

Lecture 5: Linear Gaussian Latent Variable Models

STATS305C: Applied Statistics III

Scott Linderman

April 11, 2022

Last Time...

- ▶ Directed Graphical Models
- ▶ Hierarchical Gaussian Model
- ▶ MCMC: MH and Gibbs

Today...

Outline:

- ▶ Principal Components Analysis (PCA)
- ▶ PCA as a linear autoencoder
- ▶ PCA as a linear Gaussian latent variable model
- ▶ Factor analysis
- ▶ Other continuous latent variable models

Reading:

- ▶ Required: Bishop, Ch 12

Motivating Example: More scores

Continuing on our academic theme, suppose we have not only SAT scores but entire transcripts, $\{\mathbf{x}_n\}_{n=1}^N$ where $\mathbf{x}_n \in \mathbb{R}^D$ is a vector of D grades from one student. For example, we might have grades for all their classes throughout the four years of high school. (Assume all students took the same D classes.)

We might have a few objectives in mind:

- ▶ **Dimensionality reduction:** are there a few dimensions along which the students primarily vary? E.g. do students vary along a “mathy” to “artsy” axis? Are there “late bloomers” and “senioritis sufferers”?
- ▶ **Visualization:** Like above, but how can we embed these points in 2 or 3 dimensions to best visualize them?
- ▶ **Compression:** Like above, but how can I best summarize the D scores if I am willing to sacrifice some reconstruction accuracy?

Principal Components Analysis (PCA)

Two classical definitions:

maximum variance formulation

1. An orthogonal projection of the data onto a lower dimensional linear space, known as the *principal subspace*, such that the variance of the projected data is maximized (Hotelling, 1933).
2. The linear projection that minimizes the average projection cost, defined as the mean squared distance between the data points and their projections (Pearson, 1901).

minimum-error formulation

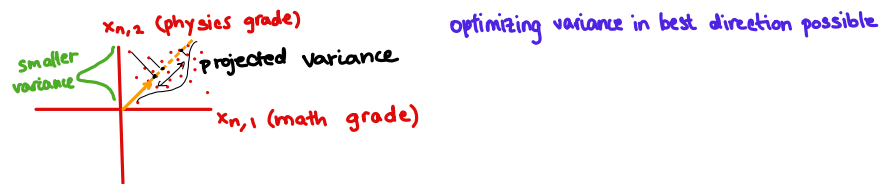
(Quoted from Bishop, Ch 12)

3. Probabilistic PCA
(Bishop, 1999)

PCA: Maximum Variance Formulation

Goal: Project data $\{\mathbf{x}_n\}_{n=1}^N$ onto a lower dimensional space of dimension $M < D$ while maximizing the variance of the projected data.

Illustration:



→ start from 1D subspace & go to M -subspace

→ use constraint that new direction is **orthogonal** to previous directions

PCA: Maximum Variance Formulation II

To start, assume $M = 1$. The principal subspace is defined by a unit vector $\mathbf{u}_1 \in \mathbb{R}^D$. This is called the first **principal component**.

Projecting a data point \mathbf{x}_n onto this subspace amounts to taking an inner product, $\mathbf{u}_1^\top \mathbf{x}_n$. These are variously called the **scores**, **embeddings**, or **signals**.

PCA: Maximum Variance Formulation III

The mean of the projected data is,

$$\frac{1}{N} \sum_{n=1}^N \overset{\text{1st principal component.}}{u_1^T} x_n = u_1^T \left(\frac{1}{N} \sum_{n=1}^N x_n \right) = u_1^T \bar{x}, \quad (1)$$

where \bar{x} is the sample mean.

The variance is

$$\frac{1}{N} \sum_{n=1}^N [u_1^T x_n - u_1^T \bar{x}]^2 = \frac{1}{N} \sum_{n=1}^N [u_1^T (x_n - \bar{x})]^2 \quad (2)$$

$$= \frac{1}{N} \sum_{n=1}^N u_1^T (x_n - \bar{x})(x_n - \bar{x})^T u_1 \quad (3)$$

$$= \underbrace{u_1^T S u_1}_{\text{want to find } u_1 \text{ to maximize this quantity}} \quad (4)$$

where $S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T \in \mathbb{R}^{D \times D}$ is the sample covariance matrix.

