

12. Support Vector Machines I

Prof. Dan Olteanu

DaST
Data • (Systems+Theory)

Nov 6, 2020



University of
Zurich

<https://lms.uzh.ch/url/RepositoryEntry/16830890400>

<https://uzh.zoom.us/j/96690150974?pwd=cnZmMTduWUtCeWoxYW85Z3RMFnpTZz09>

Support Vector Machines (SVM)

SVM is a popular **discriminative model for classification**

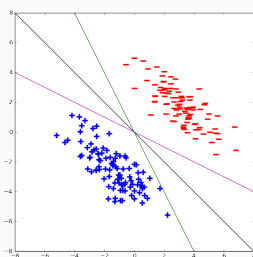
- No natural probabilistic interpretation (see one in Murphy Section 14.5.5)

SVM requires the introduction of several new concepts

- **Maximum Margin Principle**
- **Hinge Loss** optimisation
- **Primal vs Dual** Formulation
- **Kernel Methods** for non-linear classification

1

Binary Classification



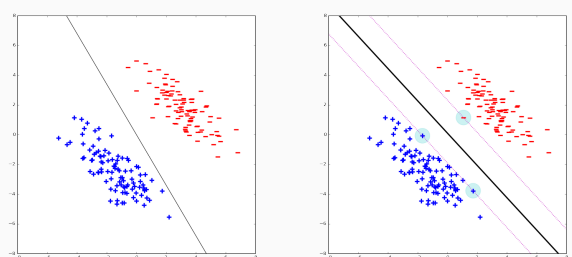
Goal: Find a linear separator

Data is **linearly separable** if there exists a linear separator that classifies all points correctly

Which separator should be picked?

2

Maximum Margin Principle



Margin for a data point = its distance to the separating hyperplane

Maximum margin principle: Pick the separating boundary that maximises the smallest margin (the least distance between data and boundary)

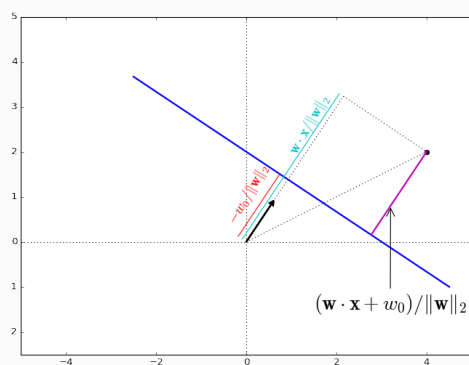
That is: Maximise the distance of the closest point from the decision boundary

Points that are closest to the decision boundary are called **support vectors**

3

Geometry Review: Distance from Point to Hyperplane

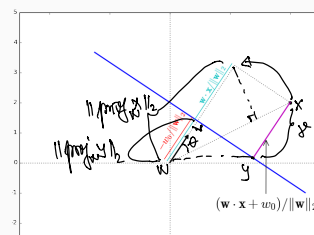
Given a hyperplane $H \equiv \mathbf{w} \cdot \mathbf{a} + w_0 = 0$ and a point $\mathbf{x} \in \mathbb{R}^D$, how far is \mathbf{x} from H ?



4

Geometry Review: Working Out the Distance

Given a hyperplane $H \equiv \mathbf{w} \cdot \mathbf{a} + w_0 = 0$ and a point $\mathbf{x} \in \mathbb{R}^D$, how far is \mathbf{x} from H ?



$$\begin{aligned} \mathbf{w} \cdot \mathbf{x} &= \|\mathbf{w}\|_2 \|\text{proj}_{\mathbf{w}} \mathbf{x}\|_2 \\ \|\text{proj}_{\mathbf{w}} \mathbf{x}\|_2 &= \frac{\mathbf{w} \cdot \mathbf{x}}{\|\mathbf{w}\|_2} \\ y &= \|\text{proj}_{\mathbf{w}} \mathbf{x}\|_2 - \|\text{proj}_{\mathbf{w}} \mathbf{y}\|_2 \\ \mathbf{w} \cdot \mathbf{y} &= \|\mathbf{w}\|_2 \|\text{proj}_{\mathbf{w}} \mathbf{y}\|_2 \\ \mathbf{w} \cdot \mathbf{y} + w_0 &= 0 \Leftrightarrow \mathbf{w} \cdot \mathbf{y} = -w_0 \\ \|\text{proj}_{\mathbf{w}} \mathbf{y}\|_2 &= -\frac{w_0}{\|\mathbf{w}\|_2} \end{aligned}$$

