



Title _____

Notes

Let's consider again a binary classification problem

We have shown that logistic regression provides a linear classification rule that maximizes the log-probability of class assignments $\sum_{i=1}^n \log P(C_i = c_i | X_i = x_i, w, b)$ for the training set

The LR solution can be found by minimizing the logistic loss

As we have mentioned, it is often useful to add a regularization term of the form $\frac{\lambda}{2} \|w\|^2$ to prevent overfitting:

$$w^*, b = \arg \min_{w,b} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-z_i(w^T x_i + b)} \right)$$

where z_i is the class label for sample x_i , encoded as

$$z_i = \begin{cases} +1 & \text{if } c_i = \mathcal{H}_T \\ -1 & \text{if } c_i = \mathcal{H}_F \end{cases}$$

We can alternatively minimize a scaled version of the objective:

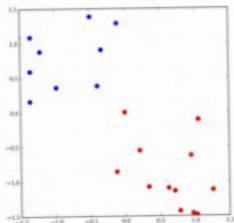
$$w^*, b = \arg \min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \log \left(1 + e^{-z_i(w^T x_i + b)} \right)$$

where $C = \frac{1}{n\lambda}$

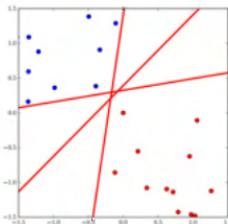
Logistic regression has a fundamental issue when classes are perfectly separable: the optimization has no true minimum.

For perfectly separate classes we could always make our model more confident by scaling up the weight vector w . As $\|w\|$ increases the model push class probabilities closer to 0 or 1, continuously decreasing the loss function. This means there's an infimum (greatest lower bound) for the loss at $\|w\|=0$ but no actual minimum value.

Assume that we have two classes that are linearly separable (i.e., we can find a linear separation hyperplane that correctly classifies all points of our training set)

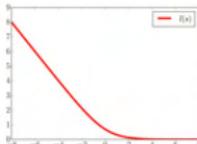


Problem: there is an infinite number of separating hyperplanes, which one should we select?



An optimal LR solution is not defined, since the objective function has an infimum for $\|w\|^2 \rightarrow +\infty$, but not a minimum

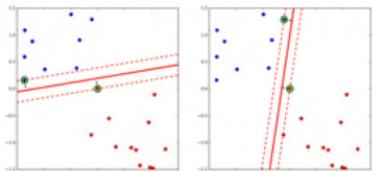
LR optimizes the Logistic Loss. As we move further along the x axis the loss tends to 0



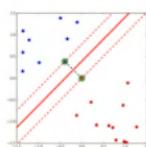
When multiple hyperplanes could separate our classes, svts answer the question. "Which hyperplane is optimal?"

The core insight is to select the hyperplane that creates the maximum margin between classes. The margin is define as the distance from the decision boundary to the nearest data point from either class.

Intuitively, we can select the hyperplane that separates the classes with the largest margin



Intuitively, we can select the hyperplane that separates the classes with the largest margin



The margin is defined as the distance of the closest point w.r.t. the separation hyperplane¹

¹Some authors separately introduce the margin for each class. This distinction has no practical effects on the derivation of the SVM objective.

Starting with our linear classifier function:

$$f(x) = w^T x + b$$

This function represent our separating hyperplane, where point with $f(x)=0$ lie directly on the boundary.

For any point in our feature space, we want to know how far it is from our decision boundary.

The distance from a point x_i to a hyperplane is given by:

$$d(x_i) = \frac{|f(x_i)|}{\|w\|}$$

The numerator $|f(x_i)|$ tells us how strongly our function classifies the point while dividing by $\|w\|$ normalize this value to give us the actual perpendicular distance.

$$s(x) = w^T x + b \stackrel{>0}{\geq 0}$$

Why this formula??



we want D , the distance

$$P_0(x_0, y_0, z_0)$$

$$P_0(x_0, y_0, z_0)$$

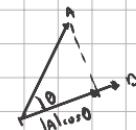
$$\vec{B} = \vec{P_0 P_2} = \langle x_2 - x_0, y_2 - y_0, z_2 - z_0 \rangle$$

\vec{p} is the normal vector (w), in order to find D , we want the scalar projection of \vec{B} on \vec{p}

negative part will be discussed in another section

$$D = \text{comp}_{\vec{p}} \vec{B} = |\vec{B}| \cos \theta = \frac{|\vec{p} \cdot \vec{B}|}{\|\vec{p}\|}$$

$$\vec{p} \cdot \vec{B} = \frac{|a|}{|a|} \frac{|b|}{|b|} |b| \cos \theta \quad [\text{projection scalar}]$$



$$D = \frac{|\vec{p} \cdot \vec{B}|}{\|\vec{p}\|} = \frac{|a(x_2 - x_0) + b(y_2 - y_0) + c(z_2 - z_0)|}{\sqrt{a^2 + b^2 + c^2}} =$$

$$= \frac{|(ax_2 + bx_2 + cx_2) - (ax_0 + bx_0 + cx_0)|}{\sqrt{a^2 + b^2 + c^2}} \quad \left. \right\} = \frac{|ax_2 + bx_2 + cx_2 + d|}{\|\vec{p}\|} = \frac{|f(P_0)|}{\|\vec{p}\|}$$

the equation of a plane is $ax + by + cz = d$

B

$$f(x) = ax + by + cz + d$$

Since we are working with binary classifications, we use z_i to denote the class label of point x_i where:

- $z_i = 1$ for the positive class

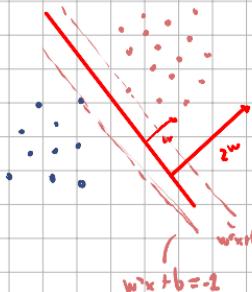
- $z_i = -1$ for the negative class

For a correctly classified point the sign of $f(x_i)$ will match the sign of z_i . This means $z_i f(x_i) > 0$ for all correctly classified points.