PCA: Maximum Variance Formulation IV

Now maximize the projected variance wrt \mathbf{u}_1 , subject to the constraint that $\mathbf{u}_1^\top \mathbf{u}_1 = 1$

$$\mathcal{L}(\mathbf{u}_1) = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1). \quad (5)$$

Taking the gradient wrt \mathbf{u}_1 and setting to zero,

$$\nabla \mathcal{L}(\mathbf{u}_1) = 2\mathbf{S} \mathbf{u}_1 - 2\lambda_1 \mathbf{u}_1 = 0 \Rightarrow \mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1. \quad (6)$$

So \mathbf{u}_1 must be an eigenvector of \mathbf{S} . Left multiplying by \mathbf{u}^\top ,

$$\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 = \lambda_1, \quad (7)$$

so the projected variance $\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$ is maximized when we choose \mathbf{u}_1 to be the eigenvector with the largest eigenvalue λ_1 .

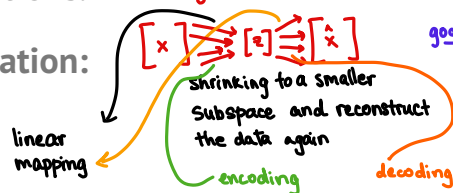
eigenvalues ≥ 0
eigenvectors can be orthonormal

More generally, to find an M dimensional principal subspace, take the M eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_M$ with the largest eigenvalues $\lambda_1, \dots, \lambda_M$.

PCA: Linear Autoencoder Formulation

Now consider Pearson's formulation of PCA as the linear projection that minimizes the average projection cost, defined as the mean squared distance between the data points and their projections. **minimizing reconstruction error.**

Illustration:



goal: minimize: $\mathcal{L}(w) = \|x_n - \hat{x}_n\|_2^2$

→ MSE gives PCA
→ other \mathcal{L} won't give PCA exactly

PCA: Linear Autoencoder Formulation II

To formalize this, let

$$\mathbf{W} = \begin{bmatrix} | & & | \\ \mathbf{w}_1 & \cdots & \mathbf{w}_M \\ | & & | \end{bmatrix} \in \mathbb{R}^{D \times M} \quad \text{with} \quad \mathbf{W}^\top \mathbf{W} = \mathbf{I} \quad (8)$$

be an **orthogonal basis** for the principal subspace.

We will **encode** each data point by subtracting the mean and projecting onto the principal subspace to obtain $\mathbf{z}_n = \mathbf{W}^\top (\mathbf{x}_n - \bar{\mathbf{x}})$.

Since \mathbf{W} is an orthogonal matrix, all we need to do to **decode** the encoded data point is multiply $\mathbf{W}\mathbf{z}_n$ and add back the mean. That gives us,

$$\hat{\mathbf{x}}_n = \mathbf{W}\mathbf{z}_n + \bar{\mathbf{x}} \Rightarrow \hat{\mathbf{x}}_n = \mathbf{W}\mathbf{W}^\top (\mathbf{x}_n - \bar{\mathbf{x}}) + \bar{\mathbf{x}}. \quad (9)$$

PCA: Linear Autoencoder Formulation III

Goal: Find an orthogonal matrix \mathbf{W} that minimizes the ^{mean} squared reconstruction error,

$$\mathcal{L}(\mathbf{W}) = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|_2^2 \quad (10)$$

We can write this in matrix notation instead. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the **centered data matrix** with rows $\mathbf{x}_n - \bar{\mathbf{x}}$. Then,

$$\hat{\mathbf{x}} = \begin{bmatrix} -(\hat{\mathbf{x}}_1 - \bar{\mathbf{x}}) \\ \vdots \\ -(\hat{\mathbf{x}}_N - \bar{\mathbf{x}}) \end{bmatrix}$$

$$\mathcal{L}(\mathbf{W}) = \frac{1}{N} \text{Tr}[(\mathbf{X} - \hat{\mathbf{X}})^\top (\mathbf{X} - \hat{\mathbf{X}})] \quad (11)$$

$$= \frac{1}{N} \text{Tr}[(\mathbf{X} - \overset{\text{X was already centered!}}{\mathbf{X}\mathbf{W}\mathbf{W}^\top})^\top (\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^\top)] \quad (12)$$

$$= \frac{1}{N} \text{Tr}[(\mathbf{X}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top))^\top (\mathbf{X}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top))] \quad (13)$$

$$= \text{Tr}[(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top \mathbf{S}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)] \quad (14)$$

where $\mathbf{S} = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$ is the sample covariance matrix.

