

Machine Learning

Lecture: Features and Kernels

Hiroshi Shimodaira and Hao Tang

2026

Ver. 1.0

Questions you should be able to answer after this week

- Feature transformation
- Meaning of linear classifiers
- Non-linear SVMs
- kernel tricks
- Mercer kernels

How to resolve a linearly non-separable case?

Feature transformation / mapping: $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$

$$h(\mathbf{x}) = \begin{cases} -1 & \text{if } \mathbf{w}^\top \mathbf{x} + b < 0 \\ +1 & \text{if } \mathbf{w}^\top \mathbf{x} + b \geq 0 \end{cases} = \text{sgn}(\mathbf{w}^\top \mathbf{x} + b) \quad (1)$$

↓

$$h(\mathbf{x}) = \begin{cases} -1 & \text{if } \mathbf{w}^\top \phi(\mathbf{x}) < 0 \\ +1 & \text{if } \mathbf{w}^\top \phi(\mathbf{x}) \geq 0 \end{cases} = \text{sgn}(\mathbf{w}^\top \phi(\mathbf{x})) \quad (2)$$

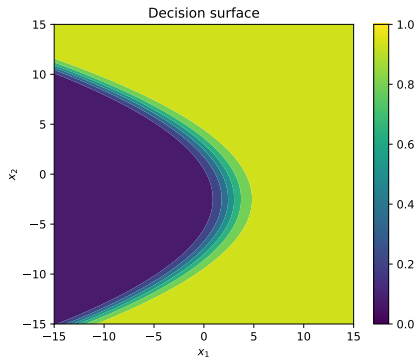
Logistic regression with feature transformation

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{1 + \exp(-y(\mathbf{w}^\top \mathbf{x} + b))} \quad (3)$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{1 + \exp(-y(\mathbf{w}^\top \phi(\mathbf{x})))} \quad (4)$$