Using this property we can rewrite the distance formula as:

$$d(x_i) = \frac{|z_i (w^T x_i + b)|}{\|w\|}$$

Visualizing how we can simplify the problem !



(w, b) The hyperplane with (w, b) and $(2w, 2b)$ is the same. When I'm thinking of a separation rule, the same rule can be represented in infinite way $(2w, 2b)$ with one I should select?

I will select only those planes(feasible planes that correctly separates the point) for which the closest points to the plane has a score of $(w^T x_i + b = \pm 1)$

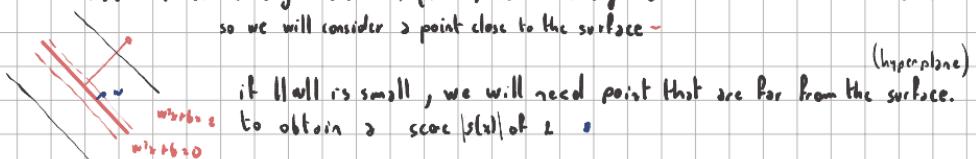
The $w^T x_i + b = \pm 1$ this will be consider an optimal solution

If the support vector have this constraint $w^T x_i + b = \pm 1$ then

$$\begin{aligned} \forall i: (w^T x_i + b) &\geq 1 \quad z_i = 1 \\ (w^T x_i + b) &\leq -1 \quad z_i = -1 \end{aligned}$$

With this limitation we can see that the range of the margin ~~$\frac{2w}{\|w\|}$~~ will be proportional to $\|w\|$

if we have $\|w\|$ is large then $x_i(w^T x_i + b)$ need to change a bit in order to obtain a score of ± 1 so we will consider a point close to the surface ~



if $\|w\|$ is small, we will need point that are far from the surface to obtain a score $|z_i|$ of ± 1

So to have a large margin we want a low $\|w\|$ or $\frac{1}{2} \|w\|^2$

Formal way

The goal of SVM is to find the hyperplane that maximizes the minimum distance to any training point.

$$w^*, b^* = \arg \max_{w, b} \min_{i \in \{1, \dots, n\}} d(x_i)$$

we select only
the support vector

Subject to the constraint that all points are correctly classified: $\forall i: (w^T x_i + b) > 0$

We can do some simplifications (for simplify the problem, keeping the same solution)

① First we can remove the absolute value thanks to the constraint

$$\forall i: (w^T x_i + b) > 0$$

② We can rewrite our optimization problem as:

$$w^*, b^* = \arg \max_{w, b} \frac{1}{\|w\|} \min_{i \in \{1, \dots, n\}} z_i (w^T x_i + b)$$

the classes are linear
separable

so far sure the minimum
of the distance will
be positive

We are looking for the hyperplane that maximizes the margin (represented by the factor $\frac{1}{\|w\|}$) while ensuring all points are correctly classified with some positive confidence value represented by $\min_i z_i (w^T x_i + b)$

Let's imagine on

We can show the two problems are equivalent:

- Hyperplanes that do not satisfy this constraint will have $\min_i z_i (w^T x_i + b) \leq 0$, and thus cannot be an optimal solution of the second formulation (classes are separable, so an optimal solution exists with $z_i (w^T x_i + b) > 0, \forall i = 1 \dots n$)
to class 0 & 1 are separable
but not w & b are same
- For hyperplanes that correctly classify all patterns, the objective functions are the same

The SVM

Focus on boundary points: The minimum distance is determined by the closest points to the boundary (the support vectors), making the solution sparse.

Scale invariance: The decision boundary depends only on the direction of w , not its magnitude, providing a unique solution unlike logistic regression.

Robustness: By maximizing the minimum distance, SVMs create the most "cautious" boundary possible, increasing resilience to small perturbations in the data.

We need to transform our function in something more tractable

Step 2

Our objective function has a special property called scale invariance. This mean that if we multiply both w and b by any positive constant α , the decision boundary remains the same.

$$\frac{1}{\|w\|_2} \min_i [z_i(w^T x_i + b)] = \frac{1}{\|\alpha w\|_2} \min_i [\alpha z_i(\alpha w^T x_i + \alpha b)]$$



Since $\|w\|_2 = \alpha \|w\|_2$ we get

$$\frac{1}{\|w\|_2} \min_i [\alpha z_i(w^T x_i + b)]$$

This mean that (w, b) and $(\alpha w, \alpha b)$ give us the same objective value for any $\alpha > 0$. They represent the same hyperplane but with different scaling of the parameters.

$y = 2x + b$ and $2y = 6x + 3b$ describe exactly the same line.

If (w, b) is an optimal solution, then also $(w, b) = (\alpha w, \alpha b)$ is optimal e viceversa.

The collections of values $\{(z_i, w^T x_i + b)\}_{z \in \mathbb{R}}$ forms an equivalent class of equivalent solutions.

