

Linear Gaussian Latent Variable Models

STATS305B: Applied Statistics II

Scott Linderman

February 19, 2025

Last Time...

- ▶ Mixture Models & the EM Algorithm
- ▶ HMMs & the Forward-Backward Algorithm

Today...

Outline:

- ▶ Principal Components Analysis (PCA)
- ▶ PCA as a linear Gaussian latent variable model
- ▶ Factor analysis
- ▶ Linear Dynamical Systems & the Kalman Filter/Smoothen

Motivating Example

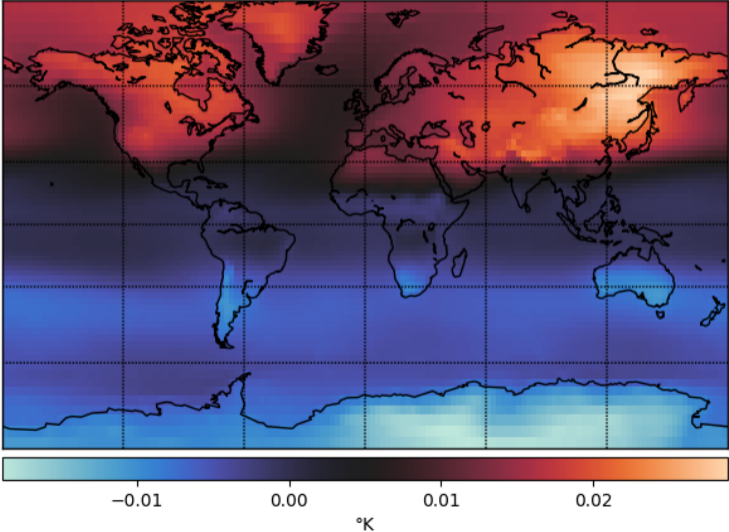
Take HW3 as an example: you have temperature measurements at 9504 locations across the globe. Are the measurements really 9504 dimensional?

We might have a few objectives in mind:

- ▶ **Dimensionality reduction:** are there a few dimensions along which the temperatures primarily vary? Maybe northern and southern hemisphere, or land and sea?
- ▶ **Visualization:** Sometimes, we want to embed high-dimensional points in 2 or 3 dimensions for visualization.
- ▶ **Compression:** How can I best summarize the data if I am willing to sacrifice some reconstruction accuracy?

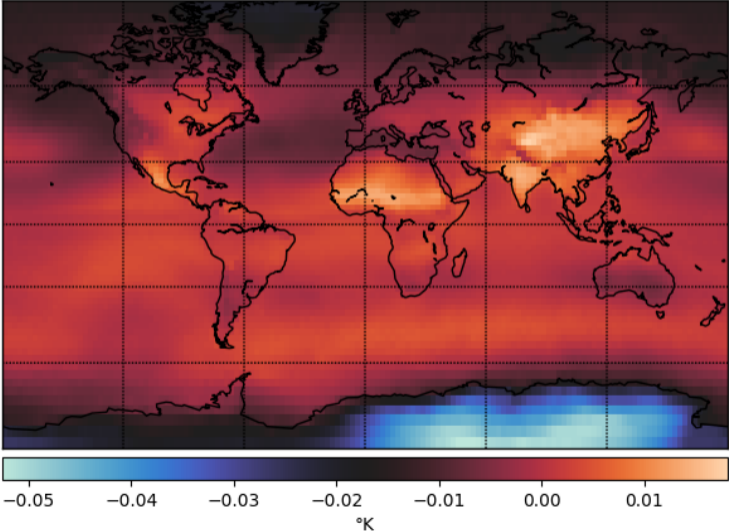
Principal Components of Global Temperature

