

Latent Variable Models

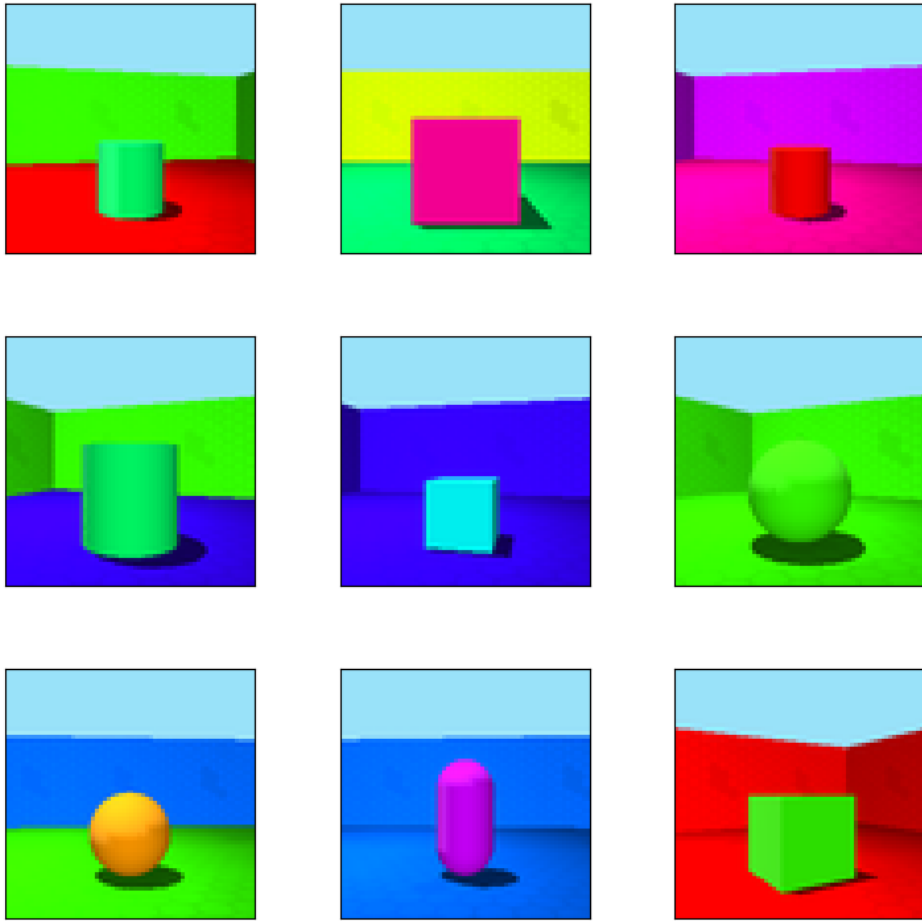
Advanced Statistical Inference

Simone Rossi

Latent variable models

Our objective is to learn the data distribution $p(\boldsymbol{x})$, but we suppose that each data point \boldsymbol{x} is associated with a **latent variable** z .

For example, image we want to learn the distribution of images of objects, like the one below:



Each image is a huge vector of pixels $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, where H is the height, W is the width and C is the number of channels (in this case, $64 \times 64 \times 3$).

Each image can be described by a set of 6 latent variables: floor, wall and object color, shape, orientation and scale.

Latent variable models

Latent variables are unobserved variables, we cannot measure them directly, but we can infer them from the observed data.

In math terms, we model our data distribution as a **marginal** of the joint distribution of the observed data and the latent variables:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad \text{or} \quad p(\mathbf{x}) = \sum_i p(\mathbf{x}, \mathbf{z}_i)$$

where $p(\mathbf{x}, \mathbf{z})$ is the **joint distribution** of the observed data and the latent variables.