For each of these equivalence classes, we are free to select any one of the equivalent solutions. In particular, we restrict our problem to solutions for which

$$\min_i z_i(w^T x_i + b) = 1$$



This mean that for the point(s) closest to the decision boundary the function marginal is equal to 1. All other point will have function marginal greater or equal to 1.

$$z_i(w^T x_i + b) \geq 1, i=1..n$$

With this new constraint (the function has the same result as before) the objective function became

$$\arg \max_{w, b} \frac{1}{\|w\|_2}$$

$$\text{s.t. } \begin{cases} z_i(w^T x_i + b) \geq 1, i=1..n \\ \min_i z_i(w^T x_i + b) = 1 \end{cases}$$

the standard distance from the closest point

became $\frac{1}{\|w\|_2}$
↓ we still want to maximize this

This is equivalent to:

$$\arg \min_{w, b} \|w\|_2$$

or for more practical way

$$\arg \min_{w, b} \frac{1}{2} \|w\|_2^2 \quad (2)$$

[the function is scale invariance]
Minimizing $\|w\|_2$ and minimizing $\|w\|_2^2$ give exactly the same result because the square function is monotonically increasing for positive value.
 $\|w\|_2 \leq \|w\|_2^2$, then $\|w\|_2^2 \leq \|w\|_2^2$

A Road Between Two Towns Analogy

Imagine you're designing a straight road between two towns. Each town has houses scattered around (these are your data points).

Town A represents the positive class (+1)

Town B represents the negative class (-1)

Your goal is to build a straight road (your decision boundary) that completely separates the two towns, with no houses from either town ending up on the wrong side of the road.

Existe un hyperplane qui permet de

Now, for safety reasons, you want this road to be as far away as possible from the houses in both towns. In other words, you want to maximize the distance from the road to the closest house from either town.

The Road Design Problem

Let's say the position and direction of your road is defined by the equation $w^T x + b = 0$.

The vector w determines the direction of the road

The term b shifts the road's position

The distance from any house to the road is given by:

$$|w^T x + b| / \|w\|$$

To maximize the safety buffer (margin) between your road and the closest houses, you need to:

Find a road that completely separates the towns (all houses are on the correct side)

Maximize the distance to the closest house

Why We Minimize w^2

Here's where mathematics helps us solve this problem efficiently:

First, we standardize our problem: we scale our road equation so that for the closest house, the value $|w^T x_i + b| = 1$. With this standardization, the actual distance from the closest house to the road becomes $1/\|w\|$.

To maximize this distance ($1/\|w\|$), we need to minimize $\|w\|$.

For mathematical convenience, we minimize $(1/2)\|w\|^2$ instead of just $\|w\|$, which gives us the same solution but with cleaner mathematics.

Why we can
remove that
constraint?

$$\begin{aligned} & \arg \min_{w \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 \\ \text{s.t. } & \begin{cases} z_i(w^T x_i + b) \geq 1, & i = 1, \dots, n \\ \min_i z_i(w^T x_i + b) = -1 \end{cases} \end{aligned}$$

SVD (110.07)
is minimized
at the margin

Indeed, an optimal solution of (2) will automatically satisfy

$$\min_i z_i(w^T x_i + b) = 1$$

If w^*, b^* was a solution for which $\min_i z_i(w^* x_i + b) = \psi > 1$, then we could build the solution $(w', b') = (\frac{w}{\psi}, \frac{b}{\psi})$ which would still satisfy all constraints, and have a lower norm, thus w^*, b^* would not be optimal. w^* was not true minimum

Let's imagine
 w^*, b^* are optimal solution
for this problem

$$\begin{aligned} & \arg \min_{w \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 \\ \text{s.t. } & z_i(w^T x_i + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

Scaling up
w reduces
norm but
satisfies all
constraints

Quadratic programming problem
VILAIN: ALSO DECISION
A VILLE GIVE A DECISION
PMLA 2

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

$$\min_i z_i(w^* x_i + b^*) = 1$$

$$\min_i z_i(w^T x_i + b) = 1$$

A half-space is one of the simplest examples of a convex set. It's the set of all points on one side of the hyperplane (a straight line in 2D, a flat plane in 3D...). In SVM each constraint is of the form

$$z_i(w^T x_i + b) \geq 1$$

For a specific data point (x_i, z_i) this constraint defines a half-space in the parameter space of (w, b) .

So to solve the problem we want to minimize $\frac{1}{2} \|w\|^2$ (to maximize the margin) but we had to ensure that all data points were correctly classified with a sufficient margin $(z_i(w^T x_i + b) \geq 1)$

This is what's known as a constrained optimization problem. We are trying to optimize one function while respecting certain boundaries or rules given by another set of functions.

We want to find the lowest point on a hill but we are allowed to walk only on certain paths.

The method of Lagrange multipliers gives us an optimal way to convert this constrained problem in an unconstrained one.

① For each constraint we introduce a new variable called a Lagrange multiplier (denoted by α_i), $\alpha_i \geq 0$

Note: Lagrange Multiplier with multiple constraints

$$\text{Optimizer: } L(w, b, \alpha)$$

$$\text{Constraints: } g_1(w, b, \alpha) = 0$$

$$g_2(w, b, \alpha) = 0$$

② We create a new function called the Lagrangian that combines our original objective function with the constraints, weighted by these multipliers.

For our SVM problem, the Lagrangian is:

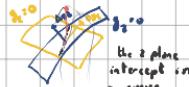
$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (z_i(w^T x_i + b) - 1)$$

the minus comes from how we formulate the constraint (following the
 - for $g_1(w, b) \geq 0$ we use $-\alpha_i g_1(w, b)$
 - for $g_2(w, b) = 0$ we use $\beta_i g_2(w, b)$ (normalizes size)

Lagrangian theory

Our constraints are $z_i(w^T x_i + b) - 1 \geq 0$

$\left\{ \begin{array}{l} \alpha_i = \lambda \alpha_1 + \lambda \alpha_2 \\ \alpha_1 = 0 \\ \alpha_2 = 0 \end{array} \right. \text{ this iteration}$
 $\alpha_1 = 0 \quad \text{curve is where the}$
 $\alpha_2 = 0 \quad \text{two constraint are true}$



Now, there are two different approaches to solve the problem that is finding the saddle point of $L(w, b, \alpha)$.

