

LECTURE 02 Machine Learning I: conventional machine learning

Dr. Suyong Eum



1) Principal Component Analysis (PCA)

- Feature selections
- Dimension reduction
- 2) Support Vector Machine (SVM)
 - Hard margin SVM: linear classification
 - Kernel trick: nonlinear classification

Principal Component Analysis (PCA)

Principal Component Analysis (PCA): definition

A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.

In Wikipedia:



- ☐ How to select principal components?
 - One that captures the largest variance of the data points
 - Intuitively speaking, you can observe more data from the direction (1) than any other direction, and then from the direction (2), you can observe the data with the least redundancy compared to the direction (1).





Copyright © 2022 OSAKA University. All right reserved.

1) Find the covariance matrix of data points.

X₁

- 2) Obtain the eigen values and vectors of the covariance matrix: eigen value decomposition.
- 3) Sort the eigen vectors in descending order in terms of their corresponding eigen values.
 - an eigen vector with the largest eigen value becomes the first principal component.

2nd principal component v_2 v_3 v_1 1st principal component component





□ Actually, there is a more convenient way of doing it, which is called "Singular Value Decomposition" or SVD.

Eigen decomposition Singular Value Decomposition (SVD) $X^T X =$ $X = U\Sigma V$ $X^{T}X = (U\Sigma V^{T})^{T}(U\Sigma V^{T})$ >> x >> [vec, val] = eig(cov(x) vec = x = -0.707110.70711 -2 0.70711 0.70711 $= V \Sigma^{\mathrm{T}} U^{\mathrm{T}} U \Sigma V^{\mathrm{T}}$ -1 1 -1 val = -1 1 1 Diagonal Matrix 1 2 2 $\Lambda = \Sigma^2$ 0.80000 Θ 4.00000 Θ >> cov(x)ans = >> [vec, val]=eig(transpose(x)*x) vec = 2.4000 1.6000 -0.707110.70711 1.6000 2.4000 0.70711 0.70711 val = Diagonal Matrix 4.0000 Θ 0 20.0000 **Eigen value** Singular value Copyright © 2022 OSAKA University. All right reserved.

□ Actually, there is a more convenient way of doing it, which is called "Singular Value Decomposition" or SVD.



Singular Value Decomposition (SVD)

 $X = U\Sigma$



Now we know how to find the principal components

- Principal Component Analysis (PCA) is nothing but finding principal components of a given data set,
 - Principal components are the directions where you look at the data set, which provides the most information of the data set.
 - They're equivalent to eigen vectors which can be found by SVD or EVD.
 - The eigen value corresponding to each eigen vector represents how widely the data set is spread along the direction which is perpendicular to the eigen vector.



- □ A data point is defined by several, let's say, features,
- The number of features to define a data point is called the dimension of the data,
- □ High dimension data implies that it contains much information,
- Sometimes, we reduce its dimension, e.g., to visualize the data or to efficiently analyze them,
- PCA can reduce the dimension without losing relatively less information of the data.

- □ The previous example shows the case of two-dimensional data
- □ How can we reduce the two-dimensional data to one dimension?
- □ Yes, just project the data points onto the eigenvector space!



Dimension Reduction



Dimension Reduction



$\mathbf{V} = \mathbf{U} \mathbf{\nabla} \mathbf{V}^{T}$	<pre>>> [U, S, V]= U =</pre>	=svd(x)					
$\Lambda - U \angle V$	-0.63246	0.00000	0.30819	-0.30819	0.28637	0.57274	
	-0.31623	-0.00000	-0.63635	0.63635	0.13426	0.26851	
	0.00000	-0.70711	0.50000	0.50000	0.00000	0.00000	
	-0.00000	0.70711	0.50000	0.50000	-0.00000	-0.00000	
	0.31623	0.00000	-0.00399	0.00399	0.94140	-0.11720	
	0.63246	0.00000	-0.00/99	0.00799	-0.11720	0.76560	
	S =						
	Diagonal Ma	Diagonal Matrix					
	4.4721	Θ					
	Θ	2.0000					
	Θ	Θ					
	Θ	Θ					
	Θ	Θ					
	Θ	Θ					
	V =						
	0.70711	-0.70711					
	0.70711	0.70711					

Copyright © 2022 OSAKA University. All right reserved.

