

# Statistical Machine Learning

## Lecture 8: Support Vector Machines

2022-23

# Linear Classification

1. Suppose there exists a function  $f : \mathcal{X} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$  such that for all  $\mathbf{x} \in \mathcal{X}$  with  $f(\mathbf{x}) > 0$  is assigned to class 1 and otherwise to class -1. We consider the case where  $f$  is a linear function of  $\mathbf{x} \in \mathcal{X}$ , so it can be expressed as:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0,$$

where  $\mathbf{W} = (w_0, \mathbf{w}) \in \mathbb{R} \times \mathbb{R}^d$  are the parameters that control the function and the decision rule  $\text{sgn}(f(\mathbf{x}))$ .

2. We usually assume that  $\|\mathbf{w}\|^2 = 1$ .
3. The parameters  $\mathbf{W}$  must be learned from the data.
4. The data set is expected to be  $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$ , where  $(\mathbf{x}_n, t_n) \in \mathcal{X} \times \{-1, +1\}$  for all  $n = 1, \dots, N$ .

# Perceptron Algorithm and Dual Representation

1. Let  $\mathbf{y} \in \{1\} \times \mathcal{X}$  such that  $\mathbf{y}_n = (1, \mathbf{x}_n^T)^T$  for all  $(\mathbf{x}_n, t_n) \in \mathcal{D}$ . The perceptron algorithm is:

$$\mathbf{W}_{l+1} = \mathbf{W}_l + \Theta(-t_{l+1} \mathbf{W}_l^T \mathbf{y}_{l+1}) t_{l+1} \mathbf{y}_{l+1}.$$

2. Assume that  $\mathbf{W}_0 = \mathbf{0}$ . Thus

$$\mathbf{W}_l = \sum_{j=1}^l \alpha_j t_j \mathbf{y}_j,$$

where  $\alpha_j$  is the number of times  $\mathbf{y}_j$  has been misclassified.

# Perceptron Algorithm and Dual Representation

1. The decision function becomes:

$$\text{sgn}(f(\mathbf{x})) = \text{sgn} \left( \sum_{j=1}^l \alpha_j t_j \mathbf{y}_j^T \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} \right),$$

where, instead of depending on  $\mathbf{W}$  it depends on  $\{\alpha_j\}$ . This is the *dual form* of the function  $f$ .

2. Observe that the points close to the boundary, and therefore harder to learn, are expected to have larger  $\alpha$ .

# Regression Problem

1. Let  $\mathbf{y} \in \{1\} \times \mathcal{X}$  such that  $\mathbf{y}_n = (1, \mathbf{x}_n^T)^T$  for all  $(\mathbf{x}_n, t_n) \in \mathcal{D} \subset (\mathcal{X} \times \mathbb{R})^N$ . The model and the least-square loss function are:

$$f(\mathbf{x}) = \mathbf{W}^T \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

$$\mathcal{E}(\mathbf{W}) = \sum_{n=1}^N \left[ t_n - \mathbf{W}^T \mathbf{y}_n \right]^2,$$

where  $\mathbf{W} = (w_0, \mathbf{w}^T)^T \in \mathbb{R} \times \mathbb{R}^d$ .

# Regression Problem: Ridge Regression

1. Let  $\mathbf{y} \in \{1\} \times \mathcal{X}$  such that  $\mathbf{y}_n = (1, \mathbf{x}_n^T)^T$  for all  $(\mathbf{x}_n, t_n) \in \mathcal{D} \subset (\mathcal{X} \times \mathbb{R})^N$ . The model and the least-square loss function are:

$$f(\mathbf{x}) = \mathbf{W}^T \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

$$\mathcal{E}_{RR}(\mathbf{W}, \lambda) = \lambda \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \left[ t_n - \mathbf{W}^T \mathbf{y}_n \right]^2,$$

where  $\mathbf{W} = (w_0, \mathbf{w}^T)^T \in \mathbb{R} \times \mathbb{R}^d$ .

2.  $\lambda$  is a parameter that controls a trade-off between the low norm of  $\mathbf{w}$  and the low quadratic error.

