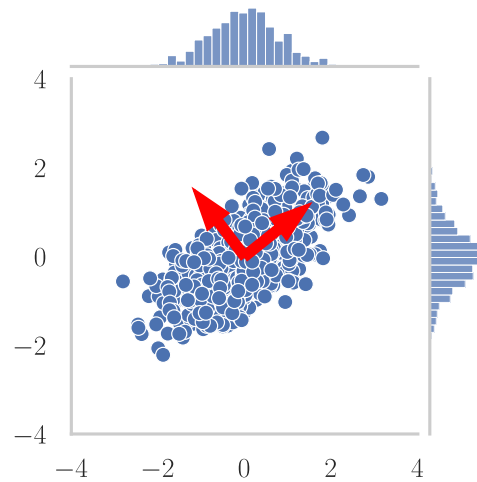
Note that the *eigenvalues* are the variances along these dimensions.

By sorting according to the eigenvalues (descending order) and choosing the top k , we get a dimensionality reduction scheme that explains $\frac{\sum_i^k \lambda_i}{\sum_i^D \lambda_i}$ of the total variance.

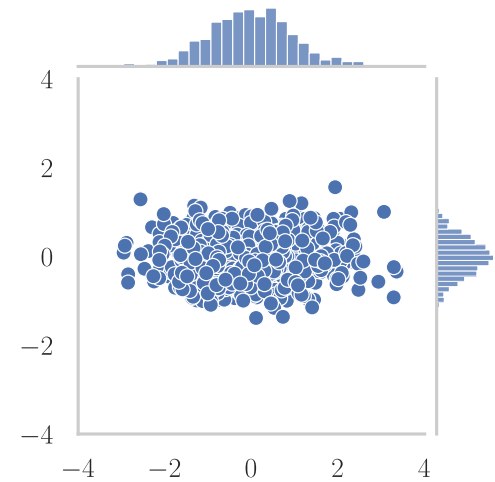
Example



$$\Sigma = \begin{bmatrix} 0.7 & 0.4 \\ 0.4 & 0.5 \end{bmatrix}$$
$$= R\tilde{\Sigma}R^T$$



$$R = \begin{bmatrix} 0.79 & -0.62 \\ 0.62 & 0.79 \end{bmatrix}$$
$$\tilde{\Sigma} = \begin{bmatrix} 1.01 & 0 \\ 0 & 0.19 \end{bmatrix}$$



$$R^T x$$

Probabilistic PCA

A probabilistic equivalent of PCA can be formulated as a latent variable model - where the latent variable is now continuous.

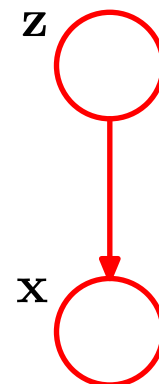
Ingredients:

Prior distribution over \mathbf{z}

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

Distribution over output variable - conditioned on \mathbf{z}

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$$



(same graphical model as
for the discrete case)

Probabilistic PCA

A probabilistic equivalent of PCA can be formulated as a latent variable model - where the latent variable is now continuous.

Ingredients:

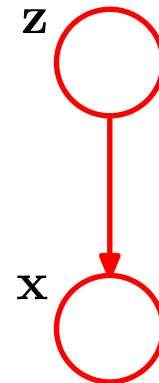
Prior distribution over \mathbf{z}

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

Distribution over output variable - conditioned on \mathbf{z}

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$$

The latent \mathbf{z} is M -dimensional, and \mathbf{W} is a $D \times M$ matrix mapping linearly from \mathbf{z} to \mathbf{x} .



(same graphical model as for the discrete case)

Probabilistic PCA

A probabilistic equivalent of PCA can be formulated as a latent variable model - where the latent variable is now continuous.

Ingredients:

Prior distribution over \mathbf{z}

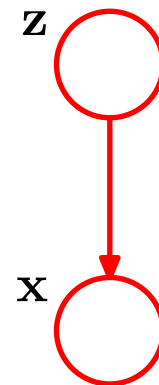
$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

Distribution over output variable - conditioned on \mathbf{z}

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$$

The latent \mathbf{z} is M -dimensional, and \mathbf{W} is a $D \times M$ matrix mapping linearly from \mathbf{z} to \mathbf{x} .

