

Machine Learning Fundamentals - LS2026

Kernel functions

Czech Technical University in Prague
V. Franc

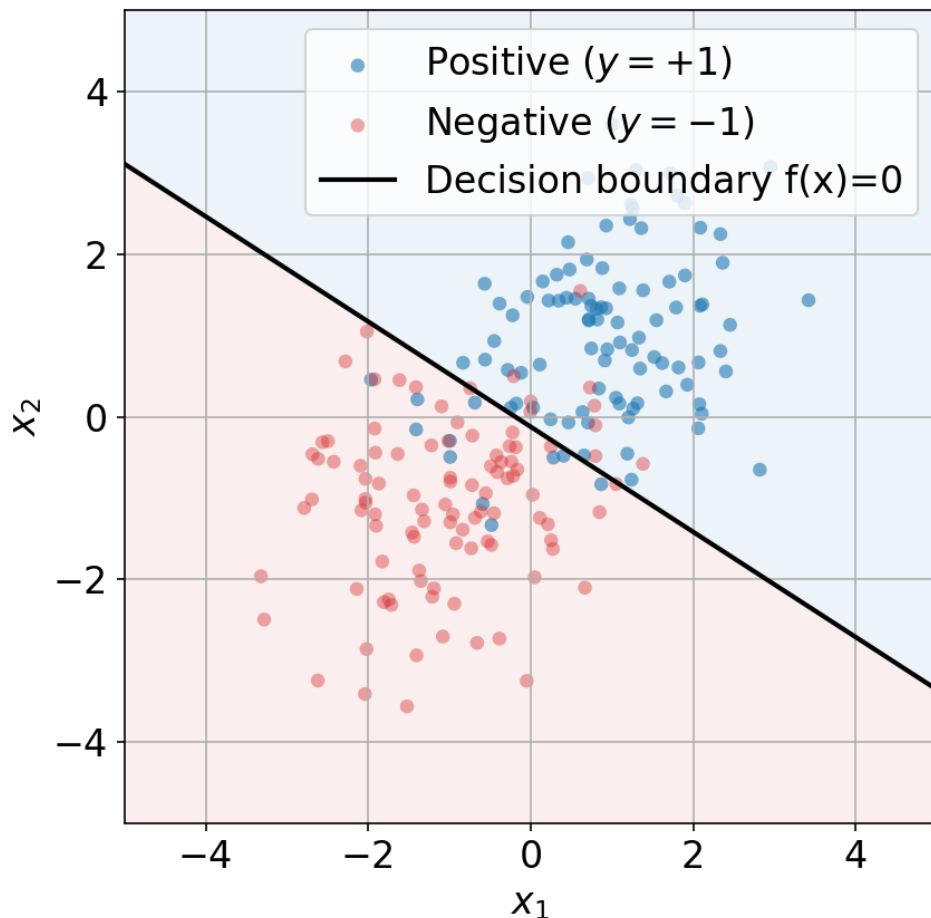
Linear models

- Linear classifier $h(\mathbf{x}; \mathbf{w}, b) = \text{sign}(f(\mathbf{x}; \mathbf{w}, b))$ uses (affine) discriminant function

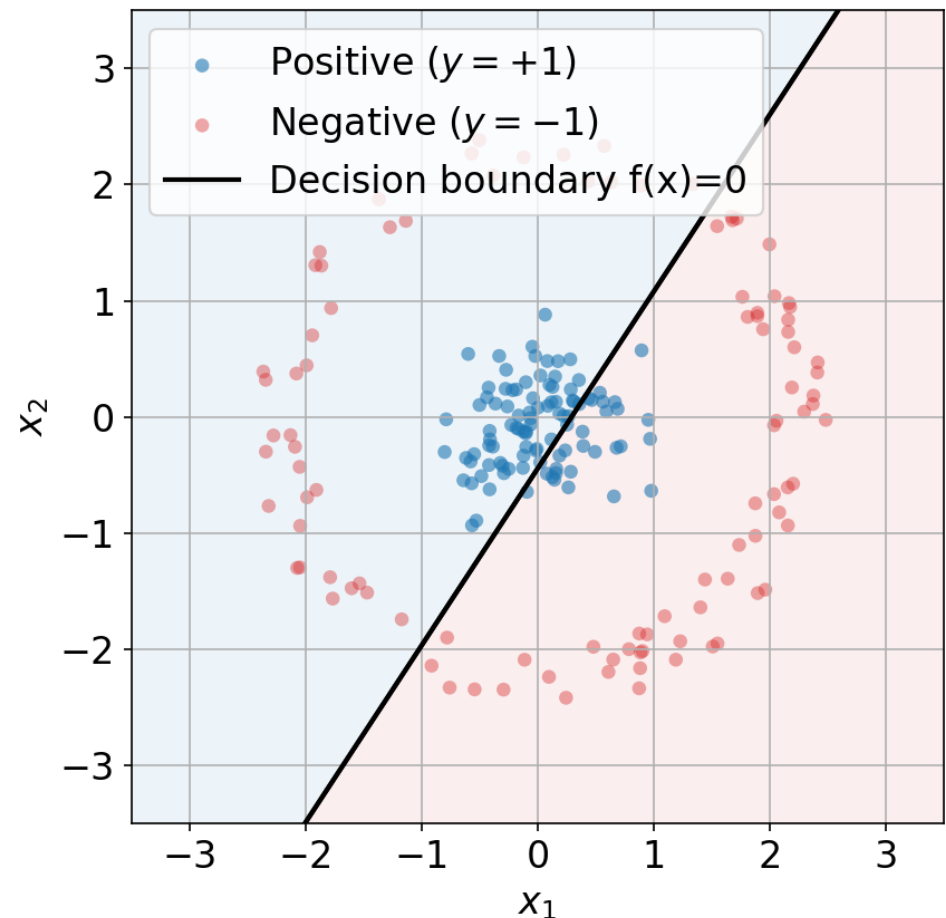
$$f(\mathbf{x}; \mathbf{w}, b) = \langle \mathbf{w}, \mathbf{x} \rangle + b = \sum_{i=1}^d w_i x_i + b$$

