

Topics in PCA

Asymptotic, Non-linear, Functional and Probabilistic PCA

Elvis Cui

MS in Biostatistics, UCLA

November 4, 2019

Overview

- 1 Ordinary PCA
 - PCA as variance maximization
 - PCA as Reduced-rank regression
- 2 Nonlinear PCA
 - Polynomial PCA
 - Kernel PCA
 - Principal Curves and Surfaces
 - Autoencoder
- 3 Functional PCA
 - Karhunen-Loeve Expansion
- 4 Probabilistic PCA
 - Gaussian Linear Model
- 5 High dimensional PCA
 - Dual Representation
 - Empirical Spectral Distribution
 - Sparse PCA

Hotelling's PCA

In 1933, Hotelling developed Principal Component Analysis (a.k.a. Karhunen-Loeve Transform).

Definition 1 (First Principal Component)

Given a random vector $\mathbf{X} \in \mathbb{R}^p$, the first principal component of \mathbf{X} is given by the following optimization problem:

$$\begin{aligned} \max_{\mathbf{v}_1} \text{Var}(\mathbf{v}_1^T \mathbf{X}) &= \mathbf{v}_1^T \boldsymbol{\Sigma}_X \mathbf{v}_1 \\ \text{s.t. } \|\mathbf{v}_1\| &= 1 \end{aligned}$$

Equivalently,

$$\mathbf{v}_1 = \arg \max_{\mathbf{v}_1} \frac{\mathbf{v}_1^T \boldsymbol{\Sigma}_X \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1}$$

The maximum is called the **operator norm** of $\boldsymbol{\Sigma}_X$.

Hotelling's PCA cont'd

Definition 2 (i^{th} Principal Component)

Given the first $(i - 1)$ principal component ($i \leq n$), the i^{th} principal component of \mathbf{X} is defined to be

$$\begin{aligned} \max_{\mathbf{v}_i} \text{Var}(\mathbf{v}_i^T \mathbf{X}) &= \mathbf{v}_i^T \boldsymbol{\Sigma}_X \mathbf{v}_i \\ \text{s.t. } \|\mathbf{v}_i\| &= 1, \langle \mathbf{v}_t, \mathbf{v}_i \rangle = 0 \quad \forall t < i \end{aligned}$$

Hotelling's PCA: Sequential Method

We could derive all n principal components sequentially:

- Introducing Lagrangian multiplier λ_1 , define

$$f(\mathbf{v}_1) = \mathbf{v}_1^T \mathbf{\Sigma}_X \mathbf{v}_1 + \lambda_1 (1 - \mathbf{v}_1^T \mathbf{v}_1)$$

- Taking derivative: $\frac{\partial f(\mathbf{v}_1)}{\partial \mathbf{v}_1} = 2(\mathbf{\Sigma}_X - \lambda_1 \mathbf{I}_p) \mathbf{v}_1 = \mathbf{0}$. The maximum of variance is attained when λ_1 is the largest eigenvalue of $\mathbf{\Sigma}_X$.

- Introducing Lagrangian multiplier λ_2 and μ , define

$$f(\mathbf{v}_2) = \mathbf{v}_2^T \mathbf{\Sigma}_X \mathbf{v}_2 + \lambda_2 (1 - \mathbf{v}_2^T \mathbf{v}_2) + \mu \mathbf{v}_2^T \mathbf{v}_1$$

- Taking derivative w.r.t. \mathbf{v}_2 :

$$\frac{\partial f(\mathbf{v}_2)}{\partial \mathbf{v}_2} = 2(\mathbf{\Sigma}_X - \lambda_2 \mathbf{I}_p) \mathbf{v}_2 + \mu \mathbf{v}_1 = \mathbf{0}$$

Hotelling's PCA: Sequential Method

Cont'd from last slide:

- Pre-multiplying \mathbf{v}_1 on both sides

$$2\mathbf{v}_1^T \mathbf{\Sigma}_X \mathbf{v}_2 + \mu = 0$$

- Pre-multiplying the equation $((\mathbf{\Sigma}_X - \lambda_1 \mathbf{I}_p) \mathbf{v}_1 = \mathbf{0})$ by \mathbf{v}_2 yields $\mathbf{v}_2^T \mathbf{\Sigma}_X \mathbf{v}_1 = 0$. Thus, $\mu = 0$
- $\mu = 0 \Rightarrow \lambda_2$ is the second largest eigenvalue of $\mathbf{\Sigma}_X$.
- Repeating the above procedure to get all n principal components.

Reduced-rank regression

PCA can also be derived from a classic multivariate analysis framework known as **Reduced-rank Regression**.

