

(1)

Spring 2014

Note Title

Kernel Trick

11/19/2006

Note that the final classifier depends on \underline{x} only by dot products.

EG. $\underline{x} \cdot \underline{x}_\mu$ in final classifier.

$\underline{x}_\mu \cdot \underline{x}_\nu$ in the dual energy.

This motivates the Kernel Trick

Compute features $\underline{\varphi}(\underline{x})$ and reformulate the problem in feature space.

→ i.e. seek a classifier of form

$$\text{Sign}(\underline{c} \cdot \underline{\varphi}(\underline{x}) + b).$$

Replace \underline{x} by $\underline{\varphi}(\underline{x})$ everywhere in the primal & dual formulation.

Then the classifier only depends on the dot product of the $\underline{\varphi}(\underline{x})$'s.

On the kernel $k(\underline{x}, \underline{x}') = \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}')$.

(2)

Why does this help?

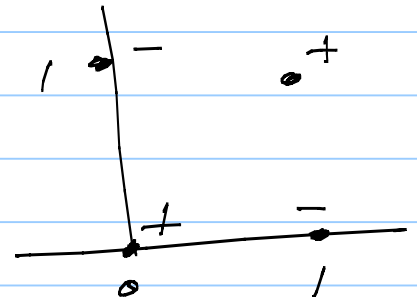
First, features can make it possible to classify data by hyperplanes.

Example Logical X-OR, $\underline{x} = (x_1, x_2)$ $x_j \in \{\pm 1\}$
 $w_j \in \{\pm 1\}$

The X-OR (exclusive or) requires a decision rule $\alpha(\underline{x})$ s.t.

$$\alpha(1, 1) = \alpha(-1, -1) = 1$$

$$\alpha(1, -1) = \alpha(-1, 1) = -1$$



Impossible to find a linear classifier to do this.

But define a feature $\phi(x_1, x_2) = x_1 x_2$.

Then we can classify the data by a linear classifier (i.e. a point) in 1-dimension

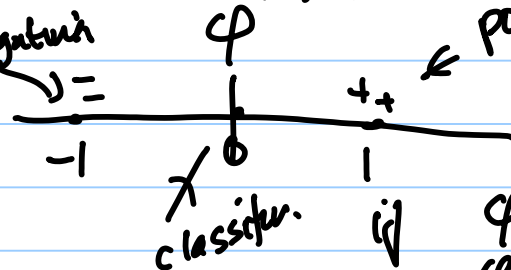
$x_1 x_2 = 1$ for +ve's

$x_1 x_2 = -1$ for -ve's

$$1 \times 1 = 1$$

$$-1 \times 1 = -1$$

negatives



positives

$$1 \times 1 = 1$$

$$-1 \times 1 = -1$$

$\phi(x_1, x_2) > 0$, +ve
 $\phi(x_1, x_2) < 0$, -ve

(3) This example is "cheating" because we guess the right feature $x_1 x_2$.

More generally, we could set

$$\underline{P}(x_1, x_2) = (x_1, x_2, x_1 x_2, x_1^2, x_2^2, x_1 x_2 x_3, \dots)$$

i.e. we specify a higher-dimensional set of features (higher than $\dim(x_1, x_2) = 2$).

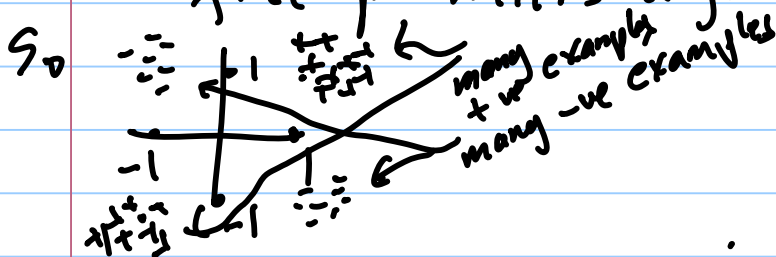
$$\underline{\lambda} = (\lambda_1, \lambda_2, \dots)$$

$$\text{Then } \underline{\lambda} \cdot \underline{P}(x) = \lambda_1 x_1 + \lambda_2 x_2 + \lambda_3 x_1 x_2 + \dots$$

The SVM algorithm selects $\underline{\lambda}$ to separate the positive and negative examples - maybe by setting $\lambda_3 = 1, \lambda_1 = \lambda_2 = \dots = 0$ which gives us our "cheating" solution.