# Regression Problem: Ridge Regression

1. The equation to be solved is

$$\nabla_{\mathbf{w}} \mathcal{E}_{RR} = \begin{pmatrix} \partial_{w_0} \\ \nabla_{\mathbf{w}} \end{pmatrix} \mathcal{E}_{RR} = \mathbf{0},$$

2. That admits the solution:

$$w_0^* = \frac{1}{N} \sum_{n=1}^N \left( t_n - \mathbf{x}_n^T \mathbf{w}^* \right) = \bar{t} - \bar{\mathbf{x}}^T \mathbf{w}^*$$

$$\mathbf{w}^* = \sum_{n=1}^N \frac{t_n - \mathbf{x}_n^T \mathbf{w}^* - w_0^*}{\lambda} \mathbf{x}_n = \sum_{n=1}^N \frac{(t_n - \bar{t}) - (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{w}^*}{\lambda} (\mathbf{x}_n - \bar{\mathbf{x}})$$

$$\alpha_n \equiv \frac{(t_n - \bar{t}) - (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{w}^*}{\lambda}$$

$$\mathbf{w}^* = \sum_{n=1}^N \alpha_n (\mathbf{x}_n - \bar{\mathbf{x}}).$$

# Ridge Regression: Dual problem

1. Let us define the matrix  $\mathbf{X} \in \mathcal{X}^N$  such that

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \quad \overline{\mathbf{X}} = \begin{pmatrix} \overline{\mathbf{x}}^T \\ \overline{\mathbf{x}}^T \\ \vdots \\ \overline{\mathbf{x}}^T \end{pmatrix} \left. \vphantom{\begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}} \right\} N \text{ times.}$$

2. Observe that  $\mathbf{w}^* = (\mathbf{X} - \overline{\mathbf{X}})^T \boldsymbol{\alpha}$ .



# Ridge Regression: Dual problem

1. Then  $\mathbf{w}^* = (\mathbf{X} - \bar{\mathbf{X}})^T \boldsymbol{\alpha}$ , where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)^T \in \mathbb{R}^N$ .  
Let  $\boldsymbol{\tau} \in \mathbb{R}^N$  such that  $[\boldsymbol{\tau}]_n = t_n - \bar{t}$ . Therefore:

$$w_0^* = \bar{t} - \bar{\mathbf{x}}^T (\mathbf{X} - \bar{\mathbf{X}})^T \boldsymbol{\alpha}$$

$$\begin{aligned} \mathcal{E}_{RR}(\boldsymbol{\alpha}, \lambda) &= \lambda \boldsymbol{\alpha}^T (\mathbf{X} - \bar{\mathbf{X}}) (\mathbf{X} - \bar{\mathbf{X}})^T \boldsymbol{\alpha} + \sum_{n=1}^N \left[ t_n - \bar{t} - \boldsymbol{\alpha}^T (\mathbf{X} - \bar{\mathbf{X}})^T \right]^2 \\ &= \lambda \boldsymbol{\alpha}^T (\mathbf{X} - \bar{\mathbf{X}}) (\mathbf{X} - \bar{\mathbf{X}})^T \boldsymbol{\alpha} + \sum_{n=1}^N \left[ (t_n - \bar{t}) - \boldsymbol{\alpha}^T (\mathbf{X} - \bar{\mathbf{X}})^T \right]^2 \\ &= \lambda \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} + \boldsymbol{\tau}^T \boldsymbol{\tau} - 2 \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\tau} + \boldsymbol{\alpha}^T \mathbf{G} \mathbf{G} \boldsymbol{\alpha}, \end{aligned}$$

$$\text{where } \mathbf{G}^T = \mathbf{G} = (\mathbf{X} - \bar{\mathbf{X}}) (\mathbf{X} - \bar{\mathbf{X}})^T.$$

## Ridge Regression: Dual problem

1. The derivatives with respect to  $\alpha$  are

$$\nabla_{\alpha} \mathcal{E}_{RR} = 2\mathbf{G}(\lambda\alpha + \mathbf{G}\alpha - \boldsymbol{\tau}) = \mathbf{0}.$$

2. The solution is

$$\alpha = (\lambda I + \mathbf{G})^{-1} \boldsymbol{\tau}.$$

3. The predictive function becomes:

$$\begin{aligned} f(\mathbf{x}) &= \left( (\mathbf{X} - \overline{\mathbf{X}})^T (\lambda I + \mathbf{G})^{-1} \boldsymbol{\tau} \right)^T (\mathbf{x} - \overline{\mathbf{x}}) + \bar{t} \\ &= \boldsymbol{\tau}^T (\lambda I + \mathbf{G})^{-1} (\mathbf{X} - \overline{\mathbf{X}}) (\mathbf{x} - \overline{\mathbf{x}}) + \bar{t} \end{aligned}$$

