

Kernel Methods

CS534

Key concepts:

Feature mapping to address non-linear separability

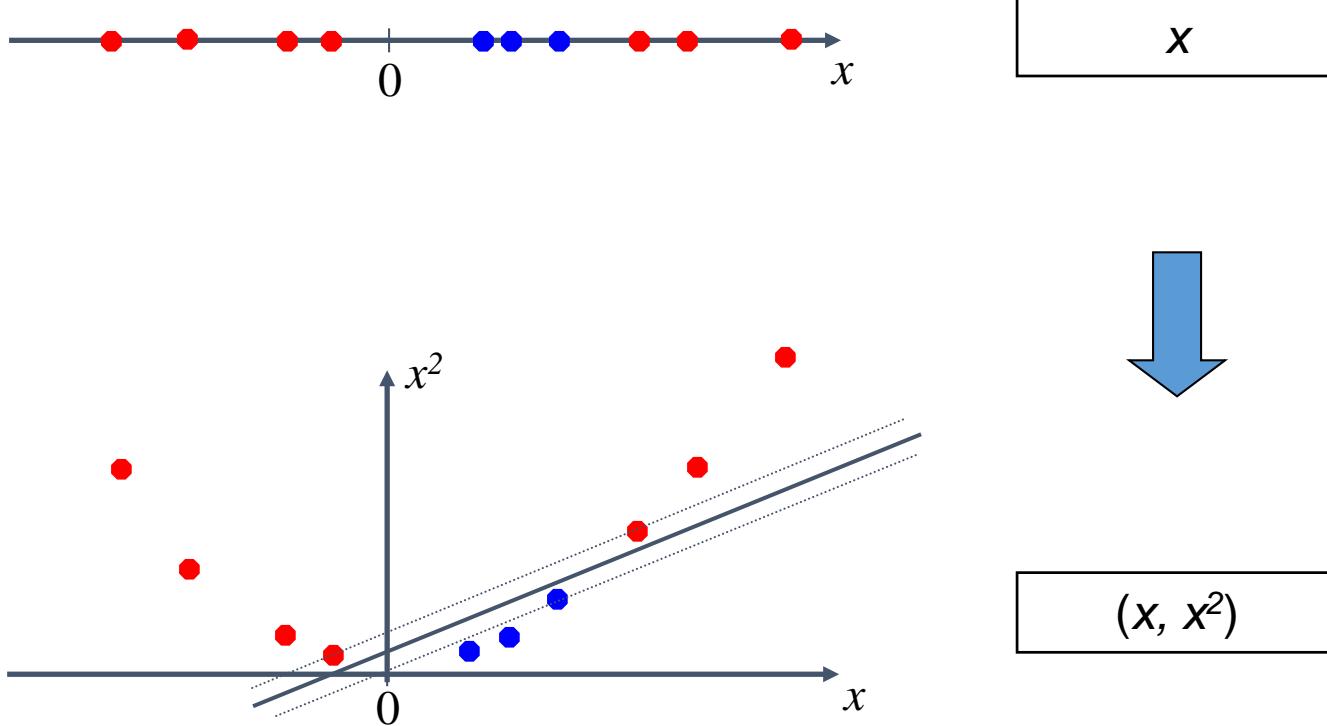
Kernel trick to avoid explicit feature mapping

