

Dimensionality Reduction

Prof. Marco Alvarez, Computer Science
University of Rhode Island

High dimensional data

- ▶ Prevalent in many areas of machine learning and data science
- ▶ Examples:
 - image data, for instance a single 224×224 RGB image has $>150k$ dimensions
 - text data, in NLP applications, transformers typically embed tokens into 768 or more dimensions
 - genomic data, gene expression data often has thousands of dimensions
 - audio signals, especially when converted to spectrograms
 - time series from sensor data or financial data
 - network traffic data
 - ...

Dimensionality reduction

- ▶ Fundamentally ...
 - unsupervised learning algorithms for extracting latent structure (potentially **low dimensional**) from **high-dimensional** data
 - can range from simple feature selection to complex nonlinear transformations
 - allows working with more compact data representation, ideally without losing information
- ▶ Examples
 - PCA, kernel PCA, t-SNE, autoencoders, matrix factorization
- ▶ Given $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ with $\mathbf{x}_i \in \mathbb{R}^d$, find a representation $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ with $\mathbf{z}_i \in \mathbb{R}^{d'}$, with $d' < d$
 - certain properties should be preserved (e.g., variance, distances, neighborhood structure)

Dimensionality reduction

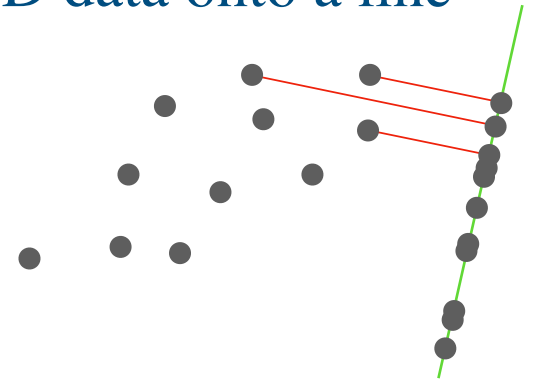
- ▶ Why?
 - visualization (2D or 3D)
 - preprocessing data before machine learning
 - focusing on important features/patterns
 - more efficient training
 - removing noise and redundant information
 - data compression

$$d' < d \quad \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{n \times d} \Rightarrow \begin{bmatrix} \mathbf{z}_1^T \\ \vdots \\ \mathbf{z}_n^T \end{bmatrix}_{n \times d'}$$

Preliminaries

Project 2D data onto a line

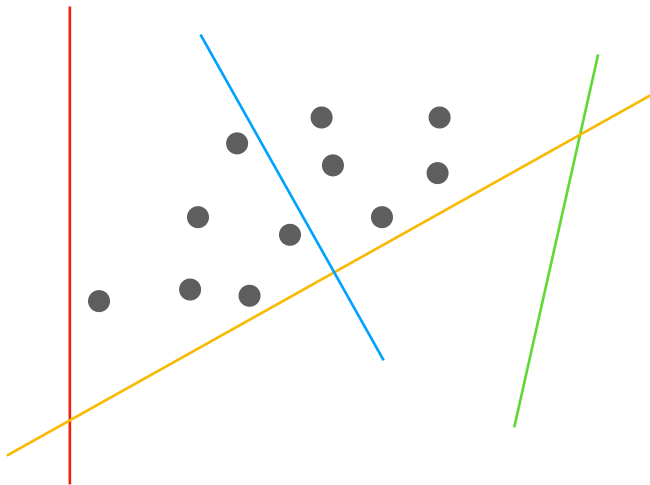
$$\begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{11 \times 2}$$



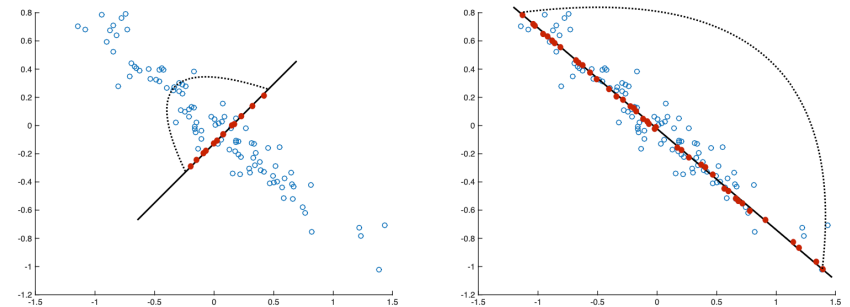
$$\begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{11 \times 1}$$



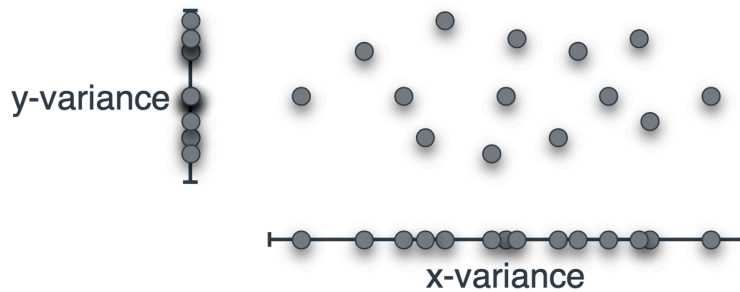
Which line is better? which is worst?