where $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ are parameters and $\mathbf{x} = [x_1, \dots, x_d] \in \mathbb{R}^d$ is an input.

Example: Linear model is good



Linear model is not good



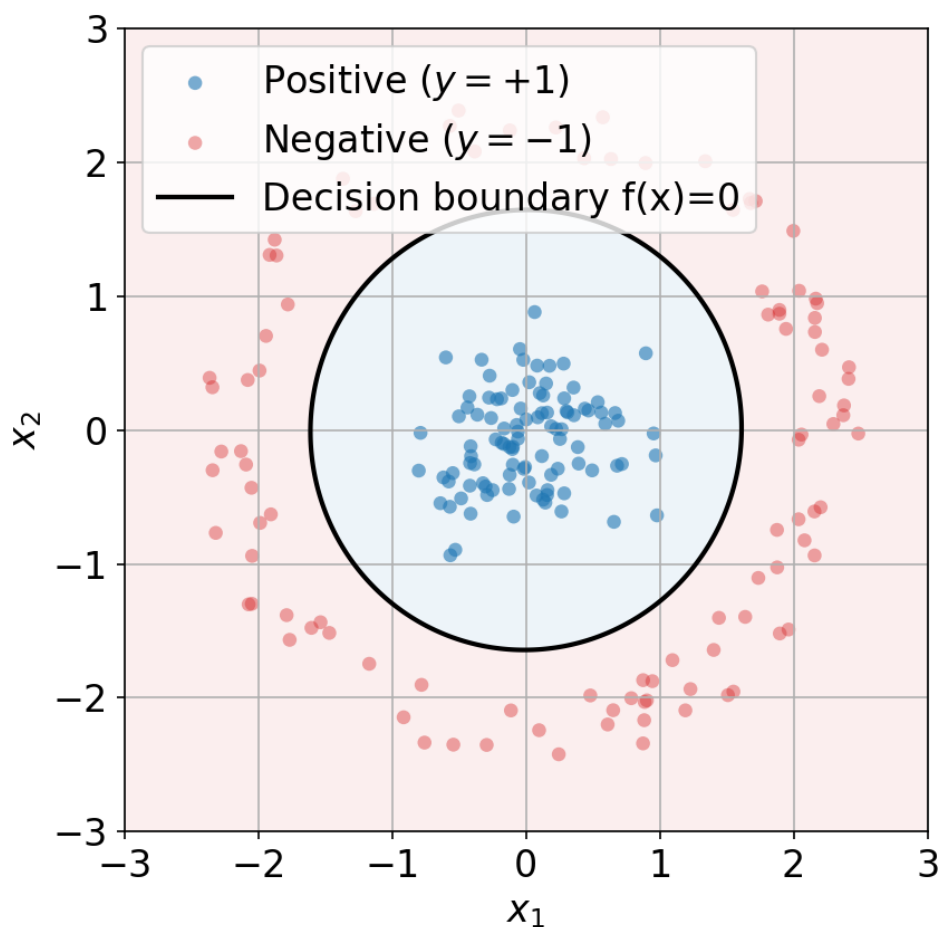
Non-linear feature mapping

◆ Linear classifier $h(\mathbf{x}; \mathbf{w}, b) = \text{sign}(f(\mathbf{x}; \mathbf{w}, b))$ using non-linear features:

$$\begin{aligned}
 f(\mathbf{x}; \mathbf{w}, b) &= w_1\phi_1(\mathbf{x}) + w_2\phi_2(\mathbf{x}) + w_3\phi_3(\mathbf{x}) + b \\
 &= w_1 x_1^2 + w_2\sqrt{2} x_1 x_2 + w_3 x_2^2 + b = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b
 \end{aligned}$$

