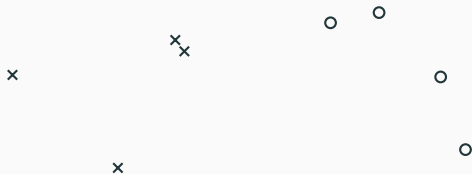# Kernel methods

UE de Master 2, AOS1
Fall 2022

S. Rousseau

## Introductory example

Let's look at one of the simplest binary classifier: the Euclidean classifier
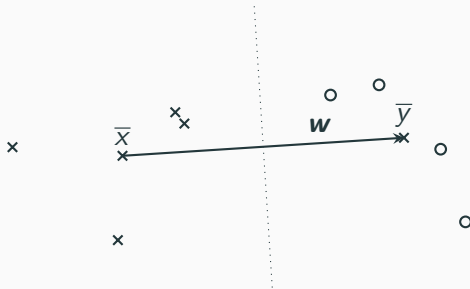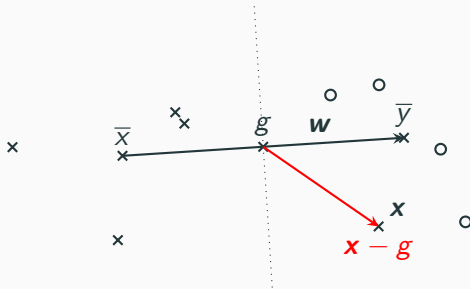
- Classification rule is: choose class whose class sample mean is the closest

Let's look at one of the simplest binary classifier: the Euclidean classifier

- Classification rule is: choose class whose class sample mean is the closest

Let's look at one of the simplest binary classifier: the Euclidean classifier

- Classification rule is: choose class whose class sample mean is the closest
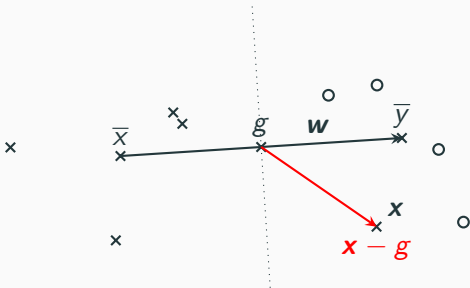
- Decision rule depends on the sign of an inner product

$$\langle \boldsymbol{x} - \boldsymbol{g}, \boldsymbol{w} \rangle$$

- Inner product leads to too simple decision boundary (hyperplane in $n$-d, line in 2-d)



How can we remedy to this?

## Making my algorithm more powerful

Two solutions

1. Pragmatic but without guarantee solution: I replace the inner product by some other similarity
   - The similarity is easy to compute, but
   - I have no guarantee that it will work
2. Theoretical but hard in practice solution: I embed the samples in a possibly infinite dimensional feature space
   - My algorithm is the same but in a different space, but
   - Computations in feature space can be computationally intensive

## Making my algorithm more powerful

Two solutions

1. Pragmatic but without guarantee solution: I replace the inner product by some other similarity
   - The similarity is easy to compute, but
   - I have no guarantee that it will work
2. Theoretical but hard in practice solution: I embed the samples in a possibly infinite dimensional feature space
   - My algorithm is the same but in a different space, but
   - Computations in feature space can be computationally intensive

   The two solutions are actually the same!

   (provided that the similarity is a **kernel**)

## Kernel: formal definition

- Arbitrary input space $\mathcal{X}$ ($\mathbb{R}^p$, strings, graphs, ...)