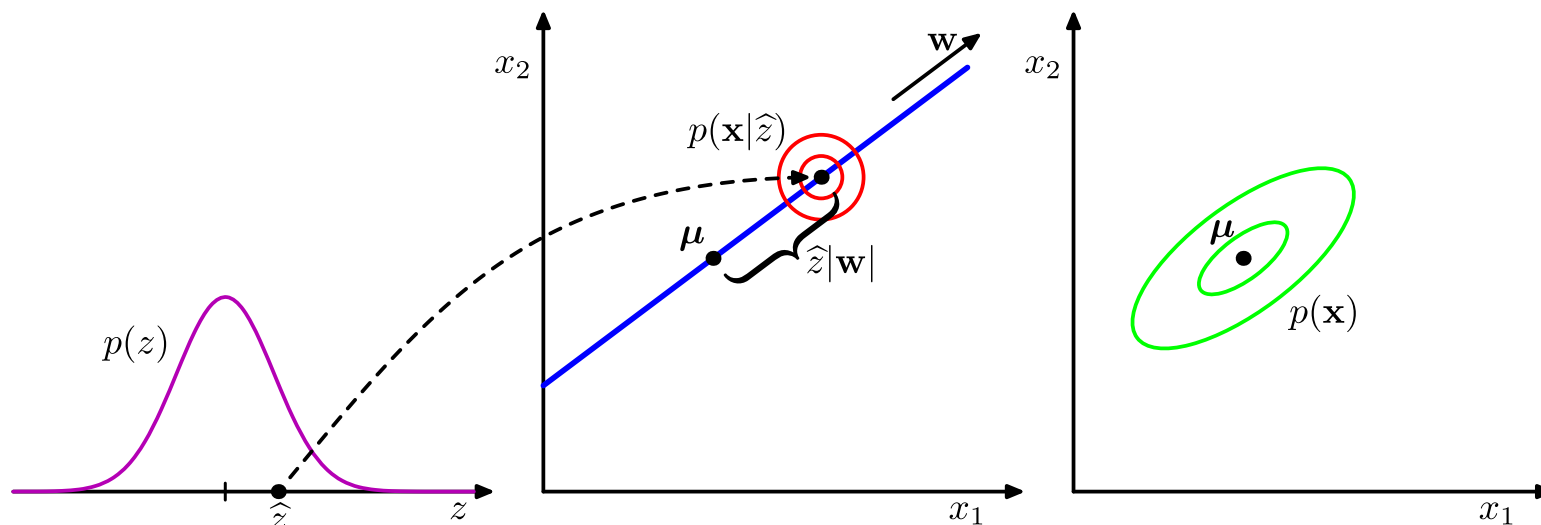
Note the simplicity of the model: the elements of \mathbf{x} are independent given \mathbf{z} , and the variances for all elements are the same.



(same graphical model as for the discrete case)