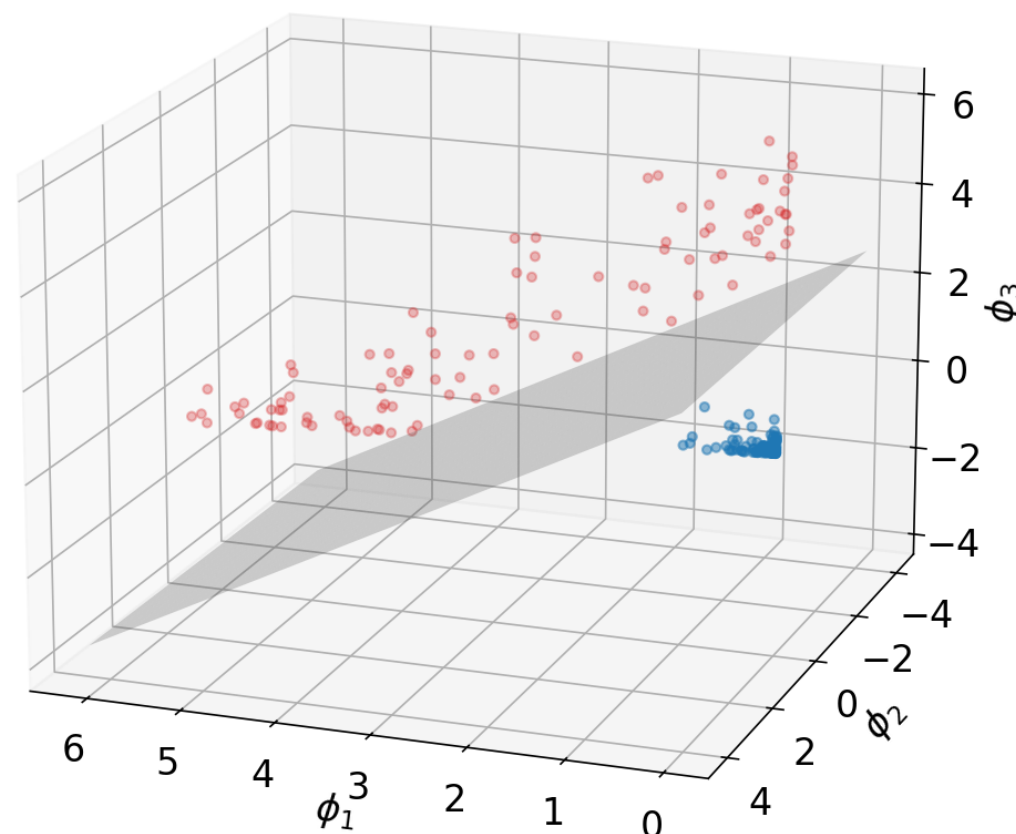
Original features

$$\mathbf{x} = [x_1, x_2] \in \mathbb{R}^2$$



New features

$$\phi(\mathbf{x}) = [x_1^2, \sqrt{2}x_1x_2, x_2^2] \in \mathbb{R}^3$$



Perceptron Learning Algorithm

Task: Given training set $T_m = ((x_i, y_i) \in \mathcal{X} \times \{+1, -1\} \mid i = 1, \dots, m)$ linearly separable under the feature map $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$, find

$$h(x; \mathbf{w}, b) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$$

with zero training error.

Perceptron Learning Algorithm (PLA)

1. $\mathbf{w} \leftarrow \mathbf{0}, b \leftarrow 0$
2. Find $(x_u, y_u) \in T_m$ such that:

$$y_y \neq \hat{y}_u = \text{sign}(\langle \mathbf{w}, \phi(x_u) \rangle + b)$$

3. If such (x_u, y_u) does not exist, return (\mathbf{w}, b) , otherwise update parameters:

$$\mathbf{w} \leftarrow \mathbf{w} + y_u \phi(x_u), \quad b \leftarrow b + y_u$$

and go to step 2.