Mixture models

Mixture models

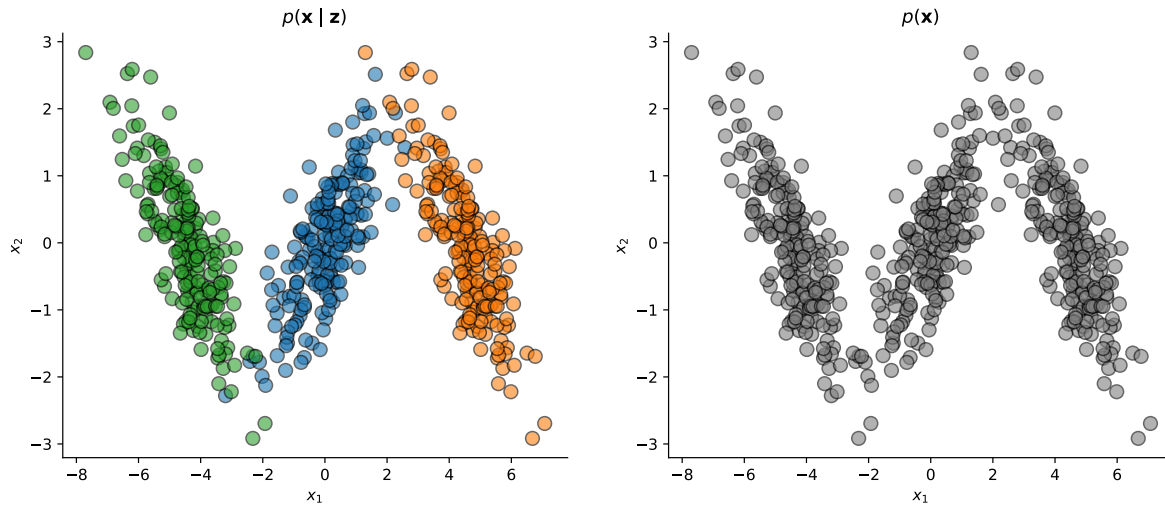
In mixture models, we assume that the data is generated from a mixture of several distributions

$$p(\mathbf{x}) = \sum_i p(\mathbf{x}, \mathbf{z}_i) = \sum_i p(\mathbf{x} | \mathbf{z}_i) p(\mathbf{z}_i)$$

where $p(\mathbf{z})$ is the **mixture distribution** and $p(\mathbf{x} | \mathbf{z})$ is the **conditional distribution** of the data given the latent variables.

- The latent variables \mathbf{z} are often called **clusters**.
- The *likelihood* is often assumed to be Gaussian, i.e., $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- The *prior* is often assumed to be a categorical distribution, i.e., $p(\mathbf{z}) = \text{Cat}(\mathbf{z}; \boldsymbol{\pi})$

Gaussian Mixture Models



Modeling the data distribution

- 1-hot encoding of the **discrete latent variable** $z_k \in \{0, 1\}$, with prior:

$$p(z_k = 1) = \pi_k \quad \text{and} \quad \sum_k \pi_k = 1$$

- The **conditional distribution** of the data given the latent variable is:

$$p(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- The **joint distribution** of the data and the latent variable is:

$$p(\mathbf{x}, z_k = 1) = p(\mathbf{x} | z_k = 1)p(z_k = 1) = \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- The **marginal distribution** of the data is:

$$p(\mathbf{x}) = \sum_k p(\mathbf{x}, z_k = 1) = \sum_k \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Posterior distribution

We are interested in the **posterior distribution** of a cluster assignment given the data:

$$\begin{aligned} p(z_k = 1 | \mathbf{x}) &= \frac{p(\mathbf{x} | z_k = 1)p(z_k = 1)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x} | z_k = 1)p(z_k = 1)}{\sum_j p(\mathbf{x} | z_j = 1)p(z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} := \gamma(z_k) \end{aligned}$$

$\gamma(z_k)$ is the **responsibility** of cluster k for data point \mathbf{x} .

Log-likelihood

Given a dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ i.i.d, the **log-likelihood** of the data is:

$$\begin{aligned} \log p(\mathbf{X} | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\boldsymbol{\Sigma}_k\}) &= \log \prod_{i=1}^N p(\mathbf{x}_i) \\ &= \sum_{i=1}^N \log p(\mathbf{x}_i | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\boldsymbol{\Sigma}_k\}) \\ &= \sum_{i=1}^N \log \sum_{k=1}^K p(\mathbf{x}_i | z_k = 1)p(z_k = 1) \\ &= \sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Warning

We cannot simplify it further, because the log and the sum do not commute. How can we optimize it?

Expectation-Maximization algorithm

We need to maximize the log-likelihood of the data w.r.t. the parameters $\boldsymbol{\pi}$, $\{\boldsymbol{\mu}_k\}$ and $\{\boldsymbol{\Sigma}_k\}$, for $k = 1, \dots, K$

$$\log p(\mathbf{X}) = \sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- The problem is not convex.
- No closed-form solution! Stationary points depend on the posterior $\gamma(\mathbf{z}_{nk})$ for each data point \mathbf{x}_n .
- ... but we will see that we can write π_k , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ in terms of $\gamma(\mathbf{z}_{nk})$.
- We can find the **local minima** by iterative algorithm: alternate update of the **expected** posterior $\gamma(\mathbf{z}_{nk})$ and the **maximized** parameters π_k , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$.