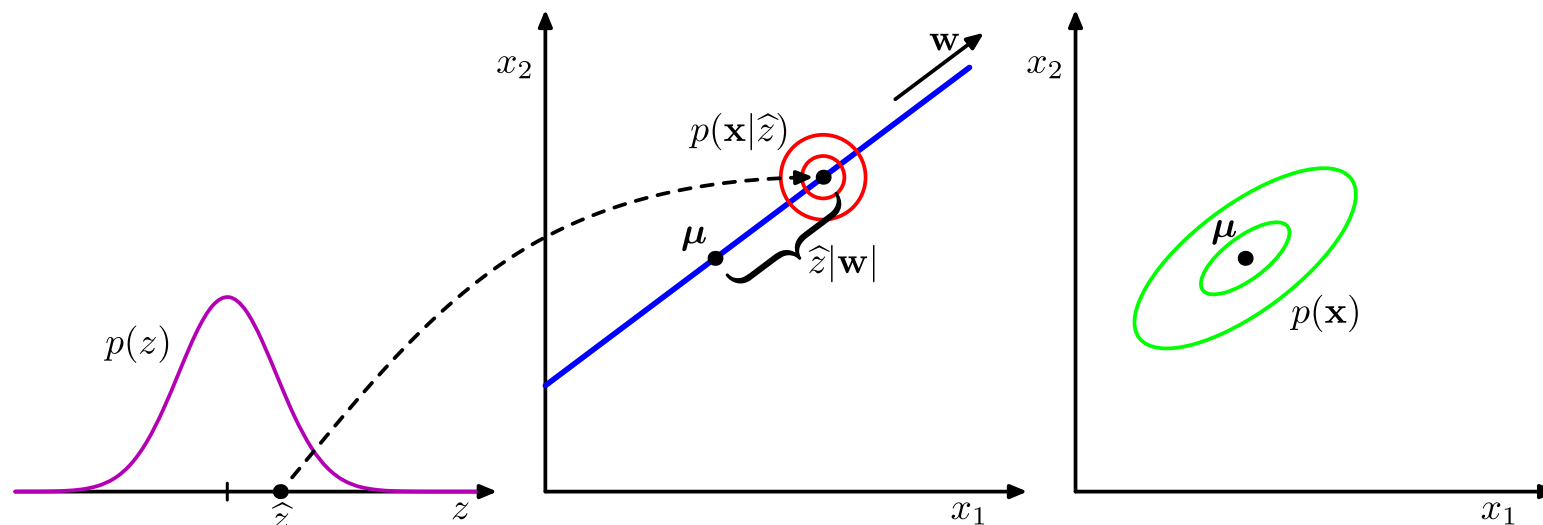
PPCA as a generative model

Similar to mixture models, we can generate samples from the model by first drawing a sample $\hat{z} \sim p(z)$, and then sampling $\mathbf{x} \sim p(\mathbf{x}|\hat{z})$:



PPCA as a generative model

Similar to mixture models, we can generate samples from the model by first drawing a sample $\hat{z} \sim p(z)$, and then sampling $\mathbf{x} \sim p(\mathbf{x}|\hat{z})$:



Note that $\mathbf{x} \sim p(\mathbf{x}|\hat{z})$ can be implemented as

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \sigma\boldsymbol{\epsilon} \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

PPCA: The joint is Gaussian

A special property of the PPCA model is that the joint is Gaussian.

$$\begin{aligned}\ln p(\mathbf{x}, \mathbf{z}) &= \ln(p(\mathbf{x}|\mathbf{z})p(\mathbf{z})) \\ &= \ln p(\mathbf{x}|\mathbf{z}) + \ln p(\mathbf{z}) = \ln \left(\mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) + \mathcal{N}(\mathbf{z}|0, \mathbf{I}) \right) \\ &= - \left(\frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu})^T (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu}) + \frac{1}{2} \mathbf{z}^T \mathbf{z} \right. \\ &\quad \left. + \frac{D}{2} \ln(2\pi\sigma^2) + \frac{M}{2} \ln(2\pi) \right)\end{aligned}$$

PPCA: The joint is Gaussian

A special property of the PPCA model is that the joint is Gaussian.

$$\begin{aligned}\ln p(\mathbf{x}, \mathbf{z}) &= \ln(p(\mathbf{x}|\mathbf{z})p(\mathbf{z})) \\&= \ln p(\mathbf{x}|\mathbf{z}) + \ln p(\mathbf{z}) = \ln \left(\mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) + \mathcal{N}(\mathbf{z}|0, \mathbf{I}) \right) \\&= - \left(\frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu})^T (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu}) + \frac{1}{2} \mathbf{z}^T \mathbf{z} \right. \\&\quad \left. + \frac{D}{2} \ln(2\pi\sigma^2) + \frac{M}{2} \ln(2\pi) \right) \\&= -\frac{1}{2\sigma^2} \left(\underbrace{\mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{W}\mathbf{z} - \mathbf{z}^T \mathbf{W}^T \mathbf{x} + \mathbf{z}^T \mathbf{W}^T \mathbf{W}\mathbf{z} + \sigma^2 \mathbf{z}^T \mathbf{z}}_{= \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} & -\mathbf{W}^T \\ -\mathbf{W} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}} \right. \\&\quad \left. - \underbrace{\mathbf{x}^T \boldsymbol{\mu} + \mathbf{z}^T \mathbf{W}^T \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{x} + \boldsymbol{\mu}^T \mathbf{W}\mathbf{z} + \boldsymbol{\mu}^T \boldsymbol{\mu}}_{= -2 \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} -\mathbf{W}^T \boldsymbol{\mu} \\ \boldsymbol{\mu} \end{pmatrix}} \right) - \frac{D}{2} \ln(2\pi\sigma^2) - \frac{M}{2} \ln(2\pi)\end{aligned}$$