Definition 3 (Reduced-rank Regression)

Suppose $\mathbf{X} \in \mathbb{R}^p$ with mean $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}_X$, $\mathbf{B} \in \mathbb{R}^{t \times p}$ with $t < p$ and $\mathbf{A} \in \mathbb{R}^{p \times t}$. The goal is to choose \mathbf{A} , \mathbf{B} and $\boldsymbol{\mu}$ to minimize a least-square criterion:

$$\min_{\mathbf{A}, \mathbf{B}, \boldsymbol{\mu}} \mathbb{E}\{\|\mathbf{X} - \boldsymbol{\mu} - \mathbf{ABX}\|_2^2\}$$

Reduced-rank Regression cont'd

Theorem 4 (Solution of Reduced-rank Regression)

The solution of the reduced-rank regression problem is given by the eigenvectors and eigenvalues of $\text{Var}(\mathbf{X}) = \mathbf{\Sigma}_X$:

$$\mathbf{A} = (\mathbf{v}_1, \dots, \mathbf{v}_t) = \mathbf{B}^T$$

$$\boldsymbol{\mu} = (\mathbf{I}_p - \mathbf{AB})\boldsymbol{\mu}_X$$

Where \mathbf{v}_j is the eigenvector associated with the j^{th} largest eigenvalue, λ_j , of $\mathbf{\Sigma}_X$.

Proof.

e.g., Chapter 7 section 2 of modern multivariate statistical techniques (Izemmann, 2013).



Polynomial PCA

- Adding quadratic, cubic terms (Gnanadesikan and Wilk, 1969)
- **Problem:** Too many terms when dimension of \mathbf{x} is large.

Example 5 (Quadratic PCA)

- $\mathbf{x} \in \mathbb{R}^p \Rightarrow$ Extended $\mathbf{x} \in \mathbb{R}^{2p+p(p-1)/2}$ (p quadratic powers and $p(p-1)/2$ cross terms).
- λ_{min} of the covariance matrix of the extended vector will be 0 if observations follow an exact quadratic curve.

Kernel PCA: Intuition

Kernel PCA can be considered as a two-step process:

- Non-linearly transform the data point $\mathbf{x} \in \mathbb{R}^r$ into Φ_i in a high dimensional **feature space** \mathcal{H} (say p dimension), where

$$\Phi_i = (\phi_1(\mathbf{x}_i), \dots, \phi_p(\mathbf{x}_i))$$

The map Φ is called a **feature map**.

- Given $\Phi_1, \dots, \Phi_n \in \mathcal{H}$ with $\sum_{i=1}^n \Phi_i = \mathbf{0}$, solve a linear PCA problem in feature space \mathcal{H} .
- Note that the dimension of \mathcal{H} can be ∞ , and in the second step, the data **MUST** be centered in \mathcal{H} (can be done via a matrix trick).

Kernel PCA: Kernel trick

Suppose Φ'_i s are centered. Define the empirical covariance matrix:

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T$$

The eigen-value and eigen-vector of \mathbf{C} is given by $\mathbf{C}\mathbf{v} = \lambda\mathbf{v}$

Lemma 6 (Equivalent form of eigen-vector)

The eigen-vector and eigen-values can also be derived from solving n linear equations:

$$\langle \Phi_i, \mathbf{C}\mathbf{v} \rangle = \lambda \langle \Phi_i, \mathbf{v} \rangle, i = 1, \dots, n$$

Or more compactly, define $\Phi^T = (\Phi_1, \dots, \Phi_n)$ (thus $\Phi \in \mathbb{R}^{n \times p}$)

$$\langle \Phi^T, \mathbf{C}\mathbf{v} \rangle = \lambda \langle \Phi^T, \mathbf{v} \rangle$$

Kernel PCA: Kernel trick cont'd

Proof.

Note that

$$\mathbf{C}\mathbf{v} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T \mathbf{v} = \frac{1}{n} \sum_{i=1}^n \Phi_i \langle \Phi_i, \mathbf{v} \rangle = \lambda \mathbf{v}$$

Thus each eigenvector is a linear combination of $\Phi_i, i = 1, \dots, n$.
So there exists coefficients $\alpha_i, i = 1, \dots, n$ s.t.

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \Phi_i$$



Kernel PCA: Kernel trick cont'd

Substituting \mathbf{C} and the representation $\mathbf{v} = \sum_{i=1}^n \alpha_i \Phi_i$ in

$$\langle \Phi^T, \mathbf{Cv} \rangle = \lambda \langle \Phi^T, \mathbf{v} \rangle$$

we get that (vectorized version, you may want to verify it):

$$\langle \Phi^T, \frac{1}{n} \Phi^T \Phi \langle \Phi, \alpha \rangle \rangle = \lambda \langle \Phi^T, \langle \Phi, \alpha \rangle \rangle$$

where $\alpha^T = (\alpha_1, \dots, \alpha_n)$. Moreover,

$$LHS = \frac{1}{n} \Phi \Phi^T \Phi \Phi^T \alpha = \lambda \Phi \Phi^T \alpha = RHS$$

Kernel PCA: Kernel trick cont'd

Now let's define $\mathbf{K} = \Phi\Phi^T$.