- A **kernel** $k$ defined on $\mathcal{X}$ is a map:

$$k : \mathcal{X} \times \mathcal{X} \to \mathbb{R},$$

that verifies

- **Symmetry**: $\forall x, y \in \mathcal{X}$, $k(x, y) = k(y, x)$, matrix $(K)_{ij} = k(x_i, x_j)$ is symmetric

- **Positive semidefiniteness**: For all $x_1, \ldots, x_n \in \mathcal{X}$, matrix $(K)_{ij} = k(x_i, x_j)$ is positive semidefinite (PSD, $K \geqslant 0$):

$$\forall \boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\alpha}^T K \boldsymbol{\alpha} = \sum_{ij} \alpha_i \alpha_j k(x_i, x_j) \geqslant 0$$

## Some examples

- Linear kernel (inner product is of course a kernel): $\mathcal{X} = \mathbb{R}^p$, $k(\boldsymbol{x}, \boldsymbol{y}) = \langle \boldsymbol{x}, \boldsymbol{y} \rangle$

- Polynomial kernel: let $\mathcal{X} = \mathbb{R}^p$, $d > 0, c \geqslant 0$

$$k_{\mathrm{poly1}}(\boldsymbol{x}, \boldsymbol{y}) = (c + \langle \boldsymbol{x}, \boldsymbol{y} \rangle)^d$$

- Gaussian kernel or RBF (radial basis function): $k_{\mathrm{gauss}}(\boldsymbol{x}, \boldsymbol{y}) = \exp\left\{ -\frac{\|\boldsymbol{x} - \boldsymbol{y}\|_2^2}{\sigma^2} \right\}$

- Sigmoid kernel: $k(\boldsymbol{x}, \boldsymbol{y}) = \tanh\left( \gamma \langle \boldsymbol{x}, \boldsymbol{y} \rangle + c \right)$

## From solution #2 to solution #1

Embedding actually defines a kernel

- Let $\Phi : \mathcal{X} \longmapsto \mathcal{F}$ some embedding of $\mathcal{X}$ in a **feature space** $\mathcal{F}$ that is equipped with an inner production $\langle \cdot, \cdot \rangle_{\mathcal{F}}$.

- Define $k$ as follows

$$k(\boldsymbol{x}, \boldsymbol{y}) = \langle \Phi(\boldsymbol{x}), \Phi(\boldsymbol{y}) \rangle_{\mathcal{F}}$$

- Define that way, $k$ is always a kernel!
  - Obviouly symmetric
  - Inner product in $\mathcal{F}$ makes $k$ PSD

## From solution #1 to solution #2: the polynomial kernel case

Suppose $\mathcal{X} = \mathbb{R}^2$, $n = 2$, $c = 1$ and $d = 2$

- Polynomial kernel is then

$$
\begin{aligned}
k_{\mathrm{poly1}}(\boldsymbol{x}, \boldsymbol{y}) &= (1 + \langle \boldsymbol{x}, \boldsymbol{y} \rangle)^2 \\
&= (\langle \boldsymbol{x}, \boldsymbol{y} \rangle + 1)^2 = (1 + x_1 y_1 + x_2 y_2)^2 \\
&= 1 + 2x_1 y_1 + 2x_2 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2
\end{aligned}
$$

- Define $\Phi(\boldsymbol{x}) = \left(1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2\right)$, we have

$$
k_{\mathrm{poly1}}(\boldsymbol{x}, \boldsymbol{y}) = \langle \Phi(\boldsymbol{x}), \Phi(\boldsymbol{y}) \rangle
$$

- The kernel $k_{\mathrm{poly1}}$ is can be interpreted as an inner product

## From solution #1 to solution #2: when $\mathcal{X}$ is finite

From the definition, $K$ is a kernel iff $K$ is symetric PSD

- $K$ is diagonalizable in an orthonormal basis $K = UDU^T$ with $U$ orthogonal and $D$ diagonal

$$
\begin{aligned}
k(\mathbf{x}_i, \mathbf{x}_j) &= K_{ij} \\
&= (UDU^T)_{ij} \\
&= \mathsf{u}_i^T D \mathsf{u}_j \\
&= \mathsf{u}_i^T \sqrt{D}\sqrt{D}\mathsf{u}_j \\
&= \left\langle \sqrt{D}\mathsf{u}_i, \sqrt{D}\mathsf{u}_j \right\rangle
\end{aligned}
$$

Define the function $\Phi$

$$
\begin{aligned}
\Phi \colon \mathcal{X} &\to \mathbb{R}^n \\
\mathbf{x}_i &\mapsto \sqrt{D}\mathsf{u}_i.
\end{aligned}
$$

we have $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$

- The kernel $k$ is can also be interpreted as an inner product

9

## Kernel trick

- Kernel trick

$$k(\boldsymbol{x}, \boldsymbol{y}) = \langle \Phi(\boldsymbol{x}), \Phi(\boldsymbol{y}) \rangle_{\mathcal{F}}$$

  No need to compute the embedding and then the inner product, the kernel is exactly doing this.