Definition of Kernel functions

Kernelized perceptron

Kernelized linear regression with L2 regularization

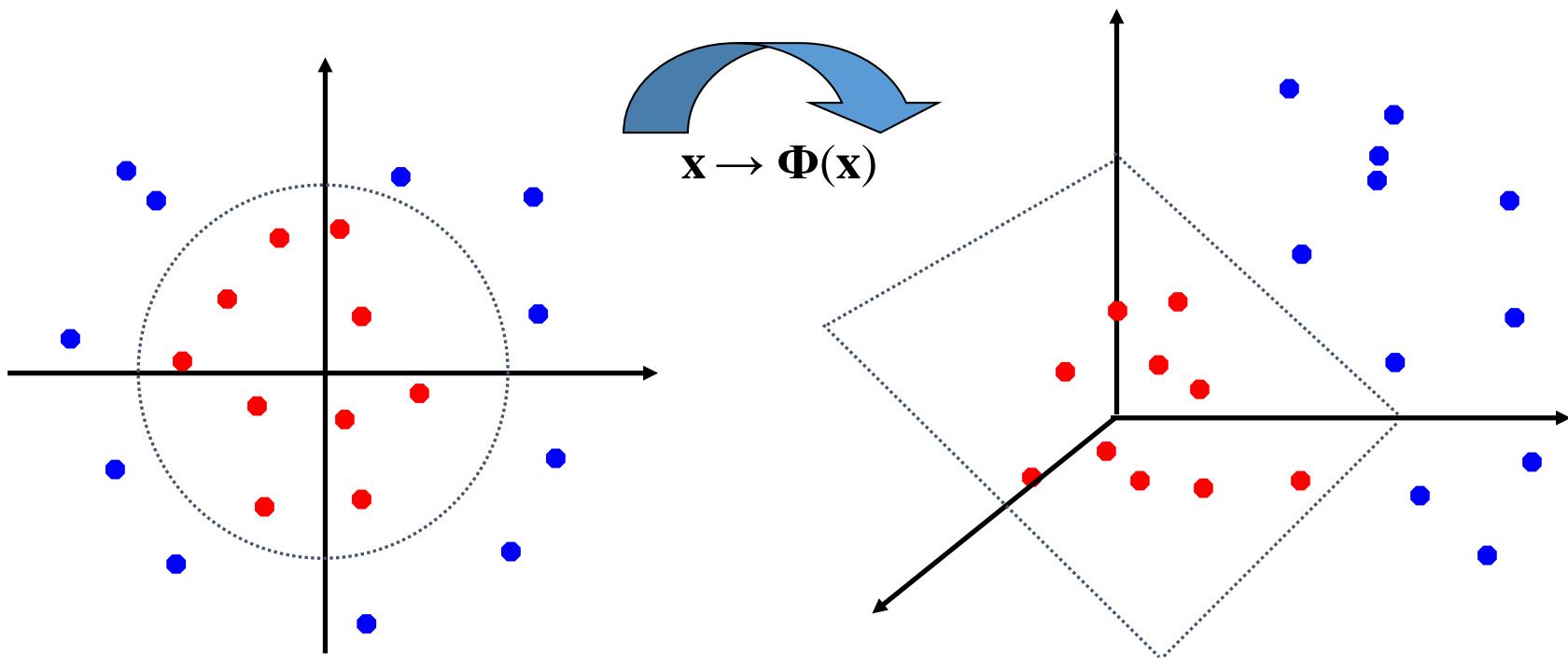
Nonlinearly separable data



Mapping the data to a higher dimensional space can introduce linear separability for data not linearly separable in the original space

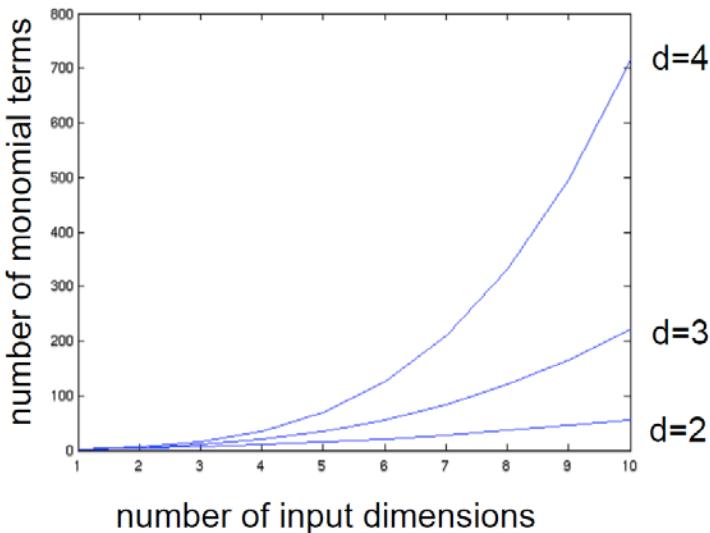
Non-linear classifier via feature mapping

- General idea: For any data set, the *original input space* can always be mapped to some higher-dimensional **feature spaces** such that the data is linearly separable:



Example: Quadratic Feature Space

- Assume m input dimensions $\mathbf{x} = (x_1, \dots, x_m)$
so you have $m \times m$ features
- Number of quadratic terms:
- $1 + m + m + m(m-1)/2 \approx O(m^2)$
- What if we want to consider even higher order features?
 - For cubic feature space: $O(m^3)$
 - The number of dimensions after mapping increase rapidly with increasing order d



$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_m \\ x_1^2 \\ x_2^2 \\ \vdots \\ x_m^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_1x_m \\ \sqrt{2}x_2x_3 \\ \vdots \\ \sqrt{2}x_{m-1}x_m \\ \vdots \\ \sqrt{2}x_{m-1}x_m \end{pmatrix}$$

Annotations for the columns of the matrix:

- Constant Term: 1 (red bracket)
- Linear Terms: $\sqrt{2}x_1, \sqrt{2}x_2, \dots, \sqrt{2}x_m$ (green bracket)
- Pure Quadratic Terms: $x_1^2, x_2^2, \dots, x_m^2$ (purple bracket)
- Quadratic Cross-Terms: $\sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_{m-1}x_m$ (blue bracket)

Revisiting perceptron

predict는 wt dot x 한다음 거기에 take sign

Given current weight w_t , predict for x by: $\hat{y}(x) = \text{sign}(w_t \cdot x)$

If mistake at time t , update the weight:

$$w_{t+1} = w_t + y_t x_t$$

If remember all the mistakes made up to t in s_t , we have:

$$w_t = \sum_{i \in s_t} y_i x_i$$

The prediction rule now becomes:

$$\hat{y}(x) = \text{sign}\left(\sum_{i \in s_t} y_i x_i \cdot x\right)$$

every previous mistake Dot x

Now data is mapped to higher dimension:

$$\hat{y}(x) = \text{sign}\left(\sum_{i \in s_t} y_i \phi(x_i) \cdot \phi(x)\right)$$

Dot product in the Quadratic Feature Space

Explicit mapping takes $O(m^2)$ time

$$\Phi(\mathbf{a}) \cdot \Phi(\mathbf{b}) = \begin{pmatrix} 1 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \vdots \\ \sqrt{2}a_m \\ a_1^2 \\ a_2^2 \\ \vdots \\ a_m^2 \\ \sqrt{2}a_1a_2 \\ \sqrt{2}a_1a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \sqrt{2}a_2a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \vdots \\ \sqrt{2}a_{m-1}a_m \end{pmatrix} \bullet \begin{pmatrix} 1 \\ \sqrt{2}b_1 \\ \sqrt{2}b_2 \\ \vdots \\ \sqrt{2}b_m \\ b_1^2 \\ b_2^2 \\ \vdots \\ b_m^2 \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \sqrt{2}b_2b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \vdots \\ \sqrt{2}b_{m-1}b_m \end{pmatrix}$$

이제 바로 커널 function 다음장에 나와영

$\kappa(a, b) = (a \cdot b + 1)^2$ computes the dot product in quadratic space in $O(m)$ time

$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) =$ 애는 $O(m^2)$

$$1 + 2 \sum_{i=1}^m a_i b_i + \sum_{i=1}^m a_i^2 b_i^2 + \sum_{i=1}^m \sum_{j=i+1}^m 2a_i a_j b_i b_j$$

Consider the following function:

$$\begin{aligned} & (a \cdot b + 1)^2 \\ &= (a \cdot b)^2 + 2(a \cdot b) + 1 \quad \text{애는 } O(m) \\ &= (\sum_{i=1}^m a_i b_i)^2 + 2 \sum_{i=1}^m a_i b_i + 1 \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i a_j b_i b_j + 2 \sum_{i=1}^m a_i b_i + 1 \\ &= \sum_{i=1}^m a_i^2 b_i^2 + 2 \sum_{i=1}^m \sum_{j=i+1}^m a_i a_j b_i b_j + 2 \sum_{i=1}^m a_i b_i + 1 \end{aligned}$$