Then $LHS = \mathbf{K}^2\alpha = n\lambda\mathbf{K}\alpha = RHS$, or as (why ?)

$$\mathbf{K}\alpha = n\lambda\alpha$$

Here \mathbf{K} is the famous **Kernel Matrix** with elements denoted as $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. Therefore, we don't need to define the mapping Φ explicitly but just the inner product (kernel) $k(\mathbf{x}_i, \mathbf{x}_j)$.

- The principal component (or direction) in the feature space \mathcal{H} is the eigenvector α_{\max} corresponding to the largest eigenvalue λ_{\max} of the kernel matrix \mathbf{K} .

Kernel PCA: Normalization

- Recall that $\mathbf{C}\mathbf{v} = \lambda\mathbf{v}$, if we require $\|\mathbf{v}_i\|_2^2 = 1, i = 1, \dots, n$, then, using the representation of \mathbf{v} by α , we have

$$\begin{aligned}
 1 &= \sum_{j=1}^n \sum_{k=1}^n \alpha_{ij} \alpha_{ik} \langle \Phi_j, \Phi_k \rangle \\
 &= \sum_{j=1}^n \sum_{k=1}^n \alpha_{ij} \alpha_{ik} K_{jk} \\
 &= \langle \alpha_i, \mathbf{K} \alpha_i \rangle \\
 &= n\lambda \langle \alpha_i, \alpha_i \rangle
 \end{aligned}$$

which determines the normalization for the vectors

$\alpha_1, \dots, \alpha_n$.

Kernel PCA: Centralization

If \mathbf{K} is not centralized, then replace it with (Scholkopf, Smola and Muller, 1998)

$$\tilde{\mathbf{K}} = \mathbf{H}\mathbf{K}\mathbf{H}$$

where $\mathbf{H} = \mathbf{I}_n - \frac{1}{n}\mathbf{J}_n$ is called the centering matrix, $\mathbf{J}_n = \mathbf{1}_n\mathbf{1}_n^T$ is an $(n \times n)$ -matrix of all ones.

- Using this trick, the mapping Φ_i in feature space \mathcal{H} is centered.

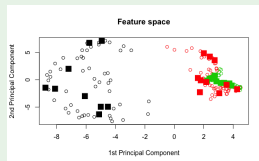
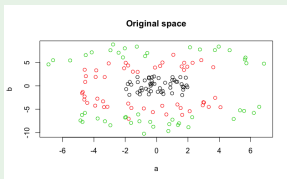
Kernel PCA: Example with toy data

Left: original data.

Right: data in first two principal component space.

Square points: predicted new feature with new inputs.

Example 7 (Three clusters (Gaussian kernel))



Principal Curves: Intuition

Another way of linear reduction is called **Principal Curves**. First, I give intuitive definitions and then we run into formal ones.

- **Principle curves and surfaces** were proposed by Hastie, 1984; Hastie and Stuetzle, 1989.
- Suppose the data observed on **X** lie close to a smooth nonlinear manifold of low dimension.
- **Principle Curve** is a smooth one-dimensional **parameterized curve f** that passes through the "middle" of the data, regardless of whether the "middle" is a straight line or a nonlinear curve.

Principle Surface: A generalization of principle curve to a high dimensional curve.

Principle Curves: A Review of Geometry

[geometry] To give formal definition of principal curves, we need to review some concepts from geometry:

- **p dimension curve/surface in \mathbb{R}^k ($p < k$):**

$$\mathbf{f}(\lambda) = (f_1(\lambda), f_2(\lambda), \dots, f_k(\lambda))^T$$

Where $\lambda = (\lambda_1, \dots, \lambda_p)$ is the parameterization.

Example 8

The unit sphere in \mathbb{R}^3 : $\{(x_1, x_2, x_3) : x_1^2 + x_2^2 + x_3^2 = 1\}$ is a two dimensional surface:

$$\mathbf{f}(\lambda) = (\cos \lambda_1 \sin \lambda_2, \sin \lambda_1 \sin \lambda_2, \cos \lambda_2)$$

Usually, it is convenient for us to assume $\mathbf{f} \in \mathcal{C}^\infty$.

Principle Curves: A Review of Geometry

For simplicity, let's consider the case $p = 1$ (or $\lambda \in \mathbb{R}$).