But there may be other equally good, or better, solutions $\underline{\lambda}$.

Note: we need more examples in this case, because we need many more examples than free parameters (eg. dimension of $\underline{\lambda}$)



Bottom-Line: By using a high-dimensional function $\underline{P}(x)$ (i.e. so that $\dim P \gg \dim x$) it is more likely that we can find a separating plane.

(4) Second, we do not need to specify the features $\underline{\varphi}(x)$ explicitly, we only need to specify the kernel

$$K(\underline{x}, \underline{x}') = \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}')$$

Remember:

the dual problem reduces to maximizing $L_d(K, \underline{\alpha}) = \sum_{\mu} \alpha_{\mu} - \frac{1}{2} \sum_{\mu, \nu} \alpha_{\mu} \alpha_{\nu} y_{\mu} y_{\nu} \underline{\varphi}(\underline{x}_{\mu}) \cdot \underline{\varphi}(\underline{x}_{\nu})$
 $= \sum_{\mu} \alpha_{\mu} - \frac{1}{2} \sum_{\mu, \nu} \alpha_{\mu} \alpha_{\nu} y_{\mu} y_{\nu} K(\underline{x}_{\mu}, \underline{x}_{\nu})$

The solution $\underline{\hat{a}} = \sum_{\mu} \hat{\alpha}_{\mu} y_{\mu} \underline{\varphi}(\underline{x}_{\mu})$

$$\begin{aligned} \underline{\hat{a}} \cdot \underline{\varphi}(\underline{x}) &= \sum_{\mu} \hat{\alpha}_{\mu} y_{\mu} \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}_{\mu}) \\ &= \sum_{\mu} \hat{\alpha}_{\mu} y_{\mu} K(\underline{x}, \underline{x}_{\mu}). \end{aligned}$$

Hence classifier is

$$\text{sign} \left\{ \sum_{\mu} \hat{\alpha}_{\mu} y_{\mu} K(\underline{x}, \underline{x}_{\mu}) + \hat{b} \right\}$$

only depends on $K(\cdot, \cdot)$

(5)

What Kernels to Use?

There are many choices of kernels. The difficulty is knowing which one to use. As always, cross-validation is useful for checking whether a kernel can generalize

$$K(\underline{x}, \underline{x}') = \langle 1 + \underline{x} \cdot \underline{x}' \rangle^d$$

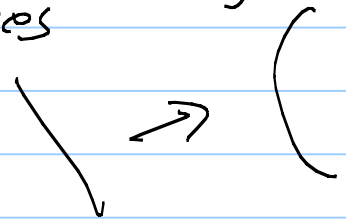
$$K(\underline{x}, \underline{x}') = e^{-\frac{1}{\sigma^2} \|\underline{x} - \underline{x}'\|^2}$$

$$K(\underline{x}, \underline{x}') = \tanh \langle C_1 \underline{x} \cdot \underline{x}' + C_2 \rangle$$

Choice of best kernel is problem dependent.

Some kernels \rightarrow eg. $\langle 1 + \underline{x} \cdot \underline{x}' \rangle^d$ naturally generalize the idea of hyperplanes

Others \rightarrow eg. $e^{-\frac{1}{\sigma^2} \|\underline{x} - \underline{x}'\|^2}$ are similar to nearest neighbours.



(6) There are two important types of Kernel.

