

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Alain Celisse

SAMM

Paris 1-Panthéon Sorbonne University

alain.celisse@univ-paris1.fr

Lecture 1: Support vector classifiers and Kernel machines

Master 2 Data Science – Centrale Lille, Lille University
Fall 2022

Outline of the lectures

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Successive topics of the coming lectures:

1. Introduction to Kernel methods
2. Support vector classifiers and Kernel methods (Today!)
3. Extending classical strategies to high dimension
 - ▶ KRR/LS-SVMs
 - ▶ KPCA
4. Duality gap and KKT conditions
5. Designing reproducing kernels
6. Maximum Mean Discrepancy (MMD)
7. Change-point detection, KCP

Outline of the lecture

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

- ▶ Optimal Separating Hyperplanes
- ▶ Support Vectors
- ▶ Kernel machines and SVM
- ▶ Towards regression

Optimal Separating Hyperplanes

Linearly separable classes

Perceptron and optimal hyperplane

Support Vectors

Reproducing Kernels

Support Vector Machines

Optimal Separating Hyperplanes

Optimal Separating Hyperplanes

Linearly separable classes

Perceptron and optimal hyperplane

Support Vectors

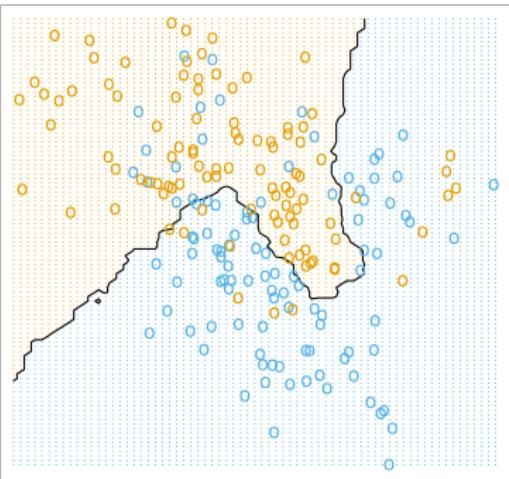
Reproducing Kernels

Support Vector Machines

Binary classification

Framework

- ▶ $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{iid}{\sim} P$: unknown
- ▶ $X_i \in \mathcal{X} \subset \mathbb{R}^d$
- ▶ $Y_i \in \{-1, 1\} = \mathcal{Y}$



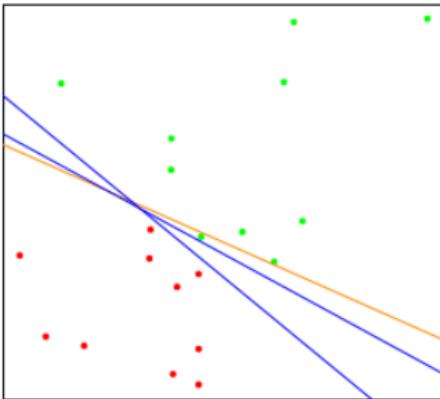
Goal and context

- ▶ Design a *data-driven* classification rule
- ▶ The classification rule achieves the smallest possible “misclassification rate”
- ▶ No parametric model for P (nonparametric)

Linearly separable classes

Kernel Machines

Alain Celisse



Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Definition (Linearly separable classes)

Two sets of points are *linearly separable* if there exists a

hyperplane $B = \left\{ x \in \mathbb{R}^d \mid \underbrace{\beta_0 + \beta^\top x}_{=f(x)} = 0 \right\}$ such that, for all i ,

$$Y_i = \text{sign}(\beta_0 + \beta^\top X_i) \Leftrightarrow Y_i = 1, \quad \underbrace{\beta_0 + \beta^\top X_i}_{=f(X_i)} > 0$$

$$Y_i = -1, \quad \underbrace{\beta_0 + \beta^\top X_i}_{=f(X_i)} < 0$$

Hyperplane and orthogonal projection

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

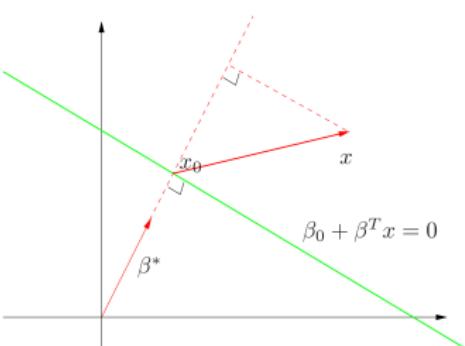
Reproducing
Kernels

Support Vector
Machines

- With linearly separable classes, several separating hyperplanes do exist
 - Need for a criterion to choose one of them

Orthogonal projection

- $x_0 \in B \Leftrightarrow \beta_0 + \langle \beta, x_0 \rangle = 0$
- $\beta^* = \beta / \|\beta\|$: normal vector to B
- x : any vector



- $p^\perp(x)$: orthogonal projection of x onto B
 $\Rightarrow \beta_0 + \langle \beta, p^\perp(x) \rangle = f(p^\perp(x)) = 0 \quad (p^\perp(x) \in B)$
- $x - p^\perp(x) = \alpha \beta^*$ for some $\alpha \in \mathbb{R}$
 $\Rightarrow \alpha = \langle x - x_0, \beta^* \rangle$ and $x - p^\perp(x) = \langle x - x_0, \beta^* \rangle \beta^*$

Distance to a hyperplane

Kernel Machines

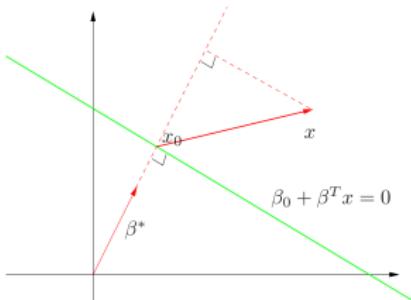
Alain Celisse

Signed distance from x to B

$$\blacktriangleright x - p^\perp(x) = \langle x - x_0, \beta^* \rangle \beta^*$$

\blacktriangleright Signed distance of x to B :

$$\begin{aligned} \langle x - p^\perp(x), \beta^* \rangle &= \langle x - x_0, \beta^* \rangle \\ &= \frac{\langle \beta, x \rangle + \beta_0}{\|\beta\|} = \frac{f(x)}{\|\beta\|} = \frac{f(x)}{\|f'(x)\|} \end{aligned}$$



Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Distance to a hyperplane

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Signed distance from x to B

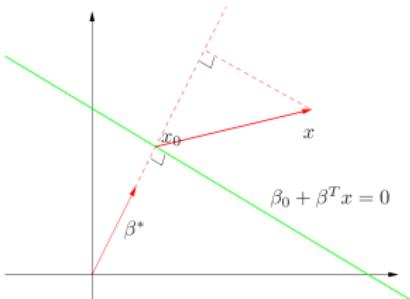
$$\triangleright x - p^\perp(x) = \langle x - x_0, \beta^* \rangle \beta^*$$

Signed distance of x to B :

$$\begin{aligned} \langle x - p^\perp(x), \beta^* \rangle &= \langle x - x_0, \beta^* \rangle \\ &= \frac{\langle \beta, x \rangle + \beta_0}{\|\beta\|} = \frac{f(x)}{\|\beta\|} = \frac{f(x)}{\|f'(x)\|} \end{aligned}$$

Remarks:

- Positive distance iff B predicts label “1” at x ($f(x) > 0$)
- Normalizing by $\|\beta\|$ mandatory because β^* is unit norm



Cumulated distance to the separating hyperplane

Definition (Misclassified points)

Any observation x_i is said *misclassified* if

$$y_i f(x_i) < 0$$

$$(f(x) = \langle \beta, x \rangle + \beta_0)$$

Cumulated distance of well-classified points to B

$$dist(\tilde{\beta}) = \sum_{i=1}^n \left(y_i \frac{f(x_i)}{\|f'(x_i)\|} \right)_+ \propto \sum_{i=1}^n (y_i f(x_i))_+$$

Perceptron learning algorithm

Perceptron algorithm

$$(\tilde{\beta} = (\beta, \beta_0))$$

- ▶ Perceptron aims at finding $\tilde{\beta} \in \mathbb{R}^d$ that minimizes

$$\sum_{i=1}^n (y_i f(\tilde{\beta}, x_i))_+$$
- ▶ Solves the optimization problem:

$$\min \sum_{i=1}^n (y_i f_{\tilde{\beta}}(x_i))_+ \quad s.t. \quad \tilde{\beta} \in \mathbb{R}^d$$