- \mathbf{f} is **closed** of **periodic** if $\mathbf{f}(\lambda + \alpha) = \mathbf{f}(\lambda) \forall \lambda$.
- Velocity at λ : $\mathbf{v} = \mathbf{f}'(\lambda)$
- Speed at λ : $\|\mathbf{v}\| = \|\mathbf{f}'(\lambda)\|$
- Acceleration: $\mathbf{a} = \mathbf{f}''(\lambda)$
- Curvature: $K(\lambda) = \|\mathbf{f}''(\lambda)\|$
- Circle of curvature (at λ): unit-speed circle with radius $r(\lambda) = \frac{1}{\|\mathbf{f}''(\lambda)\|}$
- Arc-length (from λ_0 to λ_1):

$$L(\mathbf{f}) = \int_{\lambda_0}^{\lambda_1} \|\mathbf{f}'(\lambda)\| d\lambda$$

When the speed of $\mathbf{f}(\lambda)$ is identically one, the arc-length is $\lambda_1 - \lambda_0$.

Principal Curves: Definition

Again, let's assume $\mathbf{f}(\lambda)$ is a 1-d parameterized curve in \mathbb{R}^k .

Definition 9 (Projection of \mathbf{x} on $\mathbf{f}(\lambda)$)

Suppose the minimum (or infimum) is attained, the projection is defined to be

$$\mathbf{x}_f = \arg \min_{\mathbf{f}(\mu)} \|\mathbf{x} - \mathbf{f}(\mu)\|$$

Definition 10 (Projection index)

$$\lambda_f(\mathbf{x}) = \sup_{\lambda} \{\lambda : \|\mathbf{x} - \mathbf{f}(\lambda)\|\} = \inf_{\mu} \|\mathbf{x} - \mathbf{f}(\mu)\|$$

We are taking sup because sometimes there are multiple choices of μ s.t. the infimum is attained. Note that $\lambda_f(\mathbf{x})$ can be discontinuous.

Principle Curves: Definition cont'd

Definition 11 (Reconstruction error)

The reconstruction error is the expected squared distance between \mathbf{X} (or its associated pdf) and \mathbf{f} :

$$D^2(\mathbf{X}, \mathbf{f}) = \mathbb{E}\{\|\mathbf{X} - \mathbf{f}(\lambda_{\mathbf{f}}(\mathbf{X}))\|^2\}$$

Definition 12 (Principle curve)

The 1-d curve $\mathbf{f}(\lambda)$ is called a **principle curve** if the following holds:

$$\mathbf{f}(\lambda) = \mathbb{E}\{\mathbf{X} | \lambda_{\mathbf{f}}(\mathbf{X}) = \lambda\} \text{ for a.e. } \lambda$$

In this case, $\mathbf{f}(\lambda)$ is also said to be **self-consistent**. In fact, $\mathbf{f}(\lambda)$ minimizes the reconstruction error defined above (Hastie and Stuetzle, 1989).

Note that principle curves are measurable.

Principal Curves: Generalization to surfaces

Principle curve has a natural extension to **principal surfaces** in high dimensions (w.r.t. λ) (Hastie, 1984; LeBlanc and Tibshirani, 1994).

Definition 13 (Principal surface)

A principal surface satisfies the self-consistency property,

$$\mathbf{f}(\lambda) = \mathbb{E}\{\mathbf{X} | \lambda_{\mathbf{f}}(\mathbf{X}) = \lambda\} \text{ for a.e. } \lambda$$

Principal Curves: Projection-Expectation Algorithm

- Goal: estimate \mathbf{f} from n observations $\{\mathbf{x}_i\}_1^n$.
- Idea: Empirical risk minimization (a.k.a. ERM principle in statistical learning theory)

$$D^2(\{\mathbf{x}_i\}, \hat{\mathbf{f}}) = \min_{\mathbf{f}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{f}(\lambda_{\mathbf{f}}(\mathbf{x}_i))\|^2$$

- Solution: An algorithm iterating between **projection** (estimating λ given \mathbf{f}) and **expectation** (estimating \mathbf{f} assuming a fixed λ).

Principal Curves: PE Algorithm cont'd

Definition 14 (Projection-Expectation Algorithm)

- Start with the initial curve (or line) $\mathbf{f}^{(0)}$. Then $\{\mathbf{x}_i\}$ are each projected onto this line, yielding the n points $\lambda_{\mathbf{f}^{(0)}}(\mathbf{x}_i) = \lambda_i^{(1)}$, $i = 1, 2, \dots, n$.
- The updated curve $\mathbf{f}^{(1)}$ is derived via the self-consistency property: $\mathbf{f}^{(1)}(\lambda_i^{(1)}) = \mathbb{E}\{\mathbf{X} | \lambda_{\mathbf{f}^{(0)}}(\mathbf{X}_i) = \lambda_i^{(1)}\}$, $i = 1, \dots, n$
- Given the current principal curve $\mathbf{f}^{(k-1)}$, the k^{th} iteration consists of 2 steps :

1 Projection step:

$$\lambda_{\mathbf{f}^{(k-1)}}(\mathbf{x}_i) = \lambda_i^{(k)}, \quad i = 1, 2, \dots, n$$

