

# Mathematics for Machine Learning

## Classification with Support Vector Machines

May 24, 2025

Anastasiya Kuznetsova

# Agenda

1. Refresh. Inner product
2. “Hard” SVM
3. Concept of the Margin
4. Refresh. Lagrange multiplier
5. Convex duality via Lagrange multiplier
6. Dual SVM
7. Kernels

# Refresh. Inner product

Symbol for inner product

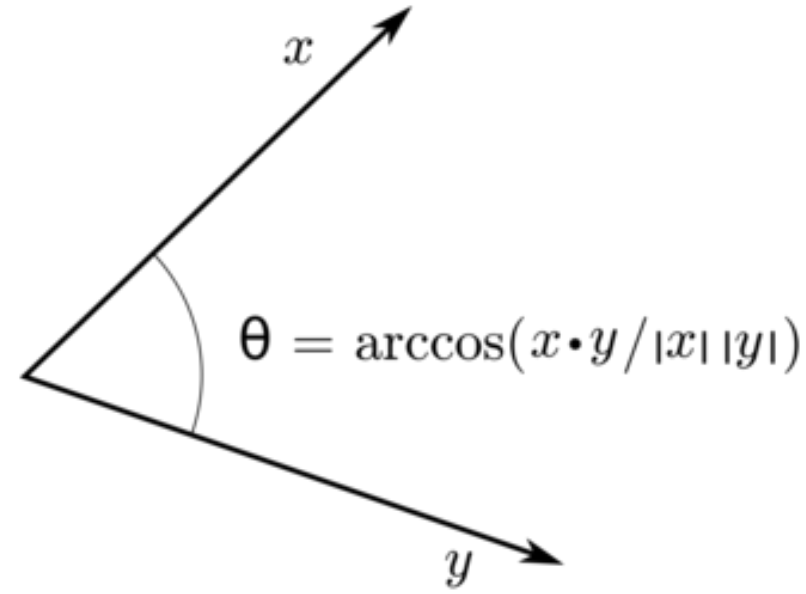
Length of vector  $\mathbf{u}, \mathbf{v}$

Angle between  $\mathbf{u}$  and  $\mathbf{v}$

$$\mathbf{u} \bullet \mathbf{v} = |\mathbf{u}| |\mathbf{v}| \cos(\theta) \quad 1$$

$$= x_1 \times x_2 + y_1 \times y_2 \quad 2$$

$$= \mathbf{u} \mathbf{v}^T \quad 3$$



$$\text{proj}_{\mathbf{v}} \mathbf{u} = \frac{\mathbf{v}}{|\mathbf{v}|} |\mathbf{u}| \cos \theta$$

Direction

Length

$$\cos \theta = \frac{|\text{proj}_{\mathbf{v}} \mathbf{u}|}{|\mathbf{u}|}$$

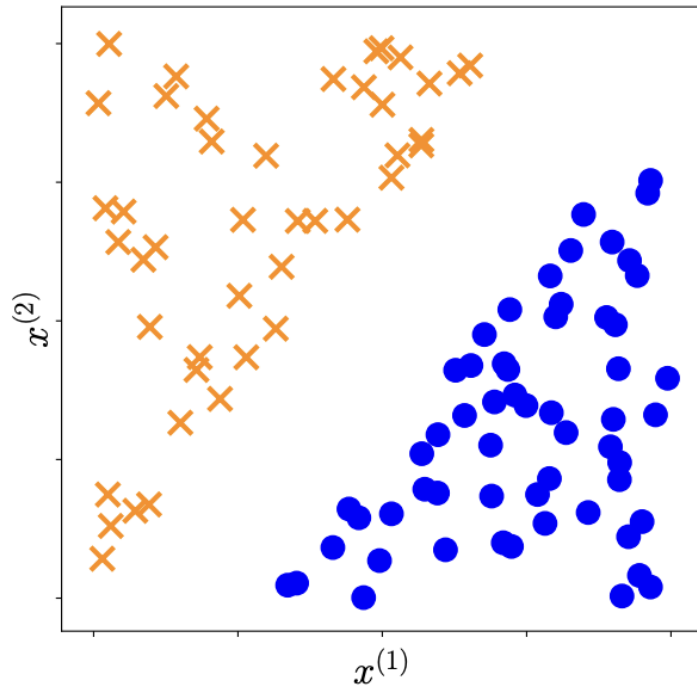
$$\begin{aligned} \text{proj}_{\mathbf{v}} \mathbf{u} &= \frac{\mathbf{v}}{|\mathbf{v}|} |\mathbf{u}| \cos \theta \\ &= \frac{\mathbf{v}}{|\mathbf{v}|} \cancel{|\mathbf{u}|} \frac{\mathbf{u} \cdot \mathbf{v}}{\cancel{|\mathbf{u}|} |\mathbf{v}|} \\ &= \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{v}|^2} \mathbf{v} \end{aligned}$$