PPCA: The joint is Gaussian (2)

$$\ln p(\mathbf{x}, \mathbf{z}) = -\frac{1}{2\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} & -\mathbf{W}^T \\ -\mathbf{W} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} + \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} -\mathbf{W}^T \boldsymbol{\mu} \\ \boldsymbol{\mu} \end{pmatrix} + \text{c}$$

PPCA: The joint is Gaussian (2)

$$\ln p(\mathbf{x}, \mathbf{z}) = -\frac{1}{2\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} & -\mathbf{W}^T \\ -\mathbf{W} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} + \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} -\mathbf{W}^T \boldsymbol{\mu} \\ \boldsymbol{\mu} \end{pmatrix} + \text{const}$$

In general, we can write the exponent of a Gaussian as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

PPCA: The joint is Gaussian (2)

$$\ln p(\mathbf{x}, \mathbf{z}) = -\frac{1}{2\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} & -\mathbf{W}^T \\ -\mathbf{W} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} + \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^T \begin{pmatrix} -\mathbf{W}^T \boldsymbol{\mu} \\ \boldsymbol{\mu} \end{pmatrix} + \text{const}$$

In general, we can write the exponent of a Gaussian as

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

We thus recognize $\ln p(\mathbf{x}, \mathbf{z})$ as a Gaussian log density $\mathcal{N}(\boldsymbol{\mu}_{z,x}, \boldsymbol{\Sigma}_{z,x})$ with:

$$\boldsymbol{\Sigma}_{z,x}^{-1} = \begin{pmatrix} \frac{1}{\sigma^2} \mathbf{W}^T \mathbf{W} + \mathbf{I} & -\frac{1}{\sigma^2} \mathbf{W}^T \\ -\frac{1}{\sigma^2} \mathbf{W} & \frac{1}{\sigma^2} \mathbf{I} \end{pmatrix} \Rightarrow \underbrace{\boldsymbol{\Sigma}_{z,x} = \begin{pmatrix} \mathbf{I} & \mathbf{W}^T \\ \mathbf{W} & \sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^T \end{pmatrix}}_{\text{Matrix Cookbook eq 399-400)}$$

and

$$\boldsymbol{\Sigma}_{z,x}^{-1} \boldsymbol{\mu}_{z,x} = \frac{1}{\sigma^2} \begin{pmatrix} -\mathbf{W}^T \boldsymbol{\mu} \\ \boldsymbol{\mu} \end{pmatrix} \Rightarrow \boldsymbol{\mu}_{z,x} = \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{pmatrix}$$

PPCA: marginal $p(\mathbf{x})$ is also Gaussian

We now know that $p(\mathbf{x}, \mathbf{z})$ is jointly Gaussian

$$p(\mathbf{x}, \mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_{z,x}, \boldsymbol{\Sigma}_{z,x}) \quad \boldsymbol{\mu}_{z,x} = \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{pmatrix} \quad \boldsymbol{\Sigma}_{z,x} = \begin{pmatrix} \mathbf{I} & \mathbf{W}^T \\ \mathbf{W} & \sigma^2 \mathbf{I} + \mathbf{W}\mathbf{W}^T \end{pmatrix}$$

That means that the marginal $p(\mathbf{x})$ is also Gaussian. We can actually read off its parameters directly from $\boldsymbol{\mu}_{z,x}$ and $\boldsymbol{\Sigma}_{z,x}$

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \quad \boldsymbol{\mu}_x = \boldsymbol{\mu} \quad \boldsymbol{\Sigma}_x = \sigma^2 \mathbf{I} + \mathbf{W}\mathbf{W}^T$$