Perceptron learning algorithm

Perceptron algorithm

$$(\tilde{\beta} = (\beta, \beta_0))$$

- ▶ Perceptron aims at finding $\tilde{\beta} \in \mathbb{R}^d$ that minimizes

$$\sum_{i=1}^n (y_i f(\tilde{\beta}, x_i))_+$$
- ▶ Solves the optimization problem:

$$\min \sum_{i=1}^n (y_i f_{\tilde{\beta}}(x_i))_+ \quad s.t. \quad \tilde{\beta} \in \mathbb{R}^d$$

SGD iterations:

- ▶ Perceptron solution computed by SGD:

$$(\tilde{\beta}^{t+1})^\top = (\tilde{\beta}^t)^\top + \rho (x_i y_i, y_i)^\top$$

for all missclassified x_i s and $\rho > 0$

Perceptron learning algorithm

Kernel Machines

Alain Celisse

Perceptron algorithm

$$(\tilde{\beta} = (\beta, \beta_0))$$

- ▶ Perceptron aims at finding $\tilde{\beta} \in \mathbb{R}^d$ that minimizes
 $\sum_{i=1}^n (y_i f(\tilde{\beta}, x_i))_+$
- ▶ Solves the optimization problem:

$$\min \sum_{i=1}^n (y_i f_{\tilde{\beta}}(x_i))_+ \quad s.t. \quad \tilde{\beta} \in \mathbb{R}^d$$

SGD iterations:

- ▶ Perceptron solution computed by SGD:

$$(\tilde{\beta}^{t+1})^\top = (\tilde{\beta}^t)^\top + \rho (x_i y_i, y_i)^\top$$

for all missclassified x_i s and $\rho > 0$

Limitations:

- ▶ Not scale-invariant (solution depends on $\|\beta\|$)
- ▶ Only converges with linearly separable classes

Optimal Separating Hyperplanes

Linearly separable classes

Perceptron and optimal hyperplane

Support Vectors

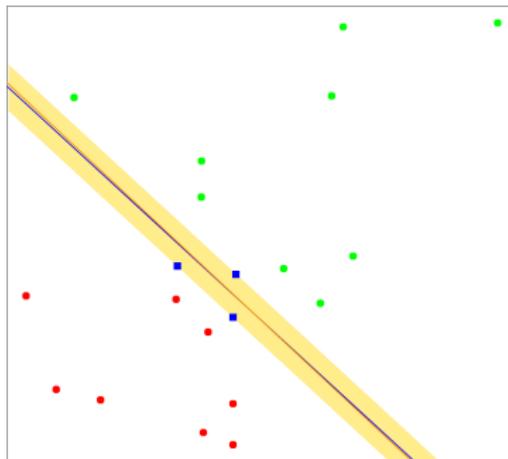
Reproducing Kernels

Support Vector Machines

Optimal Separating Hyperplane and Margin

Optimal hyperplane

- ▶ With linearly separable classes, no unique hyperplane
- ▶ **Optimal separating hyperplane:**
“The one which departs the most from each class”



Margin

Half the width of the yellow band

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Margin: Formal definition

Well-classified points and Margin $(f_{\tilde{\beta}}(x) = \langle \beta, x \rangle + \beta_0 = 0)$

- ▶ Sets of well-classified points:

$$C_1 = \{x_i \mid y_i = 1, y_i f(x_i) > 0\} \quad C_{-1} = \{x_i \mid y_i = -1, y_i f(x_i) > 0\}$$

Margin: Formal definition

Well-classified points and Margin $(f_{\tilde{\beta}}(x) = \langle \beta, x \rangle + \beta_0 = 0)$

- ▶ Sets of well-classified points:

$$C_1 = \{x_i \mid y_i = 1, y_i f(x_i) > 0\} \quad C_{-1} = \{x_i \mid y_i = -1, y_i f(x_i) > 0\}$$

- ▶ Distance from x_i to B :

$$dist(x_i, B) = y_i f_{\tilde{\beta}}(x_i) / \|f'_{\tilde{\beta}}(x_i)\|$$

Margin: Formal definition

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Well-classified points and Margin $(f_{\tilde{\beta}}(x) = \langle \beta, x \rangle + \beta_0 = 0)$

- ▶ Sets of well-classified points:

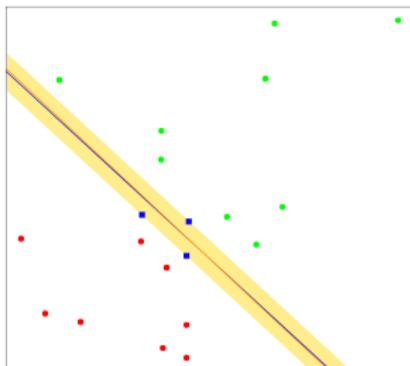
$$C_1 = \{x_i \mid y_i = 1, y_i f(x_i) > 0\} \quad C_{-1} = \{x_i \mid y_i = -1, y_i f(x_i) > 0\}$$

- ▶ Distance from x_i to B :

$$dist(x_i, B) = y_i f_{\tilde{\beta}}(x_i) / \|f'_{\tilde{\beta}}(x_i)\|$$

- ▶ Distance from a (finite) set C_ℓ to B : ($\ell \in \{-1, 1\}$)

$$dist(C_\ell, B) = \min \{ dist(x_i, B) \mid x_i \in C_\ell \}$$



Margin: Formal definition

Well-classified points and Margin $(f_{\tilde{\beta}}(x) = \langle \beta, x \rangle + \beta_0 = 0)$

- ▶ Sets of well-classified points:

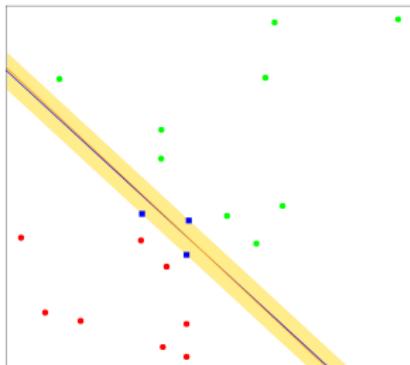
$$C_1 = \{x_i \mid y_i = 1, y_i f(x_i) > 0\} \quad C_{-1} = \{x_i \mid y_i = -1, y_i f(x_i) > 0\}$$

- ▶ Distance from x_i to B :

$$dist(x_i, B) = y_i f_{\tilde{\beta}}(x_i) / \|f'_{\tilde{\beta}}(x_i)\|$$

- ▶ Distance from a (finite) set C_ℓ to B : ($\ell \in \{-1, 1\}$)

$$dist(C_\ell, B) = \min \{ dist(x_i, B) \mid x_i \in C_\ell \}$$



Margin: Formal definition

Well-classified points and Margin $(f_{\tilde{\beta}}(x) = \langle \beta, x \rangle + \beta_0 = 0)$

- ▶ Sets of well-classified points:

$$C_1 = \{x_i \mid y_i = 1, y_i f(x_i) > 0\} \quad C_{-1} = \{x_i \mid y_i = -1, y_i f(x_i) > 0\}$$

- ▶ Distance from x_i to B :

$$dist(x_i, B) = y_i f_{\tilde{\beta}}(x_i) / \|f'_{\tilde{\beta}}(x_i)\|$$

- ▶ Distance from a (finite) set C_ℓ to B : ($\ell \in \{-1, 1\}$)

$$dist(C_\ell, B) = \min \{ dist(x_i, B) \mid x_i \in C_\ell \}$$

Definition (Margin)

The margin $M > 0$ is given by

$$M = M(B) = \min [dist(C_1, B), dist(C_{-1}, B)]$$

Remark:

Margin: function of the hyperplane

Optimization problem and Margin

Kernel Machines

Alain Celisse

$$M = M(B) = \min (\text{dist}(C_1, B), \text{dist}(C_{-1}, B))$$

Maximizing the margin

- ▶ For all well-classified points x_i , require:

$$\text{dist}(x_i, B) = y_i (\langle \beta, x_i \rangle + \beta_0) / \|\beta\| \geq M > 0$$

- ▶ M is a “confidence level” on well-classified points

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Optimization problem and Margin

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

$$M = M(B) = \min(\text{dist}(C_1, B), \text{dist}(C_{-1}, B))$$

Maximizing the margin

- ▶ For all well-classified points x_i , require:

$$\text{dist}(x_i, B) = y_i (\langle \beta, x_i \rangle + \beta_0) / \|\beta\| \geq M > 0$$