# Ridge Regression: Dual problem

1. The predictive function becomes:

$$\begin{aligned} f(\mathbf{x}) &= \boldsymbol{\tau}^T (\lambda \mathbf{I} + \mathbf{G})^{-1} \begin{pmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}}) \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}}) \\ \vdots \\ (\mathbf{x}_N - \bar{\mathbf{x}})^T (\mathbf{x} - \bar{\mathbf{x}}) \end{pmatrix} + \bar{t} \\ &= \boldsymbol{\tau}^T (\lambda \mathbf{I} + \mathbf{G})^{-1} \mathbf{z} + \bar{t}. \end{aligned}$$

2. Observe that in the dual representation the data only appears through the Gram matrix  $\mathbf{G}$ , and that in the decision function it is only the inner product of the data points with the test point ( $\mathbf{x}$ ) that are needed.

# Optimisation Theory

1. Primal Optimisation Problem: Given the functions  $f, g_i, h_j : \Omega \rightarrow \mathbb{R}$ , with  $i = 1, \dots, K$ , and  $j = 1, \dots, P$ ,

$$\begin{array}{lll} \text{minimise} & f(\mathbf{w}), & \mathbf{w} \in \Omega, \\ \text{subject to} & g_i(\mathbf{w}) \leq 0, & i = 1, \dots, K, \\ & h_j(\mathbf{w}) = 0, & j = 1, \dots, P, \end{array}$$

where  $f(\mathbf{w})$  is called the *objective function*,  $g_i(\mathbf{w})$  are the *inequality constraints* and  $h_j(\mathbf{w})$  are the *equality constraints*.

2. The region defined as

$$\mathcal{R} \equiv \{\mathbf{w} \in \Omega : \mathbf{g}(\mathbf{w}) \leq \mathbf{0}_K, \mathbf{h}(\mathbf{w}) = \mathbf{0}_P\},$$

is known as the *feasible region*.

3. The solution  $\mathbf{w}^* \in \mathcal{R}$  such that  $f(\mathbf{w}^*) \leq f(\mathbf{w})$  for all  $\mathbf{w} \in \Omega$ .

# Lagrangian Theory

1. Given an optimisation problem with domain  $\Omega \subseteq \mathbb{R}^d$ ,

$$\begin{array}{lll} \text{minimise} & f(\mathbf{w}), & \mathbf{w} \in \Omega, \\ \text{subject to} & g_i(\mathbf{w}) \leq 0, & i = 1, \dots, K, \\ & h_j(\mathbf{w}) = 0, & j = 1, \dots, P, \end{array}$$

we define the Lagrangian function as

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{w}) + \boldsymbol{\alpha}^T \mathbf{g}(\mathbf{w}) + \boldsymbol{\beta}^T \mathbf{h}(\mathbf{w}),$$

with  $\boldsymbol{\alpha} \in \mathbb{R}^K$  and  $\boldsymbol{\beta} \in \mathbb{R}^P$  are known as the Lagrange Multipliers.

# Lagrangian Theory

1. The dual problem of the primal defined above is

$$\begin{array}{ll}\text{maximise} & \theta(\boldsymbol{\alpha}, \boldsymbol{\beta}), \\ \text{subject to} & \boldsymbol{\alpha} \geq 0,\end{array}$$

where

$$\theta(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \inf_{\boldsymbol{w} \in \Omega} \mathcal{L}(\boldsymbol{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}).$$

# Weak and Strong Duality

1. WD: Let  $\mathbf{w} \in \Omega$  be a feasible solution of the primal problem and  $(\boldsymbol{\alpha}, \boldsymbol{\beta}) \in \mathbb{R}^K \times \mathbb{R}^P$  a feasible solution of the dual problem. Then  $f(\mathbf{w}) \geq \theta(\boldsymbol{\alpha}, \boldsymbol{\beta})$ . The  $\Delta(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \equiv f(\mathbf{w}) - \theta(\boldsymbol{\alpha}, \boldsymbol{\beta})$  is known as the *Duality Gap*.
2. A set  $\Upsilon$  is *convex* if for all  $\alpha, \beta \in \Upsilon$  and  $t \in [0, 1]$  the convex combination  $\gamma = t\alpha + (1 - t)\beta$  is also an element of  $\Upsilon \ni \gamma$ .
3. SD: Given an optimisation problem in a convex domain  $\Omega \subseteq \mathbb{R}^d$ ,