- Any algorithm using **only inner products** between training samples can be "kernelized"

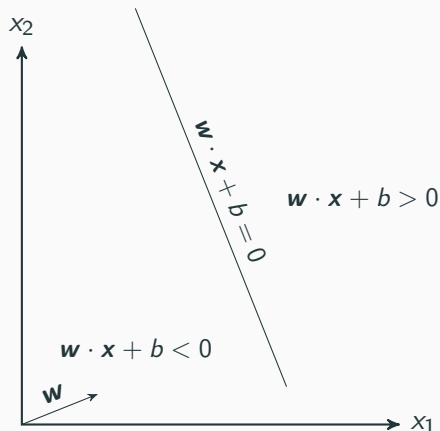- Most famous kernel method: SVM (support vector machines)

## Separating hyperplane

- In $\mathbb{R}^p$, a hyperplane is defined by $h(\boldsymbol{x}) = 0$ with $h(\boldsymbol{x}) = \boldsymbol{w} \cdot \boldsymbol{x} + b$.

- $\boldsymbol{w}$ is orthogonal to the hyperplane

- $d(\boldsymbol{x}, H) = \frac{|\boldsymbol{w} \cdot \boldsymbol{x} + b|}{\|\boldsymbol{w}\|_2}$

- Distance from the origin is $\frac{|b|}{\|\boldsymbol{w}\|_2}$

## Margin

- $H$ a separating hyperplane

- Distance from one $x_i$ to $H$ is

$$d(x_i, H) = \frac{|w \cdot x_i + b|}{\|w\|_2}$$

- Margin of a separating hyperplan is the smallest distance from $H$ to the $x_i$'s

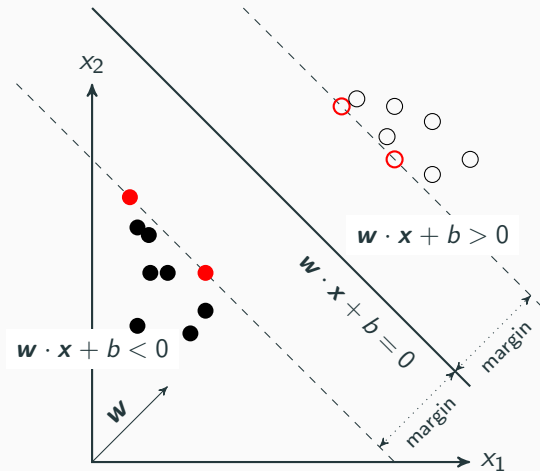$$M = \min_{i=1,\ldots,n} d(x_i, H)$$



12

## Optimal separating hyperplane

- Optimal separating hyperplane maximizes the margin

- Unique solution

- Vectors $\mathbf{x}_i$ supporting the margin, *i.e.* such that

$$M = d(\mathbf{x}_i, H)$$

  are called **support vectors**

- $H = \underset{H \in \mathcal{H}}{\arg \max} \ \underset{i=1,\dots,n}{\min} \ d(\mathbf{x}_i, H)$



$x_2$

$\mathbf{w} \cdot \mathbf{x} + b > 0$

$\mathbf{w} \cdot \mathbf{x} + b = 0$

$\mathbf{w} \cdot \mathbf{x} + b < 0$

$\mathbf{w}$

margin

margin

$x_1$

13

## One-to-one representation

- Hyperplane $\iff (\boldsymbol{w}, b)$ is not one-to-one
  If $(\boldsymbol{w}, b)$ represents $H$ so is $(c\boldsymbol{w}, cb)$ with $c \neq 0$
- Possible choice for $c$
  1. Data independant: imposing $\|\boldsymbol{w}\|_2 = 1$ (taking $c = \frac{1}{\|\boldsymbol{w}\|_2}$)

  $$(\boldsymbol{w}, b) \quad \text{with} \quad \|\boldsymbol{w}\|_2 = 1 \qquad \text{(normalization 1)}$$