2 Expectation step:

$$\mathbf{f}^{(k)}(\lambda_i^{(k)}) = \mathbb{E}\{\mathbf{X} | \lambda_{\mathbf{f}^{(k-1)}}(\mathbf{X}_i) = \lambda_i^{(k)}\}, \quad i = 1, \dots, n$$

Principal Curves: PE Algorithm cont'd

- Problem: How to derive the expectation $\mathbf{f}^{(k)}(\lambda_i^{(k)}) = \mathbb{E}\{\mathbf{X} | \lambda_{\mathbf{f}^{(k-1)}}(\mathbf{X}_i) = \lambda_i^{(k)}\}$?
- Solution: estimate the expectation via **local averaging procedure** !
- More specifically, estimation of each coordinate of $\mathbf{f}^{(k)}(\lambda_i^{(k)})$ is accomplished using a **scatterplot smoother** (e.g., kernel, cubic spline, or locally weighted running-line smoother).
- Given the "final" estimation $(\hat{\lambda}_i, \hat{\mathbf{f}}_i), i = 1, \dots, n$, the principal curve $\hat{\mathbf{f}}(\lambda)$ is then the polygon produced by joining up these n tuples (see example in next slide).
- **IMPORTANT**: Convergence of PE algorithm has **NOT** been proved. It may lead to poor "local" solution.

Principal Curves: A 2-D Example

Example 15 (Principal Curve in 2-D)

$\mathbf{X} \sim \mathcal{U}(-5, 5)$, $\mathbf{Y} = \sin \mathbf{X} + \frac{2}{\mathbf{X}} \mathcal{N}(0, 1)$, 100 data points in total.

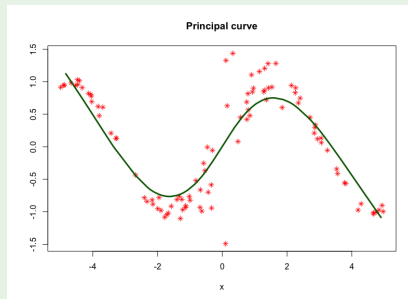


Figure: Principal Curve

Autoencoder in Deep Learning

- Peter Bickel: Imaging compression is a modern version of deriving sufficient statistic (Efron, 2016).

Example 16 (Autoencoder for MNIST)

- ➊ **Autoencoder:** Each image of MNIST dataset has 784 elements.
 - ➋ Reduce each image to 10 dimensions.
 - ➌ Reconstruct each image via its unique 10-d vector.
- Generate images from 10-d vector directly (without reduction step): This is known as **GAN** (Goodfellow, 2014).
 - See <http://cs231n.stanford.edu> for further information.

Functional PCA: Intuition

- **Intuition:** Data consisting of functions or curves.
- **Data:** Independent individuals may be recorded at different sets of time points.
- **Goal:** Characterize the main features of those curves.
- **Solution:** A functional version of PCA (e.g. Ramsay and Silverman, 1997)
- The vector-valued observations $\mathbf{X}_1, \dots, \mathbf{X}_n$ are replaced by

$$X_1(t), \dots, X_n(t), t \in [0, T]$$

- Each sample curve is considered to be an independent realization of a univariate stochastic process $X(t)$ with smooth mean function $\mathbb{E}\{X(t)\} = \mu(t)$ and covariance function

$$\text{cov}\{X(s), X(t)\} = \sigma(s, t)$$

Functional PCA: Karhunen-Loeve Expansion

By a spectral decomposition of the covariance function, we express σ as an orthogonal expansion in the L_2 sense in terms of its eigenvalues $\{\lambda_j\}$ and associated eigenfunctions $\{V_j(t)\}$, so that

$$\sigma(s, t) = \sum_{j=1}^{\infty} \lambda_j V_j(s) V_j(t)$$

where the eigenvalues quickly tend to zero and the first few eigenfunctions are slowly varying. A **random curve** can then be expressed as

$$X(t) = \mu(t) + \sum_{j=1}^{\infty} c_j V_j(t), \quad c_j = \int [X(t) - \mu(t)] V_j(t) dt$$

This is the well-known **Karhunen-Loeve expansion** in stochastic processes. The coefficient c_j is called **the j^{th} functional PC score** with $\mathbb{E}(c_j) = 0$, $\text{Var}(c_j) = \lambda_j$, $\sum_j \lambda_j < \infty$ and $\text{Cov}(c_j, c_k) = 0$.

Probabilistic PCA: Model Set-up

- The formulation of PCA as a probabilistic model was proposed independently by Tipping and Bishop in 1997 and by Roweis in 1998.
- Probabilistic PCA can be expressed as a directed graph in which each observation \mathbf{x}_i is associated with a value \mathbf{z}_i of the latent variable.
- This is a special case of **Gaussian Linear Model**.

