

# CS 584 – MACHINE LEARNING

## TOPIC: DIMENSIONALITY REDUCTION



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

# MOTIVATION

- Remove useless features
- Learn a low dimensional representation
- Visualization

# APPROACHES

- Feature selection
- Feature extraction

# FEATURE SELECTION

- Select a subset of the features
- Several approaches
  - Univariate feature selection
    - Mutual information, Chi2, ...
  - Recursive feature elimination
    - Using an external estimator, recursively remove least useful features
  - Model-based feature selection
    - L1 regularization, decision trees, ...
  - Sequential feature selection
    - Can use any model and performance metric

# FEATURE EXTRACTION

- Feature extraction, projection, latent representation learning, manifold learning, ...
- Given an input, project it into another dimension (often lower)
- Unsupervised
  - Principal component analysis, isomap, t-SNE, deep learning, ...
- Supervised
  - Linear discriminant analysis, deep learning, ...

# WE'LL COVER

- Principal component analysis
- Autoencoder

# PRINCIPAL COMPONENT ANALYSIS

- Given a dataset  $\mathbf{x}$  with  $d$  dimensions, project it into  $k$  dimensions ( $k < d$ ) with minimum loss of information

- $\mathbf{z} = \mathbf{w}^T \mathbf{x}$  —  $d$   $\begin{bmatrix} \square & \square & \dots & \square \end{bmatrix}$

$$\mathbf{z} \quad k \times \begin{bmatrix} \square & \dots & \square \end{bmatrix}$$

$\mathbf{w}^T: k \times d$

- Principal component analysis (PCA) maximizes the variance in the projected space so that objects are spread out

- $\operatorname{argmax}_{\mathbf{w}} \operatorname{Var}(\mathbf{z})$

# PCA OBJECTIVE

- $\mathbf{z} = \mathbf{w}^T \mathbf{x}$
- $\operatorname{argmax}_{\mathbf{w}} \operatorname{Var}(\mathbf{z})$
- Reminder:  $\operatorname{Var}(aX + b) = a^2 \operatorname{Var}(X)$
- $\operatorname{Var}(\mathbf{z}) = \mathbf{w}^T \operatorname{Var}(\mathbf{x}) \mathbf{w}$
- We can trivially maximize variance by multiplying  $\mathbf{w}$  by a large constant; hence, we enforce that  $\mathbf{w}$  is unit length
  - $\mathbf{w} \mathbf{w}^T = 1$
- Objective
  - maximize  $\mathbf{w}^T \operatorname{Var}(\mathbf{x}) \mathbf{w}$  subject to  $\mathbf{w} \mathbf{w}^T = 1$

$$\operatorname{Var}(\mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \cdot \operatorname{Var}(\mathbf{x}) \cdot \mathbf{w}$$



## BACKGROUND — CONSTRAINED OPTIMIZATION

Find  $\boldsymbol{\theta}$   
maximizing  $f(\boldsymbol{\theta})$   
subject to

$$c_1(\boldsymbol{\theta}) = 0$$

...

$$c_m(\boldsymbol{\theta}) = 0$$

$$\begin{aligned} & \mathbf{w} \\ & \mathbf{w}^T \mathbf{V}_a(\mathbf{x}) \mathbf{w} \\ & \mathbf{w} \mathbf{w}^T - \mathbf{I} = 0 \end{aligned}$$

Form the Lagrangian:

$$F(\boldsymbol{\theta}, \boldsymbol{\lambda}) = f(\boldsymbol{\theta}) - \sum_{j=1}^m \lambda_j c_j(\boldsymbol{\theta})$$

## EXAMPLE

- Maximize  $xy$  subject to  $x + y = 10$

$$\begin{aligned} & x + y - 10 = 0 \\ & x \cdot y - \lambda(x + y - 10) \\ & \frac{\partial}{\partial x} \quad y - \lambda = 0 \Rightarrow y = \lambda \\ & \frac{\partial}{\partial y} \quad x - \lambda = 0 \Rightarrow x = \lambda \\ & \frac{\partial}{\partial \lambda} \quad x + y - 10 = 0 \\ & \quad \lambda + \lambda - 10 = 0 \Rightarrow \lambda = 5 \end{aligned}$$