# Classification with Support Vector Machines

$$f : \mathbb{R}^D \rightarrow \{+1, -1\}$$

examples  $\mathbf{x}_n \in \mathbb{R}^D$

$$y_n \in \{+1, -1\}.$$



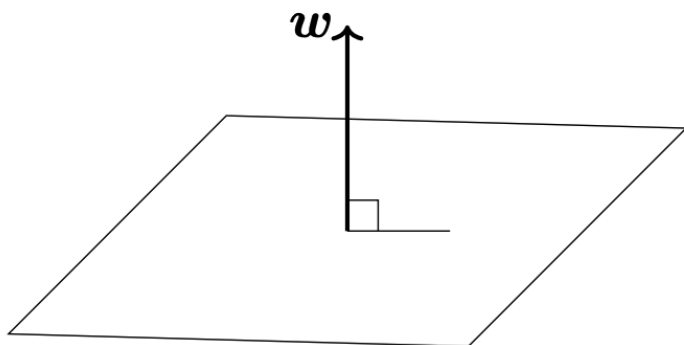
We imagine binary classification data, which can be separated by a hyperplane

Every example  $\mathbf{x}_n$  (a vector of dimension 2) is a two-dimensional location  $(x^{(1)}_n \text{ and } x^{(2)}_n)$ , and the corresponding binary label  $y_n$

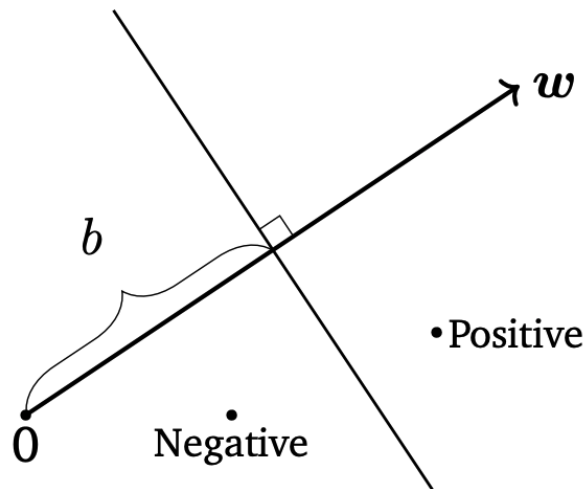
A hyperplane is an affine subspace of dimension  $D - 1$  (if the corresponding vector space is of dimension  $D$ ).

We formalize the idea of finding a linear separator of the two classes.

## Separating Hyperplanes §12.1



(a) Separating hyperplane in 3D



(b) Projection of the setting in (a) onto a plane

Represent data in  $\mathbb{R}^D$

Find a hyperplane that divides the data into two partitions such that negative examples (labelled as -1) are on the opposite side of positive examples (labelled as +1)

Vector  $w$  is a vector normal to the hyperplane and  $b$  the intercept

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto f(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + b$$

We define the hyperplane as

$$\{\mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) = 0\}$$

# Separating Hyperplanes

Therefore, to classify a test example  $\mathbf{x}_{\text{test}}$ , we calculate the value of the function  $f(\mathbf{x}_{\text{test}})$  and classify the example as +1 if  $f(\mathbf{x}_{\text{test}}) \geq 0$  and -1 otherwise.

When training the classifier, we want to ensure that the examples with positive labels are on the positive side of the hyperplane

$$\langle \mathbf{w}, \mathbf{x}_n \rangle + b \geq 0 \quad \text{when} \quad y_n = +1$$

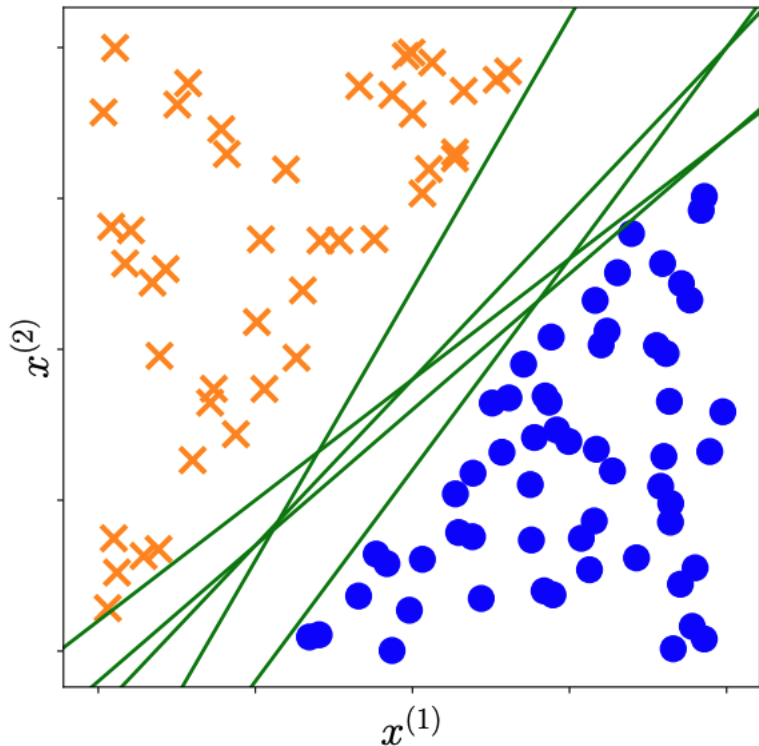
and the examples with negative labels are on the negative side, i.e.,

$$\langle \mathbf{w}, \mathbf{x}_n \rangle + b < 0 \quad \text{when} \quad y_n = -1.$$

In a single equation

$$y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 0$$

## Concept of the Margin §12.2.1



But we can create infinitely many separating hyperplanes,  
How do we choose the best?