① Minimize L with respect to the primal variables w and b , while handling the Lagrange multipliers α_i in a specific way:

For each α_i (for each point) one of two conditions must be true:

- The derivative of L with respect to α_i equals zero
- Or multiplier its self $\alpha_i = 0$

When $\hat{z}_i > 0$, $z_i(w^T x_i + b) - 1 = 0$, the datapoint x_i lies exactly on the margin boundary, they are support vector.

When $\hat{z}_i = 0$, it means this data point does not influence the final decision boundary. It's doesn't matter for our solution, we don't need it. is not a support vector now consider see also DEFINITION OF A SUPPORT VECTOR

② We first ensure that the derivative of L with respect to w and b are zero then maximize L with respect to the Lagrange multipliers $\hat{\alpha}_i$

When L is expressed in function of $\hat{\alpha}_i$ thanks to the value obtained from $\nabla_w L(w, b, \hat{\alpha})$ and $\nabla_b L(w, b, \hat{\alpha})$

We will use the second approach:

Let's start by finding the minimum with respect to w and b by setting the derivatives equal to zero.

$$\begin{aligned} \cdot \nabla_w L(w, b, \hat{\alpha}) &= w - \sum_{i=1}^n \hat{\alpha}_i z_i x_i = 0 \\ w &= \sum_{i=1}^n \hat{\alpha}_i z_i x_i \end{aligned}$$

$$\cdot \frac{\partial L(w, b, \hat{\alpha})}{\partial b} = \sum_{i=1}^n \hat{\alpha}_i z_i = 0$$

We substitute this result back into the Lagrangian to eliminate the primal variables w, b we obtain a dual objective function

$$L_0(\hat{\alpha}) = \sum_{i=1}^n \hat{\alpha}_i + \frac{1}{2} \left\| \sum_{i=1}^n \hat{\alpha}_i z_i x_i \right\|^2$$

Let's focus on the second term

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

$$\begin{aligned} \left\| \sum_{i=1}^n \hat{\alpha}_i z_i x_i \right\|^2 &= \frac{1}{2} \left(\sum_{i=1}^n \hat{\alpha}_i z_i x_i \right)^T \left(\sum_{j=1}^n \hat{\alpha}_j z_j x_j \right) = \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \hat{\alpha}_i z_i x_i^T \hat{\alpha}_j z_j x_j = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \hat{\alpha}_i \hat{\alpha}_j z_i^T z_j x_i^T x_j \end{aligned}$$

$$L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j$$

s.t.

$$\alpha_i \geq 0$$

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

↳ an optimal solution can be expressed as a linear combination of the training set sample

in matrix form

$$L(\alpha) = \alpha^T L - \frac{1}{2} \alpha^T H \alpha$$

$$, L = [1, 1, \dots, 1]$$

where H is the matrix

$$H_{ij} = z_i z_j x_i^T x_j$$

$$, x^T H x = \sum_{i=1}^n \sum_{j=1}^n \alpha_i H_{ij} \alpha_j = \sum_{i=1}^n \sum_{j=1}^n \alpha_i z_i z_j x_i^T x_j$$

[We want to find $\max_{\alpha} L(\alpha)$]

Note that the dual form depends on the data only through dot products $x_i^T x_j$. This allows to extend SVMs to non-linear classification; we can replace those dot products with more complex kernel function without changing the structure of the optimization problem

Why we want to $\max_{\alpha} ?$??

↓
is a Lagrange rule for finding the optimal solution [if you want you can see the L0 as we write before all quadratic minima (minimum in convex domain) maximum to strong minimum is called LAGRANGE POSSIBLE, THE CONVERSE CON IL VALORE MASSIMO DELLA FUNZIONE

PISSA (w, b) è così numerico nella ZONE

α MOLTIPLICATI DA IDENTIFICARE I VETTORI DI SUPPORTO, E DETERMINARE QUANDO OLTRE VETTORI DI SUPPORTO CONSIDERARE ALLA DECISIONE FINALE. (ALCUNI PUNTI E' CATTICO MA NON SÌ: GRANDE VALORE DI CUI GE NE DOVETTORE IL COTTO MA NON $\alpha > 0$) E COME UNA FUNKTION NON CONSIDERAVANO AL CALCOLO DEL $w^T b$ DOPPIO

/

Poi se vogliamo $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (z_i (w^T x_i + b) - 1)$ NON CHE I MULIPPLICATORI DI VARIANTE FUNZIONE DI LIMITAZIONI, SONO SOLO Y PER RISPECTARE I CONSTRAINT (O CUI L'OTTIMIZZAZIONE SI FA PER IL LORO MASSIMO).

When we $\max_{\alpha} L(\alpha)$ we found α_i ; $\forall i$ i è un numero di sample

Summarize we have for SVM two problem

- ① Primal Problem: Finding w and b that minimize $L_p(w, b) = \frac{1}{2} \|w\|^2$, with constraint $\sum_i (w^T x_i + b) \geq 1$ for all training points.
- ② Dual Problem: Find α that maximize $L_d(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j$, with constraint $\alpha_i \geq 0$ for all i and $\sum_i \alpha_i z_i = 0$.