- ▶ M is a “confidence level” on well-classified points

Optimization problem

Optimal hyperplane given by the unique (β, β_0) such that M is maximum

Optimization problem and Margin

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

$$M = M(B) = \min(\text{dist}(C_1, B), \text{dist}(C_{-1}, B))$$

Maximizing the margin

- ▶ For all well-classified points x_i , require:

$$\text{dist}(x_i, B) = y_i (\langle \beta, x_i \rangle + \beta_0) / \|\beta\| \geq M > 0$$

- ▶ M is a “confidence level” on well-classified points

Optimization problem

Optimal hyperplane given by the unique (β, β_0) such that M is maximum

Solve ($M = M(\beta, \beta_0)$):

$$\max_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} M, \quad \text{s.t.}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq M \|\beta\|, \quad \text{for } i = 1, \dots, n$$

Optimization problem and Margin

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

$$M = M(B) = \min(\text{dist}(C_1, B), \text{dist}(C_{-1}, B))$$

Maximizing the margin

- ▶ For all well-classified points x_i , require:

$$\text{dist}(x_i, B) = y_i (\langle \beta, x_i \rangle + \beta_0) / \|\beta\| \geq M > 0$$

- ▶ M is a “confidence level” on well-classified points

Optimization problem

Optimal hyperplane given by the unique (β, β_0) such that M is maximum

Solve ($M = M(\beta, \beta_0)$):

$$\max_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} M, \quad \text{s.t.}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq M \|\beta\|, \quad \text{for } i = 1, \dots, n$$

- ▶ Not scale invariant: β and $2 \cdot \beta$ lead to the same solution for M (Hint: $(\langle \beta, x_i \rangle + \beta_0) / \|\beta\| = \langle x - x_0, \beta^* \rangle$)

Optimization problem reformulation

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

- ▶ Add a constraint on M :

$$M \cdot \|\beta\| = 1$$

Recasting the optimization problem

$$\max_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} M, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq M \|\beta\|, \quad \text{for } i = 1, \dots, n$$

$$\Leftrightarrow \min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimization problem reformulation

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Linearly separable classes

Perceptron and optimal
hyperplane

Support Vectors

Reproducing
Kernels

Support Vector
Machines

- ▶ Add a constraint on M :

$$M \cdot \|\beta\| = 1$$

Recasting the optimization problem

$$\max_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} M, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq M \|\beta\|, \quad \text{for } i = 1, \dots, n$$

$$\Leftrightarrow \min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Remark:

Square the norm without changing the minimum location
(differentiability)

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes
Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Support Vectors

Lagrange multipliers formulation

Kernel Machines

Alain Celisse

Constrained optimization problem

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Lagrange multipliers formulation

Kernel Machines

Constrained optimization problem

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Lagrange multipliers

$$(\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}_+^n)$$

First step: Primal formulation

1. Minimize w.r.t β, β_0 :

$$\mathcal{L}_{P,\alpha}(\beta, \beta_0) = \frac{\|\beta\|^2}{2} - \sum_{i=1}^n \alpha_i (y_i (\langle \beta, x_i \rangle + \beta_0) - 1)$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Lagrange multipliers formulation

Constrained optimization problem

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Lagrange multipliers

$$(\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}_+^n)$$

First step: Primal formulation

1. Minimize w.r.t β, β_0 :

$$\mathcal{L}_{P,\alpha}(\beta, \beta_0) = \frac{\|\beta\|^2}{2} - \sum_{i=1}^n \alpha_i (y_i (\langle \beta, x_i \rangle + \beta_0) - 1)$$

2. Solution (stationary points): For $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}_+^n$

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j y_i y_j x_i^\top x_j$$

Lagrange multipliers formulation (Cont'd)

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Lagrange multipliers

Second step: Dual formulation

1. Maximize:

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j y_i y_j x_i^\top x_j$$

such that: $\alpha_i \geq 0, \quad \sum_{j=1}^n \alpha_j y_j = 0$

Lagrange multipliers formulation (Cont'd)

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Lagrange multipliers

Second step: Dual formulation

1. Maximize:

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j y_i y_j x_i^\top x_j$$

such that: $\alpha_i \geq 0, \quad \sum_{j=1}^n \alpha_j y_j = 0$

2. Solution (KKT conditions):

Solution to this problem: $(\hat{\alpha}_1, \dots, \hat{\alpha}_n)$

- ▶ If $\hat{\alpha}_i > 0$: $y_i(\langle x_i, \hat{\beta} \rangle + \hat{\beta}_0) = 1$
 $\Leftrightarrow x_i$ on the boundary of the slab

Lagrange multipliers formulation (Cont'd)

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Lagrange multipliers

Second step: Dual formulation

1. Maximize:

$$\mathcal{L}_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{1 \leq i, j \leq n} \alpha_i \alpha_j y_i y_j x_i^\top x_j$$

such that: $\alpha_i \geq 0, \quad \sum_{j=1}^n \alpha_j y_j = 0$

2. Solution (KKT conditions):

Solution to this problem: $(\hat{\alpha}_1, \dots, \hat{\alpha}_n)$

- ▶ If $\hat{\alpha}_i > 0$: $y_i(\langle x_i, \hat{\beta} \rangle + \hat{\beta}_0) = 1$
 $\Leftrightarrow x_i$ on the boundary of the slab
- ▶ If $\hat{\alpha}_i = 0$: $y_i(\langle x_i, \hat{\beta} \rangle + \hat{\beta}_0) > 1$
 $\Leftrightarrow x_i$ out of the slab

Support Vector classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

► If $\hat{\alpha}_i > 0$: $y_i(\langle x_i, \hat{\beta} \rangle + \hat{\beta}_0) = 1$

Definition (Support vectors)

The points x_i s such that $\hat{\alpha}_i > 0$ are:

► called *support vectors*. $(S = \{1 \leq i \leq n \mid \hat{\alpha}_i > 0\})$

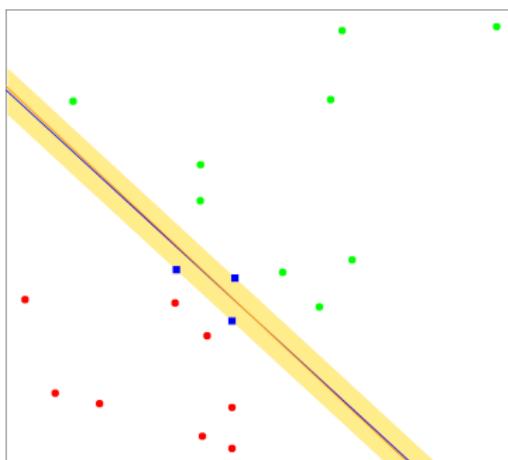
Support Vector classifier

- If $\hat{\alpha}_i > 0$: $y_i(\langle x_i, \hat{\beta} \rangle + \hat{\beta}_0) = 1$

Definition (Support vectors)

The points x_i s such that $\hat{\alpha}_i > 0$ are:

- called *support vectors*. $(S = \{1 \leq i \leq n \mid \hat{\alpha}_i > 0\})$
- located on the boundary of the slab that is, their distance to the hyperplane is equal to the margin.



Support vector classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Support Vector Classifier

From previous equality conditions:

$$\blacktriangleright \hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i y_i x_i = \sum_{i \in S} \hat{\alpha}_i y_i x_i$$

Support vector classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Support Vector Classifier

From previous equality conditions:

- ▶ $\hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i y_i x_i = \sum_{i \in S} \hat{\alpha}_i y_i x_i$
- ▶ $\hat{\beta}_0 = y_i - \left\langle x_i, \sum_{j=1}^n \hat{\alpha}_j y_j x_j \right\rangle$, for all $i \in S$

Support vector classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Support Vector Classifier

From previous equality conditions:

- ▶ $\hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i y_i x_i = \sum_{i \in S} \hat{\alpha}_i y_i x_i$
- ▶ $\hat{\beta}_0 = y_i - \left\langle x_i, \sum_{j=1}^n \hat{\alpha}_j y_j x_j \right\rangle$, for all $i \in S$
- ▶ Classifier: $\hat{f}(x) = \text{sign} \left(\left\langle \hat{\beta}, x \right\rangle + \hat{\beta}_0 \right)$

Support vector classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Support Vector Classifier

From previous equality conditions:

- $\hat{\beta} = \sum_{i=1}^n \hat{\alpha}_i y_i x_i = \sum_{i \in S} \hat{\alpha}_i y_i x_i$
- $\hat{\beta}_0 = y_i - \left\langle x_i, \sum_{j=1}^n \hat{\alpha}_j y_j x_j \right\rangle$, for all $i \in S$
- Classifier: $\hat{f}(x) = \text{sign} \left(\left\langle \hat{\beta}, x \right\rangle + \hat{\beta}_0 \right)$