$= 5$   
 $= 5$

# PCA WITH ONE DIMENSION

- maximize  $Var(z_1) = w_1^T Var(\mathbf{x}) w_1$  subject to  $w_1^T w_1 = 1$
- Let the variance of  $\mathbf{x}$  be the covariance matrix  $\Sigma$
- Objective
  - Maximize  $w_1^T \Sigma w_1$  subject to  $w_1^T w_1 = 1$
- Lagrange
  - maximize  $w_1^T \Sigma w_1 - \lambda(w_1^T w_1 - 1)$
- Take derivative with respect to  $w_1$ , set to zero
  - $2\Sigma w_1 - 2\lambda w_1 = 0$
  - $\Sigma w_1 = \lambda w_1$
- $w_1$  is an eigenvector and  $\lambda$  is an eigenvalue of  $\Sigma$
- Because we maximize variance
  - maximize  $w_1^T \Sigma w_1 = w_1^T \lambda w_1 = \lambda w_1^T w_1 = \lambda$
  - $w_1$  is the eigenvector that corresponds to the largest eigenvalue of  $\Sigma$

$d \times d$   
 $w_1^T w_1 - 1 = 0$   
 $d \times d$   
 $d \times M$   
 $d \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$

$\frac{\partial}{\partial w_1}$

# PCA WITH $k$ DIMENSIONS

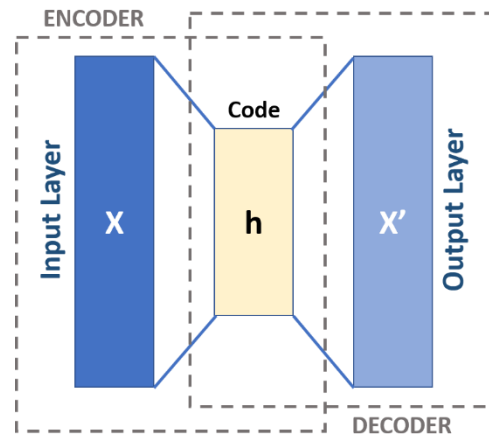
- We want the second vector  $w_2$  to be orthogonal to  $w_1$  so that the  $z_2$  is uncorrelated with  $z_1$   $w_2 \cdot w_1^T = 0$
- Skipping the derivation details
  - $w_1$  is the eigenvector for the largest eigenvalue  $\lambda_1$
  - $w_2$  is the eigenvector for the second largest eigenvalue  $\lambda_2$
  - ...
  - $w_k$  is the eigenvector for the  $k^{th}$  largest eigenvalue  $\lambda_k$
- In the end, the transformation is typically centered around zero (in the new dimensions)
  - $\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \mathbf{m})$  where  $\mathbf{m}$  is the mean of  $\mathbf{x}$

# PCA EXAMPLE

- See OneNote and Jupyter notebook

# AUTOENCODER

- A neural network architecture where the input and output are the same
  - $\mathbf{x} \rightarrow \mathbf{h} \rightarrow \mathbf{x}$
  - Input  $\mathbf{x}$  is encoded into  $\mathbf{h}$ , where  $\mathbf{h}$  is typically lower dimensional than  $\mathbf{x}$  and  $\mathbf{h}$  is decoded back to  $\mathbf{x}$ , with some error



[https://upload.wikimedia.org/wikipedia/commons/3/37/Autoencoder\\_schema.png](https://upload.wikimedia.org/wikipedia/commons/3/37/Autoencoder_schema.png)

# OTHER DL METHODS

- Transformers
- Contrastive learning
- ...
- International conference on learning representations (ICLR)
  - <https://iclr.cc/>