PCA Component 1



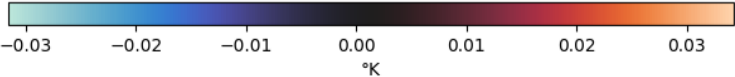
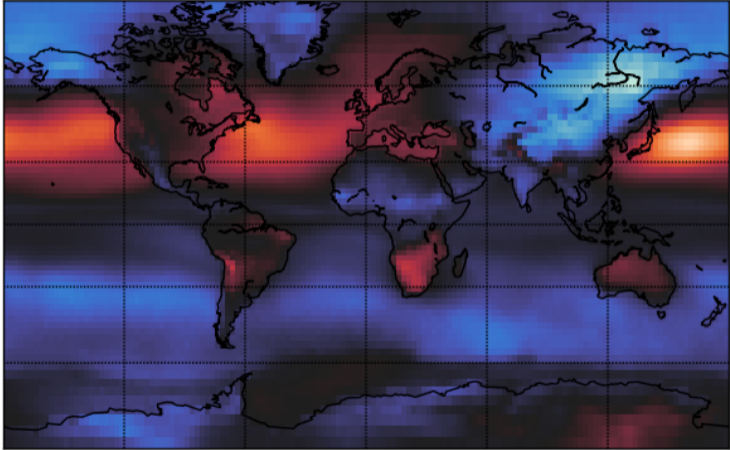
Principal Components of Global Temperature

PCA Component 2



Principal Components of Global Temperature

PCA Component 3



Principal Components Analysis (PCA)

Two classical definitions:

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).

(Quoted from Bishop, Ch 12)

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:

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 is variously called the **scores**, **embeddings**, or **signals**.

PCA: Maximum Variance Formulation III

The mean of the projected data is,

$$\frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^T \mathbf{x}_n = \mathbf{u}_1^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) = \mathbf{u}_1^T \bar{\mathbf{x}}, \quad (1)$$

where $\bar{\mathbf{x}}$ is the sample mean.

The variance is

$$\frac{1}{N} \sum_{n=1}^N [\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}]^2 = \frac{1}{N} \sum_{n=1}^N [\mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})]^2 \quad (2)$$

$$= \frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_1 \quad (3)$$

$$= \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \quad (4)$$

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

PCA: Maximum Variance Formulation IV

Now maximize the projected variance wrt \mathbf{u}_1 ,

$$\mathbf{u}_1 = \arg \max_{\mathbf{u} \in \mathbb{S}_D} \mathbf{u}^\top \mathbf{S} \mathbf{u}. \quad (5)$$

This is the variational definition of the eigenvector with maximal eigenvalue!

I.e., \mathbf{u}_1 is the eigenvector of \mathbf{S} with the largest eigenvalue, λ_1 .

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

Since \mathbf{S} is real and symmetric positive definite, the eigenvectors are orthogonal.

PCA and the Singular Value Decomposition

The first M **principal components** are the leading M **eigenvectors of the covariance matrix**. Equivalently, they are the first M **right singular vectors of the data matrix**.

Let

$$\mathbf{Y} = \frac{1}{\sqrt{N}}\mathbf{X} = \frac{1}{\sqrt{N}} \begin{bmatrix} - & (\mathbf{x}_1 - \bar{\mathbf{x}})^\top & - \\ & \vdots & \\ - & (\mathbf{x}_N^\top - \bar{\mathbf{x}})^\top & - \end{bmatrix} \quad (6)$$

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,

$$\mathbf{Y} = \mathbf{V} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^\top \Rightarrow \mathbf{Y}^\top \mathbf{Y} = \frac{1}{N} \mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{V}^\top \mathbf{V} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^\top = \frac{1}{N} \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top \quad (7)$$

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,

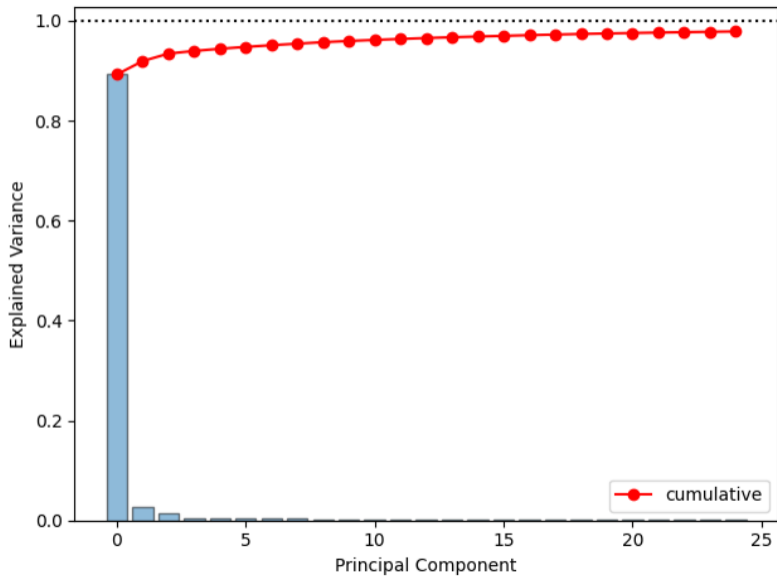
$$\text{Cov}[\mathbf{z}] = \text{Cov}[\mathbf{U}_M^\top(\mathbf{x} - \bar{\mathbf{x}})] = \mathbf{U}_M^\top \text{Cov}[\mathbf{x}] \mathbf{U}_M = \text{diag}([\lambda_1, \dots, \lambda_M]). \quad (8)$$