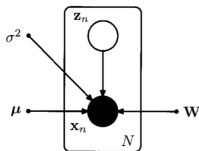


Figure: Bishop, 2006, pp574

Probabilistic PCA: Model Set-up cont'd

More formally, we write our model as

$$\mathbf{y} = \boldsymbol{\mu} + \mathbf{W}\mathbf{z} + \boldsymbol{\epsilon}$$

where $\boldsymbol{\mu} \in \mathbb{R}^p$, $\mathbf{W} \in \mathbb{R}^{p \times q}$ are fixed and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q)$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_p)$, \mathbf{z} and $\boldsymbol{\epsilon}$ are independent.

- $q \ll p$, thus the dimension of \mathbf{z} is much less than \mathbf{y} .
- \mathbf{z} is known as the latent variable and $\boldsymbol{\epsilon}$ is the noise vector.
- **Goal:** estimate \mathbf{W} and $\boldsymbol{\mu}$ so that we could recover \mathbf{z} (principal components) from observation \mathbf{y} (we assume σ) is known for simplicity.
- **Solution I:** Maximum likelihood estimation
- **Solution II:** EM Algorithm.

Probabilistic PCA: A Review of SVD

Before we go into MLE framework, let's take a quick review of **singular value decomposition**:

Definition 17 (Singular Value)

For any matrix $\mathbf{W} \in \mathbb{R}^{p \times q}$ ($p \geq q$), its singular value is defined to be

$$\sigma(\mathbf{W}) = \lambda(\mathbf{W}^T \mathbf{W})$$

where $\lambda(\cdot)$ denotes the eigen-value and $\sigma(\cdot)$ denotes the singular value.

Probabilistic PCA: SVD cont'd

Theorem 18 (Singular Value Decomposition)

For any matrix $\mathbf{W} \in \mathbb{R}^{p \times q}$, it admits the following decomposition:

$$\mathbf{W} = \tilde{\mathbf{U}}\tilde{\mathbf{L}}\mathbf{V}^T = \mathbf{U}\mathbf{L}\mathbf{V}^T$$

where $\tilde{\mathbf{U}} \in \mathbb{R}^{p \times p}$, $\mathbf{U} \in \mathbb{R}^{p \times q}$, $\mathbf{V} \in \mathbb{R}^{q \times q}$ are all **orthogonal matrices** and $\tilde{\mathbf{L}} \in \mathbb{R}^{p \times q}$, $\mathbf{L} \in \mathbb{R}^{q \times q}$ are **diagonal matrices** whose elements are known as **singular values**. Note that $\tilde{\mathbf{L}}$ has q non-zero elements lying on the diagonal.

We need both representations later in deriving MLE for \mathbf{W} and σ^2 .

MLE for Probabilistic PCA

- Knowing that \mathbf{y} has normal distribution, derive its mean and variance (left as an exercise)

$$\mathbb{E}(\mathbf{Y}) = \boldsymbol{\mu}$$

$$\text{Var}(\mathbf{Y}) = \sigma^2 \mathbf{I} + \mathbf{W}\mathbf{W}^T$$

- Write out log-likelihood function:

$$\begin{aligned} \log \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_n | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) &= \sum_{i=1}^n \log \mathbb{P}(\mathbf{y}_i | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{N}{2} (p \log(2\pi) + \log \det \mathbf{C}_y + \text{Tr}(\mathbf{C}_y^{-1} \mathbf{S})) \end{aligned}$$

where $\mathbf{C}_y = \sigma^2 \mathbf{I} + \mathbf{W}\mathbf{W}^T$ and $\mathbf{S} = \frac{1}{N} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T$.

MLE for Probabilistic PCA cont'd: 2

- A useful handbook for matrix derivatives: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>.
- Take derivative w.r.t. μ

$$-N\mathbf{C}_y^{-1}\mu = \sum_{i=1}^n \mathbf{C}_y^{-1}\mathbf{y}_i = 0$$

$$\Rightarrow \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^n \mathbf{y}_i$$

- MLE of \mathbf{W} is in general difficult to get but there is a closed form solution (Tipping and Bishop, 1999). However, in their original paper, I didn't figure out how they derive equation (A.13).
- On the other hand, I derived another equation similar to (A.13) which has the same solution. It is on the next slide.