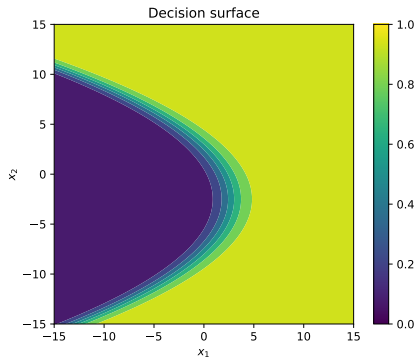
Feature transformation - examples

$$(x_1, x_2) \rightarrow (x_1, x_2, x_2^2)$$

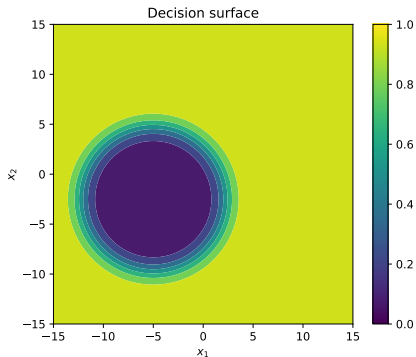


Feature transformation - examples

$$(x_1, x_2) \rightarrow (x_1, x_2, x_2^2)$$

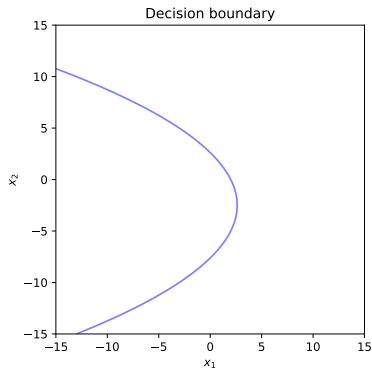


$$(x_1, x_2) \rightarrow (x_1, x_2, x_1^2, x_2^2)$$

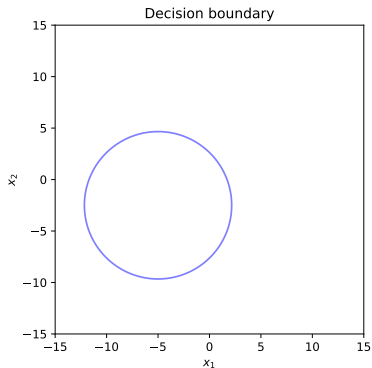


Feature transformation - examples

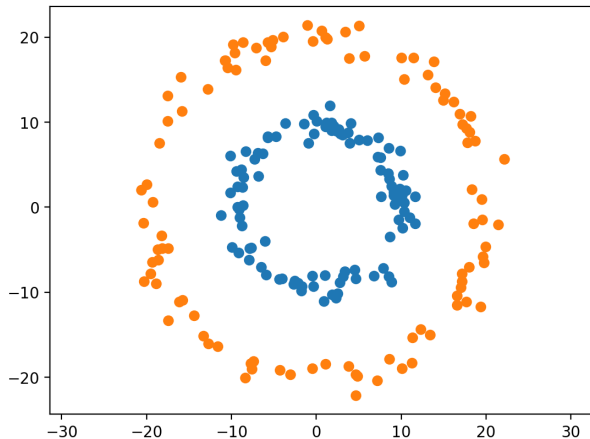
$$(x_1, x_2) \rightarrow (x_1, x_2, x_2^2)$$



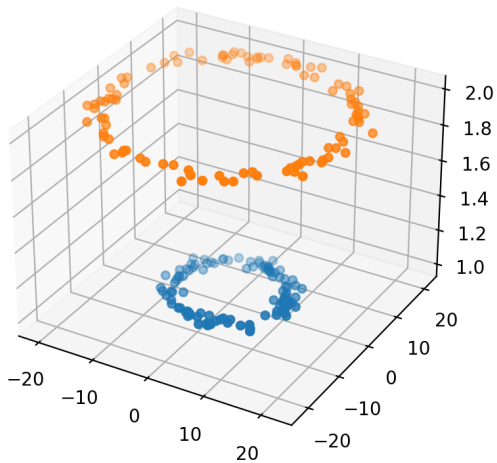
$$(x_1, x_2) \rightarrow (x_1, x_2, x_1^2, x_2^2)$$



Two-circle example



Two-circle example



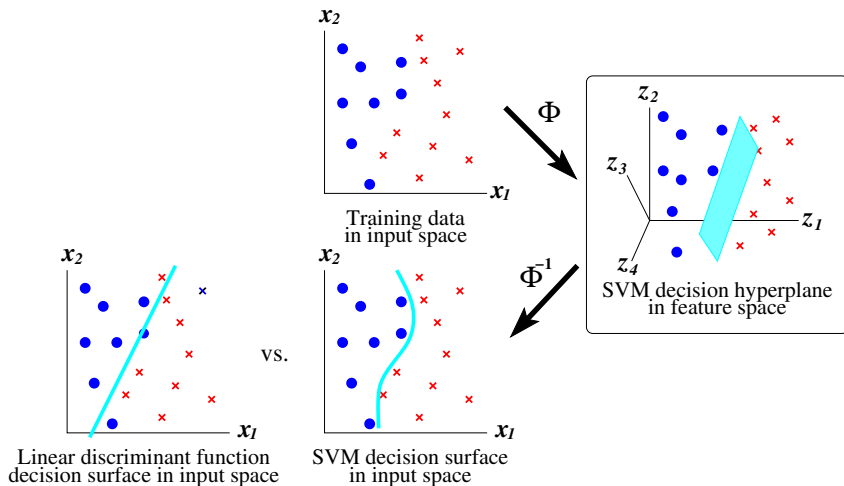
What is it meant by linear classifiers?

- A linear classifier is linear in the parameters \mathbf{w} , **not** in the features.
(i.e., classification based on a linear combination between \mathbf{w} and \mathbf{x} .)
- A linear classifier can have arbitrary nonlinear features.

Should we consider very complex transformation?

- Not necessarily so.
- Complex models may **overfit** the training data and may not **generalise** very well.
- We will come back to this in some lectures later.