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**,

$$\text{variance explained} = \frac{\text{Tr}(\text{Cov}[\mathbf{z}; M \text{ components}])}{\text{Tr}(\text{Cov}[\mathbf{z}; D \text{ components}])} = \frac{\sum_{m=1}^M \lambda_m}{\sum_{m=1}^D \lambda_m} \in [0, 1]. \quad (9)$$

Scree Plots



Outline

- ▶ Principal Components Analysis (PCA)
- ▶ **PCA as a linear Gaussian latent variable model**
- ▶ Factor analysis
- ▶ Linear Dynamical Systems & the Kalman Filter/Smoothen

Probabilistic PCA: A Continuous Latent Variable Model

We cast the principal components as the solutions to an optimization problem: maximize the projected variance.

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**, including EM.
- ▶ 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,

$$\mathbf{z}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I) \quad (10)$$

$$\mathbf{x}_n | \mathbf{z}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 I), \quad (11)$$

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,

$$\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu} + \boldsymbol{\epsilon}_n, \quad (12)$$

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

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),

$$\mathbf{W}_{\text{ML}}, \boldsymbol{\mu}_{\text{ML}}, \sigma_{\text{ML}}^2 = \arg \max \mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2), \quad (13)$$

where \mathcal{L} is the marginal likelihood,

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \log p(\mathbf{X} | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \quad (14)$$

$$= \log \int \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) p(\mathbf{z}_n) d\mathbf{z}_n \quad (15)$$

$$= \log \int \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}) d\mathbf{z}_n \quad (16)$$

Exercise: Simplify this expression.

Maximum likelihood estimation of the parameters II

The log likelihood simplifies to,

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \frac{ND}{2} \log 2\pi - \frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (17)$$

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,

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_M (\boldsymbol{\Lambda}_M - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}, \quad (18)$$

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,

$$\sigma_{\text{ML}}^2 = \frac{1}{D-M} \sum_{m=M+1}^D \lambda_m, \quad (19)$$

the average variance in the remaining dimensions.

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?

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 ?

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \propto \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \quad (20)$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{z}_n - \frac{1}{2} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \boldsymbol{\mu})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \boldsymbol{\mu}) \right\} \quad (21)$$

$$\propto \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{J}_n \mathbf{z}_n + \mathbf{h}_n^\top \mathbf{z}_n \right\} \quad (22)$$

$$(23)$$

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,

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}(\mathbf{z}_n | \mathbf{J}_n^{-1} \mathbf{h}_n, \mathbf{J}_n^{-1}). \quad (24)$$

The Posterior Distribution in the Zero Noise Limit

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

$$\lim_{\sigma^2 \rightarrow 0} \mathbb{E}[\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2] = \lim_{\sigma^2 \rightarrow 0} (\mathbf{I} + \frac{1}{\sigma^2} \mathbf{W}^\top \mathbf{W})^{-1} [\frac{1}{\sigma^2} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu})] \quad (25)$$

$$= \lim_{\sigma^2 \rightarrow 0} (\sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (26)$$

$$= (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (27)$$

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

$$\lim_{\sigma^2 \rightarrow 0} \mathbb{E}[\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \boldsymbol{\mu}, \sigma^2] = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (28)$$

$$= \boldsymbol{\Lambda}_M^{-\frac{1}{2}} \mathbf{U}_M^\top (\mathbf{x}_n - \boldsymbol{\mu}) \quad (29)$$

EM for Probabilistic PCA

The **E-step** is to compute the posterior $q(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})$, where $\boldsymbol{\theta} = (\mathbf{W}, \boldsymbol{\mu}, \sigma^2)$ are the current parameters. For simplicity, assume the data is centered so that $\boldsymbol{\mu}^* = 0$.