PPCA: the posterior is also Gaussian

Remember - joint is gaussian:

$$p(\mathbf{x}, \mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_{z,x}, \boldsymbol{\Sigma}_{z,x}) \quad \boldsymbol{\mu}_{z,x} = \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{pmatrix} \quad \boldsymbol{\Sigma}_{z,x} = \begin{pmatrix} \mathbf{I} & \mathbf{W}^T \\ \mathbf{W} & \sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^T \end{pmatrix}$$

Since the joint is Gaussian, we can evaluate the conditional $p(\mathbf{z}|\mathbf{x})$

Using eq 353 from the matrix cookbook, we have:

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \mathbf{0} + \mathbf{W}^T (\sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^T)^{-1} (\mathbf{x} - \boldsymbol{\mu}), \right. \\ \left. \mathbf{I} - \mathbf{W}^T (\sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^T)^{-1} \mathbf{W} \right)$$

We can use the Woodbury identity (cookbook eq 159) for the inversion:

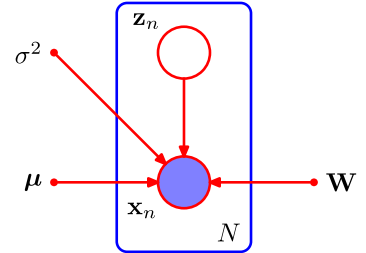
$$(\sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^T)^{-1} = \sigma^{-2} \mathbf{I} - \sigma^{-2} \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T \quad \text{where } \mathbf{M} = \mathbf{W}^T \mathbf{W} + \epsilon$$

Parameter estimation - Maximum likelihood (1)

Log likelihood function:

$$\begin{aligned} ll(\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \ln p(\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) \\ &= \sum_{n=1}^N \ln p(\mathbf{x}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \end{aligned}$$

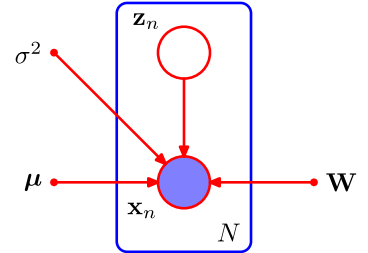
$$= -\frac{N}{2} \left(\frac{D}{2} \ln(2\pi) + \ln(\det(\boldsymbol{\Sigma}_x)) \right) - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_x^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$



Parameter estimation - Maximum likelihood (1)

Log likelihood function:

$$\begin{aligned} ll(\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \ln p(\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) \\ &= \sum_{n=1}^N \ln p(\mathbf{x}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{N}{2} \left(\frac{D}{2} \ln(2\pi) + \ln(\det(\boldsymbol{\Sigma}_x)) \right) - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_x^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$



Optimize wrt $\boldsymbol{\mu}$:

$$\begin{aligned} \frac{\partial ll}{\partial \boldsymbol{\mu}} &= -\frac{1}{2} \sum_{n=1}^N \frac{\partial ll}{\partial \boldsymbol{\mu}} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}_x^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \\ &= -\frac{1}{2} \sum_{n=1}^N -2\boldsymbol{\Sigma}_x (\mathbf{x}_n - \boldsymbol{\mu}) \quad \text{Matrix cookbook eq 86} \\ &= 0 \Rightarrow \boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \bar{\mathbf{x}} \end{aligned}$$

Parameter estimation - Maximum likelihood (2)

Optimizing $ll(\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$ wrt to \mathbf{W} and σ^2 is more difficult.

Tipping and Bishop prove in a 1999 paper that the likelihood is optimized when

$$\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

$$\sigma_{ML}^2 = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i$$

Where \mathbf{U}_M is a $D \times M$ matrix containing the M eigen vectors with largest eigenvalues, and \mathbf{L}_M is a diagonal matrix with the corresponding eigenvalues.

Parameter estimation - Maximum likelihood (2)

Optimizing $ll(\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$ wrt to \mathbf{W} and σ^2 is more difficult.

Tipping and Bishop prove in a 1999 paper that the likelihood is optimized when

$$\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

$$\sigma_{ML}^2 = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i$$

Where \mathbf{U}_M is a $D \times M$ matrix containing the M eigen vectors with largest eigenvalues, and \mathbf{L}_M is a diagonal matrix with the corresponding eigenvalues.

The \mathbf{W} thus behaves exactly as in regular PCA. We can therefore also estimate it in the same way.