The dual and primal formulations are connected through the duality gap.

$$\text{Duality Gap} = L_p(w, b) - L_d(\alpha)$$

The difference between the primal objective value and the dual objective value

For any feasible w, b (meaning they satisfy the primal constraint) and any feasible α (meaning they satisfy the dual constraint), we always have:

$$L_d(\alpha) \leq L_p(w, b)$$

Th

$$L_p(w, b, z) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n z_i (w^T x_i + b) - \frac{1}{2}$$

So for any feasible primal solution (w, b) and feasible z

↑

The term $\sum_{i=1}^n z_i (w^T x_i + b) - \frac{1}{2}$ is always positive

↑ This means

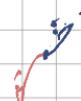
$$L_p(w, b, z) \leq \frac{1}{2} \|w\|^2 = L_p(w, b)$$

Furthermore the $L_d(\alpha)$ is the $\min L_p(w, b, z)$ so the feasible generate a value for sure

$$L_d(\alpha) \leq L_p(w, b, z) \leq \frac{1}{2} \|w\|^2$$

For convex optimization problems like SVMs, when both primal and dual problem reach their optimal solutions, the duality gap completely close.

$$L_d(\alpha^*) = L_p(w^*, b^*)$$



This is called "strong duality" and it means that solving either the primal or the dual problem will give us the same optimal value.

Observation ♦

The primal objective $\frac{1}{2} \|w\|^2$ is related to the margin of the classifier. Specifically the geometric margin is $\frac{1}{\|w\|}$, so minimising $\|w\|^2$ is equivalent to maximizing the margin.

The dual objective, $\sum_{i=1}^n z_i - \frac{1}{2} \sum_{i,j} z_i z_j x_i^T x_j$, can be interpreted as finding the most important point (support vector) that defines the margin.

The duality gap represents how far our current margin estimate (from the dual problem) is from the true maximum margin (in the primal problem). When the gap is zero we have identified exactly the right set of support vectors that define the maximum margin hyperplane.

- For a solution to be optimal, it must satisfy specific conditions called the Karush-Kuhn-Tucker (KKT) conditions. These aren't arbitrary rules - they're the mathematical requirements that any optimal solution must fulfill.

The KKT conditions provide the bridge between the primal and the dual formulations. They are both necessary and sufficient conditions for optimality in convex problems like SVMs.

Karush-Kuhn-Tucker (KKT) conditions²

Annotations on the KKT conditions:

- Primal**: w, b
- Feasible solution**: $\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i z_i x_i = 0$
- Feasible solution**: $\frac{\partial L(w, b, \alpha)}{\partial b} = -\sum_{i=1}^n \alpha_i z_i = 0$
- In margins & bias**: $\left\{ \begin{array}{l} z_i (w^T x_i + b) - 1 \geq 0 \\ \alpha_i \geq 0 \end{array} \right. \quad \forall i$
- Primal feasibility**: $\left[\alpha_i [z_i (w^T x_i + b) - 1] = 0 \right] \quad \forall i$
- SVM**: $\sum \alpha_i = 1$
- Complementary slackness**: $\alpha_i [z_i (w^T x_i + b) - 1] = 0 \quad \forall i$
- SVM**: $\sum \alpha_i = 1$

²These are both necessary and sufficient conditions for optimality of the SVM solution

The first and second equations encode that, at the optimal **primal** solution (w, b) , the gradient of the Lagrangian with respect to w and b becomes zero.

The next two equations encode that both the **primal** solution (w, b) and the corresponding **dual** solution α are feasible.

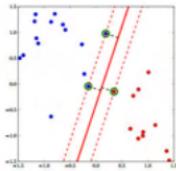
Complementary Slackness

$z_i (w^T x_i + b) - 1 = 0$ [refers to the first way of solving the problem]

For each datapoint, one of two things must be true:

- Either $\alpha_i = 0$
- Or $z_i (w^T x_i + b) = 1$ (the point lies exactly on the margin).

- For all points that do not lie on the margin (i.e. $z_i (\mathbf{w}^T \mathbf{x}_i + b) > 1$) the corresponding Lagrange multiplier is $\alpha_i = 0$
- $\alpha_i \neq 0$ implies that the corresponding point is on the margin, i.e., it is a **Support Vector**
- After we have obtained α , the KKT conditions allow us to estimate b



This means that in sum the decision boundary depends only on the support vectors. This is significant:

- When making predictions we only need to compute the dot products between the test point and the support vector not all (the training data points):

$$s(\mathbf{x}_t) = \mathbf{w}^T \mathbf{x}_t + b = \sum_{i=1}^{16384} z_i z_i^T \mathbf{x}_t + b$$

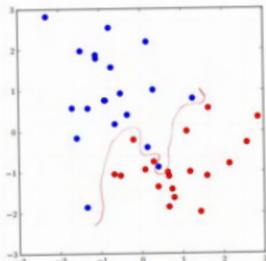
z_i is the set of support vector ($z_i \geq 0$)

The model focus only on the most informative points making it less sensitive to outliers.

To make a predictions if $s(\mathbf{x}_t) > 0$, the point belong to the positive class (H₁), if $s(\mathbf{x}_t) < 0$ it will be classified as belonging to the H₀ class.

Non linearly separable data

In many practical application, dataset aren't perfectly separable by a linear boundary. No matter how we position our hyperplane, some point will inevitably been in the wrong side or within the margin.



So No matter the value of w , some points will violate the constraint $z_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

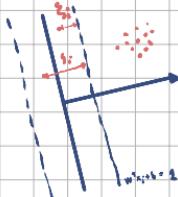
To handle this situation we introduce what are called "slack variables" ($\xi_i \geq 0$) for each training point. These variable measure how much a particular point violates the margin constraint.