The **M-step** is to maximize the expected complete data log likelihood,

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n | \mathbf{z}_n; \boldsymbol{\theta})] \\ &= \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\log \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n, \sigma^2 \mathbf{I})] \\ &= \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} \left[-\frac{D}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n)^\top (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n) \right].\end{aligned}$$

EM for Probabilistic PCA

As a function of \mathbf{W} ,

$$\begin{aligned}\mathcal{L}(\mathbf{W}) &= \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} \left[-\frac{1}{2\sigma^2} \langle \mathbf{W}^\top \mathbf{W}, \mathbf{z}_n \mathbf{z}_n^\top \rangle + \frac{1}{\sigma^2} \langle \mathbf{W}, \mathbf{x}_n \mathbf{z}_n^\top \rangle \right] \\ &= -\frac{1}{2\sigma^2} \left\langle \mathbf{W}^\top \mathbf{W}, \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\mathbf{z}_n \mathbf{z}_n^\top] \right\rangle + \frac{1}{\sigma^2} \left\langle \mathbf{W}, \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\mathbf{x}_n \mathbf{z}_n^\top] \right\rangle.\end{aligned}$$

where $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{A}^\top \mathbf{B})$ is the Frobenius inner product for matrices \mathbf{A} and \mathbf{B} .

Taking derivatives wrt \mathbf{W} and setting to zero yields,

$$\mathbf{W}^* = \left(\sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\mathbf{x}_n \mathbf{z}_n^\top] \right) \left(\sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} [\mathbf{z}_n \mathbf{z}_n^\top] \right)^{-1}.$$

It depends on sums of expected sufficient statistics!

Exercise: Derive the expected sufficient statistics and the M-step update for σ^2 .

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,

$$\mathbf{z}_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I) \quad (30)$$

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \quad (31)$$

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

Exercise: without doing any math, derive EM for this factor analysis model.

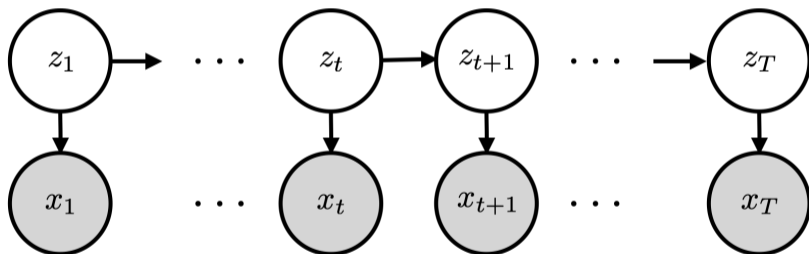
Outline

- ▶ Principal Components Analysis (PCA)
- ▶ PCA as a linear Gaussian latent variable model
- ▶ Factor analysis
- ▶ **Linear Dynamical Systems & the Kalman Filter/Smoothen**

Recap: Hidden Markov Models (HMMs)

We generalized from mixture models to HMMs by assuming that the latent states were dependent.

HMMs assume a particular factorization of the joint distribution on latent states (z_t) and observations (\mathbf{x}_t). The graphical model looks like this:



This graphical model says that the joint distribution factors as,

$$p(z_{1:T}, \mathbf{x}_{1:T}) = p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t | z_t). \quad (32)$$

State space models (SSMs)

Note that nothing above assumes that z_t is a discrete random variable!

HMM's are a special case of more general **state space models** with discrete states.

State space models assume the same graphical model but allow for arbitrary types of latent states.

For example, suppose that $\mathbf{z}_t \in \mathbb{R}^D$ are continuous valued latent states and that,

$$p(\mathbf{z}_{1:T}) = p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \quad (33)$$

$$= \mathcal{N}(\mathbf{z}_1 | \mathbf{b}_1, \mathbf{Q}_1) \prod_{t=2}^T \mathcal{N}(\mathbf{z}_t | \mathbf{A}\mathbf{z}_{t-1} + \mathbf{b}, \mathbf{Q}) \quad (34)$$

This is called a Gaussian **linear dynamical system** (LDS).