Expectation-Maximization algorithm

- **E-step:** compute the expected posterior $\gamma(\mathbf{z}_{nk})$ given the current parameters $\boldsymbol{\pi}$, $\{\boldsymbol{\mu}_k\}$ and $\{\boldsymbol{\Sigma}_k\}$:

$$\gamma(\mathbf{z}_{nk}) = \frac{\mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k}{\sum_{j=1}^K \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \pi_j}$$

- **M-step:** update the parameters $\boldsymbol{\pi}$, $\{\boldsymbol{\mu}_k\}$ and $\{\boldsymbol{\Sigma}_k\}$ to maximize the expected log-likelihood, given the current expected posterior $\gamma(\mathbf{z}_{nk})$. This can be solved in closed form:

$$\begin{aligned} N_k &= \sum_{n=1}^N \gamma(\mathbf{z}_{nk}) & \pi_k &= \frac{N_k}{N} \\ \boldsymbol{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(\mathbf{z}_{nk}) \mathbf{x}_n \\ \boldsymbol{\Sigma}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(\mathbf{z}_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \end{aligned}$$

Derivation of the M-step updates

Results of results for Gaussian distributions:

Assume $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{d}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$.

Then,

$$\frac{\partial p(\mathbf{x})}{\partial \boldsymbol{\mu}} = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})(\mathbf{x} - \boldsymbol{\mu}) \boldsymbol{\Sigma}^{-1}$$

and

$$\frac{\partial p(\mathbf{x})}{\partial \Sigma} = -\frac{1}{2} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) \left(\Sigma^{-1} - \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} \right)$$

Derivation of the M-step updates ($\boldsymbol{\mu}_k$)

$$\arg \max_{\boldsymbol{\mu}_k} p(\mathbf{X} | \boldsymbol{\pi}, \{\boldsymbol{\mu}_j\}, \{\Sigma_j\}) = \arg \max_{\boldsymbol{\mu}_k} \sum_{n=1}^N \log p(\mathbf{x}_n | \boldsymbol{\pi}, \{\boldsymbol{\mu}_j\}, \{\Sigma_j\})$$

Compute the partial derivative w.r.t. $\boldsymbol{\mu}_k$:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}_k} \log p(\mathbf{X} | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}) &= \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\})} \frac{\partial p(\mathbf{x}_n | \boldsymbol{\pi}, \{\boldsymbol{\mu}_j\}, \{\Sigma_j\})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\})} \frac{\partial \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \Sigma_k)}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_j, \Sigma_j)} (\mathbf{x}_n - \boldsymbol{\mu}_k) \Sigma_k^{-1} \\ &= \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k) \Sigma_k^{-1} = 0 \quad \Rightarrow \quad \boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \end{aligned}$$

Derivation of the M-step updates ($\Sigma_k, \boldsymbol{\pi}$)

Same for the covariance:

1. Write the derivative of the log-likelihood w.r.t. Σ_k (use the properties of the Gaussian distribution).
2. Manipulate it to make $\gamma(z_{nk})$ appear.
3. Set it to zero and solve for Σ_k .

Same for the weights $\boldsymbol{\pi}$, but we have to use the constraint $\sum_{k=1}^K \pi_k = 1$.

1. Use Lagrange multipliers to introduce a new variable λ and write the Lagrangian function:

$$\ell(\boldsymbol{\pi}, \lambda) = \log p(\mathbf{X} | \boldsymbol{\pi}, \{\boldsymbol{\mu}_k\}, \{\Sigma_k\}) + \lambda \left(1 - \sum_{k=1}^K \pi_k \right)$$

2. Compute the partial derivative w.r.t. $\boldsymbol{\pi}$ and λ
3. Set them to zero and solve for $\boldsymbol{\pi}$ and λ .

EM for Gaussian Mixture Models

Algorithm:

1. Initialize the parameters π , $\{\mu_k\}$ and $\{\Sigma_k\}$
2. Repeat until convergence (likelihood does not change significantly):
 1. E-step: compute the expected posterior $\gamma(z_{nk})$ given the current parameters π , $\{\mu_k\}$ and $\{\Sigma_k\}$:
 2. M-step: update the parameters π , $\{\mu_k\}$ and $\{\Sigma_k\}$ to maximize the expected log-likelihood, given the current expected posterior $\gamma(z_{nk})$:
 3. Compute the log-likelihood of the data given the current parameters π , $\{\mu_k\}$ and $\{\Sigma_k\}$

Important:

- EM is guaranteed to converge to a local maximum, not necessarily to the global maximum.
- EM is sensitive to the initialization of the parameters (try different initializations and choose the one that gives the highest log-likelihood)

Gaussian Mixture Models

After the EM algorithm converges, we can use the learned parameters to:

1. Generate new data points by sampling from the mixture model.
2. Compute the likelihood of new data points.
3. Compute the posterior distribution of the latent variables given new data points (e.g., clustering).

