

Principal Component Analysis (PCA)

Lecture Outline

- Motivations
- Applications
- History
- Basic Concept of PCA
- Geometric Rationale of PCA
- Assumptions of PCA

PRINCIPAL COMPONENT ANALYSIS (PCA)

MOTIVATIONS

Curse of Dimensionality — As the number of features used in the representation increases so does:

The number of samples required for training

The running time of the algorithm

The possibility of overfitting

PRINCIPAL COMPONENT ANALYSIS (PCA)

Why Principal Component Analysis?

We are interested in using machine learning to make decisions regarding phenomena that cannot be directly observed.

Our algorithms should operate on underlying latent factors rather than the raw observed data

We want machine learning to exploit hidden relationships in the data when making decisions.

PRINCIPAL COMPONENT ANALYSIS (PCA)

Applications

- Uses:
 - Data Visualization
 - Data Reduction
 - Data Classification
 - Trend Analysis
 - Noise Reduction
- Examples:
 - How many unique “sub-sets” are in the sample?
 - How are they similar / different?
 - What are the underlying factors that influence the samples?
 - Which time / temporal trends are (anti)correlated?
 - Which measurements are needed to differentiate?
 - How to best present what is “interesting”?

PRINCIPAL COMPONENT ANALYSIS (PCA)

History

Probably the most widely-used and well known of the data reduction methods

Form of unsupervised learning

Invented by Pearson (1901) and Hotelling (1933)

Also known as Principal Factor Analysis

PRINCIPAL COMPONENT ANALYSIS (PCA)

Basic Concepts

Principles

- Linear projection method to reduce the number of features
- Maps the data into a lower dimensional space
- Transfers a set of correlated variables into a new set of uncorrelated variables

Properties

- Can view as a rotation of existing axes to new positions in the original feature space.
- New axes are orthogonal and represent directions with maximum variability.
- The new variable/dimensions are linear combinations of original ones.

PRINCIPAL COMPONENT ANALYSIS (PCA)

Geometric Rationale

PCA takes a data matrix of n samples by d variables (features), which may be correlated, and summarizes it by uncorrelated axes (principal components or principal axes) that are linear combinations of the original d variables.

The first k components display as much as possible of the variation among the samples.

Samples are represented as a cloud of n points in a multidimensional space with an axis for each of the d variables (features).

PRINCIPAL COMPONENT ANALYSIS (PCA)

Geometric Rationale

The **centroid** of the points is defined by the mean of each feature.

The **variance** of each feature is the average squared deviation of its ***n*** values around the mean of that feature:

$$V_i = \frac{1}{n-1} \sum_{m=1}^n (X_{im} - \bar{X}_i)^2$$

PRINCIPAL COMPONENT ANALYSIS (PCA)

Geometric Rationale

The degree to which the features are linearly correlated is represented by their **covariances**.

$$C_{ij} = \frac{1}{n-1} \sum_{m=1}^n (X_{im} - \bar{X}_i)(X_{jm} - \bar{X}_j)^T$$

The diagram illustrates the formula for covariance C_{ij} with arrows pointing to each component:

- C_{ij}** : Covariance of features i and j
- $\frac{1}{n-1}$** : Sum over all n samples
- $\sum_{m=1}^n$** : Sum over all n samples
- X_{im}** : Value of feature i in sample m
- \bar{X}_i** : Mean of feature i
- X_{jm}** : Value of feature j in sample m
- \bar{X}_j** : Mean of feature j

PRINCIPAL COMPONENT ANALYSIS (PCA)

Assumptions

PCA assumes that the relationships among features are **linear**.

If the structure in the data is **non-linear** (the cloud of points twists and curves its way the ***d***-dimensional space), the principal axes will not be an efficient and informative summary of the data.

Questions?