  2. Data dependant: $\boldsymbol{w} \cdot \boldsymbol{x} + b$ is $\pm 1$ at margin lines
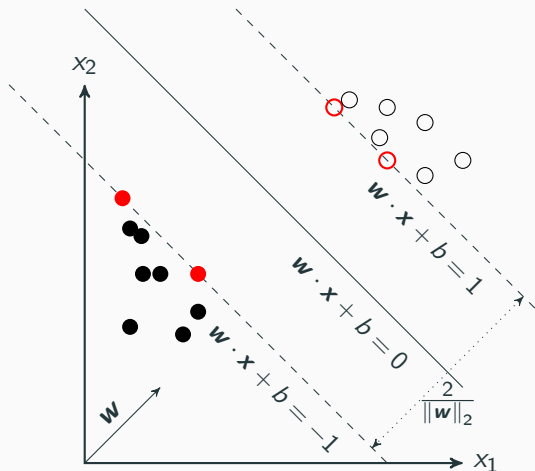
  $$(\boldsymbol{w}, b) \quad \text{with} \quad \|\boldsymbol{w}\|_2 = \frac{1}{M} \qquad \text{(normalization 2)}$$

  When $H$ is an optimal separating hyperplane, support vectors are $\pm 1$

14

Choosing data-dependent normalization

- Margin is $\frac{1}{\|\boldsymbol{w}\|_2}$
- Takes the value $\pm 1$ for support vectors



$x_2$

$\boldsymbol{w} \cdot \boldsymbol{x} + b = 1$

$\boldsymbol{w} \cdot \boldsymbol{x} + b = 0$

$\boldsymbol{w} \cdot \boldsymbol{x} + b = -1$

$\boldsymbol{w}$

$\frac{2}{\|\boldsymbol{w}\|_2}$

$x_1$

## Maximizing the margin

- Maximizing $\frac{1}{\|\boldsymbol{w}\|_2}$ $\iff$ Minimizing $\|\boldsymbol{w}\|_2^2$

- Proper labelling:
  - $\boldsymbol{w} \cdot \boldsymbol{x}_i + b \geqslant 1$ when $y_i = 1$

  - $\boldsymbol{w} \cdot \boldsymbol{x}_i + b \leqslant -1$ when $y_i = -1$

  $\left. \right\} \; y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1$

$$\underset{(\boldsymbol{w}, b)}{\arg\min} \frac{1}{2} \|\boldsymbol{w}\|_2^2 \quad \text{such that} \quad \forall i,\ y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1 \tag{1}$$

- Quadratic programming problem: quadratic objective, linear constraints

## Lagrangian formulation

- We use the Lagrangian, $n$ Lagrange multipliers $\mu_1, \ldots, \mu_n$

$$\mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\mu}) = \frac{1}{2} \|\boldsymbol{w}\|_2^2 - \sum_{i=1}^{n} \mu_i \left( y_i (\boldsymbol{w} \cdot \boldsymbol{x}_i + b) - 1 \right)$$

- Karush–Kuhn–Tucker (KKT) conditions because we have inequalities

$$\begin{cases} \dfrac{\partial \mathcal{L}}{\partial \boldsymbol{w}} = \boldsymbol{w} - \displaystyle\sum_{i=1}^{n} \mu_i y_i \boldsymbol{x}_i = 0 \\[2ex] \dfrac{\partial \mathcal{L}}{\partial b} = - \displaystyle\sum_{i=1}^{n} \mu_i y_i = 0 \end{cases} \implies \begin{cases} \boldsymbol{w} = \displaystyle\sum_{i=1}^{n} \mu_i y_i \boldsymbol{x}_i \\[2ex] \displaystyle\sum_{i=1}^{n} \mu_i y_i = 0 \end{cases}$$

$$\mu_i \geqslant 0 \quad \text{et} \quad \mu_i \left( y_i (\boldsymbol{w} \cdot \boldsymbol{x}_i + b) - 1 \right) = 0$$

