



Lecture 13: Linear Support Vector Machines

Prof. Dr. Mario Fritz

2022-06-02

<https://fritz.cispa.saarland>

Trustworthy AI

CISPA Helmholtz Center for Information Security, Saarbrücken, Germany

Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

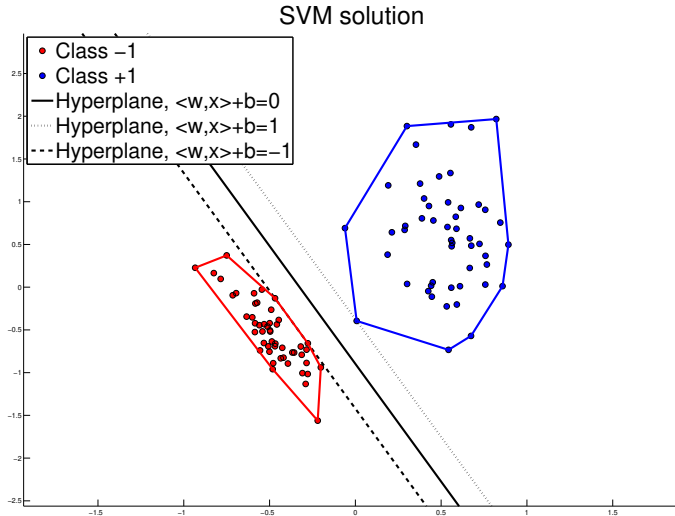
Soft-margin SVM

Summary

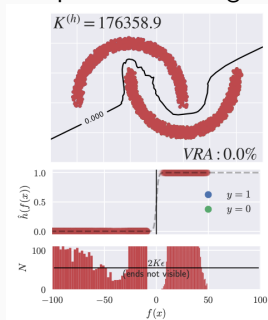
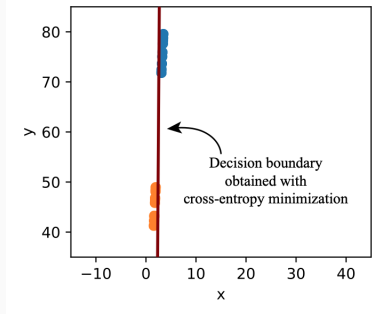
- Monday June 6th public holiday : no lecture
- Wednesday June 8th : online only lecture

Update to last lecture (v2)

Also update to this lecture (extended derivation)



In contrast, some of the recent ML methods do not optimize for margin!



e.g. Kamil Nar, Orhan Ocal, S. Shankar Sastry, Kannan Ramchandran: Cross-Entropy Loss and Low-Rank Features Have Responsibility for Adversarial Examples,

- Learning with Kernels - Chapter 7
- Bishop - Chapter 7

Software and online demos on the web:

- <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- <https://www.csie.ntu.edu.tw/~cjlin/liblinear/>
- <https://jgreitemann.github.io/svm-demo>
- <https://scikit-learn.org/stable/modules/svm.html>

Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

Soft-margin SVM

Summary

Let $\mathcal{X} = \mathbb{R}^d$ be the input space, then the classifier $\hat{y} : \mathbb{R}^d \rightarrow \{-1, 1\}$ has the form

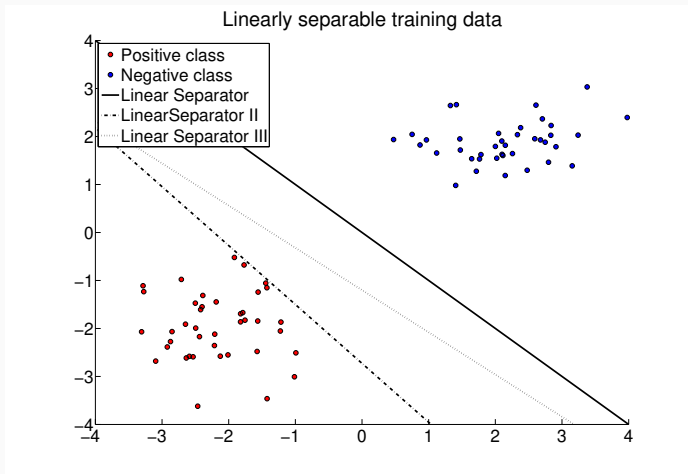
$$\hat{y}(x) = \text{sign}(f(x)) = \text{sign}(\langle w, x \rangle + b) = \begin{cases} 1 & \text{if } \langle w, x \rangle + b > 0, \\ -1 & \text{if } \langle w, x \rangle + b \leq 0. \end{cases}$$

Separation of the input space \mathbb{R}^d into two half spaces.