$$\begin{aligned} y &= \frac{\mathbf{w} \cdot \mathbf{x}}{\|\mathbf{w}\|_2} - \left(-\frac{w_0}{\|\mathbf{w}\|_2} \right) \\ &= \frac{\mathbf{w} \cdot \mathbf{x} + w_0}{\|\mathbf{w}\|_2} \quad \checkmark \end{aligned}$$

5

Distance of Point to Hyperplane

- Consider the hyperplane: $H \equiv \mathbf{w} \cdot \mathbf{x} + w_0 = 0$

- All points on the **positive** halfspace created by H satisfy

$$\mathbf{w} \cdot \mathbf{x} + w_0 > 0$$

All points on the **negative** halfspace created by H satisfy

$$\mathbf{w} \cdot \mathbf{x} + w_0 < 0$$

- We label $y = +1$ ($y = -1$) the data points in the **positive** (**negative**) halfspace

- The distance of point \mathbf{x} to H is given by:

$$\frac{|\mathbf{w} \cdot \mathbf{x} + w_0|}{\|\mathbf{w}\|_2}$$

6

Alternative Formulation as Optimisation Problem

Find distance $\|\mathbf{x} - \mathbf{x}^*\|_2$ between point \mathbf{x}^* and the hyperplane $\mathbf{w} \cdot \mathbf{x} + w_0 = 0$

The point \mathbf{x} on hyperplane that is closest to \mathbf{x}^* gives the distance

Equivalently, we seek for \mathbf{x} that optimises the following problem:

$$\begin{aligned} \text{minimise: } & \|\mathbf{x} - \mathbf{x}^*\|_2^2 \\ \text{subject to: } & \mathbf{w} \cdot \mathbf{x} + w_0 = 0 \end{aligned}$$

Lagrangian:

$$\begin{aligned} \Lambda(\mathbf{x}, \lambda) &= \|\mathbf{x} - \mathbf{x}^*\|_2^2 - 2\lambda(\mathbf{w} \cdot \mathbf{x} + w_0) \\ &= \|\mathbf{x}\|_2^2 - 2(\mathbf{x}^* + \lambda\mathbf{w}) \cdot \mathbf{x} - 2\lambda w_0 + \|\mathbf{x}^*\|_2^2 \end{aligned}$$

Find critical point of Λ by setting its gradient to 0:

$$\nabla_{\mathbf{x}} \Lambda(\mathbf{x}, \lambda) = 2\mathbf{x} - 2\mathbf{x}^* - 2\lambda\mathbf{w} = 0 \Rightarrow \mathbf{x} = \mathbf{x}^* + \lambda\mathbf{w}$$

7

Alternative Formulation as Optimisation Problem – Continued

From previous slide: We obtained the critical point $\mathbf{x} = \mathbf{x}^* + \lambda\mathbf{w}$

We next obtain an expression for λ by substituting \mathbf{x} into the hyperplane equation:

$$\mathbf{w} \cdot (\mathbf{x}^* + \lambda\mathbf{w}) + w_0 = 0 \Rightarrow \lambda = -\frac{\mathbf{w} \cdot \mathbf{x}^* + w_0}{\mathbf{w} \cdot \mathbf{w}} = -\frac{\mathbf{w} \cdot \mathbf{x}^* + w_0}{\|\mathbf{w}\|_2^2}$$

Finally, the distance between \mathbf{x}^* and \mathbf{x} becomes:

$$\|\mathbf{x} - \mathbf{x}^*\|_2 = \|\lambda\mathbf{w}\|_2 = |\lambda| \|\mathbf{w}\|_2 = \frac{|\mathbf{w} \cdot \mathbf{x}^* + w_0|}{\|\mathbf{w}\|_2}$$

8

The Linearly Separable Case

Assume that the dataset $\mathcal{D} = ((\mathbf{x}_i, y_i), \dots, (\mathbf{x}_N, y_N))$ is linearly separable