Proof.

$$\begin{aligned} y_i (\left\langle x_i, \hat{\beta} \right\rangle + \hat{\beta}_0) = 1 &\Leftrightarrow y_i \hat{\beta}_0 = 1 - y_i \left\langle x_i, \hat{\beta} \right\rangle \\ &\Leftrightarrow \hat{\beta}_0 = y_i (1 - y_i \left\langle x_i, \sum_{j=1}^n \hat{\alpha}_j y_j x_j \right\rangle) \end{aligned}$$



From non-overlap. to overlap. classes

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

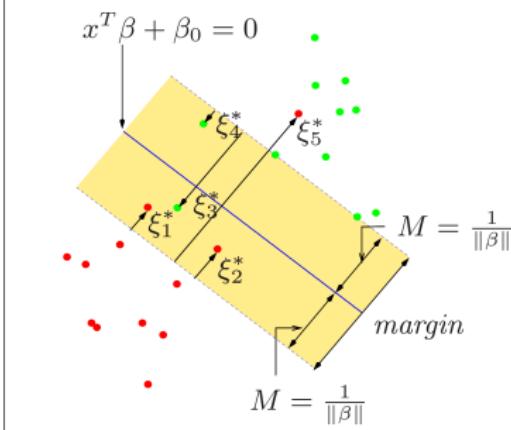
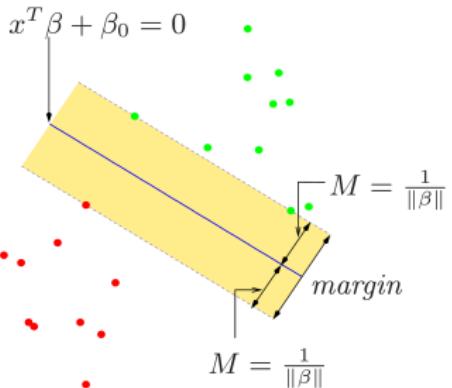
Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines



Limitation: The overlapping case

- ▶ The previous procedure addresses the non-overlapping case
- ▶ With overlapping classes, such a linear classification rule commits mistakes
- ▶ Requires to modify the optimality criterion

Misclassifications and Hinge loss

Kernel Machines

Alain Celisse

Former formulation

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Accounting for misclassifications

- ▶ x_i misclassified means:

$$y_i (\langle \beta, x_i \rangle + \beta_0) < 1 \Leftrightarrow 0 < 1 - y_i (\langle \beta, x_i \rangle + \beta_0)$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Misclassifications and Hinge loss

Kernel Machines

Alain Celisse

Former formulation

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{such that}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Accounting for misclassifications

- x_i misclassified means:

$$y_i (\langle \beta, x_i \rangle + \beta_0) < 1 \Leftrightarrow 0 < 1 - y_i (\langle \beta, x_i \rangle + \beta_0)$$

- **Minimizing misclassifications** amounts to minimize

$$\sum_{i=1}^n (1 - y_i (\langle \beta, x_i \rangle + \beta_0))_+ = \sum_{i=1}^n h(y_i, f(x_i))$$

→ Only misclassifications come into play!

Hinge loss

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

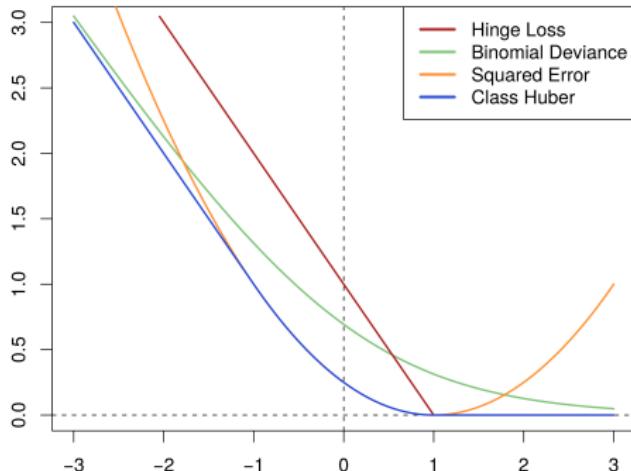
Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

$$h(y, f(x)) = (1 - y \cdot f(x))_+ = \max(0, 1 - y \cdot f(x))$$



$$x \mapsto (1 - x)_+$$

Hinge-loss and optimal classifier

Kernel Machines

Alain Celisse

Hinge-loss

$$h(y, f(x)) = (1 - y \cdot f(x))_+$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Hinge-loss and optimal classifier

Kernel Machines

Alain Celisse

Hinge-loss

$$h(y, f(x)) = (1 - y \cdot f(x))_+$$

Prediction error (Misclassification rate) for the Hinge-loss

$$PE(f) = \mathbb{E}' [h(Y', f(X'))]$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Hinge-loss and optimal classifier

Kernel Machines

Alain Celisse

Hinge-loss

$$h(y, f(x)) = (1 - y \cdot f(x))_+$$

Prediction error (Misclassification rate) for the Hinge-loss

$$PE(f) = \mathbb{E}' [h(Y', f(X'))]$$

Optimal function minimizer

$$f^* = \operatorname{Arg} \min_{f \in \mathcal{M}} \mathbb{E}' [h(Y', f(X'))]$$

yields

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Hinge-loss and optimal classifier

Kernel Machines

Alain Celisse

Hinge-loss

$$h(y, f(x)) = (1 - y \cdot f(x))_+$$

Prediction error (Misclassification rate) for the Hinge-loss

$$PE(f) = \mathbb{E}' [h(Y', f(X'))]$$

Optimal function minimizer

$$f^* = \operatorname{Arg} \min_{f \in \mathcal{M}} \mathbb{E}' [h(Y', f(X'))]$$

yields

$$\begin{aligned} f^*(x) &= \operatorname{sign} \left[\frac{2\mathbb{P}(Y = +1 | X = x) - 1}{2} \right] \\ &= \operatorname{sign} [2\mathbb{P}(Y = +1 | X = x) - 1] \\ &= \operatorname{sign} (\mathbb{E}[Y | X = x]) \end{aligned}$$

Remark:

→ Justifies the classifier expression: $\operatorname{sign}(\langle \beta, x \rangle + \beta_0)$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Optimization problem and regularization

Kernel Machines

Alain Celisse

Past formulation

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Optimization problem and regularization

Kernel Machines

Alain Celisse

Past formulation

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Relaxation

Solve for $r \geq 1$

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$\sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) \leq r$$

Optimization problem and regularization

Kernel Machines

Alain Celisse

Past formulation

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$y_i (\langle \beta, x_i \rangle + \beta_0) \geq 1, \quad \text{for } i = 1, \dots, n$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Relaxation

Solve for $r \geq 1$

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$\sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) \leq r$$

Remarks:

- ▶ Tackles misclassifications
- ▶ Addresses overlapping classes

Optimization problem and regularization

Kernel Machines

Alain Celisse

Relaxation

Solve for $r \geq 1$

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$\sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) \leq r$$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Optimization problem and regularization

Kernel Machines

Alain Celisse

Relaxation

Solve for $r \geq 1$

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$\sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) \leq r$$

Regularized formulation

Solve, for all $\lambda > 0$:

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) + \lambda \|\beta\|^2 \right\}$$

where $h(y, f(x)) = (1 - yf(x))_+$, $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Optimization problem and regularization

Kernel Machines

Alain Celisse

Relaxation

Solve for $r \geq 1$

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \|\beta\|^2 / 2, \quad \text{s.t.}$$

$$\sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) \leq r$$

Regularized formulation

Solve, for all $\lambda > 0$:

$$\min_{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(x_i)) + \lambda \|\beta\|^2 \right\}$$

where $h(y, f(x)) = (1 - yf(x))_+$, $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

Remarks:

- ▶ $\sum_{i=1}^n h(y_i, f(x_i))$ quantifies the fit to the data
- ▶ $\|\beta\|^2$: regularization term (avoids too complex β s)
- ▶ Similar idea to the one discussed for the LASSO pb!