Stability of Gaussian linear dynamical systems

Question: What is the asymptotic mean of a Gaussian LDS, $\lim_{t \rightarrow \infty} \mathbb{E}[\mathbf{z}_t]$?

Question: When is a Gaussian LDS stable? I.e. when is the asymptotic mean finite?

Message passing in HMMs

In the HMM with discrete states, we showed how to compute posterior marginal distributions using message passing,

$$p(z_t | \mathbf{x}_{1:T}) \propto \sum_{z_1} \cdots \sum_{z_{t-1}} \sum_{z_{t+1}} \cdots \sum_{z_T} p(z_{1:T}, \mathbf{x}_{1:T}) \quad (35)$$

$$= \alpha_t(z_t) p(\mathbf{x}_t | z_t) \beta_t(z_t) \quad (36)$$

where the *forward and backward messages* are defined recursively

$$\alpha_t(z_t) = \sum_{z_{t-1}} p(z_t | z_{t-1}) p(\mathbf{x}_{t-1} | z_{t-1}) \alpha_{t-1}(z_{t-1}) \quad (37)$$

$$\beta_t(z_t) = \sum_{z_{t+1}} p(z_{t+1} | z_t) p(\mathbf{x}_{t+1} | z_{t+1}) \beta_{t+1}(z_{t+1}) \quad (38)$$

The initial conditions are $\alpha_1(z_1) = p(z_1)$ and $\beta_T(z_T) = 1$.

What do the forward messages tell us?

The forward messages are equivalent to,

$$\alpha_t(z_t) = \sum_{z_1} \cdots \sum_{z_{t-1}} p(z_{1:t}, \mathbf{x}_{1:t-1}) \quad (39)$$

$$p(z_t, \mathbf{x}_{1:t-1}). \quad (40)$$

The normalized message is the *predictive distribution*,

$$\frac{\alpha_t(z_t)}{\sum_{z'_t} \alpha_t(z'_t)} = \frac{p(z_t, \mathbf{x}_{1:t-1})}{\sum_{z'_t} p(z'_t, \mathbf{x}_{1:t-1})} = \frac{p(z_t, \mathbf{x}_{1:t-1})}{p(\mathbf{x}_{1:t-1})} = p(z_t | \mathbf{x}_{1:t-1}), \quad (41)$$

The final normalizing constant yields the marginal likelihood, $\sum_{z_T} \alpha_T(z_T) = p(\mathbf{x}_{1:T})$.

Message passing in state space models

The same recursive algorithm applies (in theory) to any state space model with the same factorization, but the sums are replaced with integrals,

$$p(\mathbf{z}_t | \mathbf{x}_{1:T}) \propto \int d\mathbf{z}_1 \cdots \int d\mathbf{z}_{t-1} \int d\mathbf{z}_{t+1} \cdots \int d\mathbf{z}_T p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) \quad (42)$$

$$= \alpha_t(\mathbf{z}_t) p(\mathbf{x}_t | \mathbf{z}_t) \beta_t(\mathbf{z}_t) \quad (43)$$

where the *forward and backward messages* are defined recursively

$$\alpha_t(\mathbf{z}_t) = \int p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{z}_{t-1}) \alpha_{t-1}(\mathbf{z}_{t-1}) d\mathbf{z}_{t-1} \quad (44)$$

$$\beta_t(\mathbf{z}_t) = \int p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) \beta_{t+1}(\mathbf{z}_{t+1}) d\mathbf{z}_{t+1} \quad (45)$$

The initial conditions are $\alpha_1(\mathbf{z}_1) = p(\mathbf{z}_1)$ and $\beta_T(\mathbf{z}_T) \propto 1$.