MLE for Probabilistic PCA cont'd: 3

- Take derivative w.r.t. \mathbf{W} via the Chain rule

$$\mathbf{S}\mathbf{C}_y^{-1}\mathbf{W} = \mathbf{W}$$

- Problem:** how to derive the solution for \mathbf{W} ?
- Solution:** Use singular value decomposition !
- Singular value decomposition:** Write \mathbf{W} as

$$\mathbf{W} = \tilde{\mathbf{U}}\tilde{\mathbf{L}}\tilde{\mathbf{V}}$$

Then \mathbf{C}_y^{-1} becomes

$$(\sigma^2\mathbf{I}_p + \tilde{\mathbf{U}}\tilde{\mathbf{L}}\tilde{\mathbf{L}}^T\tilde{\mathbf{U}}^T)^{-1} = \tilde{\mathbf{U}}(\sigma^2\mathbf{I}_p + \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T)^{-1}\tilde{\mathbf{U}}^T$$

MLE for Probabilistic PCA cont'd: 4

- Using the new representation of \mathbf{C}_y^{-1} :

$$\mathbf{S}\mathbf{C}_y^{-1}\mathbf{W} = \mathbf{S}\tilde{\mathbf{U}}(\sigma^2\mathbf{I}_p + \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T)^{-1}\tilde{\mathbf{L}}\mathbf{V}^T$$

- Now here is the other form of SVD come in:

$$\mathbf{W} = \tilde{\mathbf{U}}\tilde{\mathbf{L}}\mathbf{V}^T = \mathbf{U}\mathbf{L}\mathbf{V}^T \Rightarrow$$

$$\tilde{\mathbf{U}}(\sigma^2\mathbf{I}_p + \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T)^{-1}\tilde{\mathbf{L}} = \mathbf{U}(\sigma^2\mathbf{I}_q + \mathbf{L}\mathbf{L}^T)^{-1}\mathbf{L}$$

For details, see next slide.

MLE for Probabilistic PCA cont'd: 5

- Here I made use of the fact that the inverse of the following matrix

$$(\sigma^2 \mathbf{I}_p) + \begin{bmatrix} l_1^2 & & \cdots & \\ & l_2^2 & & \\ \cdots & & & \\ & & l_q^2 & \\ & & & \mathbf{0} \end{bmatrix}$$

is nothing but

$$\begin{bmatrix} \frac{1}{l_1^2 + \sigma^2} & & & \\ & \cdots & & \\ & & \frac{1}{l_q^2 + \sigma^2} & \\ & & & \frac{1}{\sigma^2} & \cdots \end{bmatrix}$$

where l_i 's are diagonals of \mathbf{L} (singular values).

MLE for Probabilistic PCA cont'd: 6

- However, note that we could cut off all diagonals below the q^{th} row (or column) since all elements below q^{th} row of $\tilde{\mathbf{L}}$ is 0.
- This gives us the following matrix

$$\begin{bmatrix} \frac{1}{l_1^2 + \sigma^2} & & \\ & \dots & \\ & & \frac{1}{l_q^2 + \sigma^2} \end{bmatrix}$$

But this is just $(\sigma^2 \mathbf{I}_q + \mathbf{L}\mathbf{L}^T)^{-1}$.

MLE for Probabilistic PCA cont'd: 7

Recall slide 37 and 38, the equation now becomes

$$\mathbf{S}\mathbf{U}(\sigma^2\mathbf{I}_q + \mathbf{L}\mathbf{L}^T)^{-1}\mathbf{L}\mathbf{V}^T = \mathbf{U}\mathbf{L}\mathbf{V}^T$$

Since \mathbf{V} and \mathbf{L} is non-singular so they can be cancelled out on both sides. Therefore, the equation finally turns into

$$\mathbf{S}\mathbf{U} = \mathbf{U}(\sigma^2\mathbf{I}_q + \mathbf{L}\mathbf{L}^T)$$

This is equivalent to solve q eigen-value problems:

$$(\sigma^2 + l_j^2)\mathbf{u}_j = \mathbf{S}\mathbf{u}_j = \lambda_j\mathbf{u}_j, j = 1, \dots, q$$

where λ_j 's are the eigenvalues of sample covariance matrix \mathbf{S} .

MLE for Probabilistic PCA cont'd: 8

The solution of previous equations is

$$l_j = (\lambda_j - \sigma^2)^{-\frac{1}{2}}, j = 1, \dots, q$$

Thus,

$$\hat{\mathbf{W}}_{MLE} = \hat{\mathbf{U}}(\hat{\mathbf{L}}\hat{\mathbf{L}}^T - \sigma^2\mathbf{I}_q)^{\frac{1}{2}}\mathbf{V}$$

Now since the distribution of \mathbf{Y} does not depend on the rotation of \mathbf{W} (*), \mathbf{V} is chosen to be any orthogonal matrix.

- (*) That is, the mean and variance of \mathbf{Y} stays the same if we replace \mathbf{W} as \mathbf{WR} for any rotation (orthogonal matrix) \mathbf{R} .