Non-linear SVM



Non-linear SVM (*cont.*)

- Conceptual steps to construct a non-linear SVM
 - Step 1 Transform \mathbf{x} to $\phi(\mathbf{x})$ in a high-dimensional space (feature space)
 - Step 2 Train a SVM in the feature space
 - Step 3 Classify data in the feature space

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i^* y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + w_0$$

Non-linear SVM (*cont.*)

- Conceptual steps to construct a non-linear SVM

Step 1 Transform \mathbf{x} to $\phi(\mathbf{x})$ in a high-dimensional space (feature space)

Step 2 Train a SVM in the feature space

Step 3 Classify data in the feature space

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i^* y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + w_0$$

- Instead of applying the non-linear transformation and carrying out calculation in the feature space, use a kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$ such that

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (\text{cf. 'kernel trick'})$$

Non-linear SVM (cont.)

- Conceptual steps to construct a non-linear SVM

Step 1 Transform \mathbf{x} to $\phi(\mathbf{x})$ in a high-dimensional space (feature space)

Step 2 Train a SVM in the feature space

Step 3 Classify data in the feature space

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i^* y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + w_0$$

- Instead of applying the non-linear transformation and carrying out calculation in the feature space, use a kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$ such that

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (\text{cf. 'kernel trick'})$$

$$L(\alpha, \xi) = -\frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^N \alpha_i$$

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}) + w_0$$

Kernel functions for SVM

An example of kernel that maps data to a feature space explicitly

$$\begin{aligned}k(\mathbf{a}, \mathbf{b}) &\triangleq (1 + \mathbf{a}^T \mathbf{b})^2 = (1 + a_1 b_1 + a_2 b_2)^2 \\&= 1 + 2a_1 b_1 + 2a_2 b_2 + a_1^2 b_1^2 + 2a_1 b_1 a_2 b_2 + a_2^2 b_2^2 \\&= (1, \sqrt{2}a_1, \sqrt{2}a_2, a_1^2, \sqrt{2}a_1 a_2, a_2^2)^\top (1, \sqrt{2}b_1, \sqrt{2}b_2, b_1^2, \sqrt{2}b_1 b_2, b_2^2) \\&= \phi(\mathbf{a})^T \phi(\mathbf{b})\end{aligned}$$

Kernel functions for SVM

An example of kernel that maps data to a feature space explicitly

$$\begin{aligned}k(\mathbf{a}, \mathbf{b}) &\triangleq (1 + \mathbf{a}^T \mathbf{b})^2 = (1 + a_1 b_1 + a_2 b_2)^2 \\&= 1 + 2a_1 b_1 + 2a_2 b_2 + a_1^2 b_1^2 + 2a_1 b_1 a_2 b_2 + a_2^2 b_2^2 \\&= (1, \sqrt{2}a_1, \sqrt{2}a_2, a_1^2, \sqrt{2}a_1 a_2, a_2^2)^T (1, \sqrt{2}b_1, \sqrt{2}b_2, b_1^2, \sqrt{2}b_1 b_2, b_2^2) \\&= \phi(\mathbf{a})^T \phi(\mathbf{b})\end{aligned}$$

Popular kernels

Kernel	$k(\mathbf{x}_i, \mathbf{x}_j)$
Polynomial	$(1 + \langle \mathbf{x}_i, \mathbf{x}_j \rangle)^d$
Radial basis function (RBF)	$\exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right), \sigma > 0$ (bandwidth)
Hyperbolic tangent	$\tanh(\kappa_1 \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \kappa_2), \kappa_1 > 0, \kappa_2 < 0$

where $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ is an inner product (e.g. dot product) between \mathbf{x}_i and \mathbf{x}_j .

Making kernels

How can we ensure if a kernel works as an inner product in a feature space?

It should satisfy:

- $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle = k(\mathbf{x}', \mathbf{x})$
- $k(\mathbf{x}, \mathbf{x}')^2 \leq k(\mathbf{x}, \mathbf{x}) k(\mathbf{x}', \mathbf{x}')$

Making kernels

How can we ensure if a kernel works as an inner product in a feature space?