Parameter estimation - Maximum likelihood (2)

Optimizing $ll(\boldsymbol{\mu}, \mathbf{W}, \sigma^2)$ wrt to \mathbf{W} and σ^2 is more difficult.

Tipping and Bishop prove in a 1999 paper that the likelihood is optimized when

$$\mathbf{W}_{ML} = \mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

$$\sigma_{ML}^2 = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i$$

Where \mathbf{U}_M is a $D \times M$ matrix containing the M eigen vectors with largest eigenvalues, and \mathbf{L}_M is a diagonal matrix with the corresponding eigenvalues.

The \mathbf{W} thus behaves exactly as in regular PCA. We can therefore also estimate it in the same way.

\mathbf{R} is an arbitrary orthogonal matrix, which does not affect the covariance matrix:

$$\begin{aligned} \mathbf{C}_{ML} &= \sigma^2 \mathbf{I} + \underbrace{\mathbf{W}_{ML} \mathbf{W}_{ML}^T}_{=} \\ &= (\mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R})(\mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R})^T \\ &= \mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \underbrace{\mathbf{R} \mathbf{R}^T}_{=\mathbf{I}} (\mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2})^T \end{aligned}$$

i.e. an arbitrary rotation of latent space.

Interpretation of Probabilistic PCA

Consider a direction in your output space, represented by unit vector \mathbf{v}

According to the PPCA model, the variance along this direction is $\mathbf{v}^T \boldsymbol{\Sigma}_x \mathbf{v}$.

Consider two scenarios:

1. \mathbf{v} is orthogonal to the M selected principal component vectors

$$\mathbf{v}^T \boldsymbol{\Sigma}_x \mathbf{v} = \mathbf{v}^T \sigma^2 \mathbf{I} \mathbf{v} + \underbrace{\mathbf{v}^T (\mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} (\mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2})^T \mathbf{v})}_{=0} = \sigma^2$$

2. \mathbf{v} points along the i th eigenvector

$$\mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} = \mathbf{v}^T \sigma^2 \mathbf{I} \mathbf{v} + \underbrace{\mathbf{v}^T (\mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2})}_{(\lambda_i - \sigma^2)^{1/2}} \underbrace{((\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2})^T \mathbf{U}_M^T \mathbf{v}}_{(\lambda_i - \sigma^2)^{1/2}} = \lambda_i$$

This is quite nice! For all the principal components that we have included, the model correctly captures the variance of the data. For those directions that we leave out, it uses $\sigma^2 = \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i$, i.e. the average variance of the left-out principal components.

Application of PPCA

When modelling continuous data in D dimensions, you will often have to make the following choice:

- Use a full covariance matrix

 - For high dimensions, the $D \times D$ can be difficult to estimate, requiring a lot of data.

- Use a diagonal covariance matrix

 - Now, you only need to estimate D variance values - but you assume that there is no covariance among the dimensions, which is probably unrealistic.

Application of PPCA

When modelling continuous data in D dimensions, you will often have to make the following choice:

- Use a full covariance matrix

 - For high dimensions, the $D \times D$ can be difficult to estimate, requiring a lot of data.

- Use a diagonal covariance matrix

 - Now, you only need to estimate D variance values - but you assume that there is no covariance among the dimensions, which is probably unrealistic.

Probabilistic PCA provides an intermediate between these extremes - allowing you to include the most important covariances while summarizing the remaining variances with an average.

EM for Probabilistic PCA

Can't we use EM to estimate the parameters of the PPCA model?

EM for Probabilistic PCA

Can't we use EM to estimate the parameters of the PPCA model?

Yes we can! But why would we want to do that when we have a closed form solution?

EM for Probabilistic PCA

Can't we use EM to estimate the parameters of the PPCA model?

Yes we can! But why would we want to do that when we have a closed form solution?

- Can be faster in higher dimensions.

EM for Probabilistic PCA

Can't we use EM to estimate the parameters of the PPCA model?

Yes we can! But why would we want to do that when we have a closed form solution?

- Can be faster in higher dimensions.
- Is necessary for other models, like Factor analysis.

EM for Probabilistic PCA

Can't we use EM to estimate the parameters of the PPCA model?