PCA: Linear Autoencoder Formulation IV

Now apply the circular trace property,

$$\mathcal{L}(W) = \text{Tr}[S(I - WW^\top)(I - WW^\top)^\top] \quad (15)$$

Handwritten notes:
 $\text{tr}(ABC) = \text{tr}(BAC)$
 $= \text{tr}(BCA) \dots$

Question: What does $(I - WW^\top)(I - WW^\top)^\top$ equal?

Handwritten answer:
 $= \underbrace{I - WW^\top}$
a projection matrix
onto the orthogonal
complement of the subspace
spanned by W 's columns.

PCA: Linear Autoencoder Formulation V

Thus, the objective simplifies to,

$$\mathcal{L}(\mathbf{W}) = \text{Tr}[\mathbf{S}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)] = \text{Tr}[\mathbf{S}] - \text{Tr}[\mathbf{S}\mathbf{W}\mathbf{W}^\top] = \text{const} - \text{Tr}[\mathbf{W}^\top \mathbf{S} \mathbf{W}]. \quad (16)$$

one of the solvable nonlinear optimization problems!

Let $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ be the eigendecomposition of \mathbf{S} . (Since it is a covariance matrix, the eigenvectors are orthogonal.) Plugging in,
 \hookrightarrow exists as $\mathbf{S} \succeq \mathbf{0}$.

$$\mathcal{L}(\mathbf{W}) = \text{const} - \text{Tr}[\mathbf{W}^\top \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{W}] \quad (17)$$

$$= \text{const} - \text{Tr} \left[\mathbf{W}^\top \left(\sum_{d=1}^D \lambda_d \mathbf{u}_d \mathbf{u}_d^\top \right) \mathbf{W} \right] \quad (18)$$

rank 1 outer products

$$= \text{const} - \sum_{m=1}^M \sum_{d=1}^D \lambda_d \underbrace{\mathbf{w}_m^\top \mathbf{u}_d}_{\mathbb{R}} \underbrace{\mathbf{u}_d^\top \mathbf{w}_m}_{\mathbb{R}} \quad (19)$$

$$= \text{const} - \sum_{m=1}^M \sum_{d=1}^D \lambda_d (\mathbf{w}_m^\top \mathbf{u}_d)^2 \quad (20)$$

$\mathbf{w}_m = \mathbf{u}_d$

$\mathbf{w}_m^\top \mathbf{u}_d$ is the best when $\mathbf{w}_m = \mathbf{u}_d$.
(Cauchy-Schwarz)

\hookrightarrow how to find \mathbf{w}_m ?

Question: We want to minimize $\mathcal{L}(\mathbf{W})$ subject to \mathbf{W} being orthogonal. What is the solution?

PCA and the Singular Value Decomposition

We've seen two formulations of PCA, both showing us that the first M **principal components** are the leading M **eigenvectors of the covariance matrix**.

Let

$$\mathbf{Y} = \frac{1}{\sqrt{N}} \mathbf{X} = \frac{1}{\sqrt{N}} \begin{bmatrix} - & (\mathbf{x}_1 - \bar{\mathbf{x}})^\top & - \\ & \vdots & \\ - & (\mathbf{x}_N^\top - \bar{\mathbf{x}})^\top & - \end{bmatrix} \quad (21)$$

be the **centered and scaled** data matrix. Then $\mathbf{Y}^\top \mathbf{Y} = \frac{1}{N} \mathbf{X}^\top \mathbf{X} = \mathbf{S}$ is the covariance matrix.