Perceptron Learning Algorithm

Example: Use the Perceptron Learning Algorithm to learn a quadratic classifier $h(\mathbf{x}; \mathbf{w}, b) = \text{sign}(f(\mathbf{x}; \mathbf{w}, b))$, where

$$\begin{aligned} f(\mathbf{x}; \mathbf{w}, b) &= x_1^2 w_1 + \cdots + x_d^2 w_d + \sqrt{2} x_1 x_2 w_{d+1} + \cdots + \sqrt{2} x_{d-1} x_d w_n + b \\ &= \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b \end{aligned}$$

with $\mathbf{w} \in \mathbb{R}^n$, $n = \frac{d(d+1)}{2}$, and feature map

$$\phi(\mathbf{x}) = [x_1^2, \dots, x_d^2, \sqrt{2} x_1 x_2, \dots, \sqrt{2} x_{d-1} x_d].$$

- ◆ Assume the number of input features is $d = 1,000,000$.
- ◆ Storing a single input vector $\mathbf{x} \in \mathbb{R}^d$ requires $d \cdot 4$ bytes ≈ 4 MB.
- ◆ Explicitly computing or storing $\phi(\mathbf{x})$ requires $\mathcal{O}(n)$ operations and memory, where

$$n = \frac{d(d+1)}{2} = \frac{1,000,000 \cdot 1,000,001}{2} \approx 5 \times 10^{11}.$$

- ◆ This corresponds to roughly 5×10^{11} multiplications/additions and

$$n \cdot 4 \text{ bytes} \approx 1.8 \text{ TiB of memory.}$$

Polynomial kernel of degree 2

- ◆ Consider quadratic feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^n$:

$$\phi(\mathbf{x}) = [x_1^2, \dots, x_d^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{d-1}x_d]$$

which maps d dimensional feature space to $n = \frac{d(d+1)}{2}$ feature space.

- ◆ Polynomial kernel of degree 2 computes the dot product in this feature space:

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

- ◆ In case of $d = 2$, we have

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= (x_1 x'_1 + x_2 x'_2)^2 \\ &= x_1^2 x'^2_1 + 2 x_1 x_2 x'_1 x'_2 + x_2^2 x'^2_2 \\ &= \left\langle [x_1^2, \sqrt{2} x_1 x_2, x_2^2], [x'^2_1, \sqrt{2} x'_1 x'_2, x'^2_2] \right\rangle \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \end{aligned}$$

Why kernels matter: polynomial kernel of degree 2

Example: Let the input dimension be $d = 1,000,000$, and consider the quadratic feature map induced by the degree-2 polynomial kernel.

- ◆ Storing a single input vector $\mathbf{x} \in \mathbb{R}^d$ requires $d \cdot 4$ bytes ≈ 4 MB
- ◆ The corresponding quadratic feature map is

$$\phi(\mathbf{x}) = [x_1^2, \dots, x_d^2, \sqrt{2} x_1 x_2, \dots, \sqrt{2} x_{d-1} x_d]$$

Its dimension is $n = \frac{d(d+1)}{2} = \frac{1,000,000 \cdot 1,000,001}{2} \approx 5 \times 10^{11}$.

- ◆ Explicitly computing or storing $\phi(\mathbf{x})$ therefore requires $\mathcal{O}(n)$ operations and memory, that is roughly 5×10^{11} products/additions and $n \cdot 4$ bytes ≈ 1.8 TiB.
- ◆ In contrast, the degree-2 polynomial kernel can be evaluated directly as

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

which requires only $\mathcal{O}(d)$ operations, i.e. about 1,000,000 multiplications/additions.

The kernel computes the dot product in a huge feature space *without ever constructing the feature vector explicitly.*

Kernel perceptron

Task: Given training set $T_m = ((x_i, y_i) \in \mathcal{X} \times \{+1, -1\} \mid i = 1, \dots, m)$ linearly separable under the feature map $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$, find $h(x) = \text{sign}(\langle \phi(x), \mathbf{w} \rangle + b)$ with zero training error.

Key idea: Represent the parameter vector as a linear combination of the training inputs

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \phi(x_i)$$

Perceptron Learning Algorithm

1. $\mathbf{w} \leftarrow \mathbf{0}, b \leftarrow 0$
2. Find $(x_u, y_u) \in T_m$ such that:

$$y_u \neq \hat{y}_u = \text{sign}(\langle \mathbf{w}, \phi(x_u) \rangle + b)$$

3. If such (x_u, y_u) does not exist, return (\mathbf{w}, b) , otherwise update parameters:

$$\mathbf{w} \leftarrow \mathbf{w} + y_u \phi(x_u), \quad b \leftarrow b + y_u$$

and go to step 2.

PLA using only dot products

1. $\alpha \leftarrow \mathbf{0}, b \leftarrow 0$
2. Find $(x_u, y_u) \in T_m$ such that:

$$y_u \neq \hat{y}_u = \text{sign}\left(\sum_{i=1}^m \alpha_i \langle \phi(x), \phi(x_i) \rangle + b\right)$$

3. If such (x_u, y_u) does not exist return (α, b) , otherwise update parameters:

$$\alpha_u \leftarrow \alpha_u + y_u \quad b \leftarrow b + y_u$$

and go to step 2.

Kernel perceptron

Kernel Perceptron

Input: training sequence $T_m = ((x_i, y_i) \in \mathcal{X} \times \{+1, -1\} \mid i = 1, \dots, m)$; kernel function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$ under some feature map ϕ .