Kernel functions: the kernel trick

- **Definition:** A function $\kappa(\mathbf{x}, \mathbf{x}')$ is called a **kernel function** if $\kappa(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ for some mapping function ϕ
- **Implication:** we can simply replace any occurrences of dot product $\langle \mathbf{x} \cdot \mathbf{x}' \rangle$ with a kernel function κ that computes $\kappa(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- **Implication:** we do not need to explicitly compute the mapping of the features --- significant computational savings
- For example, to compute the dot product in the quadratic space:
 - With the quadratic kernel function: $O(m)$
 - With explicit mapping: $O(m^2)$

Revisiting perceptron

Given current weight w_t , predict for x by: $\hat{y}(x) = \text{sign}(w_t \cdot x)$

If mistake at time t , update the weight:

$$w_{t+1} = w_t + y_t x_t$$

If remember all the mistakes made up to t in s_t , we have:

$$w_t = \sum_{i \in s_t} y_i x_i$$

The prediction rule now becomes:

$$\hat{y}(x) = \text{sign}\left(\sum_{i \in s_t} y_i x_i \cdot x\right)$$

Now data is mapped to higher dimension:

$$\hat{y}(x) = \text{sign}\left(\sum_{i \in s_t} y_i \underbrace{\phi(x_i) \cdot \phi(x)}_{K(x_i, x)}\right) = \text{sign}\left(\sum_{i \in s_t} y_i \underbrace{K(x_i, x)}_{K(x_i, x)}\right)$$

Kernelizing perceptron

Let $\mathbf{w} \leftarrow (0, 0, \dots, 0)$

Repeat if $iter \leq iters$

for every training example $m = 1, \dots, n$

$$u_m = \mathbf{w}^T \mathbf{x}_m$$

if $y_m u_m \leq 0$ $\mathbf{w} \leftarrow \mathbf{w} + y_m \mathbf{x}_m$ update

$iter = iter + 1$

이거를 이렇게 change

Remembering all mistakes in s

$S = []$ //list of mistaken examples

Repeat if $iter \leq iters$

for every training example $m = 1, \dots, n$

$$u_m = \sum_{i \in S} y_i k(x_i, x_m)$$

if $y_m u_m \leq 0$ add m to S update

$iter = iter + 1$

Weight가 simply all mistake in s 를 더해가고 막 하는거니까
 S 를 update하면 w 가 update 되겟징

Kernelizing perceptron

$S = []$ //list of mistaken examples

Repeat if $iter \leq iters$

for every training example $m = 1, \dots, n$

$$u_m = \sum_{i \in S} y_i \kappa(x_i, x_m)$$

u_m 이 위에 sigma(yixi * x)
yi가 positive or negative sign

if $y_m u_m \leq 0$ add m to S

$iter = iter + 1$

Let $\alpha_i = 0 \ \forall i = 1, \dots, n$

Repeat if $iter \leq iters$

for every training example $m = 1, \dots, n$

$$u_m = \sum_{i=1 \text{ to } n} \alpha_i y_i \kappa(x_i, x_m)$$

if $y_m u_m \leq 0 \quad \alpha_m \leftarrow \alpha_m + 1$

$iter = iter + 1$

More efficiently, count for each example the number of times it gets misclassified.

근데 이걸 쓸 수 있는 상황이 있고 못쓰는 상황이 있겠다.

스트리밍 서비스같은 경우에는 카운트가 아니라, 그 미스테이크 자체를 기록해야한다!!

Instead remembering all s, 알파m을 써서 mistake를 count

What functions are kernels?

- A kernel function can be intuitively viewed as computing some similarity measure between examples
- In practice, we directly use the kernel functions without explicitly stating the transformation ϕ
- Given a kernel function, finding its corresponding transformation can be very cumbersome or impossible
 - RBF kernel's mapped space has infinite dimensions
- Not all functions are kernels
 - For some functions there does not exist a corresponding mapping $\phi(x)$
- If you have a good similarity measure, can we use it as a kernel?