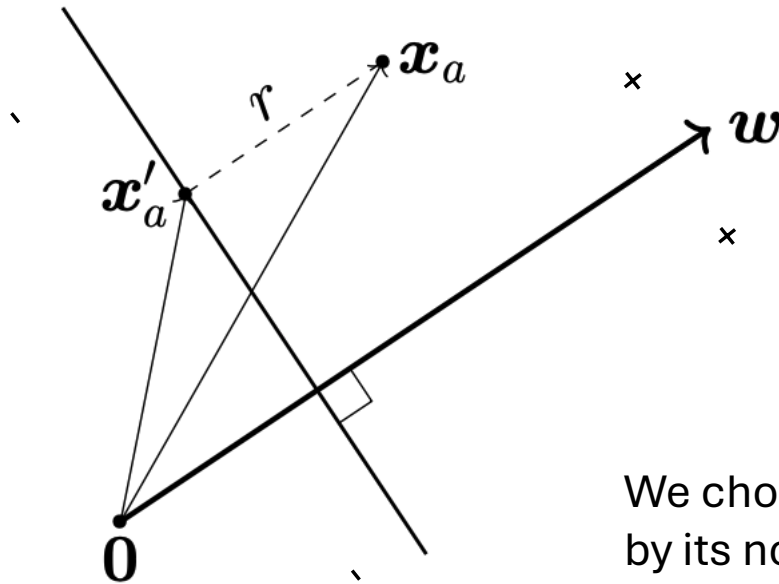
Choose the separating hyperplane that maximizes the margin  
between the positive and negative examples!

Margin is the distance of the separating hyperplane to the closest  
examples in the dataset, assuming that the dataset is linearly  
separable.

How can we calculate this distance?

# Concept of the Margin

Consider a hyperplane  $\langle \mathbf{w}, \mathbf{x} \rangle + b$ , and an example  $\mathbf{x}_a$  on the positive side of the hyperplane  $\langle \mathbf{w}, \mathbf{x}_a \rangle + b > 0$ .



We need to compute distance  $r > 0$  of  $\mathbf{x}_a$  from the hyperplane

$\mathbf{x}'_a$  orthogonal projection of the  $\mathbf{x}_a$  to the hyperplane

Since  $\mathbf{w}$  is orthogonal to the hyperplane, we know that the distance  $r$  is just a scaling of this vector  $\mathbf{w}$ .

We choose to use a vector of unit length (its norm is 1) and obtain this by dividing  $\mathbf{w}$  by its norm  $\frac{\mathbf{w}}{\|\mathbf{w}\|}$

$$\mathbf{x}_a = \mathbf{x}'_a + r \frac{\mathbf{w}}{\|\mathbf{w}\|}.$$

If  $\mathbf{x}_a$  is the closest point,  $r$  is the margin



# Concept of the Margin

We would like the positive examples to be further than  $r^+$  from the hyperplane, and the negative examples to be further than  $r^-$  in a negative direction.

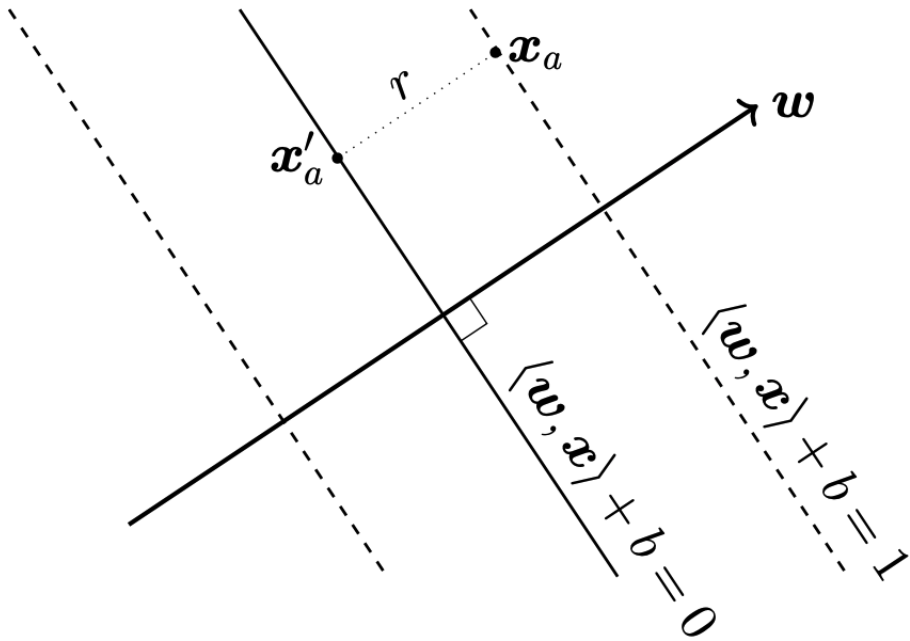
Since we are interested only in the direction, we add an assumption to our model that the parameter vector  $\mathbf{w}$  is of unit length, i.e.,  $\|\mathbf{w}\| = 1$ , where  $\|\mathbf{w}\| = \sqrt{\mathbf{w}^\top \mathbf{w}}$