$$\begin{array}{lll} \text{minimise} & f(\mathbf{w}), & \mathbf{w} \in \Omega, \\ \text{subject to} & g_i(\mathbf{w}) \leq 0, & i = 1, \dots, K, \\ & h_j(\mathbf{w}) = 0, & j = 1, \dots, P, \end{array}$$

where there exist  $\mathbf{A} \in \mathbb{R}^{K \times d}$ ,  $\mathbf{B} \in \mathbb{R}^{P \times d}$ ,  $\mathbf{a} \in \mathbb{R}^K$ , and  $\mathbf{b} \in \mathbb{R}^P$  such that  $\mathbf{g}(\mathbf{w}) = \mathbf{Aw} - \mathbf{a}$  and  $\mathbf{h}(\mathbf{w}) = \mathbf{Bw} - \mathbf{b}$ , then the duality gap is zero.

# Karush-Kuhn-Tucker Conditions

1. Given an optimisation problem in a convex domain  $\Omega \subseteq \mathbb{R}^d$ , minimise  $f(\mathbf{w})$ ,  $\mathbf{w} \in \Omega$ , subject to  $\mathbf{g}(\mathbf{w}) \leq \mathbf{0}_K$  and  $\mathbf{h}(\mathbf{w}) = \mathbf{0}_P$ , with  $f \in C^1$  (the derivative exists and it is continuous) and convex (i.e.  $f(t\mathbf{w}_1 + (1-t)\mathbf{w}_2) \leq tf(\mathbf{w}_1) + (1-t)f(\mathbf{w}_2)$ ) the necessary and sufficient conditions for a point  $\mathbf{w}^*$  to be the optimum is the existence of  $\alpha^*$  and  $\beta^*$  such that:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*, \alpha^*, \beta^*) = \mathbf{0}_d,$$

$$\nabla_{\beta} \mathcal{L}(\mathbf{w}^*, \alpha^*, \beta^*) = \mathbf{0}_P,$$

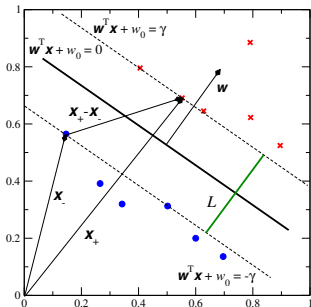
$$\alpha_i^* g_i(\mathbf{w}^*) = 0, \quad i = 1, \dots, K$$

$$g_i(\mathbf{w}^*) \leq 0, \quad i = 1, \dots, K$$

$$\alpha_i^* \geq 0, \quad i = 1, \dots, K.$$



# Maximising the Classification Margin



2. The red crosses satisfy  $w^T x_+ + w_0 > \gamma$  (classified with  $t = 1$ ) whereas the blue dots satisfy  $w^T x_- + w_0 < -\gamma$  (classified with  $t = -1$ ).
3. Observe that for all  $(x_n, t_n) \in \mathcal{D}$  we have that  $t_n(w^T x_n + w_0) - \gamma > 0$ .

# Maximising the Classification Margin

1. Observe that

$$\frac{\mathbf{w}^T}{\|\mathbf{w}\|_2}(\mathbf{x}_+ - \mathbf{x}_-) > \frac{2\gamma}{\|\mathbf{w}\|_2} = L,$$

where the left-hand side is the perpendicular distance between the margins.

2. To maximise the margin  $L$  we have to reduce

$$\|\mathbf{w}\|_2 = \sqrt{\sum_{\ell=1}^d w_{\ell}^2}.$$

# Maximising the Classification Margin

1. The loss function is then  $f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$  and the primal Lagrangian becomes:

$$\mathcal{L}_P(\mathbf{w}, w_0, \alpha') = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha'_n [\gamma - t_n(\mathbf{w}^T \mathbf{x}_n + w_0)]$$

where  $\{\alpha'_n\}$  are the Lagrange multiplier for the  $N$  constraints  $t_n(\mathbf{w}^T \mathbf{x}_n + w_0) - \gamma > 0$ .