The **singular value decomposition (SVD)** of \mathbf{Y} is,

$$\mathbf{Y} = \mathbf{V} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^\top. \quad (22)$$



I.e. the **right singular vectors** of \mathbf{Y} are the same (up to sign flips) as the **eigenvectors of \mathbf{S}** , and **singular values** of \mathbf{Y} are the **square root of the eigenvalues of \mathbf{S}** .

PCA Explained Variance

How well do the M principal components explain the data?

Let $\mathbf{z}_n = \mathbf{U}_m^\top (\mathbf{x}_n - \bar{\mathbf{x}}) \in \mathbb{R}^M$. Its covariance is,

1st m
columns of \mathbf{U}

$$\text{Cov}[\mathbf{z}] = \text{Cov}[\mathbf{U}_m^\top (\mathbf{x} - \bar{\mathbf{x}})] = \mathbf{U}_m^\top \text{Cov}[\mathbf{x}] \mathbf{U}_m = \text{diag}([\lambda_1, \dots, \lambda_M]).$$

$\mathbf{U} \mathbf{U}^\top = \mathbf{I}$ where \mathbf{U} is orthogonal

(23)

Of course, if we let $M = D$, then we have $\text{Cov}(\mathbf{z}) = \text{diag}([\lambda_1, \dots, \lambda_D])$.

One way of assessing how well M components fits the data is via the **fraction of variance explained**,

$$\text{variance explained} = \frac{\text{Tr}(\text{Cov}[\mathbf{z}; M \text{ components}])}{\text{Tr}(\text{Cov}[\mathbf{z}; D \text{ components}])} = \frac{\sum_{m=1}^M \lambda_m}{\sum_{m=1}^D \lambda_m} \in [0, 1].$$

(24)

Scree Plots

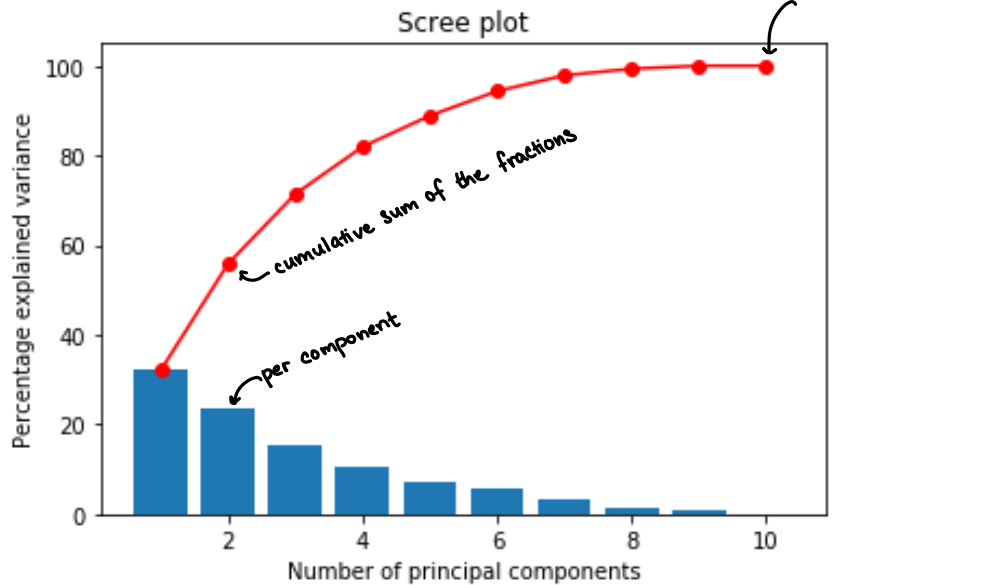


Figure: “Scree” plot showing percent variance explained per component and cumulatively.

Probabilistic PCA: A Continuous Latent Variable Model

The previous two formulations cast the principal components as the solutions to optimization problems: maximize the projected variance or minimize the reconstruction error.

A more modern view of PCA is as the **maximum likelihood estimate** of a **latent variable model**.