Key idea: maximize variance



Variance



“biased” estimator default in numpy

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Figure credit: <https://serrano.academy/unsupervised-learning/>

Different data, same variance



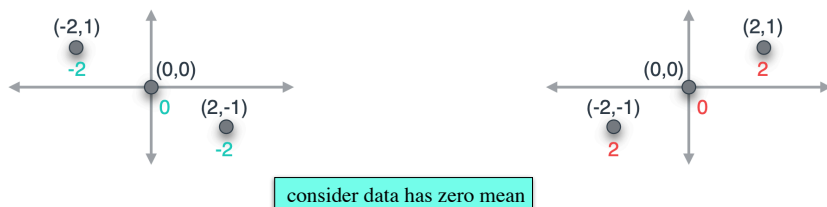
$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$V_x = \frac{2^2 + 0^2 + 2^2}{3} = \frac{8}{3}$$

$$V_y = \frac{1^2 + 0^2 + 1^2}{3} = \frac{2}{3}$$

Figure credit: <https://serrano.academy/unsupervised-learning/>

Covariance



consider data has zero mean

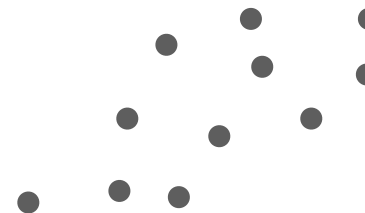
$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Cov}_x = \frac{(-2) + 0 + (-2)}{3} = \frac{-4}{3} \text{ negative covariance}$$

$$\text{Cov}_y = \frac{2 + 0 + 2}{3} = \frac{4}{3} \text{ positive covariance}$$

Figure credit: <https://serrano.academy/unsupervised-learning/>

Covariance matrix



$$X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}_{11 \times 2}$$

$$\Sigma = \begin{bmatrix} \text{Var}(i) & \text{Cov}(i, j) \\ \text{Cov}(j, i) & \text{Var}(j) \end{bmatrix}_{2 \times 2}$$

Every element (i, j) of the covariance matrix Σ represents the covariance between feature i (column) and feature j from the data matrix X

Eigenvectors and eigenvalues

- ▶ The decomposition of a **square matrix A** into **eigenvalues and eigenvectors** is known as **eigen decomposition**
 - for **real symmetric matrices** eigenvectors can be chosen real and orthonormal

$$A = V\Lambda V^T \quad A\mathbf{v} = \lambda\mathbf{v}$$

columns of V are the eigenvectors of A and Λ is a diagonal matrix whose entries are the eigenvalues of A

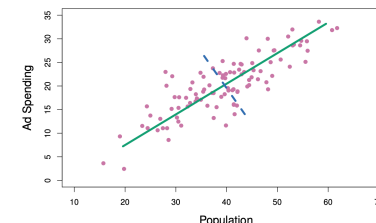
Principal Component Analysis (PCA)

PCA

- ▶ Proposed by Pearson (1901) and Hotelling (1933)
 - used for data compression, visualization, and identifying patterns/structures
- ▶ Key idea
 - high-dimensional data often lies on a lower-dimensional subspace
 - overcomplete (redundant dimensions)
 - dimensions frequently correlated

Goal of PCA

- ▶ Find projections of the data onto directions that maximize variance
 - directions are **orthogonal** to each other, and have lower intrinsic dimensionality
 - projected data points are as similar to the originals as possible



<https://www.dataschool.io/15-hours-of-expert-machine-learning-videos/>

PCA approach

- ▶ **Input:** data matrix $X_{d \times n}$
 - **center** the data (subtract the mean)
 - calculate the **covariance** matrix $\Sigma = \frac{1}{n}XX^T$
 - compute **eigendecomposition** $V\Lambda V^T$ of the covariance matrix
 - **sort** the eigenvectors by eigenvalues in decreasing order

- ▶ **Output**

- sorted orthonormal eigenvectors V and eigenvalues Λ
 - use V as the projection matrix

eigenvectors can then be used for projecting the data into lower dimensions (XV)

Remarks

- ▶ The **larger** the eigenvalue, the **more important** the corresponding eigenvector
 - that's why we **sort** eigenvalues (and corresponding eigenvectors) in **decreasing order**
- ▶ All eigenvalues of a positive semidefinite matrix are **non-negative**
 - **covariance matrix** is always symmetric and p.s.d.
- ▶ For dimensionality reduction, we can ignore eigenvectors associated with smaller eigenvalues

Explained variance

- ▶ Each eigenvalue corresponds to the amount of variance explained by its associated eigenvector
 - **explained variance** is often presented as a **percentage**, i.e., eigenvalues divided by the total sum of eigenvalues
- ▶ The sum of percentages of the top-k principal components is usually referred to as the **“cumulative explained variance”**
 - often used to select how many components to keep for a reduced dataset

Explained variance

PC	Eigenvalue	Variance (%)	Cumulative Variance
1st	23.31800	59.072%	59.072%
2nd	7.01200	17.764%	76.835%
3rd	4.61800	11.699%	88.534%
4th	1.98100	5.018%	93.553%
5th	1.00100	2.536%	96.089%
6th	0.82100	2.080%	98.168%
7th	0.64100	1.624%	99.792%
8th	0.03100	0.079%	99.871%
9th	0.02900	0.073%	99.944%
10th	0.02200	0.056%	100.000%

PCA Notebooks

<https://colab.research.google.com/drive/1MzPdVsJi8gUxhwiXFcKA8RJsw8DY1OhE>

<https://colab.research.google.com/drive/1r7JPdmmWS11yl2WVOMi0GMhlDM9s37GI>