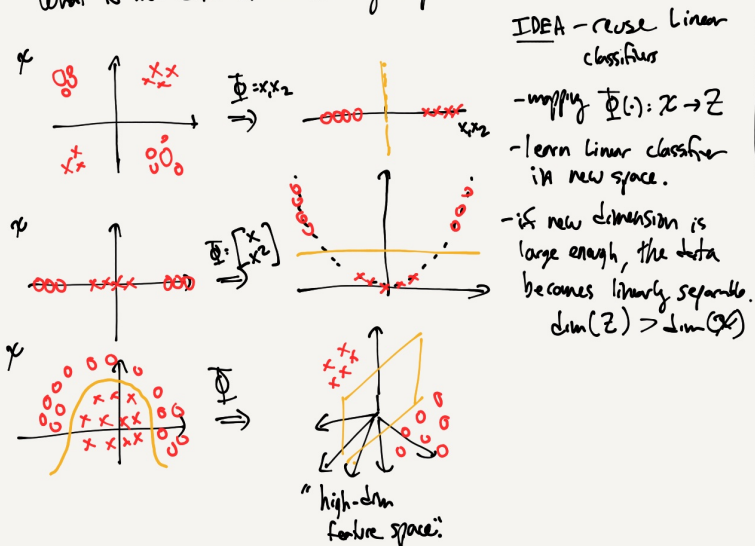


Lecture 10 Nonlinear classifiers & Kernels

Linear classifiers - SVM, LR, perceptron, etc.

What if the data is nonlinearly separable?



In the limit, $\dim(\mathcal{Z}) \rightarrow \infty$

- we are mapping each point x_i into a function

$$\Phi(x) \rightarrow \begin{bmatrix} x_1 \\ \vdots \\ x_0 \end{bmatrix} \rightarrow \phi(x; t)$$

Kernel SVM

Consider the SVM dual problem, but replace $x_i \rightarrow \Phi(x_i)$

Training $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\Phi(x_i)^T \Phi(x_j)}_{k(x_i, x_j)}$

s.t. $\sum_i \alpha_i y_i = 0, \alpha_i \geq 0$

Decision Boundary: $w^* = \sum_i \alpha_i y_i \Phi(x_i)$ (w is in same space as $\Phi(w)$)

Bias term $b = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (y_i - w^T \Phi(x_i)) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (y_i - \sum_j \alpha_j y_j \underbrace{\Phi(x_j)^T \Phi(x_i)}_{k(x_j, x_i)})$

Set of support vectors.

Decision function

$$y^* = \text{sign}(f(x^*))$$

$$f(x^*) = w^T \Phi(x^*) + b = \sum_i \alpha_i y_i \underbrace{\Phi(x_i)^T \Phi(x^*)}_{k(x_i, x^*)} + b$$

Note: entire algorithm depends only on $\Phi(x_i)^T \Phi(x_j)$.

Define function: $k(x_i, x_j) = \Phi(x_i)^T \Phi(x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$

dot product kernel

• The dual SVM can be rewritten w/ this kernel function \rightarrow nonlinear classifier

• Just define $k(x_i, x_j)$ directly w/o explicitly calculating $\Phi(x)$ ("kernel trick")

- Good: saves time/computation of calculating $\Phi(x)$

- Bad: need $O(n^2)$ terms of $k(x_i, x_j)$, store this kernel matrix.

$$K = [k(x_i, x_j)]_{i,j}$$

Example

polynomial kernel: $x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$

$$\begin{aligned} k(x, x') &= (x^T x')^2 = \left(\sum_{i=1}^d x_i x'_i \right)^2 \\ &= \sum_{i=1}^d \sum_{j=1}^d x_i x'_i x_j x'_j = \sum_{i=1}^d \sum_{j=1}^d \underbrace{(x_i x_j)}_{\Phi(x)} \underbrace{(x'_i x'_j)}_{\Phi(x')} \\ &= \underbrace{\begin{bmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_d \\ x_2 x_1 & x_2 x_2 & \dots & x_2 x_d \\ \vdots & \vdots & \ddots & \vdots \\ x_d x_1 & x_d x_2 & \dots & x_d x_d \end{bmatrix}}_{\Phi(x)^T} \underbrace{\begin{bmatrix} x'_1 x'_1 \\ x'_1 x'_2 \\ \vdots \\ x'_d x'_d \end{bmatrix}}_{\Phi(x')} \\ &= \Phi(x)^T \Phi(x') \end{aligned}$$