A training set $D = (x_i, y_i)_{i=1}^n$ is **linearly separable** if there exists a weight vector w and an offset b such that,

$$y_i f(x_i) = y_i (\langle w, x_i \rangle + b) > 0, \quad \forall i = 1, \dots, n,$$

\Rightarrow There exists a **hyperplane** $\{x \in \mathbb{R}^d \mid \langle w, x \rangle + b = 0\}$ which each separates the sets $X_+ = \{x_i \in D \mid y_i = 1\}$ and $X_- = \{x_i \in D \mid y_i = -1\}$.



A training sample of a two-class problem in \mathbb{R}^2 . The two classes are linearly separable and three different decision hyperplanes are shown which separate the two classes. (Image by Prof. Hein) 9 / 36

No distinction between the original input space $\mathcal{X} = \mathbb{R}^d$ and a possibly larger **feature space**, where we use basis functions/feature maps ϕ_i

$$\mathbf{x} \in \mathbb{R}^d \longrightarrow (\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})),$$

to the feature space \mathbb{R}^m .

Functions are linear in the parameters but not necessarily linear in the input space!

No distinction between the original input space $\mathcal{X} = \mathbb{R}^d$ and a possibly larger **feature space**, where we use basis functions/feature maps ϕ_i

$$\mathbf{x} \in \mathbb{R}^d \longrightarrow (\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})),$$

to the feature space \mathbb{R}^m .

Functions are linear in the parameters but not necessarily linear in the input space!

Definition

Let $g : \mathcal{X} \rightarrow \mathbb{R}$ be a function and $\hat{y}(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$ be the resulting classifier with output in $\mathcal{Y} = \{-1, 1\}$, then we call the set

$$\{\mathbf{x} \in \mathcal{X} \mid f(\mathbf{x}) = 0\},$$

the **decision boundary** of the classifier \hat{y} .

Three linear methods: $\hat{y}(x) = \text{sign}(f(x)) = \text{sign}(\langle w, \Phi(x) \rangle)$.

- **Linear Discriminant Analysis:**

- Loss: Squared loss, $L(y, f(x)) = (y - f(x))^2$
- Regularization: none

- **Logistic Regression:**

- Loss: Logistic loss, $L(y, f(x)) = \log(1 + \exp(-y f(x)))$
- Regularization: usually none, but there exist regularized versions.

- **Support Vector Machines** (Lecture 14).

- Loss: hinge loss, $L(y, f(x)) = \max(0, 1 - y f(x))$
- Regularization: L2-regularization, i.e., $\Omega(w) = \|w\|_2^2$

All three methods construct a **linear** classifier but all three have different **objectives**.

Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

Soft-margin SVM

Summary

The linear **support vector machine** (SVM) can be motivated from different perspectives.

Geometric Perspective: Maximum margin hyperplane

Unique hyperplane which correctly classifies the data and has maximal distance/margin to the training data.

- **hard margin** case: linearly separable data.
- **soft margin** case: all kind of data allowed.

Maximum margin hyperplane: a hyperplane which correctly classifies the data and has maximum distance/margin to the data.

Definition

A **maximum margin hyperplane** (w, b) for a **linearly separable** set of training data $(x_i, y_i)_{i=1}^n$ is defined as

$$\max_{w \in \mathbb{R}^d, b \in \mathbb{R}} \min\{\|x - x_i\| \mid \langle w, x \rangle + b = 0, x \in \mathbb{R}^d, i = 1, \dots, n\},$$

where we optimize over all (w, b) such that $y_i (\langle w, x_i \rangle + b) > 0$.

- Linear classifier is determined by the weight vector w and the offset b .

$$\hat{y}(x) = \text{sign}(\langle w, x \rangle + b).$$

- classifier and the decision boundary are not unique. For $\gamma > 0$, $\tilde{w} = \gamma w$ and $\tilde{b} = \gamma b$ gives the same classifier.