Optimal
Separating
Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

Limitations of the SV classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

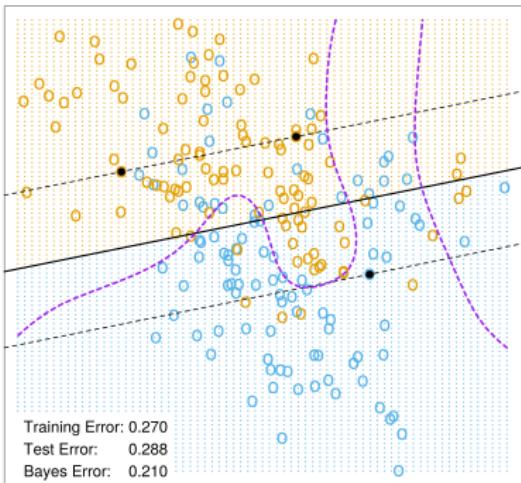
Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines

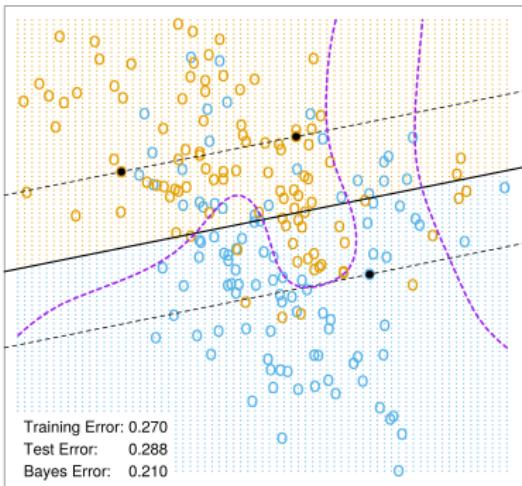


Sub-optimal linear classifier: $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

Limitations of the SV classifier

Kernel Machines

Alain Celisse



Sub-optimal linear classifier: $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

- With strongly overlapping classes, linear classifiers are less effective (straight line): lots of misclassifications

Optimal Separating Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing Kernels

Support Vector Machines

Limitations of the SV classifier

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

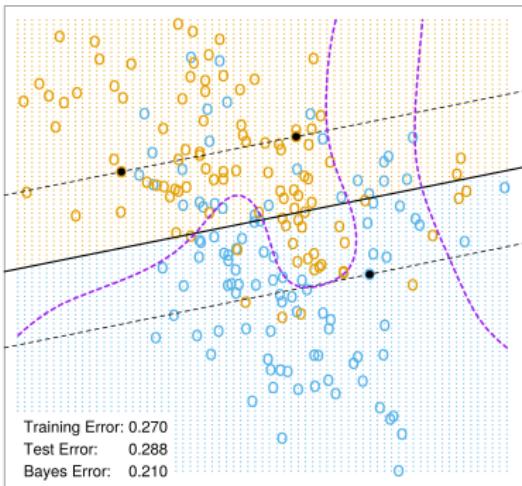
Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing
Kernels

Support Vector
Machines



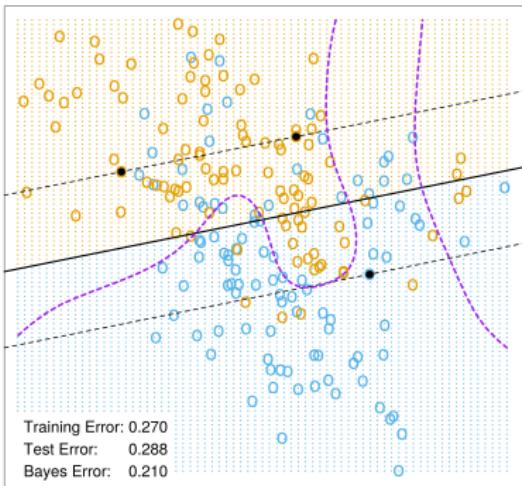
Sub-optimal linear classifier: $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

- ▶ With strongly overlapping classes, linear classifiers are less effective (straight line): lots of misclassifications
- ▶ Comparison with the Bayes classifier (purple dashed line)

Limitations of the SV classifier

Kernel Machines

Alain Celisse



Optimal Separating Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing Kernels

Support Vector Machines

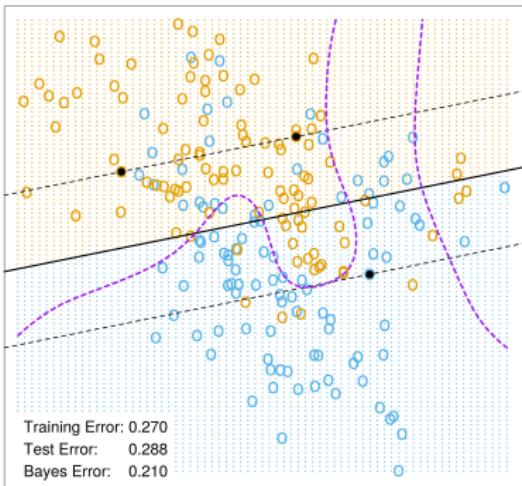
Sub-optimal linear classifier: $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

- ▶ With strongly overlapping classes, linear classifiers are less effective (straight line): lots of misclassifications
- ▶ Comparison with the Bayes classifier (purple dashed line)
- ▶ Need for alternative non-linear classifiers

Limitations of the SV classifier

Kernel Machines

Alain Celisse



Optimal Separating Hyperplanes

Support Vectors

Non-overlapping classes

Overlapping classes

Reproducing Kernels

Support Vector Machines

Sub-optimal linear classifier: $f_{\beta, \beta_0}(x) = \langle \beta, x \rangle + \beta_0$

- ▶ With strongly overlapping classes, linear classifiers are less effective (straight line): lots of misclassifications
- ▶ Comparison with the Bayes classifier (purple dashed line)
- ▶ Need for alternative non-linear classifiers
→ Where Reproducing kernels come into play...

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries
PSD Kernels
Reproducing kernels

Support Vector
Machines

Reproducing Kernels

Improvement with Reproducing Kernels

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

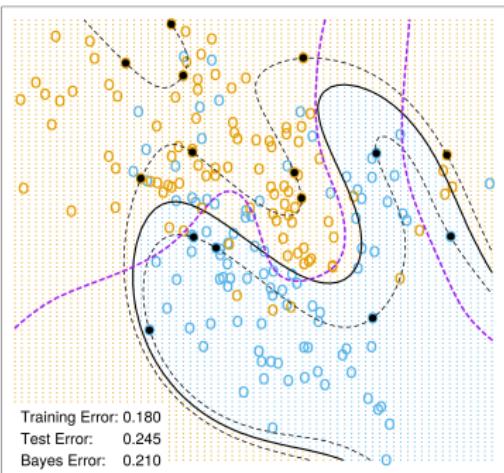
Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels
Reproducing kernels

Support Vector
Machines



Reproducing kernels provide non-linear boundaries

- ▶ The Bayes classifier: Non-linear boundary (purple dashed line)
- ▶ The Reproducing Kernel-based classifier (black line) is close to optimal

Illustration: “Polynomial extension”

Kernel Machines

Alain Celisse

Example: Degree 2-polynomial extension

- ▶ $x = (x_1, x_2) \in \mathbb{R}^2$ ($d = 2$) and $y \in \{0, 1\}$
- ▶ Create new features:

$$\phi : \quad x \in \mathbb{R}^2 \mapsto \phi(x) = (x_1^2, x_1 \cdot x_2, x_2 \cdot x_1, x_2^2)^\top \in \mathbb{R}^4$$

- ▶ $\phi(x) \in \mathbb{R}^4$: higher dimensional vector than $x \in \mathbb{R}^2$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

Illustration: “Polynomial extension”

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

Example: Degree 2-polynomial extension

- ▶ $x = (x_1, x_2) \in \mathbb{R}^2$ ($d = 2$) and $y \in \{0, 1\}$

- ▶ Create new features:

$$\phi : \quad x \in \mathbb{R}^2 \mapsto \phi(x) = (x_1^2, x_1 \cdot x_2, x_2 \cdot x_1, x_2^2)^\top \in \mathbb{R}^4$$

- ▶ $\phi(x) \in \mathbb{R}^4$: higher dimensional vector than $x \in \mathbb{R}^2$

Support Vectors and extended features

Solve, for all $\lambda > 0$:

$$\min_{\beta \in \mathbb{R}^4, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, f_{\beta, \beta_0}(\phi(x_i))) + \lambda \|\beta\|^2 \right\}$$

where $f_{\beta, \beta_0}(\phi(x_i)) = \langle \beta, \phi(x_i) \rangle_4 + \beta_0$

(see previous picture)

From extended features to PSD kernels

Kernel Machines

Alain Celisse

From extended features to kernel

- ▶ $x, x' \in \mathbb{R}^d$
- ▶ Extended features: $\phi(x) \in \mathbb{R}^p$
- ▶ Kernel: $k(\cdot, \cdot) : (x, x') \mapsto k(x, x') \in \mathbb{R}$