Kernel function or not

- Consider a finite set of m points, we define the kernel (Gram) matrix as

$$K = \begin{bmatrix} \kappa(x^1, x^1) & \kappa(x^1, x^2) & \dots & \kappa(x^1, x^m) \\ \kappa(x^2, x^1) & \kappa(x^2, x^2) & \dots & \kappa(x^2, x^m) \\ \dots & \dots & \dots & \dots \\ \kappa(x^m, x^1) & \kappa(x^m, x^2) & \dots & \kappa(x^m, x^m) \end{bmatrix}_{mxm}$$

- Kernel matrices are square and symmetric
- Mercer theorem:

A function K is a kernel function iff for any finite sample x_1, x_2, \dots, x_m , its corresponding kernel matrix is positive semi-definite (has non-negative eigenvalues) Theorem

$(x_1 * x_2)^p \Rightarrow$ 얘도 위에거랑 나중가면 equivalent해지는데 위에꺼가 easier to compute

$\exp(-\|x_i - x_j\|^2 / 2 * \text{derivation}^2)$ 이때 $2 * \text{derivation}$ kernel width라고 called.

\rightarrow relate to overfitting and potentially affect to convergence by

Closure property of kernels

If K_1 and K_2 are kernel functions, then the following are all kernel functions:

- $K(x, y) = K_1(x, y) + K_2(x, y)$ 커널function은 compatatable
 - $\phi = \phi_1 + \phi_2$
- $K(x, y) = aK_1(x, y)$, where $a > 0$ alpha가 positive면 괜찮
근데 negative면 안된다
알파가 negative면 더이상 kernel이 아니게됨
 - $\phi = \sqrt{a}\phi_1$
- $K(x, y) = K_1(x, y)K_2(x, y)$ dot product도 할 수 있당
 - If ϕ_1 has N_1 features and ϕ_2 has N_2 features
 - ϕ will have $N_1 \times N_2$ features: $\phi_{ij} = \phi_{1i} \cdot \phi_{2j}$

Key Choices in Applying Kernel

- Selecting the kernel function
 - In practice, we often try different kernel functions and use (cross-) validation to choose learning kernel: 어떤 kernel function이 좋은지
 - Linear kernel, polynomial kernels (with low degrees) and RBF kernels are popular choices
 - One can also construct a kernel using linear combinations of different kernels and learn the combination parameters (kernel learning)
 - There are many kernel functions defined for non-traditional data, e.g., graph kernel, set kernel, string kernel etc.
- Selecting the kernel parameter
 - Very strong impact on performance
 - Often the optimal range is reasonably large
 - Grid search with (cross-)validation we will talk about this

Revisiting linear regression

- Consider the L2 regularized linear regression

$$\max_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

where \mathbf{X} is the data matrix, each row contains the features of one training example

\mathbf{y} is the vector of ground truth predictions for all training examples

- Closed form solution:

$$\mathbf{w} = (\underbrace{\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X}}_{\text{regularization product}})^{-1} \mathbf{X}^T \mathbf{y}$$

where \mathbf{I} is the identity matrix

- \mathbf{w} lies in the space spanned by the training examples, i.e.,

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$$

우리는 \mathbf{w} 가 linear regression에 누워있다는 걸 알 수 있음 위에거를 통해서

Kernelizing linear regression

- Learned Function:

- Original: $\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$

$$\begin{aligned} f(\mathbf{x}) &= \text{sig } a_i * x_i * \mathbf{x} \\ &= \mathbf{w}^T \mathbf{x} \end{aligned}$$

- Kernelized: $\hat{f}(\mathbf{x}) = \sum_i \alpha_i \kappa(\mathbf{x}_i, \mathbf{x})$ (by plugging in $\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$ and replace dot product with kernel function)

- Objective:

- Original: $\frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$

- Kernelized: $\frac{1}{2} \|\mathbf{y} - \underline{\alpha K}\|^2 + \frac{\lambda}{2} \underline{\alpha}^T K \underline{\alpha}$

This is referred to as the dual of the original optimization problem

K is the kernel (gram) matrix of the training data

- Closed form solution:

$$\underline{\alpha} = (K + \lambda I)^{-1} \mathbf{y}$$

General application of Kernel methods

- Many learning tasks are framed as optimization problems and their solutions are some weighted sum of the input training examples
- By explicitly formulating the optimization in the space of the weights, we arrive at the so called “dual formulation” of the optimization problems
- The dual problem is often expressed in terms of the dot product between examples
- The kernel trick can be applied by replacing dot product with kernel functions.
- This allows the use of high dimensional feature spaces without having to pay the price of working with the high dimensional features