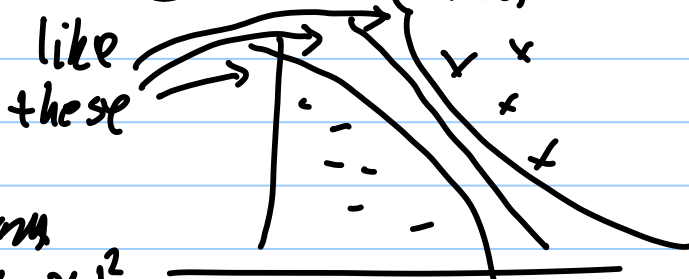
(1) Polynomial kernels

E.g. $K(\underline{x}_1, \underline{x}_2) = \underline{x}_1 \cdot \underline{x}_2$

or $K(\underline{x}_1, \underline{x}_2) = \underline{x}_1 \cdot \underline{x}_2 + |\underline{x}_1|^2 + |\underline{x}_2|^2 + (\underline{x}_1 \cdot \underline{x}_2)^2$.

(separating by plane)
previous lecture.

These give boundaries like these



(2) Radial Basis Functions

$$K(\underline{x}_1, \underline{x}_2) = e^{-\lambda |\underline{x}_1 - \underline{x}_2|^2}$$

note $K(\underline{x}_1, \underline{x}_2)$ decreases as $|\underline{x}_1 - \underline{x}_2|^2$ increases.

This kernel behaves like nearest neighbor! Only nearby data points affect classification.

Recall classifier is of form:

$$\hat{y}(\underline{x}) = \text{sign} \left\{ \sum_{\mu} \hat{\alpha}_{\mu} y_{\mu} K(\underline{x}, \underline{x}_{\mu}) + \hat{b} \right\}$$

(The $\hat{\alpha}_{\mu}$'s are non-zero only for support vectors)

Because $K(\underline{x}, \underline{x}_{\mu})$ decreases rapidly as $|\underline{x} - \underline{x}_{\mu}|$ increases, it follows that the estimate $\hat{y}(\underline{x})$ depends only on the datapoints near \underline{x} . It depends on the weighted sum of the $\hat{\alpha}_{\mu}$'s (recall $\hat{\alpha}_{\mu} > 0$).

Note: Can also specify a set of kernels $K_1(\cdot, \cdot), \dots, K_n(\cdot, \cdot)$ parameters $(\alpha_1, \dots, \alpha_n)$, Kernel = $\sum_{i=1}^n \alpha_i K_i(\cdot, \cdot)$
use cross-validation to select the best (α 's).

(7) Geometric Interpretation

The kernel can be used to specify a distance measure

Recall, Euclidean distance:

$$|\underline{x}_1 - \underline{x}_2|^2 = \underline{x}_1 \cdot \underline{x}_1 - 2 \underline{x}_1 \cdot \underline{x}_2 + \underline{x}_2 \cdot \underline{x}_2$$

Similarly $|\underline{\varphi}(\underline{x}_1) - \underline{\varphi}(\underline{x}_2)|^2 = \underline{\varphi}(\underline{x}_1) \cdot \underline{\varphi}(\underline{x}_1) + \underline{\varphi}(\underline{x}_2) \cdot \underline{\varphi}(\underline{x}_2) - 2 \underline{\varphi}(\underline{x}_1) \cdot \underline{\varphi}(\underline{x}_2)$

$$= k(\underline{x}_1, \underline{x}_1) + k(\underline{x}_2, \underline{x}_2) - 2k(\underline{x}_1, \underline{x}_2)$$

Euclidean distance in $\underline{\varphi}$ -space

Special case, if $k(\underline{x}_1, \underline{x}_2) = \underline{x}_1 \cdot \underline{x}_2$, then we recover $|\underline{x}_1 - \underline{x}_2|^2$.

Bottom Line: Think of $\underline{x} \rightarrow \underline{\varphi}(\underline{x})$ as a transformation on the geometry, which is the same as specifying a new distance - $k(\underline{x}_1, \underline{x}_1) + k(\underline{x}_2, \underline{x}_2) - 2k(\underline{x}_1, \underline{x}_2)$ - in the original \underline{x} space

Note: the RBF kernel has interesting geometry

$$k(\underline{x}_1, \underline{x}_1) + k(\underline{x}_2, \underline{x}_2) - 2k(\underline{x}_1, \underline{x}_2) = 2(1 - e^{-|\underline{x}_1 - \underline{x}_2|^2})$$

So as $|\underline{x}_1 - \underline{x}_2|$ increases, the distance between datapoints in $\underline{\varphi}$ -space tends to a maximum value of 2.