Our objective  $\max_{\mathbf{w}, b, r} \underbrace{r}_{\text{margin}}$

subject to  $\underbrace{y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq r}_{\text{data fitting}}, \underbrace{\|\mathbf{w}\| = 1}_{\text{normalization}}, \quad r > 0,$

We want to maximize the margin  $r^+$  while ensuring that the data lies on the correct side of the hyperplane

## Traditional Derivation of the Margin §12.2.2



$$\mathbf{x}_a = \mathbf{x}'_a + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

Instead of  $\|\mathbf{w}\| = 1$  assumption we choose a scale for the data

We choose this scale such that the value of the predictor  $\langle \mathbf{w}, \mathbf{x} \rangle + b$  is 1 at the closest example.

Since  $\mathbf{x}'_a$  is the orthogonal projection of  $\mathbf{x}_a$  onto hyperplane it lies exactly on the hyperplane

$$\langle \mathbf{w}, \mathbf{x}'_a \rangle + b = 0$$

$$\left\langle \mathbf{w}, \mathbf{x}_a - r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right\rangle + b = 0$$

$$\underbrace{\langle \mathbf{w}, \mathbf{x}_a \rangle + b}_{=1} - r \frac{\langle \mathbf{w}, \mathbf{w} \rangle}{\|\mathbf{w}\|} = 0$$

1 by the assumption of the scale

## Traditional Derivation of the Margin

$$1 - r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} = 0$$

$$r = \frac{1}{\|\mathbf{w}\|}$$

We derived the distance  $r$  in terms of the normal vector  $\mathbf{w}$  of the hyperplane

We want the positive and negative examples to be at least 1 away from the hyperplane

$$y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1.$$

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|}$$

subject to  $y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1$  for all  $n = 1, \dots, N$ .

## Traditional Derivation of the Margin

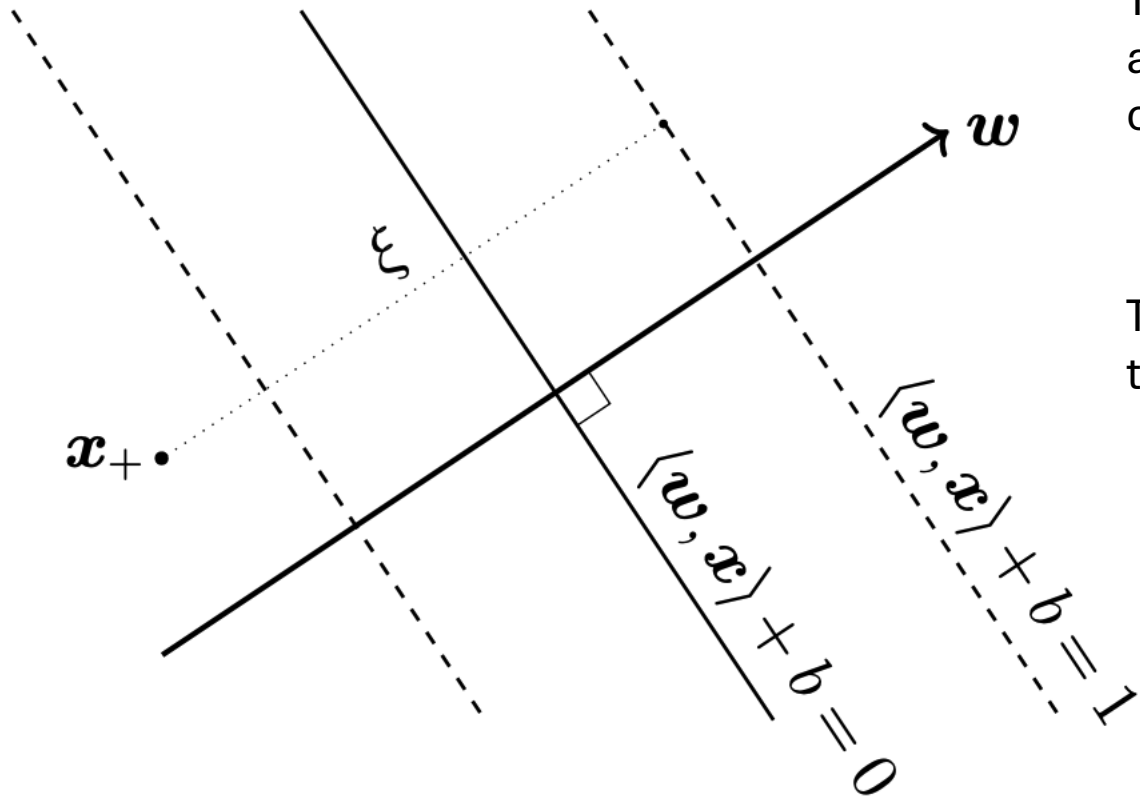
$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{1}{\|\mathbf{w}\|} \\ \text{subject to} \quad & y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 \quad \text{for all } n = 1, \dots, N. \end{aligned}$$