Hence, $\Phi(x): \mathbb{R}^d \rightarrow \mathbb{R}^{d^2}$

$$k(x, x') = \underbrace{\Phi(x)^T}_{\text{calculation: } O(d^2)} \underbrace{\Phi(x')}_{O(d)} = \underbrace{(x^T x')^2}_{\text{more efficient}}$$

Look SVM decision function:

$$\begin{aligned} f(x) &= \sum_i \alpha_i y_i k(x, x_i) + b \\ &= \sum_i \alpha_i y_i (x^T x_i)^2 + b = \sum_i \alpha_i y_i \underbrace{(x^T x_i)}_{\underbrace{x^T}_{\text{matrix}} \underbrace{x_i}_{\text{vector}}} \underbrace{(x_i^T x)}_{\text{scalar}} + b \\ &= x^T \underbrace{\left(\sum_i \alpha_i y_i x_i x_i^T \right)}_A x + b \end{aligned}$$

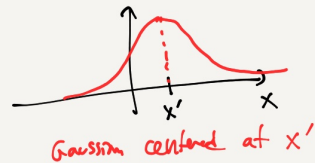
A

quadratic function.

the kernel specifies the class of functions that are used.

Gaussian kernel / Radial Basis Function (RBF) kernel

$$k(x, x') = e^{-\frac{1}{2\sigma^2} \|x - x'\|^2}$$



- Is it a dot-product kernel?
- What is the $\Phi(x)$?

Kernel Functions

kernel trick depends on $k(x_i, x_j)$ is a dot-product kernel

Defn a mapping $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a dot-product kernel iff

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

where $\Phi: \mathcal{X} \rightarrow \mathcal{H}$, \mathcal{H} is a vector space
 $\langle \cdot, \cdot \rangle$ is dot product in \mathcal{H}

How to check w/o knowing $\Phi(x)$ & $\langle \cdot, \cdot \rangle$?

Defn $k(x, x')$ is a positive definite kernel

if $\forall n$ & $\forall \{x_1, \dots, x_n\}$, $x_i \in \mathcal{X}$ (all datasets of all sizes)

the kernel matrix $K = [k(x_i, x_j)]_{i,j}$ is a positive definite matrix.

- $\left\{ \begin{array}{l} K \text{ is posdef iff } 1) y^T K y > 0, \forall y \\ 2) \text{ eigenvalues of } K > 0 \\ 3) K = X X^T, X \text{ has independent columns.} \end{array} \right.$

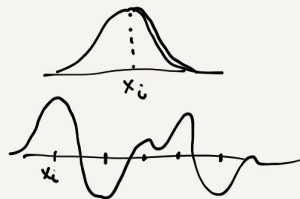
\star $k(x, x')$ is a dot-product kernel iff it is posdef kernel.

Given a posdef kernel, what is the high-dim x-form $\Phi(x)$?

Let \mathcal{H} = space of all linear combinations of function $k(\cdot, x_i)$
 $\mathcal{H} = \{ f(\cdot) \mid f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i), \forall m, \forall x_i \in \mathcal{X} \}$ function of x is fixed

e.g. use Gaussian kernel
 $k(\cdot, x_i) = e^{-\frac{1}{2\sigma^2} \|\cdot - x_i\|^2}$

$$f(\cdot) = \sum_i \alpha_i k(\cdot, x_i)$$



Let $f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i)$, $g(\cdot) = \sum_{j=1}^m \beta_j k(\cdot, x_j)$

Can show the dot-product btwn them in \mathcal{H} is:

$$\langle f, g \rangle = \sum_i \sum_j \alpha_i \beta_j k(x_i, x_j)$$

e.g. Gaussian kernel $\langle f, g \rangle = \sum_i \sum_j \alpha_i \beta_j e^{-\frac{1}{2\sigma^2} \|x_i - x_j\|^2}$

(like a non-linear similarity function in \mathcal{X})

Special case:

$$\alpha_i = 1, x_i = 0, \beta_j = 1, x_j = 0$$

$$\Rightarrow \langle \underline{k(\cdot, x_i)}, k(\cdot, x_j) \rangle = k(x_i, x_j)$$

