

Lecture 3a

Principal Component Analysis

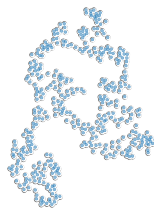
Recap Weeks 1-2

Data science

- ▶ Data science studies how one can systematically extract insights from real-world data.
- ▶ Data science addresses both the statistical and technical aspects of such process (e.g. robustness, scalability).

Visualizations:

- ▶ Visualizations can effectively convey large data distributions to the user.
- ▶ Low-dimensional embedding techniques such as MDS or T-SNE enable the visualization of large high-dimensional datasets.
- ▶ Sometimes, embeddings do not faithfully convey certain aspects of the data (e.g. distances not preserved in high dimensions, distortion of local geometry, spurious cluster structures, etc.).
- ▶ Visualizations do not give *quantifiable* insights.



Today's Lecture

- ▶ Data Dispersion
- ▶ Principal Component Analysis (PCA)
 - ▶ Dispersion maximization view
 - ▶ Error minimization view
 - ▶ PCA as a constrained optimization problem
- ▶ Lagrange Multipliers
 - ▶ Framework for solving constrained optimization problems
 - ▶ Practical examples
- ▶ Application of Lagrange Multipliers to PCA
 - ▶ Reformulation of PCA as an eigenvalue problem
- ▶ PCA explains Data Dispersion

Measuring Data Dispersion

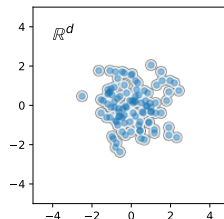
Dispersion is an important property of the data, and indicates how much variation there is in some dataset $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$. There are various possible measures of dispersion.

Examples:

- ▶ Number of distinct data points.
- ▶ Radius of minimum enclosing sphere.
- ▶ Average square Euclidean distance from the dataset mean \mathbf{m} :

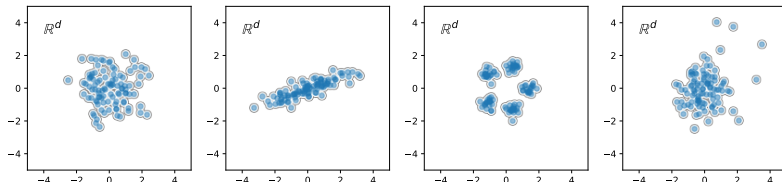
$$s(X) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{m}\|^2$$

The measure $s(X)$ can be seen as a generalization of variance from univariate to multivariate data. Other generalizations are possible.



Explaining Data Dispersion

today's lecture



Observation:

- ▶ All distributions above have the same dispersion $s(X)$. Yet, they are very different.
- ▶ Often, it is not only desirable to quantify the overall dispersion but also to explain the structure of dispersion.
- ▶ **Today's lecture** focuses on a specific type of dispersion which can be described in terms of directions in input space (aka. Principal Component Analysis).

Part 1

Principal Component Analysis

Hotelling. J Educational Psychology, 24:498–520, 1933.
Pearson. Phil. Mag. 2, 6, 1901.

Principal Component Analysis

Preliminaries:

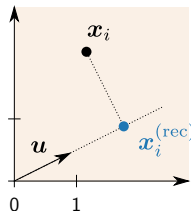
- ▶ Let $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be a dataset, where d is the number of input features and N is the number of data points (or sample size).
- ▶ Let $\mathbf{u} \in \mathbb{R}^d$ be a vector of same dimensions, that represents some direction in input space, and that is constrained to be of norm 1, i.e. $\|\mathbf{u}\| = 1$.
- ▶ Data points can be projected onto this direction by computing the dot product.

$$\forall_{i=1}^N : z_i = \mathbf{u}^\top \mathbf{x}_i$$

- ▶ These projections can be backprojected on the input space, by multiplying again by the direction \mathbf{u} .

$$\forall_{i=1}^N : \mathbf{x}_i^{(\text{rec})} = \underbrace{\mathbf{u} \mathbf{u}^\top}_{z_i} \mathbf{x}_i$$

\mathbb{R}^d (original space)



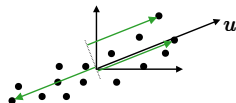
\mathbb{R} (projected space)



PCA: Two Formulations

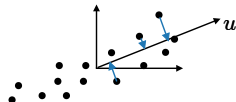
Dispersion Maximization:

- Find a projection $z = \mathbf{u}^\top \mathbf{x}$ of the data under which the dispersion (variance) is maximized.



Error Minimization:

- Find the direction that minimizes the reconstruction error (MSE) between the original data point \mathbf{x} and its backprojection $\mathbf{x}^{(\text{rec})} = \mathbf{u}\mathbf{u}^\top \mathbf{x}$.