We relax our original constraint into

$$\hat{z}_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

Geometrically:

- if $\xi_i = 0$ the point satisfies the original constraint, and is on the correct side of the margin.
- if $0 < \xi_i < 1$ the point is on the correct side of the decision boundary but inside the margin.
- if $\xi_i = 1$ The point is exactly on the decision boundary.
- if $\xi_i > 1$ The point is on the wrong side of the decision boundary (misclassified)



The point can be ξ_i inside the region

I want to keep the point that violate the constraint as few as possible
for which $\xi_i > 0$

we want to find a hyperplane that
does this

Ideally we want to minimize the number of points that violate the margin constraints, that is points with $\xi_i > 0$. This can be formulated mathematically as minimizing:

$$\sum_{i>0} I(\xi_i > 0)$$

Where I is the indicator function that equals 1 if its argument is true and 0 otherwise.
We will use

$$\Psi(\xi) = \sum_{i>0} \xi_i, \quad \xi > 0$$

For very small values of ξ , this function approximate the count of non-zero slack variables, effectively counting how many points violate the margin constraints.

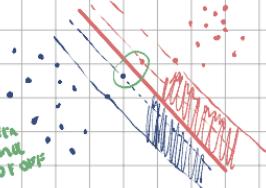
the one beyond the margin

If we remove these data points the classes become linearly separable and we can thus train a maximum margin hyperplane over the remaining points

Now combining our goal of maximizing the margin with minimizing constraint violations, we get the soft margin SVM objective:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \left(\sum_{i>0} \xi_i \right)$$

IT'S AN
HYPOTHESIS
THAT COMES
THIS TRADE OFF



If we remove all non-zero slack variables, we get a small margin

SE PARALLEL AL DAW
DE ESSERE POC
MARGINE AVIO
UN MARGINE
AMPIO MA POUCHE
POTER STONE

THAT'S LIKE CREATING

AND IMPROVE GENERALIZATION

Subject to: (z, t)

$$z_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

Where:

- F is a monotone convex function
- C is a positive constant that controls the trade-off between margin maximization and constraint violation minimization.
in fact:
 - A large C heavily penalizes constraint violations, leading to fewer margin violations but potentially smaller margin
 - A small C allows more margin violations in exchange for a potentially larger margin.

Note! Why C ?

C È UN MULTIPLOSORE DI VULCANO, DIFATTI QUANDO TRASFORMANO UN PROBLEMA VINCOLATO IN UNO NON VINCOLATO, AGGIUNGENDO DEL TERZO DI PENALITÀ MOLTIPLIATI PER COSTANTI. IL PARAMETRO C È PROPRIAMENTE QUESTA COSTANTE CHE DETERMINA QUANDO PRESE DANT'UNA PENALITÀ PER LE VIOLAZIONI RISpetto ALLA MAXIMIZZAZIONE DEL MARGINE.

QUESTO IPERPARAMESTRO È DISOUDI SCELTO DURANTE CROSS VALIDATION.

Note! Why F ?

La FUNZIONE F È INTRODOTTA CON UNA FUNZIONE CONVEXA NONLINEARE CHE APRE SULLA SOMMA DEI VARIABILI DI SVARO (ξ_i) ELEVATE ALLA POTENZA α . This function is a way to approximate il controllo delle violazioni del margine.

In TUTTA COG. OENO VOLGEMO MINIMIZZARE ESTREMALM. IL NUVLO DI PUNTI PEL CUI $\xi_i > 0$ (CHE CHIAMO IL MARGINE). TALAM SI UNA I M. QUFIA NON È DIFFERENZIABILE ^{concerno} E DA PROSPETTA A LIVELLO DI MINIMIZZAZIONE

F È UN APPROSSIMAZIONE THAT POLE DI QUESTA FUNZIONE DI CONTROLLO

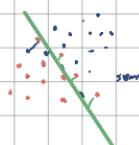
While the general formulation with arbitrary α and Function F provide theoretical flexibility, it leads to difficult optimization problem (for $\alpha \ll 1$ near 0 the problem becomes will be no more convex but...) So the standard simplification uses:

- $\alpha = 2$ (linear penalty on slack variables)
- $F(\mathbf{w}) = \|\mathbf{w}\|_2^2$ (identity function)

This gives us the practical soft margin SVM objective:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & z_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

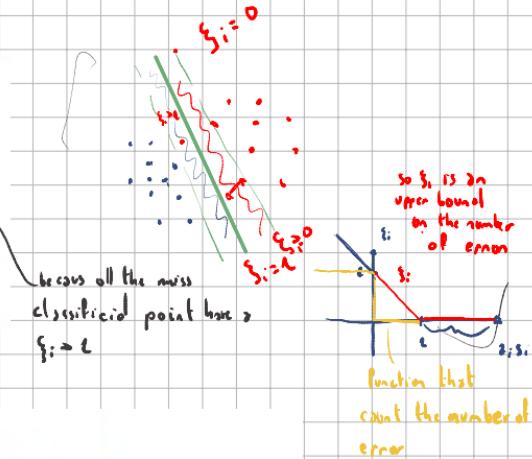
I want to minimize how much a single point is on the margin region



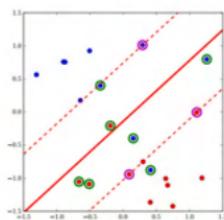
This is a convex quadratic programming problem

The terms ξ_i do not describe the number of errors, but act as a penalty

- Points inside the margin have $\xi_i > 0$, and miss-classified points have $\xi_i > 1$
- The further the point is from the hyperplane, the larger the value of ξ_i
- $\sum_{i=1}^n \xi_i$ is an upper bound on the number of miss-classified points (errors)
- C allows selecting a trade-off between margin and errors on the training set
- The solution is often called a *soft margin* hyperplane



Soft margin classification:



As we have seen with the hard-margin sum, we can solve the soft-margin sum using Lagrange multipliers. We will then find the partial derivative of $L(w, b, \xi, \mu)$ for w, b, ξ , then substitute the value to obtain $L(\alpha)$ and finally $\max_{\alpha} L(\alpha)$.