Output: kernel classifier represented by the weights $\alpha = [\alpha_1, \dots, \alpha_m] \in \mathbb{R}^m$ and bias $b \in \mathbb{R}$:

$$h(x; \alpha, b) = \text{sign} \left(\sum_{i=1}^m \alpha_i k(x, x_i) + b \right).$$

1. $\alpha \leftarrow \mathbf{0}, b \leftarrow 0$
2. Find a misclassified example $(x_u, y_u) \in T_m$ such that:

$$y_u \neq \hat{y}_u = \text{sign} \left(\sum_{i=1}^m \alpha_i k(x, x_i) + b \right)$$

3. If such (x_u, y_u) does not exist, return (\mathbf{w}, b) . Otherwise update parameters:

$$\alpha_u \leftarrow \alpha_u + y_u \quad \text{and} \quad b \leftarrow b + y_u$$

and go to step 2.

Explicit vs. kernel Perceptron

Compute and memory requirements

	Training set memory	Prediction compute	Single epoch compute
Standard PLA	$\mathcal{O}(m \cdot n)$	$\mathcal{O}(n)$	$\mathcal{O}(n \cdot m)$
Kernel PLA	$\mathcal{O}(m^2)$	$\mathcal{O}(m \cdot K)$	$\mathcal{O}(m^2 \cdot K)$

m ...number of training examples, d ..input space dimension, n ...feature space dimension, K ...kernel eval. complexity

Examples of kernel functions

	$k(x, x')$	Input space \mathcal{X}	Feature space dimension n	Evaluation compute K
2nd polynomial	$\langle \mathbf{x}, \mathbf{x}' \rangle^2$	\mathbb{R}^d	$\frac{d(d+1)}{2}$	$\mathcal{O}(d)$
RBF	$\exp\left(-\gamma \ \mathbf{x} - \mathbf{x}'\ ^2\right)$	\mathbb{R}^d	∞	$\mathcal{O}(d)$
string sub-sequence	dynamic programming	$\bigcup_{d=0}^{\infty} \Sigma^d$	$ \Sigma ^q$	$\mathcal{O}(q s t)$

Recall that the number of iterations of PLA is at most $\frac{R^2}{\gamma^2}$, where $R = \max_{i=1, \dots, m} \|\phi(x_i)\|$ and $\gamma = \max_{\|\mathbf{w}\|=1, b} \min_{i=1, \dots, m} (y_i(\langle \mathbf{w}, \phi(x_i) \rangle + b))$, i.e. it does not depend on m and n .

Radial basis function kernel

- ◆ Assume the observations are real-valued features: $\mathcal{X} = \mathbb{R}^d$
- ◆ The RBF kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

corresponds to dot product $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ in infinite dimensional space.

- ◆ In one dimensional case, $\mathcal{X} = \mathbb{R}$ the RBF kernel reads

$$\begin{aligned} k(x, x') &= \exp(-\gamma(x - x')^2) \\ &= \exp(-\gamma x^2) \exp(-\gamma x'^2) \exp(2\gamma x x') \\ &= \exp(-\gamma x^2) \exp(-\gamma x'^2) \sum_{n=0}^{\infty} \frac{(2 x x' \gamma)^n}{n!} \\ &= \exp(-\gamma x^2) \exp(-\gamma x'^2) \sum_{n=0}^{\infty} \frac{(\sqrt{2\gamma} x)^n}{\sqrt{n!}} \frac{(\sqrt{2\gamma} x')^n}{\sqrt{n!}} \end{aligned}$$

where we used the Taylor expansion $\exp(a) = \sum_{n=0}^{\infty} \frac{a^n}{n!}$