Probabilistic PCA (PPCA) has many advantages:

- ▶ It's a multivariate normal model with **low-rank plus diagonal covariance**, which takes only $O(MD)$ parameters.
- ▶ We can fit the model using a **host of inference algorithms**.
- ▶ It can handle **missing data**.
- ▶ We can obtain posterior distributions of the principal components and scores.
- ▶ It can be embedded in larger probabilistic models.

Probabilistic PCA: A Continuous Latent Variable Model

The PPCA model is quite simple,

$$\mathbf{z}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (25)$$

$$\mathbf{x}_n | \mathbf{z}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}), \quad (26)$$

where $\mathbf{z}_n \in \mathbb{R}^M$ is a latent variable, $\mathbf{W} \in \mathbb{R}^{D \times M}$ are the weights, $\boldsymbol{\mu} \in \mathbb{R}^D$ is the bias parameter, and $\sigma^2 \in \mathbb{R}_+$ is a variance.

Equivalently, we can think of \mathbf{x}_n as a linear function of \mathbf{z}_n with additive noise,

$$\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu} + \boldsymbol{\epsilon}_n, \quad (27)$$

where $\boldsymbol{\epsilon}_n \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \in \mathbb{R}^D$.

$$\begin{aligned} \mathbb{E}(\mathbf{x}_n) &= \mathbf{W} \cdot \mathbf{0} + \boldsymbol{\mu} + \mathbf{0} = \boldsymbol{\mu} \\ \text{cov}(\mathbf{x}_n) &= \mathbf{W} \cdot \text{cov}(\mathbf{z}) \cdot \mathbf{W}^T + \text{cov}(\boldsymbol{\epsilon}_n); \quad \mathbf{z}_n \perp \boldsymbol{\epsilon}_n \\ &= \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} \end{aligned}$$

Maximum likelihood estimation of the parameters

Suppose we only need a **point estimate** of the parameters \mathbf{W} , $\boldsymbol{\mu}$, and σ^2 .

A natural approach is the **maximum likelihood estimate** (MLE),

$$\mathbf{W}_{\text{ML}}, \boldsymbol{\mu}_{\text{ML}}, \sigma_{\text{ML}}^2 = \arg \max \mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2), \quad (28)$$

where \mathcal{L} is the marginal likelihood,

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \log p(\mathbf{X} | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \quad (29)$$

$$= \log \int \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}_n) d\mathbf{z}_n \quad (30)$$

likelihood *prior*

$$= \log \int \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}) d\mathbf{z}_n \quad (31)$$

Exercise: Simplify this expression.

$$= \log \prod_{n=1}^N \int \underbrace{\mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \cdot \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I})}_{= \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})} d\mathbf{z}_n$$

Maximum likelihood estimation of the parameters II

The log likelihood simplifies to,

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \frac{ND}{2} \log 2\pi - \frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (32)$$

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

Setting the derivative wrt $\boldsymbol{\mu}$ to zero and solving yields $\boldsymbol{\mu}_{\text{ML}} = \bar{\mathbf{x}}$, the sample mean.

Maximizing wrt \mathbf{W} and σ^2 is more complex but still has a closed form solution, (Bishop + Tipping, 1999)

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_M (\boldsymbol{\Lambda}_M - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}, \quad (33)$$

where $\mathbf{U}_M \in \mathbb{R}^{D \times M}$ has columns given by the leading eigenvectors of the sample covariance matrix \mathbf{S} , where $\boldsymbol{\Lambda}_M = \text{diag}([\lambda_1, \dots, \lambda_M])$, and where $\mathbf{R} \in \mathbb{R}^{M \times M}$ is an arbitrary *orthogonal* matrix.

Put differently, the MLE weights are only identifiable up to orthogonal transformation. Or, only the subspace spanned by \mathbf{U}_M is identifiable.