2. By defining  $\alpha'_n = \frac{\alpha_n}{\gamma}$  and re-scaling in  $\gamma$   $w_0$  and  $\mathbf{x}_n$ , we can write

$$\mathcal{L}_P(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n \left[ 1 - t_n \left( \mathbf{w}^T \frac{\mathbf{x}_n}{\gamma} + \frac{w_0}{\gamma} \right) \right].$$

# KKT Conditions

1. The loss function  $f(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w}$  is convex. Observe that for all  $\mathbf{w}_1, \mathbf{w}_2 \in \Omega$  and  $s \in [0, 1]$

$$\begin{aligned}f(s\mathbf{w}_1 + (1-s)\mathbf{w}_2) &= \frac{1}{2}[s\mathbf{w}_1 + (1-s)\mathbf{w}_2]^T[s\mathbf{w}_1 + (1-s)\mathbf{w}_2] \\&= \frac{1}{2}\{s^2\mathbf{w}_1^T\mathbf{w}_1 + 2s(1-s)\mathbf{w}_1^T\mathbf{w}_2 + (1-s)^2\mathbf{w}_2^T\mathbf{w}_2\} \\&= sf(\mathbf{w}_1) + (1-s)f(\mathbf{w}_2) - \frac{1}{2}s(1-s)\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \\&\leq sf(\mathbf{w}_1) + (1-s)f(\mathbf{w}_2).\end{aligned}$$

# KKT Conditions

1. The first condition produces the following equations:

$$\nabla_{\mathbf{w}} \mathcal{L}_P = \mathbf{0}_d = \mathbf{w}^* - \sum_{n=1}^N \alpha_n^* t_n \mathbf{x}_n$$

$$\frac{\partial \mathcal{L}_P}{\partial w_0} = 0 = \sum_{n=1}^N \alpha_n^* t_n$$

2. The second condition does not apply because we do not have equation constraints.
3. The last three conditions are:

$$\alpha_n^* [1 - t_n (\mathbf{x}_n^T \mathbf{w}^* + w_0^*)] = 0, \quad n = 1, \dots, N$$

$$t_n (\mathbf{x}_n^T \mathbf{w}^* + w_0^*) \geq 1, \quad n = 1, \dots, N$$

$$\alpha_n^* \geq 0, \quad n = 1, \dots, N.$$

# Dual Lagrangian

1. The dual Lagrangian of the problem is:

$$\begin{aligned}\mathcal{L}_D(\alpha) &= \inf_{\mathbf{w} \in \Omega} \mathcal{L}_P(\mathbf{w}, \alpha) \\ &= -\frac{1}{2} \left( \sum_{n=1}^N \alpha_n t_n \mathbf{x}_n \right)^T \left( \sum_{n=1}^N \alpha_n t_n \mathbf{x}_n \right) + \sum_{n=1}^N \alpha_n \\ &= -\frac{1}{2} \sum_{n,m=1}^N \alpha_n \alpha_m t_n t_m \mathbf{x}_n^T \mathbf{x}_m + \sum_{n=1}^N \alpha_n, \\ \text{subject to } 0 &= \sum_{n=1}^N \alpha_n t_n.\end{aligned}$$

2. Through the KKT conditions the *active constraints*, i.e.  $t_\ell (\mathbf{x}_\ell^T \mathbf{w}^* + w_0^*) = 1$  are the only ones that have  $\alpha_\ell > 0$ . The  $\mathbf{x}_\ell$  are the *support vectors*.

3. We can define

$$\text{sv} = \{k : 1 \leq k \leq N \text{ and } t_k (\mathbf{x}_k^T \mathbf{w}^* + w_0^*) = 1\}.$$

## Example

1. Consider the data set

$$\mathcal{D} = \{([1, 1]^T, 1); ([-1, -1]^T, -1); ([1, 2]^T, 1)\}.$$

2. The Dual Lagrangian is:

$$\mathcal{L}_D = -\frac{1}{2} \{2\alpha_1^2 + 2\alpha_2^2 + 5\alpha_3^2 + 4\alpha_1\alpha_2 + 6\alpha_1\alpha_3 + 6\alpha_2\alpha_3\} + \\ + \alpha_1 + \alpha_2 + \alpha_3$$

subject to  $0 = \alpha_1 - \alpha_2 + \alpha_3$ .



## Example

1. The system of linear equations becomes:

$$2\alpha_1 + 2\alpha_2 + 3\alpha_3 = 1$$

$$2\alpha_1 + 2\alpha_2 + 3\alpha_3 = 1$$

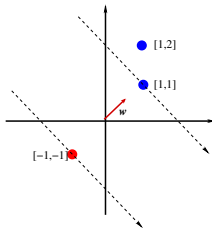
$$3\alpha_1 + 3\alpha_2 + 5\alpha_3 = 1$$

which, to be satisfied, requires that  $\alpha_3 = -1$ , which violates one of the KKT conditions. This implies that  $[1, 2]^T$  is not a support vector (does not produce an active constraint), therefore we made  $\alpha_3 = 0$ .