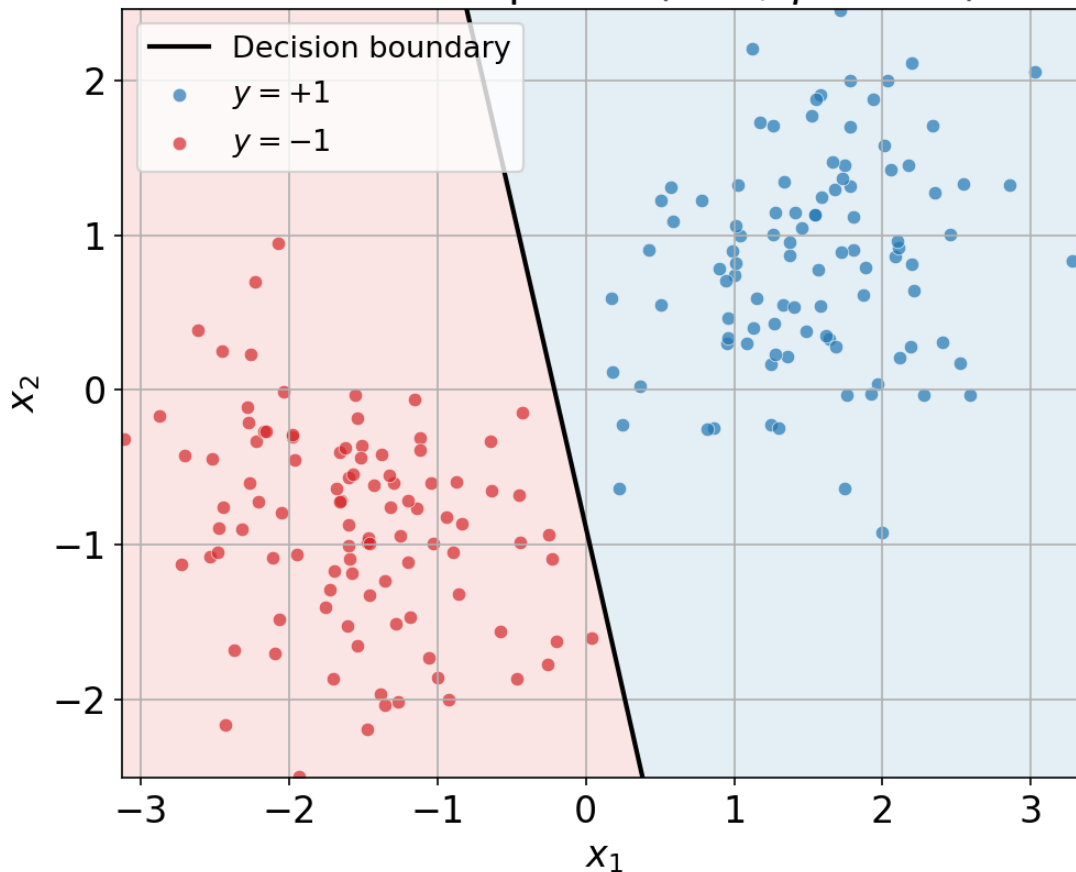
- ◆ The feature map $\phi: \mathbb{R} \rightarrow \mathbb{R}^{\infty}$ (which cannot be evaluate explicitly) then reads:

$$\phi(x) = \left[\exp(-\gamma x^2), \exp(-\gamma x^2) \sqrt{2\gamma} x, \exp(-\gamma x^2) \frac{(\sqrt{2\gamma} x)^2}{\sqrt{2!}}, \exp(-\gamma x^2) \frac{(\sqrt{2\gamma} x)^3}{\sqrt{3!}}, \dots \right]$$