In real practice we **minimize** the squared norm

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \longleftarrow \text{hard margin SVM} \\ \text{subject to} \quad & y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 \quad \text{for all } n = 1, \dots, N \end{aligned}$$

## Soft Margin SVM: Geometric View §12.2.4

In the case where data is not linearly separable, we may wish to allow some examples to fall within the margin region, or even to be on the wrong side of the hyperplane



The slack variable  $\xi_n$  to each example–label pair  $(\mathbf{x}_n, y_n)$  that allows a particular example to be within the margin or even on the wrong side of the hyperplane

To encourage correct classification of the samples, we add the sum of the  $\xi_n$  values to the objective, weighted by  $C$

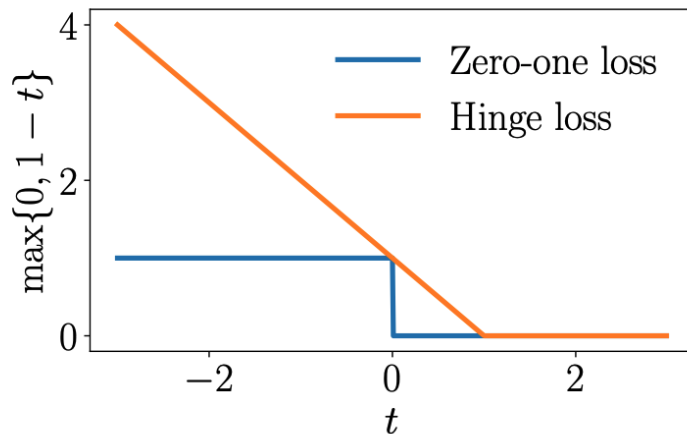
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & y_n (\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 - \xi_n \\ & \xi_n \geq 0 \end{aligned}$$

## Soft Margin SVM: Loss Function View §12.2.5

We choose the hyperplane as  $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ .

Hinge loss  $\ell(t) = \max\{0, 1 - t\}$  where  $t = yf(\mathbf{x}) = y(\langle \mathbf{w}, \mathbf{x} \rangle + b)$

1. If  $f(\mathbf{x})$  is on the correct side of the hyperplane, and further than distance 1  $\rightarrow t \geq 1 \rightarrow \ell = 0$
2. If  $f(\mathbf{x})$  is on the correct side of the hyperplane, but within the margin  $\rightarrow 0 < t < 1 \rightarrow \ell > 0$
3. If  $f(\mathbf{x})$  is on the opposite side of the hyperplane  $\rightarrow t < 0 \rightarrow \ell > 0$   
 $\ell$  larger than in 2



Using hinge loss gives us the unconstrained optimization problem

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{regularizer}} + C \underbrace{\sum_{n=1}^N \max\{0, 1 - y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b)\}}_{\text{error term}} .$$

## Lagrange multiplier refresh

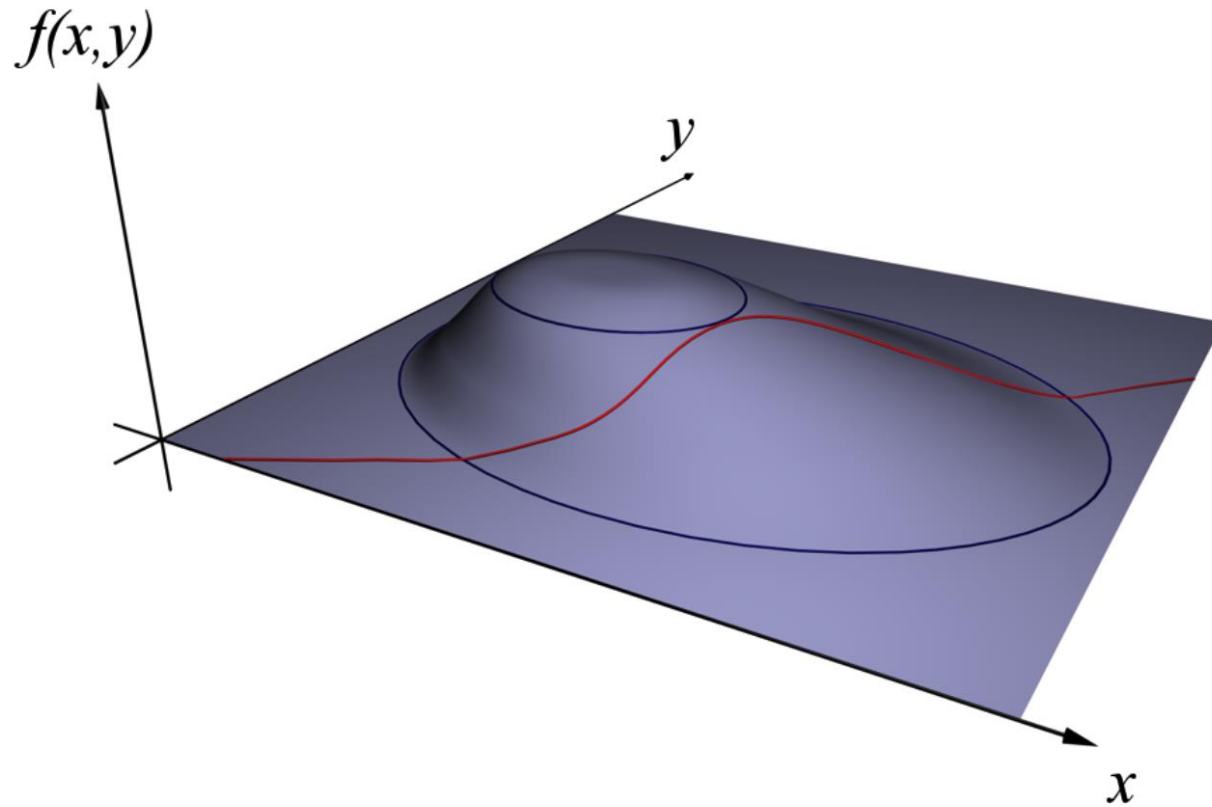
- The Lagrange multiplier technique lets you find the maximum or minimum of a multivariable function  $f(x, y, \dots)$  when there is some constraint on the input values you are allowed to use.
- This technique only applies to constraints that look something like this:

$$g(x, y, \dots) = c$$

Here,  $g$  is another multivariable function with the same input space as  $f$ , and  $c$  is some constant. [Picture](#) ✓

# Lagrange multiplier refresh

For example, if the input space is two-dimensional, the graph of  $f$  with the line representing  $g(x, y) = c$  projected onto it might look something like this:





# Lagrange multiplier refresh

- **Step 1:** Introduce a new variable  $\lambda$ , and define a new function  $\mathcal{L}$  as follows:

$$\mathcal{L}(x, y, \dots, \lambda) = f(x, y, \dots) - \lambda(g(x, y, \dots) - c)$$

This function  $\mathcal{L}$  is called the "Lagrangian", and the new variable  $\lambda$  is referred to as a "Lagrange multiplier"

- **Step 2:** Set the gradient of  $\mathcal{L}$  equal to the zero vector.

$$\nabla \mathcal{L}(x, y, \dots, \lambda) = \mathbf{0} \quad \leftarrow \text{Zero vector}$$

In other words, find the **critical points** of  $\mathcal{L}$ .

- **Step 3:** Consider each solution, which will look something like  $(x_0, y_0, \dots, \lambda_0)$ . Plug each one into  $f$ . Or rather, first remove the  $\lambda_0$  component, then plug it into  $f$ , since  $f$  does not have  $\lambda$  as an input. Whichever one gives the greatest (or smallest) value is the maximum (or minimum) point you are seeking.

## Convex Duality via Lagrange Multipliers §12.3.1

Lagrange multiplier  $\alpha_n \geq 0$  corresponding to the constraint  $y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) \geq 1 - \xi_n$

Lagrange multiplier  $\gamma_n \geq 0$  corresponding to the constraint of non-negativity of the slack variable  $\xi_n \geq 0$

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \gamma) = & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ & - \underbrace{\sum_{n=1}^N \alpha_n (y_n(\langle \mathbf{w}, \mathbf{x}_n \rangle + b) - 1 + \xi_n)}_{\text{constraint (12.26b)}} - \underbrace{\sum_{n=1}^N \gamma_n \xi_n}_{\text{constraint (12.26c)}} . \end{aligned} \quad (12.34)$$

Differentiate the Lagrangian with respect to the three parameters

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w}^\top - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n^\top \quad (12.35)$$

## Convex Duality via Lagrange Multipliers

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{n=1}^N \alpha_n y_n ,$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \alpha_n - \gamma_n .$$

By setting 12.35 to 0      $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$

Substituting the expression for  $\mathbf{w}$  into the Lagrangian

$$\begin{aligned} \mathfrak{D}(\xi, \alpha, \gamma) = & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^N y_i \alpha_i \left\langle \sum_{j=1}^N y_j \alpha_j \mathbf{x}_j, \mathbf{x}_i \right\rangle \\ & + C \sum_{i=1}^N \xi_i - b \sum_{i=1}^N y_i \alpha_i + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i \xi_i - \sum_{i=1}^N \gamma_i \xi_i . \end{aligned} \tag{12.39}$$

## Convex Duality via Lagrange Multipliers

Recall that inner products are symmetric and bilinear, then the terms in blue in 12.39 can be simplified

$$\mathfrak{D}(\xi, \alpha, \gamma) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i + \sum_{i=1}^N (C - \alpha_i - \gamma_i) \xi_i.$$

(12.40)