Gaussian Mixture Models

Pros:

1. Simple and easy to implement.
2. Can model *any* distribution as a mixture of Gaussians.
3. One model can be used for different tasks (e.g., clustering, density estimation, generation).

Cons:

1. Sensitive to the initialization of the parameters, can get stuck in local minima.
2. The number of clusters K must be specified in advance.
3. The number of clusters K must be small, otherwise the model will overfit the data.

Question: Can we be Bayesian about the number of clusters K ? **Yes**, it's called **Dirichlet Process Mixture Model (DPMM)**

Dimensionality reduction as a generative model

Linear latent variable models

Linear latent variable models are a class of generative models that assume that the data is generated from a linear combination of latent variables.

We will:

1. Review the Principal Component Analysis (PCA) algorithm.
2. Introduce the (probabilistic) PCA as a linear latent variable model.

Principal Component Analysis

1. Assume we have a dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, where $\mathbf{x}_n \in \mathbb{R}^D$.
2. We want to find a low-dimensional representation of the data, i.e., we want to find a mapping $\mathbf{x}_n \mapsto \mathbf{z}_n$ such that $\mathbf{z}_n \in \mathbb{R}^K$, where $K < D$.
3. Choose the mapping such that the variance of the projected data is maximized, i.e., we want to find the mapping that maximizes the variance of the data in the low-dimensional space.

How to find the direction of maximum variance?

- Let's start by projecting the data on 1D subspace $\mathbf{z} = \mathbf{w}^\top \mathbf{x}$, where $\mathbf{w} \in \mathbb{R}^D$.
- We only want the direction, so we can assume that $\|\mathbf{w}\| = 1$.
- Given data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, we maximize the variance of the projected data $\mathbf{z} = \mathbf{X}\mathbf{w}$ constrained to $\|\mathbf{w}\| = 1$:

$$\arg \max_{\mathbf{w}} \|\mathbf{z}\|^2 \quad \text{s.t.} \quad \|\mathbf{w}\|^2 = 1$$

Assumes that \mathbf{X} is centered (i.e., \mathbf{X} has zero mean)

Analytic solution

- Let's use the Lagrange multipliers to solve the constrained optimization problem:

$$\ell(\mathbf{w}, \lambda) = \|\mathbf{z}\|^2 - \lambda(\|\mathbf{w}\|^2 - 1) = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \lambda(\mathbf{w}^\top \mathbf{w} - 1)$$

- Compute the partial derivative w.r.t. \mathbf{w} and λ :

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{w}} &= 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\lambda \mathbf{w} = 0 \quad \Rightarrow \quad \mathbf{X}^\top \mathbf{X} \mathbf{w} = \lambda \mathbf{w} \\ \frac{\partial \ell}{\partial \lambda} &= -(\mathbf{w}^\top \mathbf{w} - 1) = 0 \end{aligned}$$

This is an eigenvalue problem, where λ is the eigenvalue and \mathbf{w} is the eigenvector of the covariance matrix $\mathbf{X}^\top \mathbf{X}$.

Analytic solution

But which eigenvalue?

- The eigenvalue λ is the variance of the projected data \mathbf{z} : $\lambda = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}$
- If we want to maximize the variance, we need to find the largest eigenvalue of the covariance matrix $\mathbf{X}^\top \mathbf{X}$.
- The corresponding eigenvector is the direction of maximum variance: $\hat{\mathbf{w}}$

Probabilistic PCA

- PCA can be seen as a linear continuous latent variable model, where the latent variable \mathbf{z} is a linear combination of the observed data \mathbf{x} .
- Question:** Can we approach PCA from a probabilistic point of view? **Answer:** Yes!

Recall the continuous latent variable model:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z}) d\mathbf{z}$$

PPCA Modeling Assumptions

- The **generative model** works as follows:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$$

where:

1. $\mathbf{W} \in \mathbb{R}^{D \times K}$ is the weight matrix
2. $\mathbf{z} \in \mathbb{R}^K$ is the latent variable
3. $\boldsymbol{\mu} \in \mathbb{R}^D$ is the mean of the data
4. $\boldsymbol{\varepsilon} \in \mathbb{R}^D$ is the noise

How can we treat these variables?