Maximum likelihood estimation of the parameters III

Finally, the MLE of the variance is,

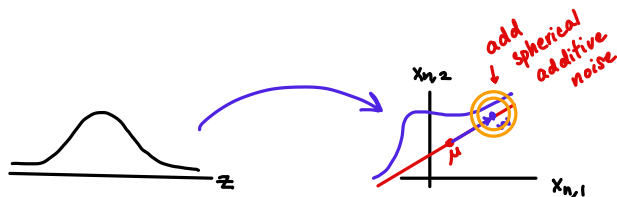
$$\sigma_{\text{ML}}^2 = \frac{1}{D-M} \sum_{m=M+1}^D \lambda_m, \quad (34)$$

the average variance in the remaining dimensions. $WW^T + \sigma^2 \mathbf{I} = \mathbf{U}_m \mathbf{\Lambda}_m \mathbf{U}_m^T$

Question: What is the marginal covariance \mathbf{C} using the MLE \mathbf{W}_{ML} and σ_{ML}^2 ?

Question: Intuitively, why is the marginal covariance invariant to rotations of the weights?

→ $\mathbf{R}\mathbf{R}^T$ cancels out.



The Posterior Distribution on the Latent Variables

Now fix \mathbf{W} , $\boldsymbol{\mu}$, and σ^2 (e.g. to their maximum likelihood values). What is the posterior of \mathbf{z}_n ?

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \propto \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \quad (35)$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{z}_n - \frac{1}{2} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \boldsymbol{\mu})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \boldsymbol{\mu}) \right\} \quad (36)$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{J}_n \mathbf{z}_n + \mathbf{h}_n^\top \mathbf{z}_n \right\} \quad (37)$$

$$(38)$$

where $\mathbf{J}_n = \mathbf{I} + \frac{1}{\sigma^2} \mathbf{W}^\top \mathbf{W}$ and $\mathbf{h}_n = \frac{1}{\sigma^2} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu})$

Completing the square,

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{z}_n | \mathbf{J}_n^{-1} \mathbf{h}_n, \mathbf{J}_n^{-1}). \quad (39)$$

The Posterior Distribution in the Zero Noise Limit

In the limit where $\sigma^2 \rightarrow 0$, the posterior mean of \mathbf{z}_n is,

$$\lim_{\sigma^2 \rightarrow 0} \mathbb{E}[\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2] = \lim_{\sigma^2 \rightarrow 0} \underbrace{\left(I + \frac{1}{\sigma^2} \mathbf{W}^\top \mathbf{W} \right)^{-1}}_{\mathbf{J}_n^{-1}} \underbrace{\left[\frac{1}{\sigma^2} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \right]}_{\mathbf{h}_n} \quad (40)$$

$$= \lim_{\sigma^2 \rightarrow 0} (\sigma^2 I + \mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (41)$$

$$= (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (42)$$

Now suppose $\mathbf{W} = \mathbf{W}_{\text{ML}} = \mathbf{U}_M (\boldsymbol{\Lambda}_M - \sigma^2 I)^{\frac{1}{2}} \mathbf{R}$ and set $\mathbf{R} = \mathbf{I}$. This goes to $\mathbf{W} = \mathbf{U}_M \boldsymbol{\Lambda}_M^{\frac{1}{2}}$. Then,

$$\lim_{\sigma^2 \rightarrow 0} \mathbb{E}[\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2] = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (43)$$

$$= \boldsymbol{\Lambda}_M^{-\frac{1}{2}} \mathbf{U}_M^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (44)$$

This is the same as \mathbf{z}_n from the linear autoencoder formulation, except here the \mathbf{z}_n 's are scaled by $\boldsymbol{\Lambda}^{-\frac{1}{2}}$ to be unit variance in all dimensions.

Gibbs Sampling for Probabilistic PCA I

For simplicity, assume the data is centered and fix $\boldsymbol{\mu} = \mathbf{0}$. Now let's put a prior on the parameters,

$$p(\mathbf{W}, \sigma^2) = \chi^{-2}(\sigma^2 \mid \nu_0, \sigma_0^2) \prod_{d=1}^D \mathcal{N}(\mathbf{w}_d \mid \mathbf{0}, \frac{\sigma^2}{\kappa_0} \mathbf{I}), \quad (45)$$

where $\mathbf{w}_d \in \mathbb{R}^M$ are the *rows* of \mathbf{W} .