- $\alpha_i \geq 0$ for the constraint $z_i(w^T x_i + b) = 1 - \xi_i$
- $\mu_i \geq 0$ " " $\xi_i \geq 0$

$$L(w, b, \xi, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [z_i(w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^n \mu_i \xi_i.$$

The KKT conditions are

The necessary and sufficient KKT conditions are

$$\nabla_w L(w, b, \alpha, \mu) = w - \sum_{i=1}^n \alpha_i z_i x_i = 0$$

$$\frac{\partial L(w, b, \alpha, \mu)}{\partial b} = -\sum_{i=1}^n \alpha_i z_i = 0$$

$$\frac{\partial L(w, b, \alpha, \mu)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad \forall i$$

$$z_i(w^T x_i + b) - 1 + \xi_i \geq 0 \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

$$\alpha_i \geq 0 \quad \forall i$$

$$\mu_i \geq 0 \quad \forall i$$

$$\alpha_i [z_i(w^T x_i + b) - 1 + \xi_i] = 0 \quad \forall i \quad \text{non-zero unit test vs weight vector}$$

$$\mu_i \xi_i = 0 \quad \forall i$$

if $\alpha_i \neq 0$ then

$$z_i(w^T x_i + b) - 1 + \xi_i = 0$$

so we have two option

$$\alpha_i \neq 0 \rightarrow \alpha_i < C \Rightarrow \alpha_i > 0 \Rightarrow \xi_i = 0$$

$$\alpha_i = C \Rightarrow \xi_i = 0$$

$$w^T x_i + b = 1$$

$$\alpha_i = 0 \Rightarrow \xi_i \geq 0$$

From the KKT conditions (1)-(3)

$$w = \sum_{i=1}^n \alpha_i z_i x_i , \quad \sum_{i=1}^n \alpha_i z_i = 0 , \quad \alpha_i = C - \mu_i$$

Note that $\mu_i \geq 0$ implies that $\alpha_i \leq C$

We can then optimize the Lagrangian w.r.t. w , b and ξ to obtain the dual problem

$$\begin{aligned} L_0(\alpha) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j + \sum_{i=1}^n \alpha_i \xi_i + \sum_{i=1}^n \mu_i \xi_i + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n 2 \alpha_i \xi_i - \sum_{i=1}^n \mu_i \xi_i \\ &\quad - \sum_{i=1}^n \alpha_i z_i \left(\sum_{j=1}^n \alpha_j z_j x_j \right)^T x_i = \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j \end{aligned}$$

So the Dual problem is:

$$\max_{\alpha} L_0(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j$$

s.t

$$0 \leq \alpha_i \leq C \quad i=1, \dots, n$$

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

The only difference with the hard-margin is the upper bound on the Lagrange multipliers $\alpha_i \leq C$. The constraint on α_i 's in the soft-margin case have important implications.

① if $\alpha_i = 0$ the point is correctly classified and outside the margin. it doesn't influence the decision boundary

② if $0 < \alpha_i < C$ For complementary slackness, we know that:

- $z_i(w^T x_i + b) - \xi_i = 0$ (the point lies exactly on the margin)
- $\mu_i \xi_i = 0$ and since $\mu_i = C - \alpha_i > 0$ we must have $\xi_i = 0$

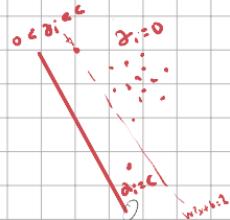
This point are the support vector

③ if $\alpha_i = C$ The point either:

- Inside the margin but correctly classified ($0 < \xi_i < \epsilon$)
- Exactly on the decision boundary ($\xi_i = 0$) ($w^T x_i + b = 0$)
- Misclassified ($\xi_i > \epsilon$)

These are also support vector but they are bounded or saturated at the upper limit.

inside the margin region
or misclassified



Once we have solved for the optimal α values, we make prediction using the score.

$$s(x_t) = w^T x_t + b = \sum_{i=1}^n \alpha_i z_i x_i^T x_t + b$$

the boundary depends
only on the support vector

Again depends only on the dot products $x_i \cdot x_j$

Several methods to solve the *dual* problem

Packages are available for several solvers

If we are interested in linear separation we can directly solve the *primal SVM* problem

Let's recall the primal formulation

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

we use only the primal formulation

s.t.

$$z_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

We can observe that for any data point, the optimal value of the slack variable ξ_i is:

- $\xi_i = 0$ if the point is correctly classified and outside the margin

- $\xi_i = 1 - z_i(w^T x_i + b)$, if the point violates the margin constraint

Putting these two cases together, for each data point, the optimal value of ξ_i is:

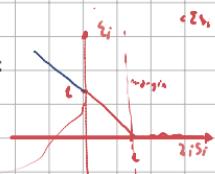
$$\xi_i = \max(0, 1 - z_i(w^T x_i + b))$$

The objective function then becomes:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(0, 1 - z_i(w^T x_i + b))$$

cost that increases linearly in

that way



This is an unconstrained optimization problem! We've effectively absorbed the constraint into the objective function through the max function, which is precisely the hinge loss

$$l(s) = \max(0, 1 - s) \quad \text{cost that I pay for those points inside of the margin}$$

$$l(s) = \{s\}$$

The hinge loss has a clear geometric interpretation:

- When a point is correctly classified outside the margin ($z_i f(x_i) > 1$), the loss is 0. The point doesn't contribute to the objective function.
- When a point is inside the margin or misclassified ($z_i f(x_i) \leq 1$), the loss grows linearly with the distance from the margin boundary. The further a point is on the wrong side, the more it "pulls" on the decision boundary.