It should satisfy:

- $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle = k(\mathbf{x}', \mathbf{x})$
- $k(\mathbf{x}, \mathbf{x}')^2 \leq k(\mathbf{x}, \mathbf{x}) k(\mathbf{x}', \mathbf{x}')$
- The Gram matrix $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))$, which is a N -by- N matrix, is positive definite.

→ Mercer kernels or positive-definite kernel

Making kernels

How can we ensure if a kernel works as an inner product in a feature space?

It should satisfy:

- $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \langle \phi(\mathbf{x}'), \phi(\mathbf{x}) \rangle = k(\mathbf{x}', \mathbf{x})$
- $k(\mathbf{x}, \mathbf{x}')^2 \leq k(\mathbf{x}, \mathbf{x}) k(\mathbf{x}', \mathbf{x}')$
- The Gram matrix $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))$, which is a N -by- N matrix, is positive definite.

→ Mercer kernels or positive-definite kernel

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ & \ddots & \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} \quad (5)$$

$$\mathbf{a}^T \mathbf{K} \mathbf{a} > 0 \quad \forall \mathbf{a} \in \mathbb{R}^N \setminus \{\mathbf{0}\} \quad (6)$$

$$\text{i.e., } \sum_{i=1}^N \sum_{j=1}^N a_i a_j k(\mathbf{x}_i, \mathbf{x}_j) > 0 \quad \forall a_i \in \mathbb{R} \setminus \{0\} \quad (7)$$

Mercer's theorem

Suppose k is a continuous symmetric non-negative definite kernel (i.e. Mercer kernel), then k can be expressed as:

$$k(\mathbf{x}, \mathbf{z}) = \sum_{i=1}^{\infty} \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{z})$$

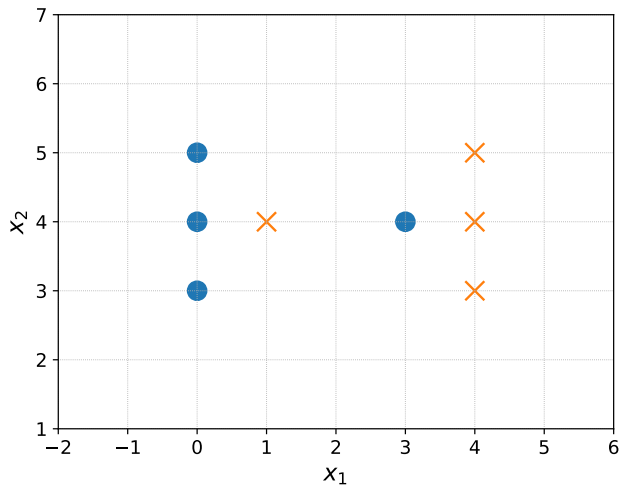
where $\{\phi_i\}$ are eigen-functions, $\|\phi_i\| = 1$, and $\{\lambda_i\}$ are positive eigenvalues $\lambda_i > 0$.

Making kernels from kernels

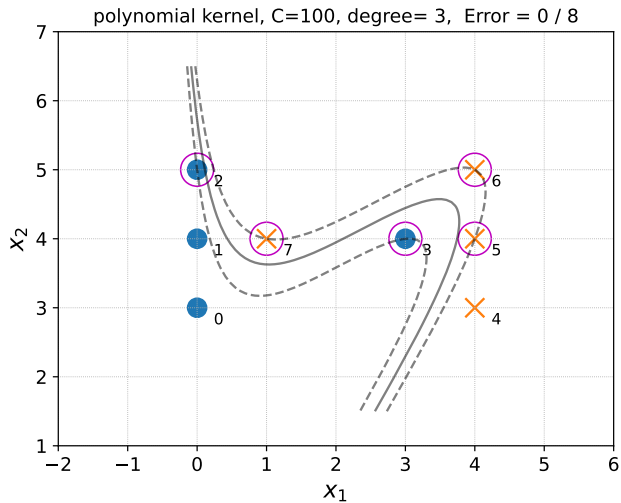
Letting k_1 , k_2 , and k_3 are kernels, we can create a new kernel k .