Definition (Geometrical margin)

For a hyperplane $\{x \mid \langle w, x \rangle + b = 0\}$, the **geometrical margin** of a point (x, y) is:

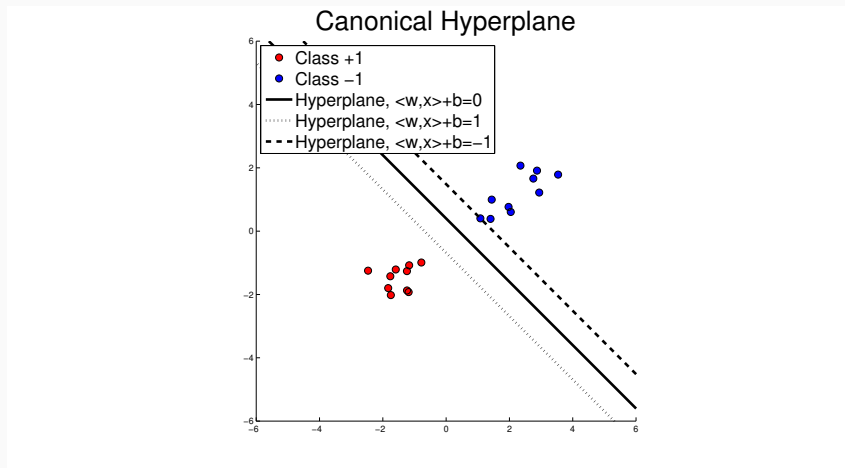
$$\rho_{w,b}(x, y) = y(\langle w, x \rangle + b) / \|w\|.$$

Definition (Canonical hyperplane)

The pair $(w, b) \in \mathbb{R}^d \times \mathbb{R}$ is said to be in **canonical** form with respect to $x_1, \dots, x_n \in \mathbb{R}^d$, if it is scaled such that

$$\min_{i=1, \dots, n} |\langle w, x_i \rangle + b| = 1,$$

which implies that the point closest to the hyperplane $h = \{x \mid \langle w, x \rangle + b = 0\}$ has distance $\rho = \frac{1}{\|w\|}$.



Canonical hyperplane for a set of training points $(x_i)_{i=1}^n$.
(Image by Prof. Hein)

Formulation:

$$\begin{aligned} & \max_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{\|w\|} \\ & \text{subject to: } y_i(\langle w, x_i \rangle + b) \geq 1, \quad \forall i = 1, \dots, n \end{aligned}$$

Second equivalent formulation:

$$\begin{aligned} & \min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|w\|^2 \\ & \text{subject to: } y_i(\langle w, x_i \rangle + b) \geq 1, \quad \forall i = 1, \dots, n \end{aligned}$$

Observation: convex optimization problem – quadratic program

Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

Soft-margin SVM

Summary

Lagrange function: Let $w \in \mathbb{R}^d$ and $\alpha \in \mathbb{R}^n$

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i \left[1 - y_i (\langle w, x_i \rangle + b) \right],$$

where $\alpha_i \geq 0$, $\forall i = 1, \dots, n$, are the **Lagrange multipliers**.

Dual Lagrange function:

$$q(\alpha) = \inf_{w \in \mathbb{R}^d, b \in \mathbb{R}} L(w, b, \alpha).$$

Observations:

- L is convex!
- Slater condition fulfilled if data is linearly separable \Rightarrow strong duality
- We can solve primal problem via the dual problem.

Derivatives:

$$\nabla_{\mathbf{w}} L(\mathbf{w}, b, \alpha) = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, \quad \frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = - \sum_{i=1}^n \alpha_i y_i.$$

Conditions for global minimum:

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, \quad \sum_{i=1}^n \alpha_i y_i = 0.$$

Plugging these expressions into $L(\mathbf{w}, b, \alpha)$ we get **the dual Lagrangian**:

$$q(\alpha) = -\frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i,$$

where $\alpha_i \geq 0, \quad \forall i = 1, \dots, n.$

- ▶ develop equation further:

$$\begin{aligned} L_p &\equiv \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i (y_i (w^T \cdot x_i + b) - 1) \\ &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i w^T x_i - \sum_{i=1}^N \alpha_i y_i b + \sum_{i=1}^N \alpha_i \end{aligned}$$

use following constraint:

$$\sum_{i=1}^N \alpha_i y_i = 0$$
$$L_p \equiv \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i w^T x_i + \sum_{i=1}^N \alpha_i$$