2. This implies that the new system of linear equations is  $2\alpha_1 + 2\alpha_2 = 1$  subject to  $\alpha_1 = \alpha_2$  ( $= \frac{1}{4}$ ).

# Results

1. The weight  $\mathbf{w}^* = \frac{1}{2}[1, 1]^T$  and bias  $w_0 = 1 - [1, 1]^T \mathbf{w}^* = 0$  can be immediately computed.



2.

## Remarks

1. The bias cannot be computed from the equations derived from the dual problem. It has to be computed from the primal constraints:

$$w_0^* = -\frac{\min_{t_\ell=1}\{\mathbf{x}_\ell^T \mathbf{w}^*\} + \max_{t_\ell=-1}\{\mathbf{x}_\ell^T \mathbf{w}^*\}}{2}.$$

2. The optimal hyperplane can be expressed as:

$$h(\mathbf{x}, \boldsymbol{\alpha}^*, w_0^*) = \sum_{\ell \in \text{sv}} t_\ell \alpha_\ell^* \mathbf{x}_\ell^T \mathbf{x} + w_0^*.$$

# Remarks

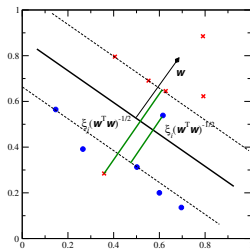
1. Observe that for all  $k \in \text{sv}$

$$t_k h(\mathbf{x}_k, \boldsymbol{\alpha}^*, w_0^*) = t_k \left( \sum_{\ell \in \text{sv}} t_\ell \alpha_\ell^* \mathbf{x}_\ell^T \mathbf{x}_k + w_0^* \right) = 1$$

$$\begin{aligned} (\mathbf{w}^*)^T \mathbf{w}^* &= \left( \sum_{\ell \in \text{sv}} t_\ell \alpha_\ell^* \mathbf{x}_\ell \right)^T \sum_{k \in \text{sv}} t_k \alpha_k^* \mathbf{x}_k \\ &= \sum_{k \in \text{sv}} \alpha_k^* t_k \sum_{\ell \in \text{sv}} t_\ell \alpha_\ell^* \mathbf{x}_\ell^T \mathbf{x}_k \\ &= \sum_{k \in \text{sv}} \alpha_k^* (1 - t_k w_0^*) = \sum_{k \in \text{sv}} \alpha_k^* - w_0^* \sum_{k \in \text{sv}} \alpha_k^* t_k \end{aligned}$$

$$\|\mathbf{w}^*\|_2 = \sqrt{\sum_{k \in \text{sv}} \alpha_k^*}$$

# Soft Margins



Consider the situation presented in the figure, where the problem is not linearly separable (due to noise in the data set).

2. We can soften the margins by considering the following constraints:

$$t_k(\mathbf{w}^T \mathbf{x}_k + w_0) \geq 1 - \xi_k$$
$$\xi_k \geq 0,$$

where  $\{\xi_k\}$  is a set of suitable slack variables.

# Soft Margins- Optimisation Problem

1. The solution to the constraint optimisation problem requires now to reduce the size of  $\mathbf{w}$  and  $\xi$  while satisfying the constraints  $t_k(\mathbf{w}^T \mathbf{x}_k + w_0) \geq 1 - \xi_k$  ( $\xi_k \geq 0$  is not needed, allow  $\xi_k < 0$  will not change the optimal solution).
2. The problem becomes:

$$\begin{aligned} \text{minimise} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} \xi^T \xi, & \mathbf{w} \in \Omega_{\mathbf{w}}, \quad \xi \in \Omega_{\xi}, \\ \text{subject to} \quad & t_k(\mathbf{w}^T \mathbf{x}_k + w_0) \geq 1 - \xi_k, & k = 1, \dots, N. \end{aligned}$$

3. The primal Lagrangian is

$$\mathcal{L}_P(\mathbf{w}, w_0, \xi, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} \xi^T \xi - \sum_{n=1}^N \alpha_n [t_n(\mathbf{w}^T \mathbf{x}_n + w_0) - 1 + \xi_n].$$