Forward pass in a linear dynamical system

Consider an linear dynamical system (LDS) with Gaussian emissions,

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) = p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \quad (46)$$

$$= \mathcal{N}(\mathbf{z}_1 | \mathbf{b}_1, \mathbf{Q}_1) \prod_{t=2}^T \mathcal{N}(\mathbf{z}_t | \mathbf{A}\mathbf{z}_{t-1} + \mathbf{b}, \mathbf{Q}) \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t | \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \quad (47)$$

Let's derive the forward message $\alpha_{t+1}(\mathbf{z}_{t+1})$. Assume $\alpha_t(\mathbf{z}_t) \propto \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1})$.

$$\alpha_{t+1}(\mathbf{z}_{t+1}) = \int p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_t | \mathbf{z}_t) \alpha_t(\mathbf{z}_t) d\mathbf{z}_t \quad (48)$$

$$= \int \mathcal{N}(\mathbf{z}_{t+1} | \mathbf{A}\mathbf{z}_t + \mathbf{b}, \mathbf{Q}) \mathcal{N}(\mathbf{x}_t | \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) d\mathbf{z}_t \quad (49)$$

The update step

The first step is the **update step**, where we **condition on** the emission \mathbf{x}_t ,

Exercise: Expand the densities, collect terms, and complete the square to compute the mean $\boldsymbol{\mu}_{t|t}$ and covariance $\boldsymbol{\Sigma}_{t|t}$ after the update step,

$$\mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \propto \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}). \quad (50)$$

The update step II

Write the joint distribution,

$$p(\mathbf{z}_t, \mathbf{x}_t \mid \mathbf{x}_{1:t-1}) = \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) \quad (51)$$

$$= \mathcal{N}\left(\begin{bmatrix} \mathbf{z}_t \\ \mathbf{x}_t \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_{t|t-1} \\ \mathbf{C}\boldsymbol{\mu}_{t|t-1} + \mathbf{d} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{t|t-1} & \boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top \\ \mathbf{C}\boldsymbol{\Sigma}_{t|t-1} & \mathbf{C}\boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top + \mathbf{R} \end{bmatrix}\right) \quad (52)$$

Since $(\mathbf{z}_t, \mathbf{x}_t)$ are jointly Gaussian, \mathbf{z}_t must be conditionally Gaussian as well,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) = \mathcal{N}(\boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}). \quad (53)$$

Exercise: Now use the **Schur complement** from Week 1 to solve for $\boldsymbol{\mu}_{t|t}$ and $\boldsymbol{\Sigma}_{t|t}$

The update step III

Exercise: Write $\mu_{t|t}$ and $\Sigma_{t|t}$ in terms of the **Kalman gain**,

$$K_t = \Sigma_{t|t-1} C^T (C \Sigma_{t|t-1} C^T + R)^{-1} \quad (54)$$

What is the Kalman gain doing?

The predict step

The predict step yields $p(\mathbf{z}_t | \mathbf{x}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t})$. To complete the forward pass, we need the **predict step**,

$$\alpha_{t+1}(\mathbf{z}_{t+1}) = \int p(\mathbf{z}_{t+1} | \mathbf{z}_t) p(\mathbf{x}_t | \mathbf{z}_t) \alpha_t(\mathbf{z}_t) d\mathbf{z}_t \quad (55)$$

$$= \int \mathcal{N}(\mathbf{z}_{t+1} | \mathbf{A}\mathbf{z}_t + \mathbf{b}, \mathbf{Q}) \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) d\mathbf{z}_t \quad (56)$$

$$= \mathcal{N}(\mathbf{z}_{t+1} | \boldsymbol{\mu}_{t+1|t}, \boldsymbol{\Sigma}_{t+1|t}) \quad (57)$$

Exercise: Solve for the mean $\boldsymbol{\mu}_{t+1|t}$ and covariance $\boldsymbol{\Sigma}_{t+1|t}$ after the predict step.

Completing the recursions

That wraps up the recursions! All that's left is the base case, which comes from the initial state distribution,

$$\mu_{1|0} = \mathbf{b}_1 \quad \text{and} \quad \Sigma_{1|0} = \mathbf{Q}_1. \quad (58)$$

Computing the marginal likelihood

Like in the discrete HMM, we can compute the marginal likelihood along the forward pass.

$$p(\mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{x}_{1:t-1}) \quad (59)$$

$$= \prod_{t=1}^T \int p(\mathbf{x}_t \mid \mathbf{z}_t) p(\mathbf{z}_t \mid \mathbf{x}_{1:t-1}) d\mathbf{z}_t \quad (60)$$

$$= \prod_{t=1}^T \int \mathcal{N}(\mathbf{x}_t \mid \mathbf{C}\mathbf{z}_t + \mathbf{d}, \mathbf{R}) \mathcal{N}(\mathbf{z}_t \mid \boldsymbol{\mu}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}) d\mathbf{z}_t \quad (61)$$

Exercise: Obtain a closed form expression for the integrals.