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^N \alpha_i$$

$$\text{subject to} \quad \sum_{i=1}^N y_i \alpha_i = 0$$

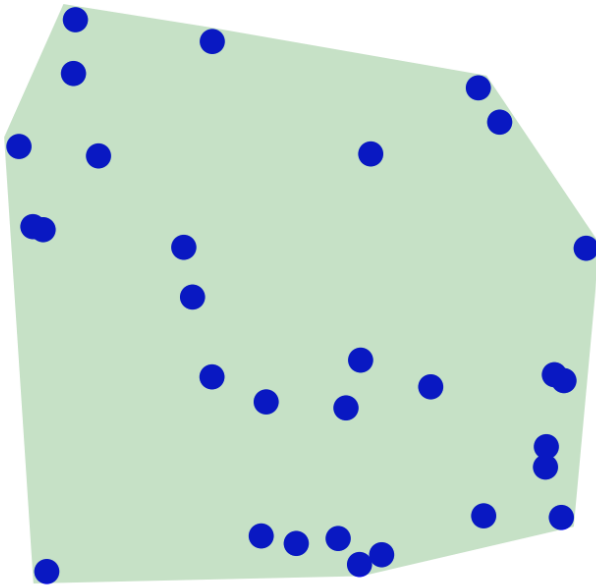
$$0 \leq \alpha_i \leq C \quad \text{for all } i = 1, \dots, N.$$

## Dual SVM: Convex Hull View §12.3.2

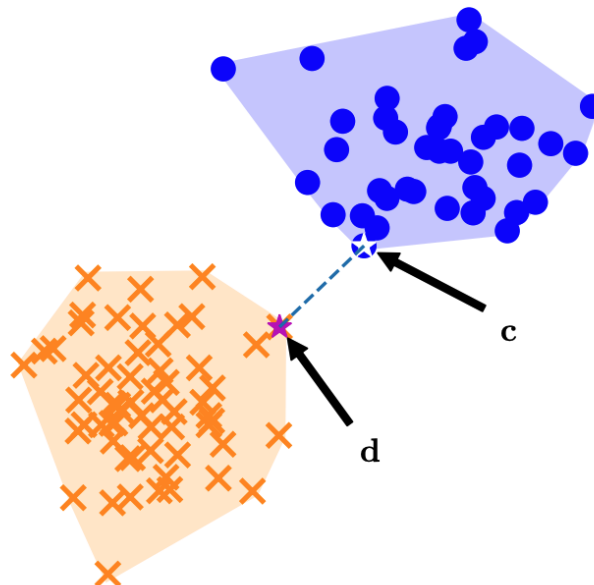
Consider the set of examples  $x_n$  with the same label. We would like to build a convex set that contains all the examples such that it is the smallest possible set. This is called the convex hull.

Building a convex hull can be done by introducing non-negative weights  $\alpha_n \geq 0$  corresponding to each example  $x_n$

Then the convex hull can be described as the set  $\text{conv}(\mathbf{X}) = \left\{ \sum_{n=1}^N \alpha_n \mathbf{x}_n \right\}$  with  $\sum_{n=1}^N \alpha_n = 1$  and  $\alpha_n \geq 0$ , (12.43)



(a) Convex hull.



(b) Convex hulls around positive (blue) and negative (orange) examples. The distance be-

We create convex hulls for positive and negative classes, respectively.

## Dual SVM: Convex Hull View

We pick a point  $\mathbf{c}$ , which is in the convex hull of the set of positive examples, and is closest to the negative class distribution.

We pick a point  $\mathbf{d}$  in the convex hull of the set of negative examples and is closest to the positive class distribution;

We define a difference vector between  $\mathbf{d}$  and  $\mathbf{c}$  as  $\mathbf{w} := \mathbf{c} - \mathbf{d}$ .

Corresponding optimization problem  $\arg \min_{\mathbf{w}} \|\mathbf{w}\| = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$ .

$\mathbf{c}$  can be expressed as a convex combination of the positive examples  $\mathbf{c} = \sum_{n:y_n=+1} \alpha_n^+ \mathbf{x}_n$ .

Similarly for the  $\mathbf{d}$   $\mathbf{d} = \sum_{n:y_n=-1} \alpha_n^- \mathbf{x}_n$ .

Objective  $\min_{\alpha} \frac{1}{2} \left\| \sum_{n:y_n=+1} \alpha_n^+ \mathbf{x}_n - \sum_{n:y_n=-1} \alpha_n^- \mathbf{x}_n \right\|^2$ .

## Dual SVM: Convex Hull View

Let  $\alpha$  be the set of all coefficients, i.e., the concatenation of  $\alpha^+$  and  $\alpha^-$ .

$$\sum_{n:y_n=+1} \alpha_n^+ = 1 \quad \text{and} \quad \sum_{n:y_n=-1} \alpha_n^- = 1.$$

This implies the constraint  $\sum_{n=1}^N y_n \alpha_n = 0$ .

Could be solved using Lagrangian!

But we won't solve it here 😊

## Kernels §12.4

Recall the dual SVM formulation. The inner products are only between examples, not parameters. This means you can modify the representation of the data and only need to replace the inner product values.