Dimension Reduction



- Let's say, we have one image representing one data point as shown below,
- Then, we decide to present the data by all pixels which are 64 in this case, in other words, it is 64-dimensional data,
- □ What happens if we reduce its dimension to 2 dimension?



- Let's say, we have one image representing one data point as shown below,
- Then, we decide to present the data by all pixels which are 64 in this case, in other words, it is 64-dimensional data,
- □ What happens if we reduce its dimension to 2 dimension?



Copyright © 2022 OSAKA University. All right reserved.



- Well, now we have a new set of data which have two dimension, so they can be presented in the two-dimensional space. Data visualization!
- □ Also, we may be able to classify those data by drawing a line???



Support Vector Machine (SVM)

Why Support Vector Machine?

- □ Most widely used classification approach (practical)
 - Linearly separable data set
 - Non-linearly separable data set

- ❑ Supported by well defined mathematical theories
 - Geometry
 - Optimization

Which line is better to split two data sets?







Terminology used in SVM





$$\mathbf{x}^{c} = \mathbf{x}^{b} + \| r \| \frac{\mathbf{w}}{\| \mathbf{w} \|}$$

D Let's multiply w^T and add w_0 in both sides.

 $W^{T}X^{c} + w_{0} = W^{T}X^{b} + w_{0} + W^{T} ||r|| \frac{W}{||w||}$ $y(x^{c}) = W^{T} ||r|| \frac{W}{||w||}$ $||r|| = \frac{y(x^{c})}{||w||}$ $U^{Let's say}||y(x^{c})| = 1$ We use it later ...



Finding a decision boundary which maximizes the margin.

$$\max \parallel r \parallel = \frac{1}{\parallel \mathbf{w} \parallel}$$

s.t.



Every data points are classified correctly.

$$\begin{cases} t_n = +1, & y(\mathbf{x}_n) > 0 \\ t_n = -1, & y(\mathbf{x}_n) < 0 \end{cases}$$



Problem formulation





Copyright © 2022 OSAKA University. All right reserved.

How about non-linearly separable case?





Kernel Trick

Lagrange method for an optimization problem with inequality constraints



$$\min_{x} \max_{\lambda} x^{2} - \lambda(x-b)$$

s.t. $\lambda \ge 0$

- □ If $b \le 0$, the minima is 0 ... so $\lambda = 0$
- If b > 0, the minima is $b^2 \dots$ so x=b
- **C** So, either λ or (x b) becomes zero, in other words,
 - $\lambda(x-b) = 0$ (complementary slackness)
- □ Since $x \ge b$, maximizing λ minimizes the objective value
 - $\lambda \ge 0$



$$\min_{\mathbf{w}} \max_{\lambda} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1)$$

s.t. $\lambda_{n} \ge 0$

Proof begins



$$\begin{array}{l} \min_{\mathbf{w}} \max_{\lambda} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1) \\ S.t. \quad \lambda_{n} \geq 0 \end{array}$$

We would like to convert again the optimization problem above into another form, which provides same results.

- Because we want to solve the optimization problem in term of "lagrange multiplier (λ_n) ".

$$\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1)$$

s.t. $\lambda_{n} \ge 0$

Primal

problem

probl

$$\begin{array}{l} \underset{\mathbf{w}}{\min} \max_{\lambda} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1) \\
s.t. \quad \lambda_{n} \ge 0
\end{array}$$

- We would like to convert again the optimization problem above into another form, which provides same results.
 - Because we want to solve the optimization problem in term of "lagrange multiplier (λ_n)".

Dual
Toblem
$$\begin{array}{c}
\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1) \\
s.t. \quad \lambda_{n} \ge 0
\end{array}$$

Karush–Kuhn–Tucker conditions

KKT conditions

 $\lambda_n \geq 0$

1) Stationarity condition

 $\frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} - \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{\mathrm{T}} x_{n} + w_{0}) - 1) = 0$

Primal

problem

Dual

problem

- 2) Complementary slackness condition $\lambda_n(t_n(\mathbf{w}^{\mathrm{T}}x_n + w_0) - 1) = 0$
- 3) Duality feasibility condition

$$\begin{array}{l} \underset{\mathbf{w}}{\min} \max_{\lambda} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1) \\
s.t. \quad \lambda_{n} \ge 0
\end{array}$$