Under this prior, the complete conditional distribution of the parameters is,

$$p(\mathbf{w}_d \mid \{\mathbf{x}_n, \mathbf{z}_n\}_{n=1}^N, \sigma^2, \kappa_0) \propto \mathcal{N}(\mathbf{w}_d \mid \mathbf{0}, \frac{\sigma^2}{\kappa_0} \mathbf{I}) \prod_{n=1}^N \mathcal{N}(x_{n,d} \mid \mathbf{w}_d^\top \mathbf{z}_n, \sigma^2) \quad (46)$$

$$= \mathcal{N}(\mathbf{w}_d \mid \mathbf{J}_d^{-1} \mathbf{h}_d, \mathbf{J}_d^{-1}) \quad (47)$$

where $\mathbf{J}_d = \frac{\kappa_0}{\sigma^2} \mathbf{I} + \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{z}_n \mathbf{z}_n^\top$ and $\mathbf{h}_d = \frac{1}{\sigma^2} \sum_{n=1}^N x_{n,d} \mathbf{z}_n$.

Gibbs Sampling for Probabilistic PCA II

For the variance,

$$p(\sigma^2 \mid \{\mathbf{x}_n, \mathbf{z}_n\}_{n=1}^N, \mathbf{W}, \kappa_0, \nu_0, \sigma_0^2) \propto \chi^{-2}(\sigma^2 \mid \nu_0, \sigma_0^2) \prod_{d=1}^D \mathcal{N}(\mathbf{w}_d \mid \mathbf{0}, \frac{\sigma^2}{\kappa_0} \mathbf{I}) \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n \mid \mathbf{W} \mathbf{z}_n, \sigma^2 \mathbf{I}) \quad (48)$$

$$\propto \chi^{-2}(\sigma^2 \mid \nu_N, \sigma_N^2) \quad (49)$$

where

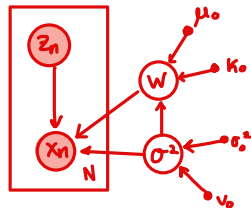
$$\nu_N = \nu_0 + DM + DN \quad (50)$$

$$\sigma_N^2 = \nu_N^{-1} \left[\nu_0 \sigma_0^2 + \kappa_0 \sum_{d=1}^D \overbrace{\mathbf{w}_d^\top \mathbf{w}_d}^{\text{outer products!}} + \sum_{n=1}^N (\mathbf{x}_n - \mathbf{W} \mathbf{z}_n)^\top (\mathbf{x}_n - \mathbf{W} \mathbf{z}_n) \right] \quad (51)$$

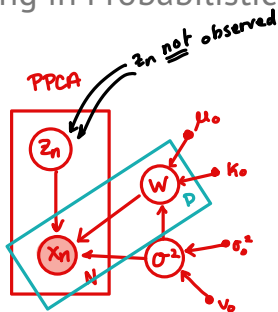
Note: It's a little strange to put a prior on the rows of \mathbf{W} ; it's more natural to put a prior on the columns. We only did this for simplicity. It turns out you can put a conjugate prior on (\mathbf{W}, σ^2) that specifies the conditional variance of the columns. It's called a **matrix normal inverse Wishart** prior.

Connection to Bayesian Linear Regression

Question: How does Gibbs sampling in Probabilistic PCA relate to Bayesian linear regression from HW1?



Bayesian
linear
regression



Factor Analysis

Factor analysis is another continuous latent variable model. In fact, it's almost the same as probabilistic PCA!

The difference is that FA allows σ^2 to vary across output dimensions. The generative model is,

$$\mathbf{z}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (52)$$

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \quad (53)$$

where $\boldsymbol{\sigma}^2 = [\sigma_1^2, \dots, \sigma_D^2]^\top$.

Put a similar prior on the parameters as before,

$$p(\mathbf{W}, \boldsymbol{\sigma}^2) = \prod_{d=1}^D \left[\chi^{-2}(\sigma_d^2 \mid \nu_0, \sigma_0^2) \mathcal{N}(\mathbf{w}_d \mid \mathbf{0}, \frac{\sigma_d^2}{\kappa_0} \mathbf{I}) \right]. \quad (54)$$

Exercise: without doing any math, derive a Gibbs sampler for this factor analysis model.