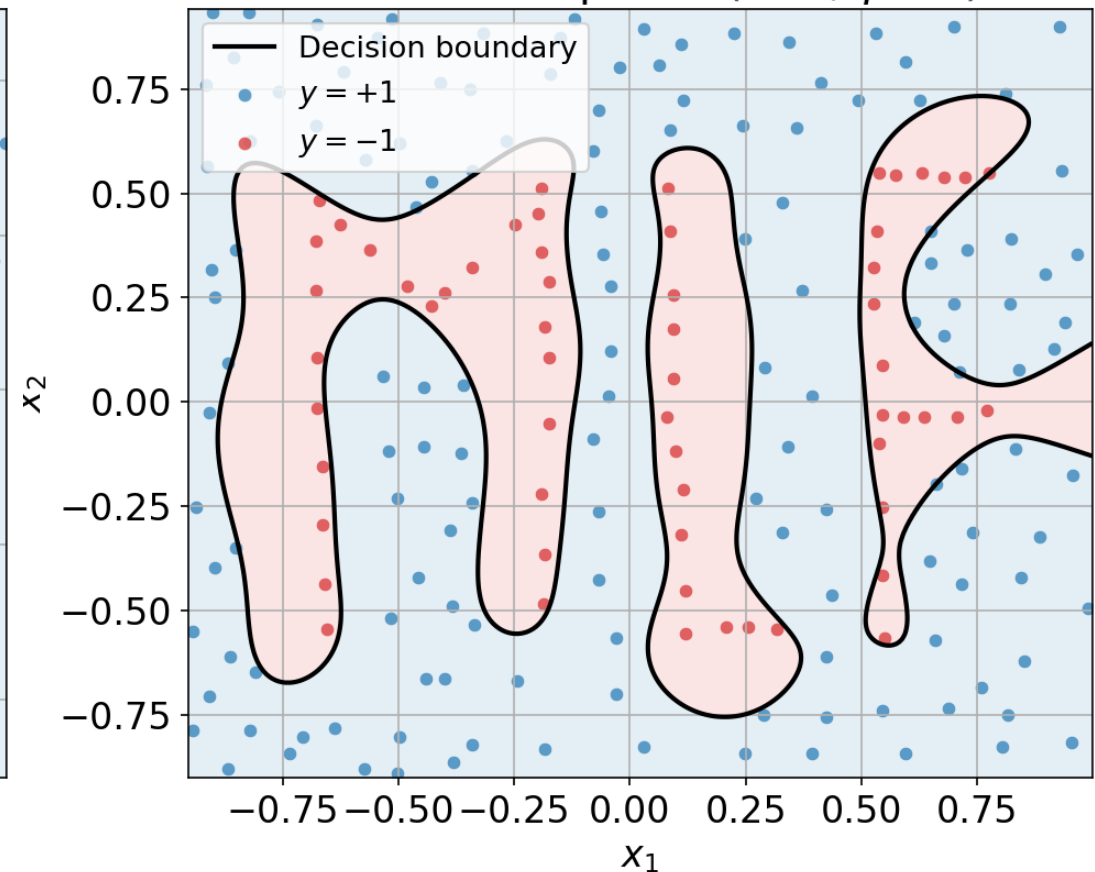
Radial basis function kernel

- ◆ Universal: can approximate arbitrary continuous decision functions on compact domains.
- ◆ Highly flexible, with a single scale parameter γ :
 - small γ : smoother / less complex boundary
 - large γ : more intricate boundary, higher overfitting risk

Kernel Perceptron (RBF, $\gamma = 0.01$)



Kernel Perceptron (RBF, $\gamma = 5$)



String Subsequence kernel

- ◆ Input space $\mathcal{X} = \cup_{d=1}^{\infty} \Sigma^d$ contains all strings on a finite alphabet Σ
- ◆ The features: occurrence+continuity of sub-sequences of length q :

$$\phi: \mathcal{X} \rightarrow \mathbb{R}^{|\Sigma|^q} \quad \text{and} \quad \phi_u(s) = \sum_{i: u=s[i]} \lambda^{l(i)}, \quad \forall u \in \Sigma^q$$

where $s[i]$ denotes a sub-sequence of s and $\lambda \in (0, 1]$ is a decay factor.

- ◆ Example for strings "cat", "car", "bat" and "bar" and $q = 2$:

	ca	ct	at	ba	bt	cr	ar	br
$\phi(\text{"cat"})$	λ^2	λ^3	λ^2	0	0	0	0	0
$\phi(\text{"car"})$	λ^2	0	0	0	0	λ^3	λ^2	0
$\phi(\text{"bat"})$	0	0	λ^2	λ^2	λ^3	0	0	0
$\phi(\text{"bar"})$	0	0	0	λ^2	0	0	λ^2	λ^3

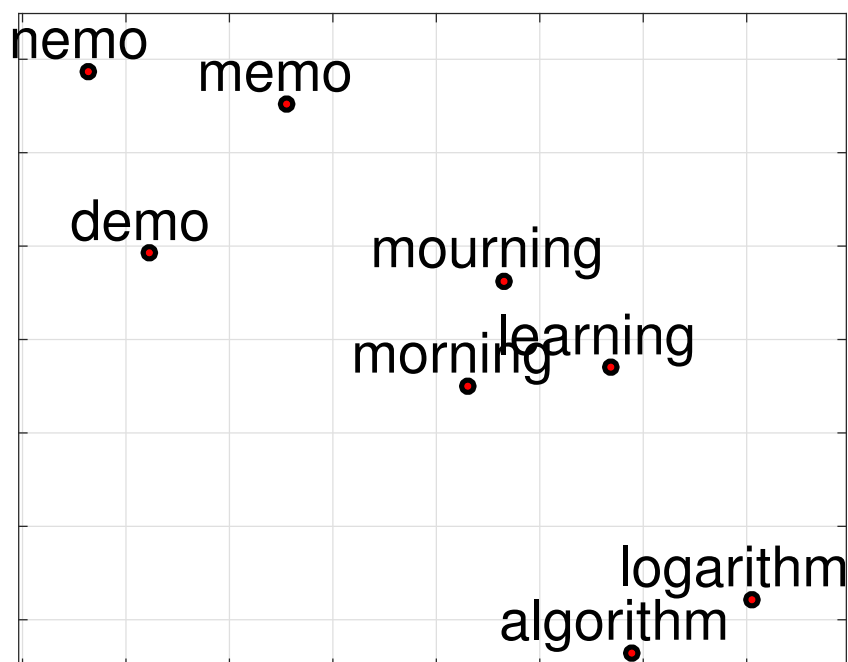
$$k(\text{"cat"}, \text{"car"}) = \lambda^4, \quad k(\text{"cat"}, \text{"bat"}) = \lambda^4, \quad k(\text{"cat"}, \text{"bar"}) = 0, \dots$$

- ◆ The kernel value $k(s, t)$ can be computed by dynamic programming in time $\mathcal{O}(q|s||t|)$.