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_p$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

From extended features to PSD kernels

Kernel Machines

Alain Celisse

From extended features to kernel

- ▶ $x, x' \in \mathbb{R}^d$
- ▶ Extended features: $\phi(x) \in \mathbb{R}^p$
- ▶ Kernel: $k(\cdot, \cdot) : (x, x') \mapsto k(x, x') \in \mathbb{R}$

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_p$$

Definition (Psd kernel)

A kernel $k(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is positive semi-definite if

- ▶ $k(x, y) = k(y, x)$, for all $x, y \in \mathbb{R}^d$ (symmetric)
- ▶ $\forall n \in \mathbb{N}^*, \forall x_1, \dots, x_n \in \mathbb{R}^d, \forall a_1, \dots, a_n \in \mathbb{R}$,

$$\sum_{1 \leq i, j \leq n} a_i a_j k(x_i, x_j) \geq 0$$

which is equivalent to

$\forall n \in \mathbb{N}^*, \forall x_1, \dots, x_n \in \mathbb{R}^d$,

the matrix $K = \{k(x_i, x_j)\}_{1 \leq i, j \leq n} \in \mathcal{S}_n^+(\mathbb{R})$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

PSD kernel: Examples

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

How to build a PSD kernel? First example

- ▶ $\{\psi_a(\cdot)\}_{1 \leq a \leq A}$: a set (basis) of real-valued functions
- ▶ $\phi(x) = (\psi_1(x), \psi_2(x), \dots, \psi_A(x))^\top \in \mathbb{R}^A$
- ▶ PSD kernel:

$$k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathbb{R}^A}$$

Other classical examples

- ▶ Linear kernel:

$$k(x, y) = \langle x, y \rangle_{\mathbb{R}^d}$$

- ▶ Polynomial kernel: $(c \geq 0, d > 0)$

$$k(x, y) = (\langle x, y \rangle_{\mathbb{R}^d} + c)^d$$

- ▶ Gaussian (Radial Basis Function) kernel:

$$k(x, y) = e^{-\frac{(x-y)^2}{2}}$$

Reproducing kernel and RKHS (Cont'd)

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

Definition (Reproducing kernel)

$k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a reproducing kernel on a Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ if

- ▶ \mathcal{H} contains all functions $k_x : x \mapsto k(x, \cdot)$, for all $x \in \mathcal{X}$
- ▶ $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$,

$$f(x) = \langle f, k_x \rangle_{\mathcal{H}} \quad (\text{Reproducing property})$$

Then, $\mathcal{H} = \mathcal{H}_k$: Reproducing Kernel Hilbert Space (RKHS) associated with k .

Definition (Gram matrix)

The so-called Gram matrix is the $n \times n$ matrix denoted by K and defined by

$$K = \{k(x_i, x_j)\}_{1 \leq i, j \leq n}$$

Finite-dimensional examples

- ▶ If $\text{Card}(\mathcal{X}) < +\infty$, then $\mathcal{H} = \mathbb{R}^{\text{Card}(\mathcal{X})}$
- ▶ If $k(x, y) = \langle x, y \rangle_d$, then $\mathcal{H} = \mathbb{R}^d$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
KernelsNon-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines

Infinite-dimensional examples

- ▶ If $k(x, y) = \min(x, y)$ ($x, y \in [0, 1]$), then

$$\mathcal{H} = \left\{ g : \mathcal{X} \rightarrow \mathbb{R} \mid g \in L^2(\mathcal{X}), g(0) = 0, \int_{\mathcal{X}} (g'(t))^2 dx < +\infty \right\}$$

→ function space...

Remark:

RKHS is an example where high-dimensional tools are involved (High-dimensional statistics) although data can belong to \mathbb{R}^d with $d = 1$!

Canonical feature map and kernel trick

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

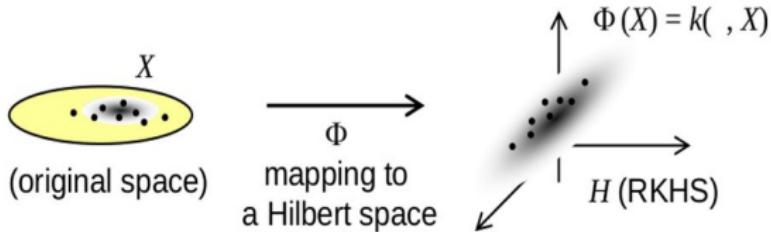
Reproducing
Kernels

Non-linear decision
boundaries

PSD Kernels

Reproducing kernels

Support Vector
Machines



Problem:

\mathcal{H} is high-dimensional (function space)

Kernel trick

$$\langle k_x, k_y \rangle_{\mathcal{H}} = k(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathbb{R}^d$$

- ▶ Avoids manipulating data that are high-dimensional in \mathcal{H}

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Support Vector Machines

From SV to SVM classifiers

Support Vectors with extended feature vectors

- Extended feature vector: $\phi(x) \in \mathbb{R}^p$
- Solve, for all $\lambda > 0$:

$$\hat{f} = \operatorname{Arg} \min_{\beta \in \mathbb{R}^p, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle \beta, \phi(x_i) \rangle_p + \beta_0) + \lambda \|\beta\|^2 \right\}$$

Support vectors with general feature map

- Kernel-based feature vector: $\phi(x) = k(x, \cdot) \in \mathcal{H}$
- Key ingredients: $(g \in \mathcal{H})$

$$\begin{array}{lll} \text{Replace:} & \langle \beta, \phi(x_i) \rangle_p & \text{by} & \langle g, \phi(x_i) \rangle_{\mathcal{H}} \\ & \|\beta\| & \text{by} & \|g\|_{\mathcal{H}} \end{array}$$

Support Vector Machine (SVM)

Solve, for all $\lambda > 0$:

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

where $\hat{f}(x) = \langle \hat{g}, \phi(x) \rangle_{\mathcal{H}} + \hat{\beta}_0$, for all $x \in \mathcal{X}$

Support Vector Machine (SVM)

Solve, for all $\lambda > 0$:

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

where $\hat{f}(x) = \langle \hat{g}, \phi(x) \rangle_{\mathcal{H}} + \hat{\beta}_0$, for all $x \in \mathcal{X}$

Remark:

- ▶ $h(y, f(x)) = (1 - y \cdot f(x))_+$
- ▶ $\sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0)$ measures how well g fits the data
- ▶ $\lambda \|g\|_{\mathcal{H}}$: regularization term
 - penalizes for too large (complex) $g \in \mathcal{H}$
- ▶ $\lambda > 0$: Regularization parameter
 - ▶ The larger λ , the smaller $\|\hat{g}\|_{\mathcal{H}}$, the less complex \hat{g}
 - ▶ Tradeoff between: fitting the data vs model complexity

Optimal Separating Hyperplanes

Support Vectors

Reproducing Kernels

Support Vector Machines

Binary classification

Primal/Dual formulations

Support vectors

Representer theorem

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

Theorem (Representer theorem)

$\Psi : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathbb{R}$, nondecreasing w.r.t. its $n+1$ th argument

$$\operatorname{Arg} \min_{g \in \mathcal{H}} \{\Psi [g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}}]\}$$

Any solution \hat{g} to the above optimization problem can be written as

$$\hat{g}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x), \quad \forall x \in \mathcal{X}$$

where $\hat{\alpha}_i \in \mathbb{R}$, for all $1 \leq i \leq n$

Proof of the Representer theorem

Kernel Machines

Alain Celisse

Proof.

- ▶ $V = \text{Vect}(k_{x_1}, \dots, k_{x_n})$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Proof of the Representer theorem

Kernel Machines

Alain Celisse

Proof.

- ▶ $V = \text{Vect}(k_{x_1}, \dots, k_{x_n})$
- ▶ Any $f \in V$ satisfies

$$x \mapsto f(x) = \sum_{i=1}^n \alpha_i k(x_i, x) \in \mathbb{R}$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Proof of the Representer theorem

Kernel Machines

Alain Celisse

Proof.

► $V = \text{Vect}(k_{x_1}, \dots, k_{x_n})$

► Any $f \in V$ satisfies

$$x \mapsto f(x) = \sum_{i=1}^n \alpha_i k(x_i, x) \in \mathbb{R}$$

Let g be a solution to this problem. Then



$$g(x_j) = g_V(x_j) + g_{V^\perp}(x_j) = g_V(x_j) + \underbrace{\langle g_{V^\perp}, k_{x_j} \rangle_{\mathcal{H}}}_{=0}$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Proof of the Representer theorem

Kernel Machines

Alain Celisse

Proof.