## Dual formulation

- Plugging the KKT conditions back in the objective function we have

$$\mathcal{L}(\boldsymbol{\mu}, b) = \frac{1}{2} \left\| \sum_{j=1}^{n} \mu_j y_j \mathbf{x}_j \right\|_2^2 - \sum_{i=1}^{n} \mu_i \left( y_i \left( \sum_{j=1}^{n} \mu_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + b \right) - 1 \right)$$

$$\mathcal{L}(\boldsymbol{\mu}) = \frac{1}{2} \left\| \sum_{j=1}^{n} \mu_j y_j \mathbf{x}_j \right\|_2^2 - \sum_{i,j=1}^{n} \mu_i \mu_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_{i=1}^{n} \mu_i$$

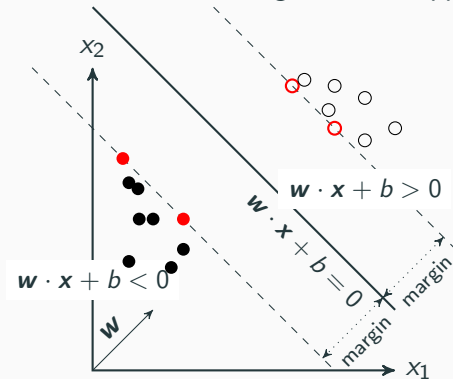$$= -\frac{1}{2} \left\| \sum_{j=1}^{n} \mu_j y_j \mathbf{x}_j \right\|_2^2 + \sum_{i=1}^{n} \mu_i$$

## Dual formulation

SVM minimization problem (dual form)

$$\arg\min_{\boldsymbol{\mu}} \frac{1}{2} \sum_{i=1}^{n} \mu_i \mu_j y_i y_j \boldsymbol{x}_i \cdot \boldsymbol{x}_j - \sum_{i=1}^{n} \mu_i \quad \text{such that} \quad \begin{cases} \forall i,\ \mu_i \geqslant 0 \\ \sum_{i=1}^{n} \mu_i y_i = 0 \end{cases}$$

- Only depends on $\boldsymbol{x}_i \cdot \boldsymbol{x}_j$: SVM algorithm is kernelizable

- $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) > 1$: $\mu_i = 0$, $\mathbf{x}_i$ is not used
- $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$: $\mu_i > 0$, $\mathbf{x}_i$ is on the margin, it is a *support vector*



- Only support vectors are used to classify

## Classification

- Classification of a new example $\mathbf{x}$:

$$y(\mathbf{x}) = \text{sgn}\left(\mathbf{w} \cdot \mathbf{x} + b\right)$$
$$= \text{sgn}\left(\sum_{i=1}^{n} \mu_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$$

- Only the support vectors are used to classify
  - We only need to compute $\mathbf{x}_i \cdot \mathbf{x}$ with $\mathbf{x}_i$ a support vector
  - Robust to outliers in training data
- How to determine $b$:
  - If $\mathbf{x}_i$ is a support vector, $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$ hence

$$b = y_i - \mathbf{w} \cdot \mathbf{x}_i$$

  - More stable numerically

$$\frac{1}{\#\{i, \mu_i > 0\}} \sum_{i,\, \mu_i > 0} y_i - \mathbf{w} \cdot \mathbf{x}_i$$

## Extension to the non-separable case

- Until now, classes were was supposed to be linearly separable but:
  - it is never the case in practice
  - minimization problem (1) or dual form do not even have a solution: set of admissible solution is empty
- Typical solution: relaxing the constraints

## Slack variables

- Original minimization problem

$$\underset{(\boldsymbol{w},b)}{\arg\min} \frac{1}{2} \|\boldsymbol{w}\|_2^2 \quad \text{such that} \quad \forall i,\, y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1$$