The two views coincide (Pearson 1901).

PCA as Dispersion Maximization

Objective:

- ▶ Find a direction \mathbf{u} (with $\|\mathbf{u}\| = 1$) so that the data projected onto this direction (i.e. z_1, \dots, z_N) has **maximum variance**.
- ▶ This can be cast into an optimization problem:

$$\arg \max_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\underbrace{\mathbf{u}^\top \mathbf{x}_i}_{z_i} - \tilde{m})^2 \right]$$

where $\tilde{m} = \frac{1}{N} \sum_{i=1}^N z_i$ is the dataset mean in projected space.

- ▶ Making sure we first center the data, i.e. $\mathbf{m} = \mathbf{0}$, it implies that $\tilde{m} = 0$, and the optimization problem simplifies to:

$$\arg \max_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{u}^\top \mathbf{x}_i)^2 \right]$$

PCA as Error Minimization

- ▶ Find a direction \mathbf{u} (with $\|\mathbf{u}\| = 1$) so that the data projected on the corresponding subspace best reconstructs the original data, specifically, has **minimal squared distance** to the original data.
- ▶ Recall that the reconstruction model is given by:

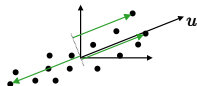
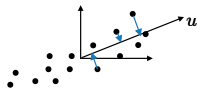
$$\mathbf{x}^{(\text{rec})} = \mathbf{u}\mathbf{u}^\top \mathbf{x}$$

- ▶ This can be cast into the optimization problem:

$$\arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \underbrace{\mathbf{u}\mathbf{u}^\top \mathbf{x}_i}_{\mathbf{x}_i^{(\text{rec})}} \right\|^2 \right]$$

Connecting the two Approaches

$$\begin{aligned} & \arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{u} \mathbf{u}^\top \mathbf{x}_i\|^2 \right] \\ &= \arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{u} \mathbf{u}^\top \mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{u} \mathbf{u}^\top \mathbf{x}_i) \right] \\ &= \arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N \underbrace{\mathbf{x}_i^\top \mathbf{x}_i}_{\text{cst.}} - 2\mathbf{x}_i^\top \mathbf{u} \mathbf{u}^\top \mathbf{x}_i + (\mathbf{u} \mathbf{u}^\top \mathbf{x}_i)^\top (\mathbf{u} \mathbf{u}^\top \mathbf{x}_i) \right] \\ &= \arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N -2(\mathbf{x}_i^\top \mathbf{u})^2 + \mathbf{x}_i^\top \underbrace{\mathbf{u} \mathbf{u}^\top \mathbf{u}}_{=1} \mathbf{u}^\top \mathbf{x}_i \right] \\ &= \arg \min_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N -2(\mathbf{x}_i^\top \mathbf{u})^2 + (\mathbf{x}_i^\top \mathbf{u})^2 \right] \\ &= \arg \max_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{u})^2 \right] \end{aligned}$$



Both approaches give the same result!

Finding Principal Components

- ▶ We first recall the PCA optimization problem of slide 8:

$$\arg \max_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{u})^2 \right] \quad \text{s.t. } \|\mathbf{u}\| = 1$$

and observe it can be rewritten as:

$$\begin{aligned} &= \arg \max_{\mathbf{u}} \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{u}^\top \mathbf{x}_i)(\mathbf{x}_i^\top \mathbf{u}) \right] \\ &= \arg \max_{\mathbf{u}} \left[\mathbf{u}^\top \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \right)}_{\Sigma} \mathbf{u} \right] \quad \text{s.t. } \|\mathbf{u}\| = 1 \end{aligned}$$

where Σ is the covariance matrix (remember that the data is assumed to be centered). The covariance matrix does not depend on \mathbf{u} and can be precomputed.

- ▶ Optimization problems with equality constraints can be solved using the method of Lagrange Multipliers.

Part 2

Method of Lagrange Multipliers

Method of Lagrange Multipliers

The method of Lagrange multipliers is a general framework for finding solutions of constrained optimization problems of the type:

$$\arg \max_{\theta} f(\theta) \quad \text{subject to} \quad g(\theta) = 0$$

It consists of applying the following two steps:

- **Step 1:** Construct the 'Lagrangian':

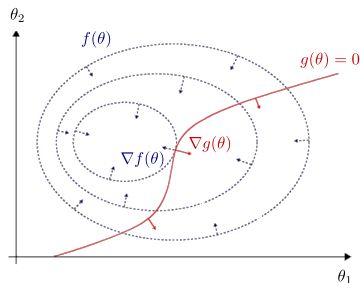
$$\mathcal{L}(\theta, \lambda) = f(\theta) + \lambda \cdot g(\theta)$$

where λ is called the Lagrange multiplier.