► $V = \text{Vect}(k_{x_1}, \dots, k_{x_n})$

► Any $f \in V$ satisfies

$$x \mapsto f(x) = \sum_{i=1}^n \alpha_i k(x_i, x) \in \mathbb{R}$$

Let g be a solution to this problem. Then



$$g(x_j) = g_V(x_j) + g_{V^\perp}(x_j) = g_V(x_j) + \underbrace{\langle g_{V^\perp}, k_{x_j} \rangle_{\mathcal{H}}}_{=0}$$



$$\|g\|_{\mathcal{H}}^2 = \|g_V\|_{\mathcal{H}}^2 + \|g_{V^\perp}\|_{\mathcal{H}}^2 \geq \|g_V\|_{\mathcal{H}}^2$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Proof of the Representer theorem

Kernel Machines

Alain Celisse

Proof.

► $V = \text{Vect}(k_{x_1}, \dots, k_{x_n})$

► Any $f \in V$ satisfies

$$x \mapsto f(x) = \sum_{i=1}^n \alpha_i k(x_i, x) \in \mathbb{R}$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Let g be a solution to this problem. Then

►

$$g(x_j) = g_V(x_j) + g_{V^\perp}(x_j) = g_V(x_j) + \underbrace{\langle g_{V^\perp}, k_{x_j} \rangle_{\mathcal{H}}}_{=0}$$

►

$$\|g\|_{\mathcal{H}}^2 = \|g_V\|_{\mathcal{H}}^2 + \|g_{V^\perp}\|_{\mathcal{H}}^2 \geq \|g_V\|_{\mathcal{H}}^2$$

Hence,

$$\Psi[g(x_1), \dots, g(x_n), \|g\|_{\mathcal{H}}] \geq \Psi[g_V(x_1), \dots, g_V(x_n), \|g_V\|_{\mathcal{H}}]$$

Representer thm and finite dimension

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

$$\hat{g}(x) = \sum_{j=1}^n \hat{\alpha}_j k(x_j, x), \quad \forall x \in \mathcal{X}$$

From infinite to finite dimension

Solving

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

Representer thm and finite dimension

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

$$\hat{g}(x) = \sum_{j=1}^n \hat{\alpha}_j k(x_j, x), \quad \forall x \in \mathcal{X}$$

From infinite to finite dimension

Solving

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

equivalent to solving

(K: Gram matrix)

$$\hat{\alpha} = \operatorname{Arg} \min_{\alpha \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, [K\alpha]_i + \beta_0) + \lambda \alpha^\top K \alpha \right\}$$

Representer thm and finite dimension

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

$$\hat{g}(x) = \sum_{j=1}^n \hat{\alpha}_j k(x_j, x), \quad \forall x \in \mathcal{X}$$

From infinite to finite dimension

Solving

$$\hat{f} = \operatorname{Arg} \min_{g \in \mathcal{H}, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \|g\|_{\mathcal{H}}^2 \right\}$$

equivalent to solving

(K: Gram matrix)

$$\hat{\alpha} = \operatorname{Arg} \min_{\alpha \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, [K\alpha]_i + \beta_0) + \lambda \alpha^\top K \alpha \right\}$$

Remarks:

- ▶ Representer thm \Rightarrow Finite dimensional space
- ▶ Dimension of the solution space at most $n!$
- ▶ Not yet any closed-form expression for the solution...

Introducing Slack variables ξ_i s

Solving

(fixed $\lambda > 0$)

$$\min_{\alpha \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \alpha^T K \alpha \right\}$$

is equivalent to solving

Introducing Slack variables ξ_i s

Solving

(fixed $\lambda > 0$)

$$\min_{\alpha \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \alpha^T K \alpha \right\}$$

is equivalent to solving

$$\min_{\alpha, \xi \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n \xi_i + \lambda \alpha^T K \alpha \right\}$$

such that $\xi_i \geq h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0)$, $\forall 1 \leq i \leq n$

is equivalent to solving

Introducing Slack variables ξ_i s

Solving

(fixed $\lambda > 0$)

$$\min_{\alpha \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0) + \lambda \alpha^T K \alpha \right\}$$

is equivalent to solving

$$\min_{\alpha, \xi \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n \xi_i + \lambda \alpha^T K \alpha \right\}$$

such that $\xi_i \geq h(y_i, \langle g, \phi(x_i) \rangle_{\mathcal{H}} + \beta_0), \quad \forall 1 \leq i \leq n$

is equivalent to solving

$$\min_{\alpha, \xi \in \mathbb{R}^n, \beta_0 \in \mathbb{R}} \left\{ \sum_{i=1}^n \xi_i + \lambda \alpha^T K \alpha \right\}$$

such that $\xi_i \geq 1 - y_i \cdot \left[\sum_{j=1}^n \alpha_j k(x_j, x_i) + \beta_0 \right], \quad \forall 1 \leq i \leq n$

$$\xi_i \geq 0, \quad \forall 1 \leq i \leq n$$

Lagrangian formulation

Kernel Machines

Alain Celisse

Lagrangian formulation ($\mu, \nu \geq 0$: Langrange multipliers)

$$L_P(\xi, \alpha; \mu, \nu) = \sum_{i=1}^n \xi_i + \lambda \alpha^\top K \alpha - \sum_{i=1}^n \mu_i \left[\xi_i - 1 + y_i \cdot \left[\sum_{j=1}^n \alpha_j k(x_j, x_i) + \beta_0 \right] \right] - \sum_{i=1}^n \nu_i \xi_i$$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Lagrangian formulation

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Lagrangian formulation ($\mu, \nu \geq 0$: Langrange multipliers)

$$L_P(\xi, \alpha; \mu, \nu) = \sum_{i=1}^n \xi_i + \lambda \alpha^\top K \alpha - \sum_{i=1}^n \mu_i \left[\xi_i - 1 + y_i \cdot \left[\sum_{j=1}^n \alpha_j k(x_j, x_i) + \beta_0 \right] \right] - \sum_{i=1}^n \nu_i \xi_i$$

Primal problem

- ▶ Minimizing $L_P(\xi, \alpha; \mu, \nu)$ w.r.t. α, ξ (computing derivatives)
- ▶ Leads to the dual $L_D(\mu, \nu) = L_P(\hat{\xi}(\mu, \nu), \hat{\alpha}(\mu, \nu); \mu, \nu)$

Dual optimization problem

Optimizing the Dual

$$\max_{0 \leq \mu_i \leq 1/n} \underbrace{\left\{ \sum_{i=1}^n \mu_i - \frac{1}{4\lambda} \sum_{1 \leq i,j \leq n} y_i y_j \mu_i \mu_j k(x_i, x_j) \right\}}_{=L_D(\mu, \nu)}$$

solved for $\mu_i = 2\lambda\alpha_i y_i$, for all $1 \leq i \leq n$ $(\mu_i + \nu_i = 1/n)$

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Dual optimization problem

Optimizing the Dual

$$\max_{\substack{0 \leq \mu_i \leq 1/n}} \underbrace{\left\{ \sum_{i=1}^n \mu_i - \frac{1}{4\lambda} \sum_{1 \leq i,j \leq n} y_i y_j \mu_i \mu_j k(x_i, x_j) \right\}}_{=L_D(\mu, \nu)}$$

solved for $\mu_i = 2\lambda\alpha_i y_i$, for all $1 \leq i \leq n$ ($\mu_i + \nu_i = 1/n$)

Towards a simplified Primal expression

$$\begin{aligned} & \max_{\alpha \in \mathbb{R}^n} \left\{ 2 \sum_{i=1}^n \alpha_i y_i - \sum_{1 \leq i,j \leq n} \alpha_i \alpha_j k(x_i, x_j) \right\} \\ &= \max_{\alpha \in \mathbb{R}^n} \left\{ 2\alpha^\top Y - \alpha^\top K\alpha \right\} \end{aligned}$$

such that $0 \leq \alpha_i y_i \leq \frac{1}{2n\lambda}, \quad \forall i = 1, \dots, n$

Remark:

Easier to solve numerically (No closed-form expression)