► so far:

$$L_p \equiv \frac{1}{2} ||w||^2 - \sum_{i=1}^N \alpha_i y_i w^T x_i + \sum_{i=1}^N \alpha_i$$

use following constraint:

$$w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$L_p \equiv \frac{1}{2} ||w||^2 - \sum_{i=1}^N \alpha_i y_i \sum_{j=1}^N \alpha_j y_j x_j^T x_i + \sum_{i=1}^N \alpha_i$$

$$L_p \equiv \frac{1}{2} ||w||^2 - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_j^T x_i) + \sum_{i=1}^N \alpha_i$$

$$\frac{1}{2} \|w\|^2 = \frac{1}{2} w^T w$$

use following constraint:

$$w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\frac{1}{2} w^T w = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_j^T x_i)$$

which results in the so called (Wolfe dual)

$$L_D \equiv \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (x_j^T x_i)$$

Dual problem:

$$\max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle ,$$

subject to: $\alpha_i \geq 0, \quad i = 1, \dots, n,$

$$\sum_{i=1}^n y_i \alpha_i = 0.$$

Observations:

- The dual problem is solved in practice using SMO (Sequential minimal optimization) method.
- Complexity is in the worst case cubic in n but often much faster.

Karush-Kuhn-Tucker (KKT) conditions: The most important one is the complementary slackness condition:

$$\begin{aligned} & \left[1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \right] = 0 \quad \text{if} \quad \alpha_i > 0 \\ \text{and} \quad & \alpha_i = 0 \quad \text{if} \quad \left[1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \right] < 0. \end{aligned}$$

or more compactly

$$\alpha_i \left[1 - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \right] = 0.$$

The offset b can thus be determined by averaging the value $y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle$ over all points with $\alpha_i > 0$ (as these are supposed to have margin = 1):

$$b = \frac{1}{\sum_{i=1}^n \mathbb{1}_{\alpha_i > 0}} \sum_{i=1}^n \mathbb{1}_{\alpha_i > 0} \left(y_i - \sum_{j=1}^n \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right).$$

Final weight vector:

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i.$$

Only the closest points to the decision boundary contribute to solution, i.e.,

$$\alpha_i > 0 \quad \Leftrightarrow \quad \left[1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \right] = 0,$$

The points \mathbf{x}_i for which $\alpha_i > 0$ are called **support vectors**. The area between the two supporting hyperplanes $\{\mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 1\}$ and $\{\mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = -1\}$ is called the **margin**.

Observations:

1. The weight vector of the support vector machine is typically **sparse** in terms of α .
2. Modifications of the training points matter only if they move into the margin.

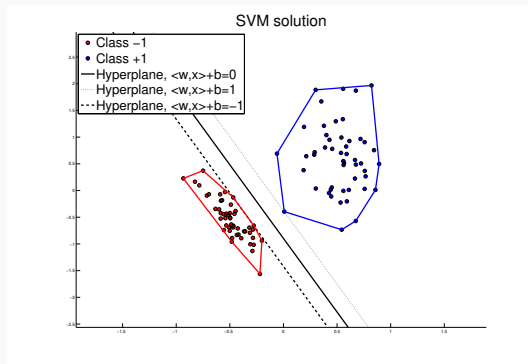
Equivalent reformulation of the dual problem:

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} & \left\| \sum_{i=1, y_i=1}^n \alpha_i x_i - \sum_{j=1, y_j=-1}^n \alpha_j x_j \right\|^2, \\ \text{subject to: } & \alpha_i \geq 0, \quad i = 1, \dots, n, \\ & \sum_{i=1, y_i=1}^n \alpha_i = \sum_{j=1, y_j=-1}^n \alpha_j = 1. \end{aligned}$$

Observations:

- It can be shown that the above problem maximizes the distance between the convex hulls of the positive and negative class.
- The maximum margin hyperplane is the one bisecting the shortest line orthogonally connecting both hulls.

Example: linearly separable case



A linearly separable problem. The hard margin solution of the SVM is shown together with the convex hulls of the positive and negative class. The points on the margin, that is $\langle w, x \rangle + b = \pm 1$, are called **support vectors**. (Image by Prof. Hein)

Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

Soft-margin SVM

Summary

Problems of the hard margin case:

- in general, data is not linearly separable,
- the **hard margin** case is often too strict since it is sensitive to outliers.

Relaxation of the constraints:

$$y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i$$

where $\xi_i \geq 0$ are the **slack variables**.