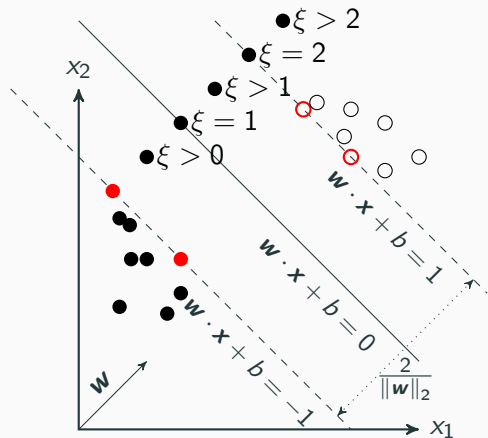- Some constraint might be violated; we give them some slack by introducing nonnegative **slack variables** $\xi_1, \ldots, \xi_n$ so that

$$\forall i,\, y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1 \quad \text{becomes} \quad \forall i,\, y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1 - \xi_i$$

- The $\xi_i$'s catch up on mistakes

New constraints: $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geqslant 1 - \xi_i$

- If $\xi_i = 0$ the original constraint is met
- If $0 < \xi_i < 2$ the original constraint is violated and the sample $\mathbf{x}_i$ belongs to the margin
- If $\xi_i \geqslant 2$ the original constraint is violated and the sample $\mathbf{x}_i$ belong to the wrong side

## Soft margin minimization problem

- Still maximizing the margin but regularizing by the amount of slack $\sum_{i=1}^{n} \xi_i$

$$\underset{(\boldsymbol{w},b)}{\arg\min} \frac{1}{2} \|\boldsymbol{w}\|_2^2 + C \sum_{i=1}^{n} \xi_i \quad \text{such that} \quad \begin{cases} \forall i,\ y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1 - \xi_i \\ \xi_i \geqslant 0 \end{cases}$$

- We introduce a new hyperparameter $C$: cost of adding some slack
  - Hard margin if $C \to +\infty$: no slack

## Hinge loss

- Equivalence of constraints

$$\begin{cases} \forall i, \, y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b) \geqslant 1 - \xi_i \\ \xi_i \geqslant 0 \end{cases} \iff \xi_i \geqslant \max\left(0, 1 - y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b)\right)$$

- Smaller objective function if $\xi_i = \max\left(0, 1 - y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b)\right)$

- If we plug the $\xi_i$'s in the objective function we have

$$\underset{(\boldsymbol{w}, b)}{\arg\min} \, \frac{1}{2} \|\boldsymbol{w}\|_2^2 + C \sum_{i=1}^{n} \max\left(0, 1 - y_i(\boldsymbol{w} \cdot \boldsymbol{x}_i + b)\right)$$

- Defining hinge loss function as: $L_{\text{hinge}}(x, y) = \max\left(0, 1 - xy\right)$

$$\underset{(\boldsymbol{w}, b)}{\arg\min} \, \sum_{i=1}^{n} L_{\text{hinge}}(y_i, \boldsymbol{w} \cdot \boldsymbol{x}_i + b) + \frac{1}{2C} \|\boldsymbol{w}\|_2^2$$
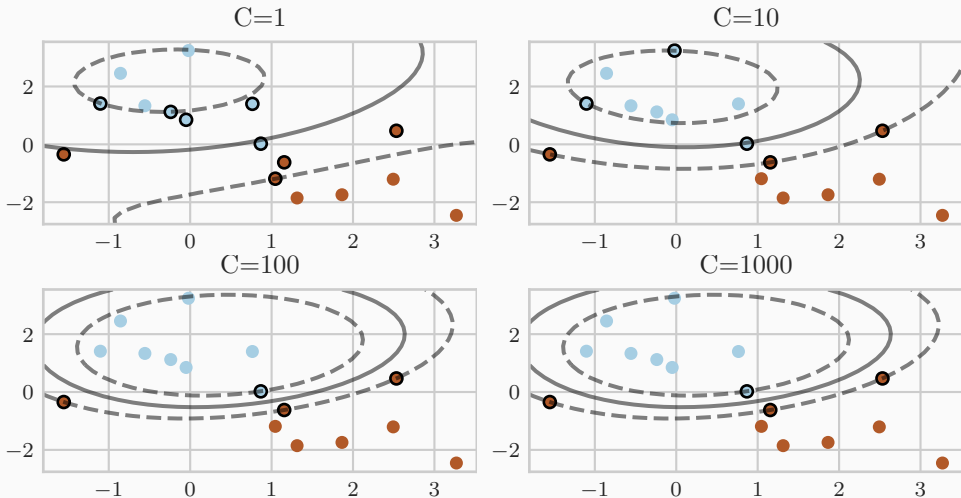
26