# Soft Margins- Optimisation Problem

1. The equations to be solved are:

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L} &= \mathbf{w}^* - \sum_{n=1}^N \alpha_n^* t_n \mathbf{x}_n = \mathbf{0}_d \\ \nabla_{\xi} \mathcal{L} &= C \xi^* - \alpha^* = \mathbf{0}_N \\ \partial_{w_0} \mathcal{L} &= \mathbf{t}^T \alpha^* = 0.\end{aligned}$$

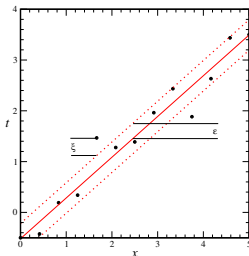
2. The dual Lagrangian becomes:

$$\begin{aligned}\text{maximise} \quad \mathcal{L}_D(\alpha) &= \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N t_m t_n \alpha_m \alpha_n \left( \mathbf{x}_m^T \mathbf{x}_n + \frac{\delta_{m,n}}{C} \right) \\ \text{subject to} \quad \mathbf{t}^T \alpha &= 0, \quad \alpha_n \geq 0, \quad 1 \leq n \leq N.\end{aligned}$$

3. The corresponding KKT complementary conditions are:  
 $\alpha_i [t_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i] = 0.$

# SVM for Regression

1. Relies on a loss function that ignores errors within a margin  $\varepsilon$  of the expected value. The variable  $\xi$  measures the cost of errors.



2. We define the slack variables  
$$\xi_n = \max\{0, |t_n - \mathbf{w}^T \mathbf{x}_n - w_0| - \varepsilon\}.$$



## Quadratic $\varepsilon$ -Insensitive Loss

1. We define the quadratic  $\varepsilon$ -insensitive loss as:

$$L_2^\varepsilon(\mathbf{x}, t, f) = (\max\{0, |t - f(\mathbf{x})| - \varepsilon\})^2.$$

2. The objective function for the problem is:

$$R\mathbf{w}^T\mathbf{w} + \sum_{n=1}^N L_2^\varepsilon(\mathbf{x}_n, t_n, f),$$

for a suitable  $R \in \mathbb{R}^+$ .

3. The primal problem is defined as:

$$\begin{aligned} &\text{minimise} && \frac{1}{2}\mathbf{w}^T\mathbf{w} + \frac{C}{2} \left( \boldsymbol{\xi}^T\boldsymbol{\xi} + \hat{\boldsymbol{\xi}}^T\hat{\boldsymbol{\xi}} \right), && \mathbf{w} \in \Omega_{\mathbf{w}}, \boldsymbol{\xi}, \hat{\boldsymbol{\xi}} \in \Omega_{\boldsymbol{\xi}}, \\ &\text{subject to} && \mathbf{w}^T\mathbf{x}_k + w_0 - t_k \leq \varepsilon + \xi_k, && k = 1, \dots, N, \\ &&& t_k - \mathbf{w}^T\mathbf{x}_k + w_0 \leq \varepsilon + \hat{\xi}_k, && k = 1, \dots, N, \\ &&& \xi_k, \hat{\xi}_k \geq 0, && k = 1, \dots, N. \end{aligned}$$

# Quadratic $\varepsilon$ -Insensitive Loss

1. The primal Lagrangian is

$$\begin{aligned}\mathcal{L}_P(\mathbf{w}, w_0, \xi, \hat{\xi}, \alpha, \hat{\alpha}) = & \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{2} (\xi^T \xi + \hat{\xi}^T \hat{\xi}) + \\ & + \sum_{n=1}^N \alpha_n (\mathbf{w}^T \mathbf{x}_n + w_0 - t_n - \varepsilon - \xi_n) + \\ & + \sum_{n=1}^N \hat{\alpha}_n (t_n - \mathbf{w}^T \mathbf{x}_n - w_0 - \varepsilon - \hat{\xi}_n)\end{aligned}$$