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

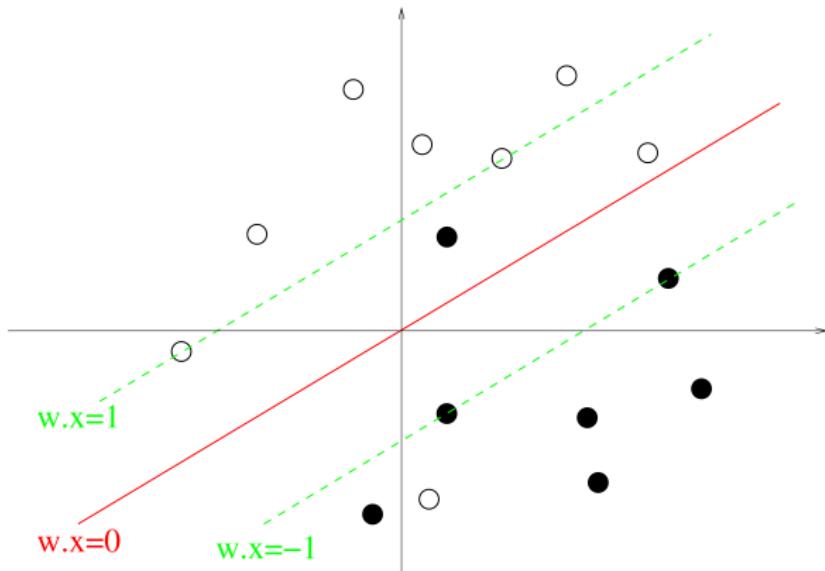
Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

Illustration: Where are support vectors?



Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

KKT and support vectors

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

KKT and support vectors

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

Interpretation

1. $\mu_i = 0$:

- ▶ $\nu_i = 1/n$ and $\xi_i = 0 \Rightarrow 1 < y_i f(x_i)$
- ▶ x_i : well-classified

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

Interpretation

1. $\mu_i = 0$:

- ▶ $\nu_i = 1/n$ and $\xi_i = 0 \Rightarrow 1 < y_i f(x_i)$
- ▶ x_i : well-classified

2. $0 < \mu_i < 1/n$:

- ▶ $0 < \nu_i < 1/n$ and $\xi_i = 0 = 1 - y_i f(x_i) \Rightarrow 1 = y_i f(x_i)$

KKT and support vectors

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

Interpretation

1. $\mu_i = 0$:

- ▶ $\nu_i = 1/n$ and $\xi_i = 0 \Rightarrow 1 < y_i f(x_i)$
- ▶ x_i : well-classified

2. $0 < \mu_i < 1/n$:

- ▶ $0 < \nu_i < 1/n$ and $\xi_i = 0 = 1 - y_i f(x_i) \Rightarrow 1 = y_i f(x_i)$
- ▶ x_i on the slab boundary

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

KKT and support vectors

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

Interpretation

1. $\mu_i = 0$:

- ▶ $\nu_i = 1/n$ and $\xi_i = 0 \Rightarrow 1 < y_i f(x_i)$
- ▶ x_i : well-classified

2. $0 < \mu_i < 1/n$:

- ▶ $0 < \nu_i < 1/n$ and $\xi_i = 0 = 1 - y_i f(x_i) \Rightarrow 1 = y_i f(x_i)$
- ▶ x_i on the slab boundary

3. $\mu_i = 1/n$:

- ▶ $\nu_i = 0$ and $\xi_i \geq 0 \Rightarrow 1 \geq y_i f(x_i)$
- ▶ x_i within the slab or misclassified

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

KKT and support vectors

KKT conditions (when optimizing the Primal problem)

For all $1 \leq i \leq n$,

$$\mu_i [\xi_i - (1 - y_i f(x_i))] = 0 \quad \text{and} \quad \nu_i \xi_i = 0$$

with $0 \leq \mu_i \leq 1/n$ and $\mu_i + \nu_i = 1/n$

Interpretation

1. $\mu_i = 0$:

- ▶ $\nu_i = 1/n$ and $\xi_i = 0 \Rightarrow 1 < y_i f(x_i)$
- ▶ x_i : well-classified

2. $0 < \mu_i < 1/n$:

- ▶ $0 < \nu_i < 1/n$ and $\xi_i = 0 = 1 - y_i f(x_i) \Rightarrow 1 = y_i f(x_i)$
- ▶ x_i on the slab boundary

3. $\mu_i = 1/n$:

- ▶ $\nu_i = 0$ and $\xi_i \geq 0 \Rightarrow 1 \geq y_i f(x_i)$
- ▶ x_i within the slab or misclassified

Definition (SV)

Support vectors are the x_i s such that $\mu_i = 2\lambda y_i \alpha_i > 0$

Illustration: Support vectors locations

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

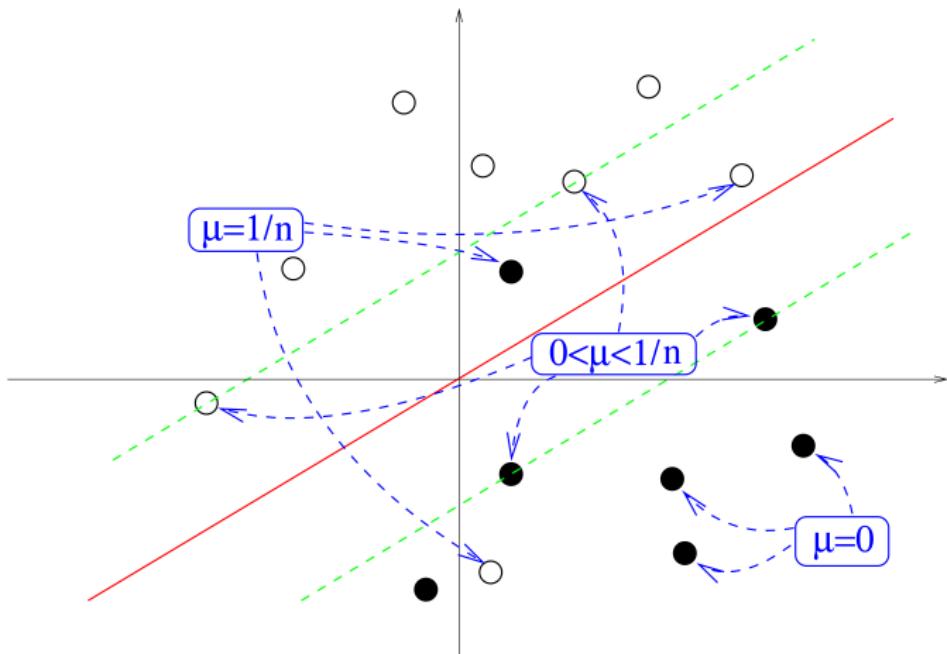
Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors



(from J.-Ph. Vert's slides)

Remark:

- ▶ $f(x) = \sum_{i=1}^n \alpha_i k(x_i, x)$, for all $x \in \mathcal{X}$
- ▶ Only a few α_i s (μ_i s) are non zero! (Support Vectors)

Optimal
Separating
Hyperplanes

Support Vectors

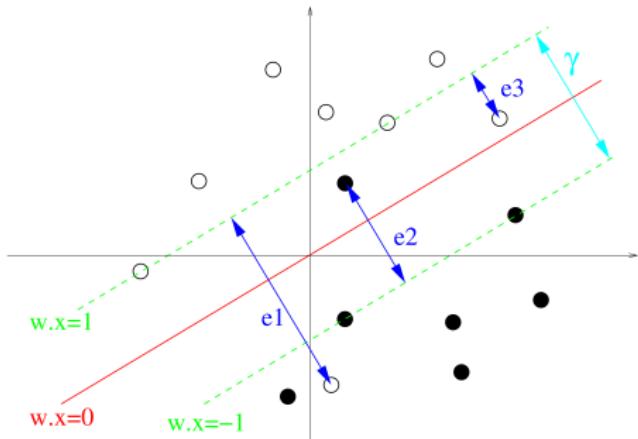
Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors



Definition (Margin)

In the SVM context, the margin $\gamma > 0$ denotes the width of the region within the RKHS between the hyperplanes of equations

$$\langle g, \phi(x) \rangle + \beta_0 = 1, \quad \text{and} \quad \langle g, \phi(x) \rangle + \beta_0 = -1$$

Remark:

$$\gamma = 2 \frac{y_i \left(\langle \hat{g}, \phi(x_i) \rangle + \hat{\beta}_0 \right)}{\|\hat{g}\|_{\mathcal{H}}} = \frac{2}{\|\hat{g}\|_{\mathcal{H}}}$$

SVM improvement

Kernel Machines

Alain Celisse

Optimal
Separating
Hyperplanes

Support Vectors

Reproducing
Kernels

Support Vector
Machines

Binary classification

Primal/Dual
formulations

Support vectors