Primal problem of the soft-margin case:

$$\begin{aligned} \min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \quad & \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ \text{subject to: } & y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad \forall i = 1, \dots, n, \\ & \xi_i \geq 0, \quad \forall i = 1, \dots, n \end{aligned}$$

At the optimum: (note that $\xi_i \geq 0$)

$$\xi_i = \max \left(0, 1 - y_i (\langle w, x_i \rangle + b) \right),$$

where we recall that $\max \left(0, 1 - y_i f(x_i) \right)$ is the **hinge loss**.

Soft Margin SVM is RERM with Hinge loss and L_2 -regularization:

$$\min_{w \in \mathbb{R}^d, b \in \mathbb{R}} C \frac{1}{n} \sum_{i=1}^n \max \left(0, 1 - y_i (\langle w, x_i \rangle + b) \right) + \|w\|^2,$$

Error parameter C is the inverse of the regularization parameter $\lambda = \frac{1}{C}$.

Lagrangian of the soft margin problem:

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i \left[1 - \xi_i - y_i (\langle w, x_i \rangle + b) \right] - \sum_{i=1}^n \beta_i \xi_i$$

where $\alpha_i \geq 0$, $i = 1, \dots, n$ and $\beta_i \geq 0$, $i = 1, \dots, n$.

Conditions for a stationary point: ($\mathbf{1}$ is an n -dimensional vector of ones)

$$w = \sum_{i=1}^n \alpha_i y_i x_i, \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \beta = \frac{C}{n} \mathbf{1} - \alpha.$$

The last equation can be used to get rid of β . Due to the positivity of β we get the new constraint for α as

$$0 \leq \alpha_i \leq \frac{C}{n}, \quad i = 1, \dots, n.$$

Dual Lagrangian of the soft margin problem:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle, \\ \text{subject to: } \quad & 0 \leq \alpha_i \leq \frac{C}{n}, \quad i = 1, \dots, n, \quad \sum_{i=1}^n y_i \alpha_i = 0. \end{aligned}$$

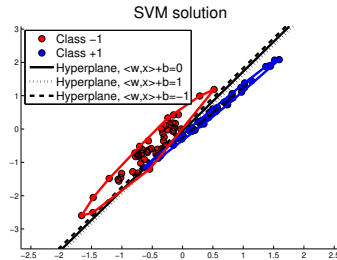
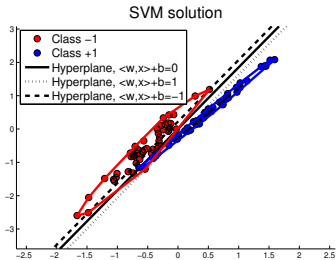
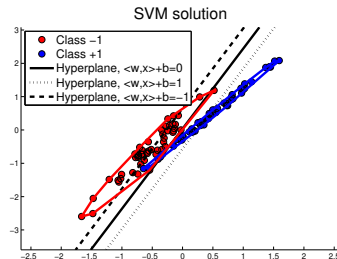
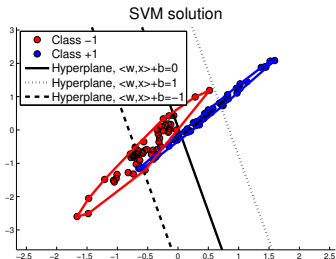
Complementary slackness conditions (part of KKT conditions) of the original problem:

$$\alpha_i \left[1 - \xi_i - y_i (\langle w, x_i \rangle + b) \right] = 0 \quad \text{and} \quad \beta_j \xi_j = 0, \quad \text{for } i, j = 1, \dots, n.$$

Three classes of points:

- $\alpha_i = 0$, outside the margin and all correctly classified.
- $0 < \alpha_i < \frac{C}{n}$, lie exactly on the margin and are all correctly classified.
- $\alpha_i = \frac{C}{n}$, inside the margin and may be misclassified.

Comparison of different C



Bibliography

Recap on linear classification

Support Vector Machines

Dual formulation

Soft-margin SVM

Summary

- Linear SVMs find the hyperplane that maximizes the margin between two classes in linearly separable datasets. A solution is found using the dual optimization problem.
- The resulting classifier (hyperplane) is computed only using the support vectors, i.e., those datapoints that lie exactly at the margin.
- Thus, the SVM classifier only varies across datasets if the support vectors change. This is due to the robust Hinge loss.
- For nonlinearly separable datasets, we relax the formulation and allow a subset of observations to lie inside the margin. The parameter C controls the proportion of observations that can lie inside the margin.
- We can generalize SVM to non-linear problems using specific basis functions that result from a similarity function over pairs of data points, known as **kernels** (next lecture!).