- We would like to convert again the optimization problem above into another form, which provides same results.
 - Because we want to solve the optimization problem in term of "lagrange multiplier (λ_n) ".

$$\max_{\lambda} \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^{T} \mathbf{w} - \sum_{n=1}^{n} \lambda_{n} (t_{n} (\mathbf{w}^{T} x_{n} + w_{0}) - 1)$$

s.t. $\lambda_{n} \ge 0$

$$\max_{\lambda} \min_{\mathbf{w}, \mathbf{w}_0} L(\mathbf{w}, \mathbf{w}_0, \lambda) = \frac{1}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} - \sum_{n=1}^{N} \lambda_n (t_n(\mathbf{w}^{\mathrm{T}} x_n + w_0) - 1)$$



$$\frac{\partial L}{\partial w} = \mathbf{w} - \sum_{n=1}^{N} \lambda_n t_n x_n = 0$$

- The first one is called stationarity condition.
 - when we partial differentiate the problem with respect to its parameter "w", each of them should be zero.

$$\max_{\lambda} \min_{w,w_0} L(w,w_0,\lambda) = \frac{1}{2} w^{\mathsf{T}} w - \sum_{n=1}^{N} \lambda_n (t_n (w^{\mathsf{T}} x_n + w_0) - 1)$$



The first one is called stationarity condition.

 \succ Again, this time in terms of "w₀"

$$L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n t_n w_0 + \sum_{n=1}^{N} \lambda_n t_n w_0$$



$$\max_{\lambda} L(\lambda) = \sum_{n=1}^{N} \lambda_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m$$

s.t. $\lambda_n \ge 0$, $\sum_{n=1}^{N} \lambda_n t_n = 0$

- □ Let's change it to a quadratic programming again.
- As mentioned previously, a quadratic programming problem needs to be minimized

$$\max_{\lambda} L(\lambda) = \sum_{n=1}^{N} \lambda_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m$$

$$s.t. \quad \lambda_n \ge 0, \qquad \sum_{n=1}^{N} \lambda_n t_n = 0$$

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda_n \ge 0, \qquad \sum_{n=1}^{N} \lambda_n t_n = 0$$

□ Again, the optimization problem becomes a quadratic programming problem.

Let's summarize

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

s.t. $\lambda \ge 0$, $t^T \lambda = 0$

The solution from the quadratic programming is "lagrange multipliers"(λ_n)
 Many of the solutions (lagrange multipliers) are zero
 Complementary slackness (one of KKT conditions) should be satisfied.

 $\lambda_n(t_n(\mathbf{w}^{\mathrm{T}}x_n+w_0)-1)=0$

□ In other words, if λ_n are not zero, $(t_n(w_tx_n+w_0)-1)$ should be zero where corresponding data points should be support vectors.



Let's summarize

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

s.t. $\lambda \ge 0$, $t^T \lambda = 0$

□ The solution from the quadratic programming is "lagrange multipliers"(λ_n) □ Many of the solutions (lagrange multipliers) are zero

□ Complementary slackness (one of KKT conditions) should be satisfied.

 $\lambda_n(t_n(\mathbf{w}^{\mathrm{T}}x_n+w_0)-1)=0$

In other words, if λ_n are not zero, (t_n(w_tx_n+w₀)-1) should be zero where corresponding data points should be support vectors.
 With the non-zero λ_n, w and w₀ can be calculated using t_n(w_tx_n+w₀)=1

$$\mathbf{w} = \sum_{n=1}^{N} \lambda_n t_n x_n \qquad \mathbf{w}_0 = t_n - \sum_{n=1}^{N} \lambda_n t_n x_n x_n$$

 $W^{T}X + W_{0} = -1$

 $W' X + W_0 = 0$

Let's summarize

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^{N} \lambda_n$$

s.t. $\lambda \ge 0$, $t^T \lambda = 0$
 $\mathbf{x} = \mathbf{x}_n \mathbf{x$

The solution from the quadratic programming is "lagrange multipliers"(λ_n)
 Many of the solutions (lagrange multipliers) are zero
 Complementary slackness (one of KKT conditions) should be satisfied.

