Slides adapted from Prof Carpuat and Duraiswami

CMSC 422 introduction to Machine Learning Lecture 17 Unsupervised Learning – Principal Component Analysis

Furong Huang / furongh@cs.umd.edu





UNIVERSITY OF MARYLAND

Unsupervised Learning

Discovering hidden structure in data

What algorithms do we know for unsupervised learning?

K-Means Clustering

Today: how can we learn better representations of our data points?



Dimensionality Reduction

Goal: extract hidden lower-dimensional structure from high dimensional datasets

Why? To visualize data more easily To remove noise in data To lower resource requirements for storing/processing data To improve classification/clustering





Examples of data points in D dimensional space that can be effectively represented in a ddimensional subspace (d < D)



Principal Component Analysis

Goal: Find a **projection** of the data onto directions that **maximize variance** of the original data set

Intuition: those are directions in which most information is encoded

Definition: **Principal Components** are orthogonal directions that capture most of the variance in the data



PCA: finding principal components

1st PC



Projection of data points along 1st PC discriminates data most along any one direction

2nd PC

next orthogonal direction of greatest variability

And so on...



PCA: notation

Data points

Represented by matrix X of size NxD X_i is the i-th row, i.e., the i-th example X_{ij} is the value of j-th feature for example i Let's assume data is centered, i.e., $\sum_i X_i = \vec{0}$

Principal components are d vectors: $v_1, v_2, ..., v_d$ $v_i^{\mathsf{T}} v_j = 0, i \neq j \text{ and } v_i^{\mathsf{T}} v_i = 1, \forall i$

The sample variance data projected on vector v is $\sum_{i=1}^{n} (x_i^T v)^2 = (Xv)^T (Xv)$





Finding vector that maximizes sample variance of projected data: $argmax_v v^T X^T X v$ such that $v^T v = 1$

A constrained optimization problem(method of Lagrange multipliers)

-Lagrangian folds constraint into objective: $argmax_v v^T X^T X v - \lambda(v^T v - 1)$

•Solutions are vectors v such that $X^T X v = \lambda v$

•i.e. eigenvectors of $X^T X$ (sample covariance matrix)



Lagrange Multiplier

- A strategy for finding the local maxima and minima of a function subject to equality constraints
- Consider the optimization problem

 $\max_{v} f(v)$

Subject to g(v) = 0

- Introduce a new variable λ called a Lagrange multiplier
- Study the Lagrange (Lagrangian) function $\mathcal{L}(v,\lambda) = f(v) \lambda g(v)$



Relationship between PCA and Eigen Value Decomposition

$$argmax_{v} v^{T} X^{T} X v - \lambda (v^{T} v - 1) \qquad (*)$$

where X is the N by D data matrix

- For this optimization problem, taking derivative with respect to v, set the gradient to 0 to solve for the stationary point $X^T X v - \lambda v = 0$
- Therefore the eigenvector of covariance matrix X^T X is a stationary point of the optimization problem (*).



PCA formally

- The eigenvalue λ denotes the amount of variability captured along dimension vSample variance of projection $v^T X^T X v = \lambda$
- If we rank eigenvalues from large to small The 1st PC is the eigenvector of $X^T X$ associated with largest eigenvalue The 2nd PC is the eigenvector of $X^T X$ associated with 2nd largest eigenvalue



Alternative interpretation of PCA

PCA finds vectors v such that projection on to these vectors minimizes reconstruction error

$$\frac{1}{n}\sum_{i=1}^{n} \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i)\mathbf{v}\|^2$$





Resulting PCA algorithm

Algorithm 36 PCA(X, K)

 1: $\mu \leftarrow MEAN(X)$ // compute data mean for centering

 2: $\mathbf{D} \leftarrow (\mathbf{X} - \mu \mathbf{1}^{\top})^{\top} (\mathbf{X} - \mu \mathbf{1}^{\top})$ // compute covariance, 1 is a vector of ones

 3: $\{\lambda_k, u_k\} \leftarrow$ top K eigenvalues/eigenvectors of D
 // project data using U



How to choose the hyperparameter K?

i.e. the number of dimensions



We can ignore the components of smaller significance



An example: Eigenfaces





PCA pros and cons

Pros

Eigenvector method No tuning of the parameters No local optima

Cons

Only based on covariance (2nd order statistics) Limited to linear projections



What you should know

Principal Components Analysis

Goal: Find a **projection** of the data onto directions that **maximize variance** of the original data set

PCA optimization objectives and resulting algorithm

Why this is useful!





Furong Huang 3251 A.V. Williams, College Park, MD 20740 301.405.8010 / furongh@cs.umd.edu