Yes we can! But why would we want to do that when we have a closed form solution?

- Can be faster in higher dimensions.
- Is necessary for other models, like Factor analysis.
- Natural way to deal with missing data.

EM for Probabilistic PCA (2)

Same procedure as for mixture models

In iteration t :

Expectation step:

Calculate the expectation of the complete-data likelihood $p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)$, under the distribution $p(z|x, \boldsymbol{\mu}_{t-1}, \mathbf{W}_{t-1}, \sigma_{t-1}^2)$

Maximization step

Optimize this completed likelihood to obtain new values $\boldsymbol{\mu}_t, \mathbf{W}_t, \sigma_t^2$

EM for Probabilistic PCA - E-step

Complete-data log-likelihood:

$$\begin{aligned} ll_c &= \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N (\ln p(\mathbf{x}_n | z_n) + \ln p(z_n)) \\ &= \sum_{n=1}^N (\ln (\mathcal{N}(\mathbf{x}_n | \mathbf{W} \mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) + \mathcal{N}(z_n | 0, \mathbf{I}))) \\ &= - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \mathbf{W} \mathbf{z}_n - \boldsymbol{\mu}\|^2 + \frac{M}{2} \ln(2\pi) + \right. \\ &= - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \underbrace{(\mathbf{W} \mathbf{z}_n)^T (\mathbf{W} \mathbf{z}_n)}_{= \mathbf{z}_n^T \mathbf{W}^T \mathbf{W} \mathbf{z}_n = \text{tr}(\mathbf{z}_n \mathbf{z}_n^T)} \right. \\ &\quad \left. - \frac{1}{\sigma^2} \underbrace{(\mathbf{W} \mathbf{z}_n)^T}_{= \mathbf{z}_n^T \mathbf{W}^T} (\mathbf{x}_n - \boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi) + \frac{1}{2} \underbrace{\mathbf{z}_n^T \mathbf{z}_n}_{= \text{tr}(\mathbf{z}_n \mathbf{z}_n^T)} \right) \end{aligned}$$

EM for Probabilistic PCA - E-step (2)

Taking the expectation wrt \mathbf{z} :

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}}[l_c] = & - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \text{tr}(\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right. \\ & \left. - \frac{1}{\sigma^2} (\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \mathbf{W}^T)(\mathbf{x}_n - \boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi) + \frac{1}{2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^T) \right)\end{aligned}$$

EM for Probabilistic PCA - E-step (2)

Taking the expectation wrt \mathbf{z} :

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}}[l_c] = & - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \text{tr}(\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right. \\ & \left. - \frac{1}{\sigma^2} (\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \mathbf{W}^T) (\mathbf{x}_n - \boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi) + \frac{1}{2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^T) \right)\end{aligned}$$

Remember that all these expectations are taken wrt

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}_{t-1}^{-1} \mathbf{W}_{t-1}^T (\mathbf{x} - \boldsymbol{\mu}_{t-1}), \sigma_{t-1}^2 \mathbf{M}_{t-1}^{-1}) \quad \text{where } \mathbf{M}_{t-1} = \mathbf{W}_{t-1}^T \mathbf{W}_{t-1}.$$

EM for Probabilistic PCA - E-step (2)

Taking the expectation wrt z :

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}}[l_c] = & - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \text{tr}(\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right. \\ & \left. - \frac{1}{\sigma^2} (\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \mathbf{W}^T)(\mathbf{x}_n - \boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi) + \frac{1}{2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^T) \right)\end{aligned}$$

Remember that all these expectations are taken wrt

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}_{t-1}^{-1} \mathbf{W}_{t-1}^T (\mathbf{x} - \boldsymbol{\mu}_{t-1}), \sigma_{t-1}^2 \mathbf{M}_{t-1}^{-1}) \quad \text{where } \mathbf{M}_{t-1} = \mathbf{W}_{t-1}^T \mathbf{W}_{t-1}.$$

So, in the E-step, all we need to calculate is:

$$\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] = \mathbf{M}_{t-1}^{-1} \mathbf{W}_{t-1}^T (\mathbf{x} - \boldsymbol{\mu}_{t-1})$$

$$\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma_{t-1}^2 \mathbf{M}_{t-1}^{-1} + \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \quad \text{since } \text{cov}(\mathbf{z}_n) = \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] - \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T$$