High dimensional PCA: Computational difficulty

Now suppose the dimension of each data point is high (larger the the sample size). The standard PCA will solve the following eigen-vector equation:

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

Where here $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the (centered) data matrix. Since we know there are (at least) $p - n - 1$ eigen-values are 0, we could pre-multiply both sides by \mathbf{X} to give

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} \mathbf{X} \mathbf{v}_i = \lambda_i \mathbf{X} \mathbf{v}_i$$

High dimensional PCA: Computational difficulty cont'd

Define $\mathbf{u}_i = \mathbf{X}\mathbf{v}_i$, we obtain

$$\frac{1}{n} \mathbf{X} \mathbf{X}^T \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

Note that solving this eigen-vector problem only needs $O(n^3)$ instead of $O(p^3)$ for solving the standard one. Besides, there is a relationship between \mathbf{v}_i and \mathbf{u}_i :

$$\mathbf{v}_i = \frac{1}{(n\lambda_i)^{1/2}} \mathbf{X}^T \mathbf{u}_i$$

In fact, the technique used above can be derived from a higher point of view known as **Woodbury identity** (see next slide). Such techniques are common in machine learning to reduce the computation of algorithms.

High dimensional PCA: Woodbury identity

Lemma 19 (Woodbury)

Provided all matrix multiplications are compatible, we have:

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

Proof.

Left as an exercise.

Hint 1: ETS the inverse of LHS multiplies RHS is the identity.

Hint 2: $\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} = \mathbf{B}\mathbf{D}^{-1}\mathbf{D}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}$



High dimensional PCA: Another useful matrix lemma

Lemma 20 (Another Useful Formula)

Provided all matrix multiplications are compatible, we have:

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1}$$

It can be verified by right multiplying both sides by $(\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})$.

Example 21 (Dual Representation of Ridge Regression)

For ridge regression, we have

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_n)^{-1} \mathbf{X}^T \mathbf{y}$$

By the above lemma, the dual representation is

$$\hat{\beta} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_p)^{-1} \mathbf{y}$$

Exact Distribution of the Eigenvalues

Before we dive into the high dimensional case where $p \gg n$, let's recall a classic result in multivariate analysis.

Definition 22 (Wishart Distribution)

Suppose $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, $i = 1, \dots, n$ and $n > p$, then

$$\mathbf{X}\mathbf{X}^T \sim \mathcal{W}_p(n, \mathbf{I}_p)$$

where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ and $\mathcal{W}_p(n, \mathbf{I}_p)$ is called **the Wishart distribution** with parameter n and \mathbf{I}_p , and the pdf is

$$p_{\mathbf{X}}(\mathbf{X}) = \frac{(\det \mathbf{X})^{(n-p-1)/2} e^{-\text{Tr}(\mathbf{X})/2}}{2^{\frac{np}{2}} \Gamma_p(\frac{n}{2})}$$

Exact Distribution of the Eigenvalues cont'd

Theorem 23 (Dist. of the Eigenvalues of Wishart Dist.)

Now suppose $\mathbf{XX}^T \sim \mathcal{W}_p(n, \mathbf{I}_p)$, then the exact distribution of the eigenvalues of this random matrix has the form:

$$p(\lambda_1, \dots, \lambda_p) = c_{p,n} \prod_{j=1}^p [w(\lambda_j)]^{1/2} \prod_{j < k} (\lambda_j - \lambda_k)$$

where $w(x) = x^{n-r-1} e^{-x}$, $c_{p,n}$ is the normalizing constant and $\prod_{j < k} (\lambda_j - \lambda_k)$ is known as the **Vandermonde determinant**

Proof: Hsu, 1939.

e.g. Chapter 11 section 5 and chapter 13 section 3 of Anderson, 1984. □

Empirical Spectral Distribution

If p is larger than n , the previous theorem fails. However, random matrix theory gives us the asymptotic results:

Theorem 24 (Marcenko-Pastur Law)

Suppose $\mathbf{X}\mathbf{X}^T \sim \mathcal{W}_p(n, \mathbf{I}_p)$. Define the **empirical spectral distribution**

$$G_p(k) = \frac{1}{p} \# \{ \hat{\lambda}_j \leq k \}$$

If $\frac{p}{n} \rightarrow \gamma \in (0, \infty)$, then, $G_r(k) \rightarrow G(k)$ almost surely, where the limiting distribution $G(k)$ has density $g(k) = G'(k)$:

$$g(k) = \frac{\sqrt{(b_+ - k)(k - b_-)}}{2\pi\gamma k}, \quad b_{\pm} = (1 \pm \sqrt{\gamma})^2$$

This is the famous **Marcenko-Pastur Law** in statistical physics.

Sparse PCA

- Due to the previous theorem, we know that in high dimensional case, neither the eigen-values nor eigen-vectors of sample covariance matrix will be consistent. Therefore, we need to modify the standard PCA in such case. This is the so-called **Sparse PCA**.
- This part is quite technical (concentration inequalities) so we can omit it.

END !

Next time:

Radon-Nikodym Theorem,
Conditional Expectation
and Martingale !