Recall

$$\mathbf{w} \cdot \mathbf{x}_i + w_0 > 0 \text{ and } y_i = +1 \quad \text{or} \quad \mathbf{w} \cdot \mathbf{x}_i + w_0 < 0 \text{ and } y_i = -1$$

More compactly,

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) > 0$$

Since \mathcal{D} is finite, we can always find $\epsilon > 0$ such that

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq \epsilon$$

Alternative formulation (divide both sides by ϵ):

$$y_i \left(\underbrace{\frac{\mathbf{w}}{\epsilon}}_{\text{new } \mathbf{w}} \cdot \mathbf{x}_i + \underbrace{\frac{w_0}{\epsilon}}_{\text{new } w_0} \right) \geq 1$$

9

Margin Maximisation

Recall the margin of data point \mathbf{x}_i to hyperplane:

$$\frac{|\mathbf{w} \cdot \mathbf{x}_i + w_0|}{\|\mathbf{w}\|_2} = \frac{y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0)}{\|\mathbf{w}\|_2} \geq \frac{1}{\|\mathbf{w}\|_2}$$

Maximum margin principle: pick the boundary that maximises the margin

$$\begin{aligned} \text{maximise: } & \frac{1}{\|\mathbf{w}\|_2} \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \text{ for } 1 \leq i \leq N \end{aligned}$$

Equivalently, we **minimise the squared ℓ_2 norm of \mathbf{w} subject to the constraints**

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \text{ for } 1 \leq i \leq N \end{aligned}$$

The constant factor $\frac{1}{2}$ is added without loss of generality.

10

Margin Maximisation is a Convex Quadratic Optimisation Problem

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1, \text{ for } 1 \leq i \leq N \end{aligned}$$

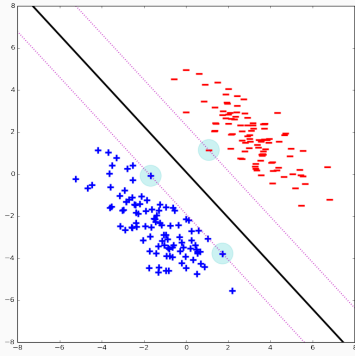
- The objective is a convex quadratic function
- The constraints are convex linear functions
- The feasible set as defined by the constraints is convex as well

\Rightarrow Our optimisation problem is convex quadratic

\Rightarrow Solvable using generic convex optimisation methods

11

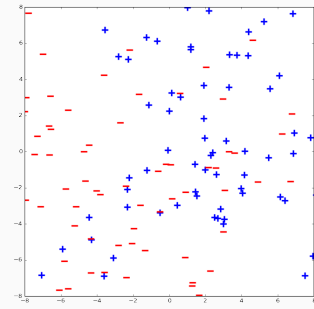
SVM in the Linearly Separable Case



We find an SVM classifier with no classification error on the training set

12

The Non-Separable Case



- The quadratic program from the previous slides has **no feasible** solution
- What would be a good relaxation of the previous optimisation problem?

13

Relaxation of the SVM Formulation

Original SVM formulation for linearly-separable data:

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 \quad \text{for } 1 \leq i \leq N \end{aligned}$$

Relaxations:

1. Find a separator that makes the least mistakes on the training error
 - Minimising the number of misclassifications is NP-hard :(
2. Proxy for least mistakes: Satisfy as many of the N constraints as possible
3. Alternative: Allow all constraints to be satisfied **with some slack**

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \zeta_i \quad \text{for } 1 \leq i \leq N \\ & \zeta_i \geq 0 \quad \text{for } 1 \leq i \leq N \end{aligned}$$

14

SVM Formulation in the Non-Separable Case

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \zeta_i \quad \text{for } 1 \leq i \leq N \\ & \zeta_i \geq 0 \quad \text{for } 1 \leq i \leq N \end{aligned}$$

- Feasible solution always exists thanks to slack variables ζ_i
 - ζ_i can be as large (constraints can be violated as much) as necessary
- The constraints $\zeta_i \geq 0$ are important!
 - Assume $\zeta_i \ll 0$ and \mathbf{x}_i correctly classified by a huge margin
 - This would compensate the penalty incurred on misclassified points
 - $\zeta_i \geq 0$ ensures no bonus for correct classification by a huge margin