## Other losses

Losses when true class is $y_i = 1$

- 0/1 loss: penalises wrong sign of prediction $\widehat{y}_i$
- Square loss penalises the gap whatever the direction
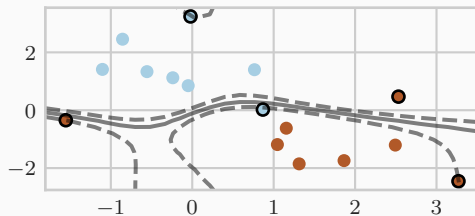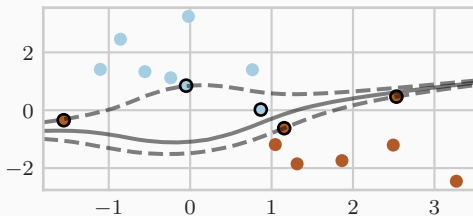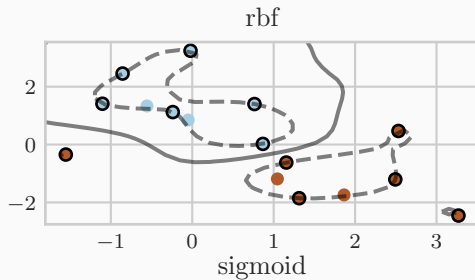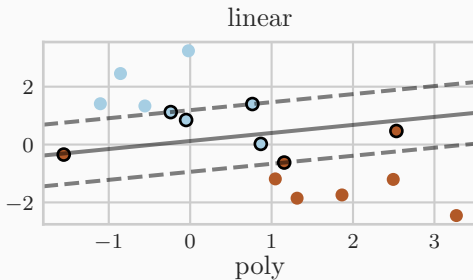- hinge loss: relaxing 0/1 loss

Gaussian kernel with varying $C$

# Illustrations: changing the kernel

## Kernel PCA

The idea is to apply the principal component analysis **on the feature space**

- New points are $\boldsymbol{y}_i = \Phi(\boldsymbol{x}_i)$
- Design matrix changes from

$$X = \begin{bmatrix} \boldsymbol{x}_1^T \\ \vdots \\ \boldsymbol{x}_n^T \end{bmatrix} \quad \text{to} \quad X_\Phi = \begin{bmatrix} \boldsymbol{y}_1^T \\ \vdots \\ \boldsymbol{y}_n^T \end{bmatrix}$$

- Suppose for now that $X_\Phi$ is **centered**
- The kernel matrix $K$ gathers the inner products in feature space: $K = X_\Phi X_\Phi^T$

## Kernel PCA

Kernel PCA on feature space

- Let $V_\Phi = \frac{1}{n} X_\Phi^T X_\Phi$ the sample variance–covariance matrix, $v_1, \ldots, v_q$ the eigen vectors and $\lambda_1 \geqslant \ldots \geqslant \lambda_q$ the corresponding eigenvalues
- We have seen that the $i$-th principal component is $X_\Phi v_i$

We don't want to compute $X_\Phi v_i$

- It is easy to see that $K(X_\Phi v_i) = n\lambda_i(X_\Phi v_i)$; PC are just eigenvectors (properly rescaled) of the kernel matrix $K$

## Centering the kernel

We supposed that $X_\Phi$ was centered but what if it's not?