# Quadratic $\varepsilon$ -Insensitive Loss

1. The stability conditions are:

$$\nabla_{\mathbf{w}} \mathcal{L}_P = \mathbf{w} - \sum_{n=1}^N (\hat{\alpha}_n - \alpha_n) \mathbf{x}_n = \mathbf{0}_d,$$

$$\nabla_{\boldsymbol{\xi}} \mathcal{L}_P = C\boldsymbol{\xi} - \boldsymbol{\alpha} = \mathbf{0}_N,$$

$$\nabla_{\hat{\boldsymbol{\xi}}} \mathcal{L}_P = C\hat{\boldsymbol{\xi}} - \hat{\boldsymbol{\alpha}} = \mathbf{0}_N,$$

$$\partial_{w_0} \mathcal{L}_P = \sum_{n=1}^N (\hat{\alpha}_n - \alpha_n) = 0.$$

# Quadratic $\varepsilon$ -Insensitive Loss

1. The dual Lagrangian is:

$$\begin{aligned}\mathcal{L}_D(\boldsymbol{\alpha}, \hat{\boldsymbol{\alpha}}) = & \sum_{n=1}^N t_n(\hat{\alpha}_n - \alpha_n) - \varepsilon \sum_{n=1}^N (\alpha_n + \hat{\alpha}_n) - \\ & - \frac{1}{2} \sum_{m,n=1}^N (\hat{\alpha}_m - \alpha_m)(\hat{\alpha}_n - \alpha_n) \mathbf{x}_m^T \mathbf{x}_n - \\ & - \frac{1}{2C} (\boldsymbol{\alpha}^T \boldsymbol{\alpha} + \hat{\boldsymbol{\alpha}}^T \hat{\boldsymbol{\alpha}})\end{aligned}$$

$$\text{subject to } \sum_{n=1}^N (\hat{\alpha}_n - \alpha_n) = 0.$$

# Quadratic $\varepsilon$ -Insensitive Loss

1. The KKT complementary conditions are:

$$\alpha_n \left( \mathbf{w}^T \mathbf{x}_n + w_0 - t_n - \varepsilon - \xi_n \right) = 0, \quad n = 1, \dots, N,$$

$$\hat{\alpha}_n \left( t_n - \mathbf{w}^T \mathbf{x}_n - w_0 - \varepsilon - \hat{\xi}_n \right) = 0, \quad n = 1, \dots, N,$$

$$\alpha_n \hat{\alpha}_n = 0, \quad i = 1, \dots, N,$$

$$\xi_n \hat{\xi}_n = 0, \quad i = 1, \dots, N,$$

the last two conditions indicate that it is not possible a data point, at the same time, to be wrong for excess and defect.

# Quadratic $\varepsilon$ -Insensitive Loss

1. The dual becomes:

$$\begin{aligned}\mathcal{L}_D(\boldsymbol{\alpha}, \hat{\boldsymbol{\alpha}}) &= \sum_{n=1}^N t_n(\hat{\alpha}_n - \alpha_n) - \varepsilon \sum_{n=1}^N (\alpha_n + \hat{\alpha}_n) - \\ &\quad - \frac{1}{2} \sum_{m,n=1}^N (\hat{\alpha}_m - \alpha_m)(\hat{\alpha}_n - \alpha_n) \mathbf{x}_m^T \mathbf{x}_n - \\ &\quad - \frac{1}{2C} \left( \boldsymbol{\alpha}^T \boldsymbol{\alpha} + \hat{\boldsymbol{\alpha}}^T \hat{\boldsymbol{\alpha}} \right) + \frac{1}{C} \sum_{n=1}^N \alpha_n \hat{\alpha}_n \\ &= \sum_{n=1}^N t_n(\hat{\alpha}_n - \alpha_n) - \varepsilon \sum_{n=1}^N (\alpha_n + \hat{\alpha}_n) - \\ &\quad - \frac{1}{2} \sum_{m,n=1}^N (\hat{\alpha}_m - \alpha_m)(\hat{\alpha}_n - \alpha_n) \left( \mathbf{x}_m^T \mathbf{x}_n + \frac{\delta_{m,n}}{C} \right)\end{aligned}$$

## Quadratic $\varepsilon$ -Insensitive Loss

1. By making the substitution  $\beta = \hat{\alpha} - \alpha$  and by using that  $\hat{\alpha}_n \alpha_n = 0$  for all  $1 \leq n \leq N$ ,

$$\mathcal{L}_D(\beta) = \sum_{n=1}^N t_n \beta_n - \varepsilon \sum_{n=1}^N |\beta_n| - \frac{1}{2} \beta_m \beta_n \left( \mathbf{x}_m^T \mathbf{x}_n + \frac{\delta_{m,n}}{C} \right)$$

subject to  $\sum_{n=1}^N \beta_n = 0.$