String Subsequence kernel

Example: Consider substring kernel for sub-sequence length $q = 2$ and decay factor $\lambda = 0.4$:

Feature space
(2D embedding via tSNE)



Kernel matrix

algorithm	0.24	0.15	0.01	0.04	0.02	0.00	0.00	0.00
logarithm	0.15	0.24	0.04	0.02	0.01	0.00	0.00	0.00
learning	0.01	0.04	0.23	0.13	0.13	0.00	0.00	0.00
morning	0.04	0.02	0.13	0.19	0.17	0.03	0.03	0.03
mourning	0.02	0.01	0.13	0.17	0.23	0.03	0.03	0.03
demo	0.00	0.00	0.00	0.03	0.03	0.09	0.06	0.06
memo	0.00	0.00	0.00	0.03	0.03	0.06	0.09	0.06
nemo	0.00	0.00	0.00	0.03	0.03	0.06	0.06	0.09
	algorithm	logarithm	learning	morning	mourning	demo	memo	nemo

- Kernel $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ can be thought of as a similarity measure between the inputs embedded to a feature space:

$$k(A, B) = \text{similarity between input A and B}$$

Positive definite kernel

Definition 1. Let \mathcal{X} be a non-empty set. The function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a positive definite kernel if it is symmetric and for any finite set of inputs x_1, \dots, x_m , the kernel matrix $\mathbf{K} \in \mathbb{R}^{m \times m}$ with elements $K_{i,j} = k(x_i, x_j)$ is positive semi-definite.

- ◆ The kernel matrix $\mathbf{K} \in \mathbb{R}^{m \times m}$ represents similarities between each pair of inputs (x_1, \dots, x_m) :

$$\mathbf{K} = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_m) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_m) \\ \vdots & & & \\ k(x_m, x_1) & k(x_m, x_2) & \dots & k(x_m, x_m) \end{bmatrix}$$

- ◆ A matrix $\mathbf{K} \in \mathbb{R}^{m \times m}$ is PSD if for every $\alpha \in \mathbb{R}^m$, $\alpha^T \mathbf{K} \alpha \geq 0$.

Every kernel defines a feature space

Theorem 1. *For every positive definite kernel $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ there exists a Hilbert space \mathcal{H} and a feature map $\phi: \mathcal{X} \rightarrow \mathcal{H}$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$.*

Proof for a kernel defined on a finite input space \mathcal{X} :

The kernel matrix $\mathbf{K} \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ is symmetric and PSD hence the spectral decomposition exists $\mathbf{K} = \mathbf{V}\mathbf{D}\mathbf{V}^T$ where $\mathbf{V} \in \mathbb{R}^{m \times m}$ is orthogonal and $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_{|\mathcal{X}|})$ is diagonal matrix of non-negative eigenvalues.

Therefore $\mathbf{K} = \mathbf{\Phi}^T \mathbf{\Phi}$ where $\mathbf{\Phi}^T = \mathbf{V}\mathbf{D}^{\frac{1}{2}}$.

- ◆ A Hilbert space is a vector space, i.e. a set of objects that you can add together and multiply by a number, equipped with a dot product, which lets you measure lengths and angles between vectors.

It generalizes ordinary Euclidean space to possibly infinite dimensions.

Construction of valid kernels

Let $k_1: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k_2: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be p.d. kernels, $\lambda \in [0, 1]$, $\alpha \geq 0$, $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$ and $k_3: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ a p.d. kernel, \mathbf{K} a symmetric positive definite matrix. Then the following functions are also p.d. kernels:

$$k(x, z) = (1 - \lambda) k_1(x, z) + \lambda k_2(x, z)$$

$$k(x, z) = \alpha k_1(x, z)$$

$$k(x, z) = k_1(x, z) k_2(x, z)$$

$$k(x, z) = k_3(\phi(x), \phi(z))$$

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

Summary

- ◆ **Motivation:** Linear models may fail in the original input space, but can become linear after a nonlinear feature map $f(x) = \langle \mathbf{w}, \phi(x) \rangle + b$.
- ◆ **Kernel trick:** avoid constructing the high-dimensional feature vector explicitly by using $k(x, x') = \langle \phi(x), \phi(x') \rangle$.
- ◆ **Kernel perceptron:** represent the classifier by training examples

$$h(x) = \text{sign} \left(\sum_{i=1}^m \alpha_i k(x, x_i) + b \right),$$

so learning and prediction use only kernel evaluations.

- ◆ **Examples of kernels:**
 - polynomial kernel: finite nonlinear feature space,
 - RBF kernel: implicit infinite-dimensional feature space,
 - string subsequence kernel: similarity for structured inputs such as strings.
- ◆ **Key takeaway:** Kernels enable nonlinear decision boundaries while keeping algorithms formulated only through similarities between examples.