Independent Components Analysis (PCCA)

Independent Components Analysis (ICA) is yet another linear latent variable model. It aims to find *independent* factors of variation in the data. One probabilistic formulation is,

$$\mathbf{z}_n \stackrel{\text{iid}}{\sim} p(\mathbf{z}) \quad \leftarrow \text{non Gaussian prior} \quad (55)$$

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \quad (56)$$

where the prior distribution on \mathbf{z} assumes the coordinates are independent,

$$p(\mathbf{z}) = \prod_{m=1}^M p(z_m). \quad (57)$$

The success of this approach requires that $p(\mathbf{z})$ be non-Gaussian; otherwise, we could always reduce it to factor analysis by moving any correlation in $p(\mathbf{z})$ into the weights.

Instead, we typically choose priors that have **heavy tails**, like a **Laplace distribution**.

Probabilistic Canonical Correlations Analysis

Now consider a slightly different setting in which we have two types of observations, $\mathbf{x} \in \mathbb{R}^{D_x}$ and $\mathbf{y} \in \mathbb{R}^{D_y}$.

Goal: find *shared* latent variables $\mathbf{z}_n^{(s)}$ that capture common factors of variation across domains, as well as *private* latent variables $\mathbf{z}_n^{(x)}$ and $\mathbf{z}_n^{(y)}$ that capture domain-specific variation.

Model:

$$\mathbf{z}_n^{(s)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I_{M_s}) \quad \mathbf{z}_n^{(x)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I_{M_x}) \quad \mathbf{z}_n^{(y)} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I_{M_y}) \quad (58)$$

and

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}_{xx}\mathbf{z}_n^{(x)} + \mathbf{W}_{xs}\mathbf{z}_n^{(s)} + \boldsymbol{\mu}_x, \sigma^2 I_{D_x}) \quad (59)$$

$$\mathbf{y}_n \sim \mathcal{N}(\mathbf{W}_{yy}\mathbf{z}_n^{(y)} + \mathbf{W}_{ys}\mathbf{z}_n^{(s)} + \boldsymbol{\mu}_y, \sigma^2 I_{D_y}) \quad (60)$$

where M_s , M_x , and M_y are the dimensions of $\mathbf{z}_n^{(s)}$, $\mathbf{z}_n^{(x)}$, and $\mathbf{z}_n^{(y)}$, respectively, and the parameters consist of \mathbf{W}_{xx} , \mathbf{W}_{xs} , $\boldsymbol{\mu}_x$, \mathbf{W}_{yy} , \mathbf{W}_{ys} , $\boldsymbol{\mu}_y$, and σ^2 .

Probabilistic Canonical Correlations Analysis II

Exercise: Draw the graphical model for probabilistic CCA.

Probabilistic Canonical Correlations Analysis III

Just as the MLE for probabilistic PCA yields the classical PCA solution, Bach and Jordan [2005] showed that the MLE for probabilistic CCA yields the classical CCA solution,

$$\mathbf{w}_{xs,1}, \mathbf{w}_{ys,1} = \arg \max \text{corr}(\mathbf{w}_{xs,1}^\top \mathbf{x}, \mathbf{w}_{ys,1}^\top \mathbf{y}), \quad (61)$$

where $\mathbf{w}_{xs,1}$ and $\mathbf{w}_{ys,1}$ are the first columns of \mathbf{W}_{xs} and \mathbf{W}_{ys} respectively; aka the first **pair of canonical variables**.

Subsequent pairs of canonical variables are found by maximizing the same objective, subject to being orthogonal to previous pairs.

As with PCA, the classical CCA solution can be found with an SVD. Here, \mathbf{W}_{xs} and \mathbf{W}_{ys} are the left and right singular vectors of the sample correlation matrix $\text{diag}(\mathbf{S}_{xx})^{-\frac{1}{2}} \mathbf{S}_{xy} \text{diag}(\mathbf{S}_{yy})^{-\frac{1}{2}}$.

See Witten et al. [2009] for (non-Bayesian) sparse extensions to PCA and CCA.

References I

Francis R Bach and Michael I Jordan. A probabilistic interpretation of canonical correlation analysis. 2005.

Daniela M Witten, Robert Tibshirani, and Trevor Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3): 515–534, 2009.