(7)

When do kernels correspond to features?

i.e. if we specify $K(\underline{x}, \underline{x}')$, is it equal to $\underline{\phi}(\underline{x}) \cdot \underline{\phi}(\underline{x}')$ for some features $\underline{\phi}(\underline{x})$?

There are theoretical results \rightarrow Mercer's Theorem.

Here we do two things:

(1) give some partial understanding, by showing that this relates to the Spectral Theorem in Linear Algebra.

(2) explaining why this result matters in practice, giving alternatives to SVM.

(part i)

Consider $K(\underline{x}, \underline{y})$

allow $\underline{x}, \underline{y}$ to take a finite set of possible values

$\{z_1, z_2, \dots, z_n\}$, evenly spaced. K^z values

e.g. grid $\begin{matrix} \uparrow x_2 & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{matrix}$

Then consider

$$K(z_i, z_j) = K_{ij}$$

This is an $n \times n$ matrix \underline{K} with positive terms ($K_{ij} > 0$)

By the Spectral theorem (linear analysis)

we can express it as $K_{ij} = \sum_{\mu=1}^n \lambda_{\mu} e_{\mu}^i e_{\mu}^j$

where λ, e are the eigenvalues & eigenvectors of \underline{K} ,

$$\text{i.e. } \sum_j K_{ij} e_{\mu}^j = \lambda_{\mu} e_{\mu}^i$$

(9) The eigenvalues $\lambda^n > 0$ (because $K_{ij} > 0$ for all i, j)
 so $K_{ij} = \sum_{\mu=1}^N (\lambda_{\mu}^{\frac{1}{2}} e_{\mu}^i) (\lambda_{\mu}^{\frac{1}{2}} e_{\mu}^j)$

i.e. $K(\underline{z}_i, \underline{z}_j) = \underline{\varphi}(\underline{z}_i) \cdot \underline{\varphi}(\underline{z}_j)$
 where $\underline{\varphi}(\underline{z}_i) = (\lambda_1^{\frac{1}{2}} e_i^1, \lambda_2^{\frac{1}{2}} e_i^2, \dots, \lambda_N^{\frac{1}{2}} e_i^N)$

This answers the question - when can we find the φ s -
 if we restrict the $\underline{x}, \underline{y}$ to only take values $\underline{z}_1, \dots, \underline{z}_N$

Instead, we have to consider $K(\underline{x}, \underline{y})$ for
 continuous values of \underline{x} and \underline{y} .

This requires - functional analysis, eigenvalue
 finding eigenfunctions and eigenvalues ← eigenfunction

$$\int K(\underline{x}, \underline{y}) \underline{e}(\underline{y}) d\underline{y} = \lambda \underline{e}(\underline{x})$$

Functional analysis extends Linear Algebra to functions.
 Out of scope of this course. But Mercer's theorem is like a
 generalization of the spectral theorem to Functional Analysis

(part 2) why does this matter? New Algorithms

Consider the spectral theorem again. Order the
 eigenvalues so that $\lambda_1 > \lambda_2 > \dots > \lambda_N$.

Then - depending on $K(\underline{z}_i, \underline{z}_j)$ - we may get a
 very good approximation to K by using only a few terms in
 the spectral theorem. Those with the biggest λ 's

i.e. replace $K(\underline{z}_i, \underline{z}_j)$ by $\sum_{\mu=1}^n \lambda_{\mu} e_{\mu}^i e_{\mu}^j$
where $n \ll N$

This approximation can be very useful for some kernel methods. (It means we need only
 evaluate $\sum_{\mu=1}^n \lambda_{\mu}^{-1/2} e_{\mu}^i$. (Full method requires Functional Analysis))
(so that $|\lambda_{\mu}^{-1/2}|$ is small, $\mu > n$)