- Define the centering operator

$$Q_n = I_n - \frac{1}{n}\mathbb{1}_n = \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \ddots & \vdots \\ \vdots & \ddots & \ddots & -\frac{1}{n} \\ -\frac{1}{n} & \cdots & -\frac{1}{n} & 1 - \frac{1}{n} \end{pmatrix}.$$
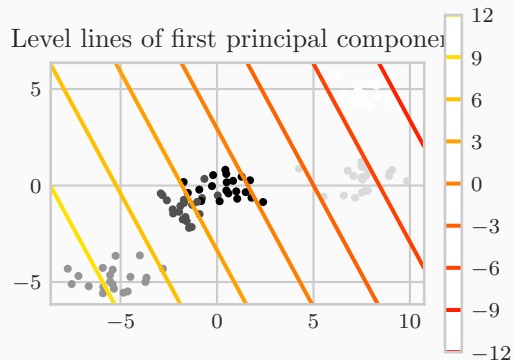
so that $X_\Phi^0 = Q_n X_\Phi$ is the centering of $X_\Phi$

- The corresponding kernel is $K^0 = Q_n X_\Phi (Q_n X_\Phi)^T = Q_n X_\Phi X_\Phi^T Q_n = Q_n K Q_n$
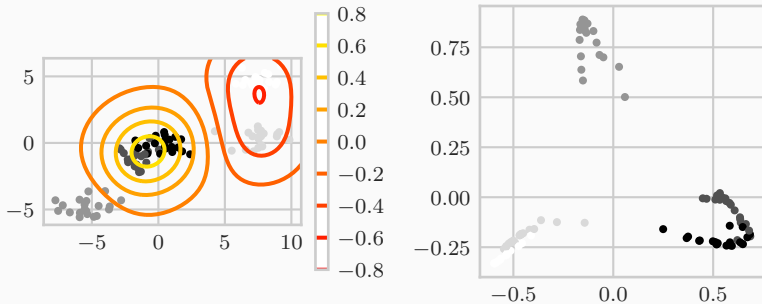- No need to compute $X_\Phi$, just center column-wise and row-wise the kernel

Classic PCA (linear kernel)

- Level lines of first principal component are straight lines
- And they are orthogonal to first principal direction



Level lines of first principal component

Kernel PCA with an RBF kernel (gaussian kernel)

## Kernel ridge regression

- Matrix inversion trick: from $X\left(I + X^T X\right) = \left(I + XX^T\right) X$ we get

$$\left(I + XX^T\right)^{-1} X = X\left(I + X^T X\right)^{-1}$$

where $XX^T$ is $n \times n$ and $X^T X$ is $p \times p$

- Using this in $\widehat{\boldsymbol{y}}^{\text{ridge}} = X\left(X^T X + \lambda I_p\right)^{-1} X^T \boldsymbol{y}$ we get

$$\widehat{\boldsymbol{y}}^{\text{ridge}} = \left(\lambda I_p + XX^T\right)^{-1} XX^T \boldsymbol{y}$$

and then $\widehat{\boldsymbol{y}}^{\text{ridge}} = \left(\lambda I_p + K\right)^{-1} K \boldsymbol{y}$

## Kernel $k$-means / kernel $k$-nearest neighbors

Those algorithms only depend on distances between samples

- Since interdistance can be expressed as inner products

$$\|\mathbf{x}_i - \mathbf{x}_j\|^2 = \langle \mathbf{x}_i, \mathbf{x}_i \rangle - 2 \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \langle \mathbf{x}_j, \mathbf{x}_j \rangle$$

- $K$-means is kernelizable. The "kernel distance" is then

$$d(\mathbf{x}_i, \mathbf{x}_j)^2 = k(\mathbf{x}_i, \mathbf{x}_i) - 2k(\mathbf{x}_i, \mathbf{x}_j) + k(\mathbf{x}_j, \mathbf{x}_j)$$