 $\lambda_n(t_n(\mathbf{w}^{\mathrm{T}}x_n+w_0)-1)=0$

In other words, if λ_n are not zero, (t_n(w_tx_n+w₀)-1) should be zero where corresponding data points should be support vectors.
 With the non-zero λ_n, w and w₀ can be calculated using t_n(w_tx_n+w₀)=1

$$\mathbf{w} = \sum_{n=1}^{N} \lambda_n t_n x_n \qquad \qquad \mathbf{w}_0 = t_n - \sum_{n=1}^{N} \lambda_n t_n x_n x_n$$

 $w_0 = t_n - \sum_{n=1}^N \lambda_n t_n x_n x_n$

Proof ends

Kernel trick



 \Box If data x_n are not linearly separable, what should we do?



Kernel trick

□ The idea of Kernel trick begins from here: to find the scalar values (the inner product of two vectors: z_n and z_m) and so we can formulate the quadratic problem which can be linearly separable.



Kernel trick: Kernel function

□ Kernel function K() is a function which returns the scalar values (the inner product of two vectors: z_n and z_m in Z space) when the data points (x_n and x_m in X space) are given.

$$K(\mathbf{x}_n^T, \mathbf{x}_m) = \mathbf{z}_n^T \mathbf{z}_m$$





Kernel trick: Kernel function

□ With the Kernel function defined previously, we want to change the quadratic problem as follows:

- Because the Kernel function is a function of data points $(x_n \text{ and } x_m)$ which we already have.

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^{N} \lambda_n \\
s.t. \quad \lambda \ge 0, \quad t^T \lambda = 0$$

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{K}(\mathbf{x}_n^T \mathbf{x}_m) - \sum_{n=1}^{N} \lambda_n \\
s.t. \quad \lambda \ge 0, \quad t^T \lambda = 0$$

Kernel trick: Kernel function

□ With the Kernel function defined previously, we want to change the quadratic problem as follows:

- Because the Kernel function is a function of data points $(x_n \text{ and } x_m)$ which we already have.

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \ge 0, \quad t^T \lambda = 0$$

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} t_n t_m \lambda_n \lambda_m \mathbf{K}(\mathbf{x}_n^T \mathbf{x}_m) - \sum_{n=1}^{N} \lambda_n$$

$$s.t. \quad \lambda \ge 0, \quad t^T \lambda = 0$$

$$\min_{\lambda} L(\lambda) = \frac{1}{2} \lambda^{T} \begin{bmatrix} t_{1}t_{1}K(\mathbf{x}_{1},\mathbf{x}_{1}) & t_{1}t_{2}K(\mathbf{x}_{1}^{T},\mathbf{x}_{2}) & \cdots & t_{1}t_{N}K(\mathbf{x}_{1}^{T},\mathbf{x}_{N}) \\ t_{2}t_{1}K(\mathbf{x}_{2},\mathbf{x}_{1}) & t_{2}t_{2}K(\mathbf{x}_{2}^{T},\mathbf{x}_{2}) & \cdots & t_{2}t_{N}K(\mathbf{x}_{2}^{T},\mathbf{x}_{N}) \\ \cdots & \cdots & \cdots & \cdots \\ t_{N}t_{1}K(\mathbf{x}_{N}\mathbf{x}_{1}) & t_{N}t_{2}K(\mathbf{x}_{N}^{T},\mathbf{x}_{2}) & \cdots & t_{N}t_{N}K(\mathbf{x}_{N}^{T},\mathbf{x}_{N}) \end{bmatrix} \lambda + (-1^{T})\lambda$$

Polynomial kernel of degree 2



$$K(\mathbf{x}, y) = (\mathbf{x}y)^{2}$$

= $((x_{1}, x_{2}) \cdot (y_{1}, y_{2}))^{2}$
= $(x_{1}y_{1} + x_{2}y_{2})^{2}$
= $x_{1}^{2}y_{1}^{2} + 2x_{1}x_{2}y_{1}y_{2} + x_{2}^{2}y_{2}^{2}$

 $\phi(\mathbf{x})\phi(y) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2)$

$$= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2$$

- PCA and SVM are probably the most representative conventional machine learning algorithms.
- □ PCA helps you to manipulate a set of data in a way that
 - determining which features are important,
 - reducing its dimension, so that the data can be processed or visualized more efficiently.
- SVM is a classification method founded on well defined mathematical framework, which can handle linear or nonlinear classification problems.