- $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = ak_1(\mathbf{x}, \mathbf{x}')$, $a > 0$
- $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$
- $k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$
- $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T B \mathbf{x}'$, where B is a N -by- N positive-definite matrix

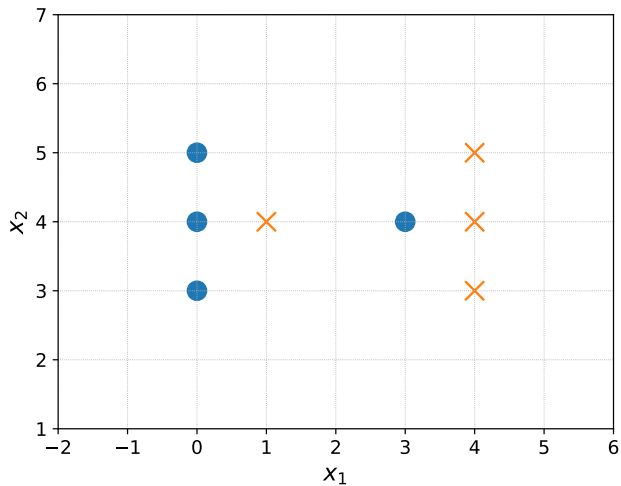
Example – SVM with a polynomial kernel



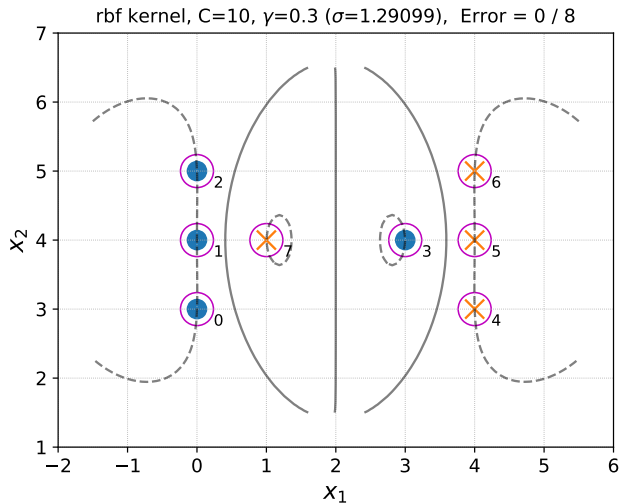
Example – SVM with a polynomial kernel



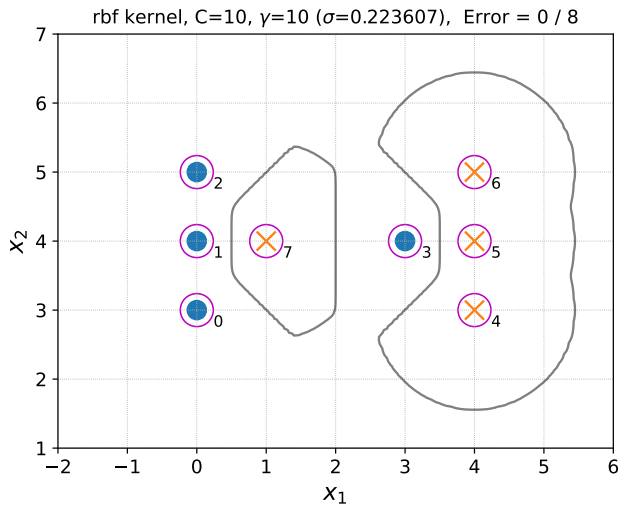
Example – SVM with an RBF kernel



Example – SVM with an RBF kernel



Example – SVM with an RBF kernel



Quizzes

- Considering a kernel $k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^2$, find the corresponding transformation function $\phi(\mathbf{a})$ when $\mathbf{a} \in \mathbb{R}^2$.
- Discuss how the dimensionality of the input vectors and the size of the training set affect the time and memory required to train a linear SVM. Do you expect the same for an SVM with the RBF kernel?
- How many dimensions does the feature space induced by the RBF kernel have? Explain your answer.