- **Step 2:** Solve the equation

$$\nabla \mathcal{L}(\theta, \lambda) = 0$$

which is a necessary condition for the solution.

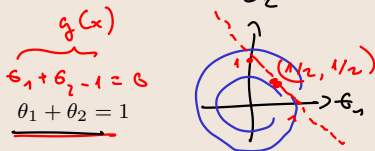


Intuition for step 2: the equation includes the equation $\nabla f(\theta) = -\lambda \nabla g(\theta)$, i.e. the gradient of objective and constraint are aligned, but point in opposite directions (cf. 2d plot).

Method of Lagrange Multipliers

Example 1: Solve the optimization problem:

$$\arg \max_{\theta} \underbrace{[1 - (\theta_1^2 + \theta_2^2)]}_{g(\theta)} \quad \text{s.t.} \quad \underline{\theta_1 + \theta_2 = 1}$$



$$\mathcal{L}(\theta, \lambda) = 1 - \theta_1^2 - \theta_2^2 + \lambda \cdot (\theta_1 + \theta_2 - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = -2\theta_1 + \lambda \cdot 1 = 0 \Rightarrow \underline{2\theta_1 = \lambda}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = -2\theta_2 + \lambda \cdot 1 = 0 \Rightarrow \underline{2\theta_2 = \lambda}$$
$$\underbrace{\frac{\lambda}{2}}_{\theta_1} + \underbrace{\frac{\lambda}{2}}_{\theta_2} = 1 \Rightarrow \lambda = 1$$

$$\theta_1 = \frac{\lambda}{2} = \frac{1}{2} \quad \theta_2 = \frac{\lambda}{2} = \frac{1}{2}$$

Method of Lagrange Multipliers

Example 2: Let $\theta, m, b \in \mathbb{R}^d$ and $\|b\| = 1$. Solve the optimization problem:

$$\arg \min_{\theta} \underbrace{\|\theta - m\|}_{f(\theta)} \quad \text{s.t.} \quad \underbrace{\theta^T b = 0}_{g(\theta)}$$



$$\mathcal{L}(\theta, \lambda) = \frac{1}{2} \|\theta - m\|^2 + \lambda \theta^T b$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = (\theta - m) + \lambda b = 0$$

$$\begin{pmatrix} b^T \theta - b^T m + \lambda b^T b = 0 \\ 0 - b^T m + \lambda 1 = 0 \Rightarrow \lambda = b^T m \end{pmatrix}$$

$$(\theta - m) + b b^T m = 0$$

$$\theta = m - b b^T m$$

Part 3

Back to PCA

The Solution of PCA

Recall that our PCA optimization problem has the form:

$$\arg \max_{\mathbf{u}} [\mathbf{u}^\top \Sigma \mathbf{u}] \quad \text{s.t.} \quad \|\mathbf{u}\| = 1$$

We rewrite the constraint as $\|\mathbf{u}\|^2 = 1$ and look for a solution by applying the method of Lagrange multipliers.

- **Step 1:** Build the Lagrangian

$$\mathcal{L}(\mathbf{u}, \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} + \lambda \cdot (1 - \|\mathbf{u}\|^2)$$

- **Step 2:** Set gradient of Lagrangian to zero:

$$\nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \lambda) = \mathbf{0} \quad \Rightarrow \quad \Sigma \mathbf{u} = \lambda \mathbf{u}$$

$$\nabla_{\lambda} \mathcal{L}(\mathbf{u}, \lambda) = 0 \quad \Rightarrow \quad \|\mathbf{u}\|^2 = 1$$

! PCA solution is an eigenvector of Σ

- Note that the eigenvector is determined up to a sign flip (i.e. in practice, our interpretation of the PCA result should not depend on that sign flip).

Which Eigenvector?

- ▶ Start with the eigenvalue problem

$$\Sigma \mathbf{u} = \lambda \mathbf{u}$$

- ▶ Multiply the eigenvalue problem by \mathbf{u}^\top on both sides, and identify the terms of the resulting equation:

$$\underbrace{\mathbf{u}^\top \Sigma \mathbf{u}}_{\text{objective}} = \lambda \underbrace{\mathbf{u}^\top \mathbf{u}}_1$$

- ▶ In other words, for the objective to be maximized, we should choose the eigenvector in a way that the corresponding eigenvalue λ is maximum. In other words, one should choose the leading eigenvector.