15

SVM Formulation: Loss Function

$$\begin{aligned} \text{minimise: } & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i \\ \text{subject to: } & y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \zeta_i \quad \text{for } 1 \leq i \leq N \\ & \zeta_i \geq 0 \quad \text{for } 1 \leq i \leq N \end{aligned}$$

Optimal solution must satisfy

- either $\zeta_i = 0$: This is ideal, no slack is needed to classify \mathbf{x}_i
- or $\zeta_i = 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 0$: Minimum ζ_i that satisfies the constraint

Equivalent formulation of the optimal solution for ζ_i :

$$\zeta_i = \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0)) \stackrel{\text{def}}{=} \ell_{\text{hinge}}(\mathbf{w}, w_0; \mathbf{x}_i, y_i)$$

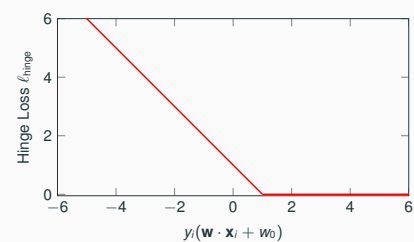
Our SVM formulation becomes equivalent to minimising the objective function:

$$\mathcal{L}_{\text{SVM}}(\mathbf{w}, w_0 | \mathbf{X}, \mathbf{y}) = \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2}_{\ell_2 \text{ regulariser}} + C \underbrace{\sum_{i=1}^N \ell_{\text{hinge}}(\mathbf{w}, w_0; \mathbf{x}_i, y_i)}_{\text{hinge loss}}$$

16

Hinge Loss Function

$$\ell_{\text{hinge}}(\mathbf{w}, w_0; \mathbf{x}_i, y_i) = \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0))$$



17

Logistic Loss Function

$y_i \in \{0, 1\}$ for logistic regression and $y_i \in \{-1, +1\}$ for SVM

We resolve the mismatch by using $z_i = 2y_i - 1$ to map from $\{0, 1\}$ to $\{-1, +1\}$.

$$\text{NLL}(y_i | \mathbf{w}, \mathbf{x}_i) = - \left(y_i \log \left(\frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}_i}} \right) + (1 - y_i) \log \left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}_i}} \right) \right)$$

If $y_i = 1$, then $z_i = 1$ and:

$$\text{NLL}(y_i = 1 | \mathbf{w}, \mathbf{x}_i) = - \log \left(\frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}_i}} \right) = \log \left(1 + e^{-\mathbf{w} \cdot \mathbf{x}_i} \right) = \log \left(1 + e^{-z_i \mathbf{w} \cdot \mathbf{x}_i} \right)$$

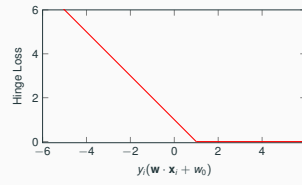
If $y_i = 0$, then $z_i = -1$ and obtain the **logistic loss function**:

$$\text{NLL}(y_i = 0 | \mathbf{w}, \mathbf{x}_i) = - \log \left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}_i}} \right) = \log \left(1 + e^{\mathbf{w} \cdot \mathbf{x}_i} \right) = \log \left(1 + e^{-z_i \mathbf{w} \cdot \mathbf{x}_i} \right)$$

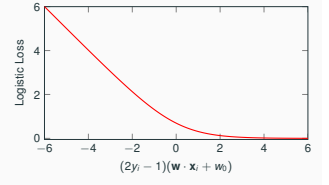
By replacing $z_i = (2y_i - 1)$, we obtain: $\text{NLL}(y_i | \mathbf{w}, \mathbf{x}_i) = \log \left(1 + e^{-(2y_i - 1) \mathbf{w} \cdot \mathbf{x}_i} \right)$

18

Loss Functions: Hinge vs Logistic



$$\max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + w_0))$$



$$\log \left(1 + e^{-(2y_i - 1)(\mathbf{w} \cdot \mathbf{x}_i + w_0)} \right)$$

19