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^N \alpha_i \\ \text{subject to} \quad & \sum_{i=1}^N y_i \alpha_i = 0 \\ & 0 \leq \alpha_i \leq C \quad \text{for all } i = 1, \dots, N. \end{aligned} \quad (12.41)$$

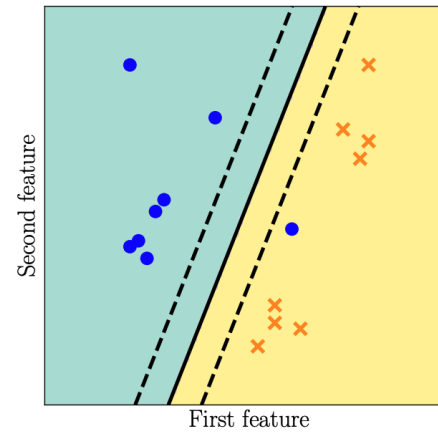
Since  $\phi(\mathbf{x})$  could be a non-linear function, we can use the SVM (which assumes a linear classifier) for separating non-linear data.

We will define a similarity function  $k$ , between samples, without explicitly defining the non-linear feature map  $\phi(\mathbf{x})$

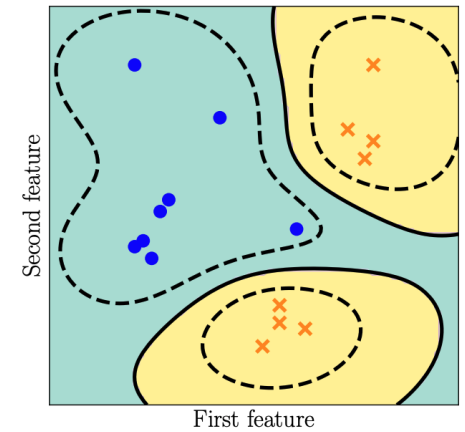
Kernels are by definition functions  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

for which there is a Hilbert space  $\mathcal{H}$  and  $\phi : \mathcal{X} \rightarrow \mathcal{H}$

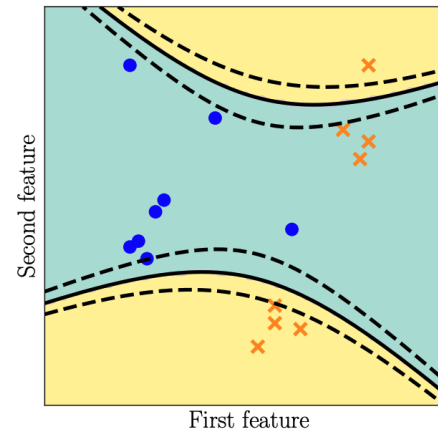
Is a feature map such that  $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}}$ .



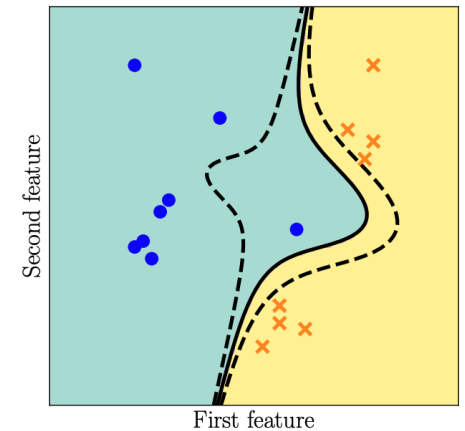
(a) SVM with linear kernel



(b) SVM with RBF kernel



(c) SVM with polynomial (degree 2) kernel



(d) SVM with polynomial (degree 3) kernel



## Numerical Solution §12.5

The loss view of SVM is a convex unconstrained optimization problem

But the hinge loss (12.28) is not differentiable, so the subgradient approach is used to solve it

Both the primal and dual SVM result in a convex quadratic programming problem

- The primal SVM has optimization variables that have the size  $D$  of the dimension of samples
- The dual SVM has optimization variables that have the size  $N$  of the number of samples

The primal soft margin SVM in matrix form can be expressed as:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix}^\top \begin{bmatrix} \mathbf{I}_D & \mathbf{0}_{D, N+1} \\ \mathbf{0}_{N+1, D} & \mathbf{0}_{N+1, N+1} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} + [\mathbf{0}_{D+1, 1} \quad C\mathbf{1}_{N, 1}]^\top \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} \\ \text{subject to} \quad & \begin{bmatrix} -\mathbf{Y}\mathbf{X} & -\mathbf{y} & -\mathbf{I}_N \\ \mathbf{0}_{N, D+1} & & -\mathbf{I}_N \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \\ \xi \end{bmatrix} \leq \begin{bmatrix} -\mathbf{1}_{N, 1} \\ \mathbf{0}_{N, 1} \end{bmatrix}. \end{aligned} \quad (12.56)$$

The dual SVM can be written as:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^\top \mathbf{Y} \mathbf{K} \mathbf{Y} \alpha - \mathbf{1}_{N, 1}^\top \alpha \\ \text{subject to} \quad & \begin{bmatrix} \mathbf{y}^\top \\ -\mathbf{y}^\top \\ -\mathbf{I}_N \\ \mathbf{I}_N \end{bmatrix} \alpha \leq \begin{bmatrix} \mathbf{0}_{N+2, 1} \\ C\mathbf{1}_{N, 1} \end{bmatrix}. \end{aligned} \quad (12.57)$$