Computing the smoothing distributions

- ▶ The forward pass gives us the filtering distributions $p(\mathbf{z}_t | \mathbf{x}_{1:t})$. How can we compute the smoothing distributions, $p(\mathbf{z}_t | \mathbf{x}_{1:T})$?
- ▶ In the discrete HMM we essentially ran the *same algorithm in reverse* to get the backward messages, starting from $\beta_T(\mathbf{z}_T) \propto 1$.
- ▶ We can do the same sort of thing here, but it's a bit funky because we need to start with an improper Gaussian distribution $\beta_T(\mathbf{z}_T) \propto \mathcal{N}(\mathbf{0}, \infty I)$.
- ▶ Instead, we'll derive an alternative way of computing the smoothing distributions.

Bayesian Smoothing

Note: \mathbf{z}_t is conditionally independent of $\mathbf{x}_{t+1:T}$ given \mathbf{z}_{t+1} , so

$$p(\mathbf{z}_t \mid \mathbf{z}_{t+1}, \mathbf{x}_{1:T}) = p(\mathbf{z}_t \mid \mathbf{z}_{t+1}, \mathbf{x}_{1:t}) \quad (62)$$

$$= \frac{p(\mathbf{z}_t, \mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})} \quad (63)$$

$$= \frac{p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) p(\mathbf{z}_{t+1} \mid \mathbf{z}_t)}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})} \quad (64)$$

Question: what rules did we apply in each of these steps?

Bayesian Smoothing II

Now we can write the joint distribution as,

$$p(\mathbf{z}_t, \mathbf{z}_{t+1} \mid \mathbf{x}_{1:T}) = p(\mathbf{z}_t \mid \mathbf{z}_{t+1} \mid \mathbf{x}_{1:T})p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T}) \quad (65)$$

$$= \frac{p(\mathbf{z}_t \mid \mathbf{x}_{1:t})p(\mathbf{z}_{t+1} \mid \mathbf{z}_t)p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T})}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})}. \quad (66)$$

Marginalizing over \mathbf{z}_{t+1} gives us,

$$p(\mathbf{z}_t \mid \mathbf{x}_{1:T}) = p(\mathbf{z}_t \mid \mathbf{x}_{1:t}) \int \frac{p(\mathbf{z}_{t+1} \mid \mathbf{z}_t)p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:T})}{p(\mathbf{z}_{t+1} \mid \mathbf{x}_{1:t})} d\mathbf{z}_{t+1} \quad (67)$$

Question: Can we compute each of these terms?

The Rauch-Tung-Striebel Smoother, aka Kalman Smoother

These recursions apply to any Markovian state space model. For the special case of a Gaussian linear dynamical system, the smoothing distributions are again Gaussians,

$$p(\mathbf{z}_t | \mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}_t | \boldsymbol{\mu}_{t|T}, \boldsymbol{\Sigma}_{t|T}) \quad (68)$$

where

$$\boldsymbol{\mu}_{t|T} = \boldsymbol{\mu}_{t|t} + \mathbf{G}_t(\boldsymbol{\mu}_{t+1|T} - \boldsymbol{\mu}_{t+1|t}) \quad (69)$$

$$\boldsymbol{\Sigma}_{t|T} = \boldsymbol{\Sigma}_{t|t} + \mathbf{G}_t(\boldsymbol{\Sigma}_{t+1|T} - \boldsymbol{\Sigma}_{t+1|t})\mathbf{G}_t^\top \quad (70)$$

$$\mathbf{G}_t \triangleq \boldsymbol{\Sigma}_{t|t}\mathbf{A}^\top \boldsymbol{\Sigma}_{t+1|t}^{-1}. \quad (71)$$

This is called the **Rauch-Tung-Striebel (RTS) smoother** or the **Kalman smoother**.

References