Hence, $\underline{k(x_i, x_j)} = \langle \Phi(x_i), \Phi(x_j) \rangle$

where $\Phi: \mathcal{X} \rightarrow \mathcal{H}$

$$x_i \rightarrow \Phi(x_i) = k(\cdot, x_i)$$

x-formulation is from x to a very-high-dim space (infinite-dim space).

e.g. Gaussian kernel

$$x \rightarrow e^{-\frac{1}{2\sigma^2} \|\cdot - x\|^2}$$

polynomial kernel

$$x \rightarrow \cdot^T (xx^T) \cdot$$

Final Note

$$f(\cdot) = \sum_i \alpha_i k(\cdot, x_i)$$

$$\langle k(\cdot, x), f(\cdot) \rangle = \sum_i \alpha_i k(x, x_i) = f(x)$$

"reproducing property" - like convolving a signal w/ delta function, which gives back the signal.

\mathcal{H} is called a "Reproducing Kernel Hilbert Space" (RKHS)
(vector space + dot product)

RKHS uniquely specifies the kernel function & vice versa.

(sometimes called a Mercer kernel)

Representer Thm

Empirical Risk: $R_{\text{emp}} = \sum_i L(y_i, f(x_i))$

Regularizer: $\Omega(\|f\|_p)$, $\Omega \geq 0$ & strictly monotonically increasing
 $k(x, x') \in \text{RKHS}$

Then for:

$$f^* = \underset{f}{\operatorname{argmin}} R_{\text{emp}}(f) + \lambda \Omega(\|f\|_p)$$

$$\Rightarrow f^* \text{ has the form } f^*(x) = \sum_{i=1}^n \alpha_i k(x, x_i), \quad f^* \in \text{RKHS}$$

- optm over inf. dim space of functions \rightarrow finite dim space of α_i 's only @
- many ML algorithms fit this framework \rightarrow they can be kernelized

Kernel functions

Kernels on \mathbb{R}^d :

linear kernel: $k(x, x') = x^T x'$
 $k(x, x') = x^T A x'$, A is posdef

poly kernel $k(x, x') = (x^T x')^d \Rightarrow$ purely d order ter
 $k(x, x') = (x^T x' + 1)^d \Rightarrow$ get all terms with order $\leq d$.

Gaussian/RBF: $k(x, x') = e^{-\frac{1}{2\sigma^2} \|x - x'\|^2}$

exponential: $k(x, x') = e^{-\frac{1}{2\sigma^2} \|x - x'\|}$

Combine kernels: $k(x, x') = \underbrace{x^T x'}_{\text{linear}} + \underbrace{e^{-\frac{1}{2\sigma^2} \|x - x'\|^2}}_{\text{RBF}}$
wiggly, straight line.

Kernels on histogram

correlation kernel: $k(x, x') = x^T x' = \sum_{i=1}^d x_i x'_i$

Bhattacharyya kernel: $k(x, x') = \sum_{i=1}^d \sqrt{x_i} \sqrt{x'_i}$

χ^2 -RBF kernel: $k(x, x') = e^{-\frac{1}{2} \chi^2(x, x')}$
 $\chi^2(x, x') = \sum_i \frac{(x_i - x'_i)^2}{\frac{1}{2}(x_i + x'_i)}$

Histogram intersection: $k(x, x') = \sum_i \min(x_i, x'_i)$



Kernels on sets: $X = \{x_1, \dots, x_n\}$
 $X' = \{x'_1, \dots, x'_m\}$ $n \neq m$

Intersection kernel: $k(X, X') = \frac{|X \cap X'|}{2}$ \leftarrow # of common elements

pairwise distance: $k(X, X') = e^{-\frac{1}{2\sigma^2} \sum_{i,j} d(x_i, x'_j)}$

pyramid watch kernel: approx to sum of min. distances btwn points.

Kernels on strings/trees/graphs

$k(x, x') = \sum_s w_s \phi_s(x) \phi_s(x')$, $\phi_s(x) = \begin{matrix} \# \text{ of times} \\ \text{substring } s \text{ appears} \\ \text{in } x. \end{matrix}$
 $w_s > 0$, coefficient (weight)

Kernels on probability densities: $p(x), q(x)$

corr. kernel: $k(p, q) = \int p(x) q(x) dx$

prob product kernel: $k(p, q) = \int p(x)^2 q(x)^2 dx$

Fisher kernel, ...