! PCA solution is the leading eigenvector of Σ

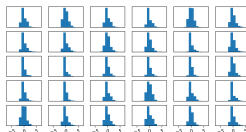
Application: Breast Tumor Analysis



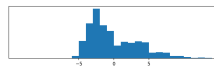
Breast Cancer Wisconsin (Diagnostic)

- 569 instances
- 30 features
- meta-data (benign/
malignant)

statistics of each features

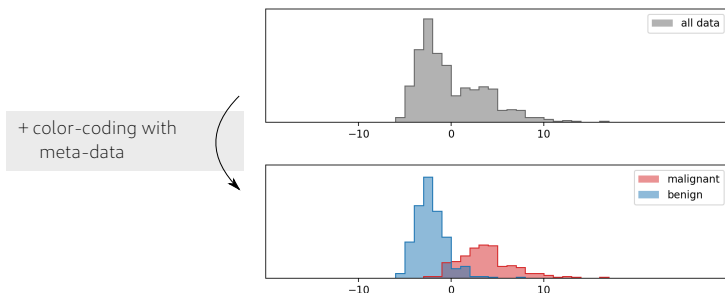


statistics in PCA space



- PCA conveys in a more concise way the main variations in a dataset. Instead of looking at as many histograms as there are input features, one only needs to look at a single one.

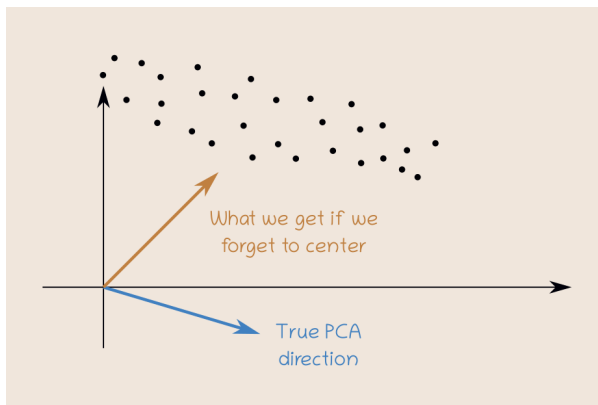
Application: Breast Tumor Analysis



- PCA analysis can be color-coded based on meta-data. Here, for example, one can see that the two tumor types are well-differentiated with the principal component, i.e. features collectively represent well variations associated to tumor malignancy.

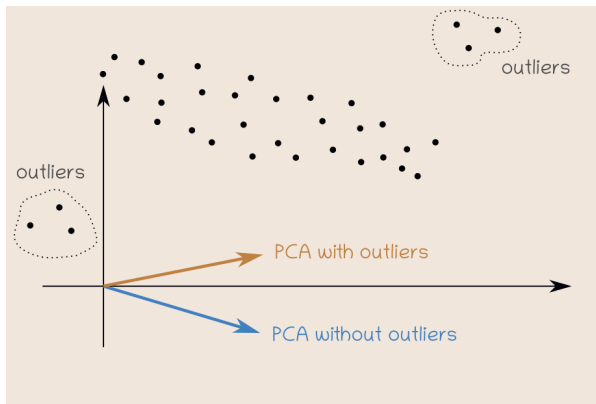
Further Remarks about PCA

Further Remarks



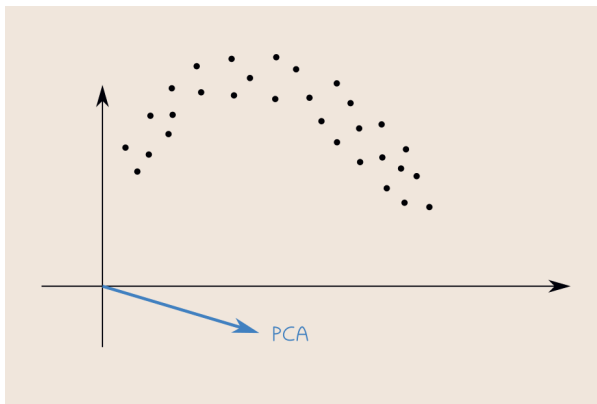
- Don't forget to center the data before applying PCA.

Further Remarks



- ▶ PCA is not very robust to outliers. For PCA to be meaningful, data needs to be first cleaned from outliers. Alternatively, robust variants of PCA need to be used.

Further Remarks



- ▶ PCA does not describe well the data when it is strongly non-Gaussian. It fails to account for the fact that data may vary locally in different directions.

Summary

Summary

- ▶ **Principal Component Analysis** is a dimensionality reduction technique that implements [Pearson 1901]'s principle of minimizing *noise* and maximizing *signal*. (It does both simultaneously!).
- ▶ PCA reduces high-dimensional data to one dimension (the principal component). The latter represents the **main trend** in the data, which often captures interesting information such as class membership.
- ▶ PCA is a well-known method, available in most ML libraries. No hyperparameters to select. It is furthermore is an **exact method**. It always finds the optimum of the objective (defined up to a sign flip).