EM for Probabilistic PCA - E-step (2)

Taking the expectation wrt \mathbf{z} :

$$\begin{aligned}\mathbb{E}_{\mathbf{Z}}[ll_c] = & - \sum_{n=1}^N \left(\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \text{tr}(\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right. \\ & \left. - \frac{1}{\sigma^2} (\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \mathbf{W}^T)(\mathbf{x}_n - \boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi) + \frac{1}{2} \text{tr}(\mathbf{z}_n \mathbf{z}_n^T) \right)\end{aligned}$$

Remember that all these expectations are taken wrt

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}_{t-1}^{-1} \mathbf{W}_{t-1}^T (\mathbf{x} - \boldsymbol{\mu}_{t-1}), \sigma_{t-1}^2 \mathbf{M}_{t-1}^{-1}) \quad \text{where } \mathbf{M}_{t-1} = \mathbf{W}_{t-1}^T \mathbf{W}_{t-1}.$$

So, in the E-step, all we need to calculate is:

$$\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] = \mathbf{M}_{t-1}^{-1} \mathbf{W}_{t-1}^T (\mathbf{x} - \boldsymbol{\mu}_{t-1})$$

$$\mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma_{t-1}^2 \mathbf{M}_{t-1}^{-1} + \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T \quad \text{since } \text{cov}(\mathbf{z}_n) = \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n \mathbf{z}_n^T] - \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n] \mathbb{E}_{\mathbf{Z}}[\mathbf{z}_n]^T$$

For the M-step, we can now differentiate ll_c wrt \mathbf{W} and σ^2 and set to zero.

Factor analysis

Factor analysis is similar to Probabilistic PCA, but with slightly more complex output covariance

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \mathbf{I}\boldsymbol{\sigma}^2)$$

Remember that in Probabilistic PCA we had simply:

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \mathbf{I}\sigma^2)$$

In Factor analysis, we thus estimate variances for all data dimensions separately.

Factor analysis

Factor analysis is similar to Probabilistic PCA, but with slightly more complex output covariance

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \mathbf{I}\boldsymbol{\sigma}^2)$$

Remember that in Probabilistic PCA we had simply:

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \mathbf{I}\sigma^2)$$

In Factor analysis, we thus estimate variances for all data dimensions separately.

With this small change, the model no longer has a closed-form maximum likelihood solution. But EM is still applicable.

Probabilistic PCA - Summary

Probabilistic PCA is a simple continuous latent variable model that brings PCA functionality into the rich framework of probabilistic modelling

Probabilistic PCA - Summary

Probabilistic PCA is a simple continuous latent variable model that brings PCA functionality into the rich framework of probabilistic modelling

Due to the particular choice of output variance, there is a closed-form maximum likelihood solution to the \mathbf{W} , μ , σ parameters. We can also use EM, which can be more efficient in higher dimensions.

Probabilistic PCA - Summary

Probabilistic PCA is a simple continuous latent variable model that brings PCA functionality into the rich framework of probabilistic modelling

Due to the particular choice of output variance, there is a closed-form maximum likelihood solution to the \mathbf{W} , μ , σ parameters. We can also use EM, which can be more efficient in higher dimensions.

In Factor Analysis, we change this assumption (now having separate variances for each dimension). This no longer has a closed form maximum likelihood solution.

Probabilistic PCA - Summary

Probabilistic PCA is a simple continuous latent variable model that brings PCA functionality into the rich framework of probabilistic modelling

Due to the particular choice of output variance, there is a closed-form maximum likelihood solution to the \mathbf{W} , μ , σ parameters. We can also use EM, which can be more efficient in higher dimensions.

In Factor Analysis, we change this assumption (now having separate variances for each dimension). This no longer has a closed form maximum likelihood solution.

For both models, we should remember that the maximum likelihood only determines the latent space up to an arbitrary rotation.

Perspectives

The techniques covered today are important steps towards understanding the variational autoencoder.

- The considerations about bounds in the discussion of EM set the stage for variational inference.
- We will see that variational autoencoders arise as a natural generalization of probabilistic PCA.