PPCA Modeling Assumptions

- The **latent variable** \mathbf{z} is assumed to be a standard Gaussian distribution:

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

- The **noise** $\boldsymbol{\varepsilon}$ is assumed to be a Gaussian distribution with zero mean and covariance matrix $\sigma^2 \mathbf{I}$:

$$p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon}; \mathbf{0}, \sigma^2 \mathbf{I})$$

...

- The **observed data** \mathbf{x} conditioned on the latent variable \mathbf{z} is a Gaussian distribution

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

- The **marginal distribution** of the data \mathbf{x} is obtained by integrating out the latent variable \mathbf{z} :

$$p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z}) \, d\mathbf{z} = \int \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) \, d\mathbf{z}$$

It's a Gaussian distribution $\mathcal{N}(?, ?)$.

PPCA Modeling Assumptions

1. For the mean, we have:

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}] \\ &= \mathbf{W}\mathbb{E}[\mathbf{z}] + \boldsymbol{\mu} + \mathbb{E}[\boldsymbol{\varepsilon}] \\ &= \mathbf{W} \cdot \mathbf{0} + \boldsymbol{\mu} + 0 = \boldsymbol{\mu}\end{aligned}$$

...

2. For the covariance, we have:

$$\begin{aligned}\text{Cov}(\mathbf{x}) &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \\ &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})^\top] \\ &= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^\top\mathbf{W}^\top + \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top + \mathbf{W}\mathbf{z}\boldsymbol{\varepsilon}^\top + \boldsymbol{\varepsilon}\mathbf{z}^\top\mathbf{W}] \\ &= \mathbf{W}\mathbb{E}[\mathbf{z}\mathbf{z}^\top]\mathbf{W}^\top + \mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top] + \mathbf{W}\mathbb{E}[\mathbf{z}]\mathbb{E}[\boldsymbol{\varepsilon}^\top] + \mathbb{E}[\boldsymbol{\varepsilon}]\mathbb{E}[\mathbf{z}^\top]\mathbf{W} \\ &= \mathbf{W} \cdot \mathbf{I} \cdot \mathbf{W}^\top + \sigma^2\mathbf{I} + 0 + 0 \\ &= \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}\end{aligned}$$

PPCA Solution

Now, we need to find the parameters \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 that maximize the log-likelihood of the data:

$$\begin{aligned}\log p(\mathbf{X}; \mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \sum_{n=1}^N \log p(\mathbf{x}_n; \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}) \\ &= -\frac{NK}{2} \log(2\pi) - \frac{N}{2} \log(\det(\mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})) - \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$

Again, compute the partial derivative w.r.t. \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 and set them to zero.

PPCA Solution

Let's define the sample covariance matrix $\mathbf{S} = \mathbf{X}\mathbf{X}^\top$. Solutions for \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 can be found in closed form, no need to use the EM algorithm.

...

1. For the mean, we have: $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$

...

2. For the noise variance, we have: $\sigma^2 = \frac{1}{D-K} \sum_{j=K+1}^D \lambda_j$, where λ_j are the eigenvalues of \mathbf{S}

...

3. For the weight matrix, we have: $\mathbf{W} = \mathbf{U}_K(\boldsymbol{\Lambda}_K - \sigma^2\mathbf{I})^{1/2}$, where \mathbf{U}_K are the first K eigenvectors \mathbf{S} and $\boldsymbol{\Lambda}_K$ is the diagonal matrix of the first K eigenvalues.

...

Because everything is Gaussian, we can also compute the posterior distribution of the latent variable \mathbf{z} given the data \mathbf{x} :

$$p(\mathbf{z}_n | \mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n; \mathbf{C}^{-1}\mathbf{W}^\top(\mathbf{x}_n - \boldsymbol{\mu}), \sigma^2\mathbf{C}^{-1})$$

with $\mathbf{C} = \mathbf{W}^\top\mathbf{W} + \sigma^2\mathbf{I}$

Is the PPCA solution the same as PCA?

- The PPCA solution is more general than PCA, because it assumes that the data is generated from a Gaussian distribution with noise.
- However, when the noise variance σ^2 goes to zero, the PPCA solution converges to the PCA solution.

Note: The principal components found by PCA and PPCA span the same subspace, but the actual components (i.e., the directions) may differ by an (un-identifiable) rotation matrix \mathbf{R}