Only the points that are close to or on the wrong side of the margin exert (nonzero) influence.

The hinge loss formulation reveals that SVMs fit into the broader framework of regularized empirical risk minimization:

Note C is obtained by cross Validation
I can multiply for λ in some way

$$\min_{w,b} \frac{\lambda}{2} \|w\|^2 + \underbrace{\frac{1}{n} \sum_{i=1}^n l(z_i, f(x_i))}_{\text{Regularization term}} + \underbrace{\frac{1}{n} \sum_{i=1}^n l(z_i, f(x_i))}_{\text{Empirical risk}}$$

Where $l(z_i, f(x_i)) = \max(0, 1 - z_i f(x_i))$ is the hinge loss

- The first term regularize the complexity of the model (favoring large margin)
- The second term ensure the model fit the training data.

$$\min_{w,b} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max[0, 1 - z_i(w^T x_i + b)]$$

Compare with Logistic Regression:

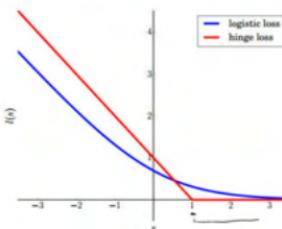
$$\min_{w,b} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \log \left[1 + e^{-z_i(w^T x_i + b)} \right]$$

LR: Logistic Loss

$$l(s) = \log(1 + e^{-s})$$

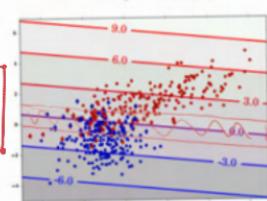
SVM: Hinge Loss

$$l(s) = \max(0, 1 - s)$$



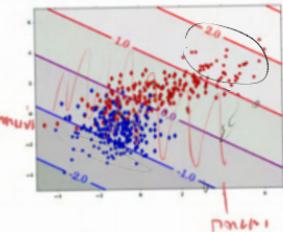
Binary classification:

We are trying to minimize the number of errors



$$C = 1.0 (\lambda = 0.0025)$$

$$C = 0.001 (\lambda = 2.5)$$



One of the most powerful aspect of the support vector machines is their ability to create non-linear decision boundaries through what's called "Kernel trick".

We have seen that we can obtain the separation hyperplane by solving either the primal or the dual problem.

Primal Formulation

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max[0, 1 - z_i(w^T x_i + b)]$$

Dual Formulation

$$\begin{aligned} \max_{\alpha} & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j z_i z_j x_i^T x_j \\ & = \alpha^T 1 - \frac{1}{2} \alpha^T H \alpha \end{aligned}$$

s.t.

$$0 \leq \alpha_i \leq C \quad \forall i$$

$$\sum \alpha_i z_i = 0 \quad \forall i$$

At the optimal solution this formulation yield the same result, with the duality gap closing to 0

$$L_p(w^*, b^*) = L_d(\alpha^*)$$

Parameters $w \in \mathbb{R}^D$ and $b \in \mathbb{R}$ (where D is the feature dimension)

Making prediction $s(x_t) = w^T x_t + b$ with complexity $O(D)$

Embedding a non-linear transformation requires explicitly defining the mapping from \mathbb{R}^D to the Hilbert space \mathcal{H}

- ANDROID FEATURE

$\Phi: \mathbb{R}^D \rightarrow \mathcal{H}$
 $w^T \Phi(x) + b$ as logistic regression
 we train the model in the expanded feature function

Parameters $\alpha \in \mathbb{R}^n$ (where n is the number of training samples)

Making predictions:

$$\begin{aligned} s(x_t) &= \sum_{i=1}^n \alpha_i z_i x_i^T x_t + b \\ &= \sum_{i: \alpha_i > 0} \alpha_i z_i x_i^T x_t + b \end{aligned}$$

with complexity $O(nv)$, where nv is the number of support vector

Embedding a non-linear transformation only requires dot-products in the expanded space

$$\Phi(x_i)^T \Phi(x_j)$$

The fundamental observation is that in the dual formulation, both training and prediction only required to compute the dot products between data points. We never need the transformed features explicitly - we only need their dot products.

If we have a function that efficiently computes the dot-products in the expanded space

$$k(x_1, x_2) = \Phi(x_1)^T \Phi(x_2)$$

then we can replace all the dot product in the dual problem [Remember that $H_{ij} = z_i z_j \Phi(x_i)^T \Phi(x_j)$ for mapping Φ]

Training: The Matrix H became: $H_{ij} = z_i z_j k(x_i, x_j)$

Prediction: The scoring function becomes: $s(x_t) = \sum_{i: \alpha_i > 0} \alpha_i z_i k(x_i, x_t) + b = \sum_{i: \alpha_i > 0} \alpha_i z_i k(x_i, x_t) + b$

Function k is called **kernel function**

A kernel function allows training a SVM in a large (even infinite) dimensional Hilbert space \mathcal{H} , without requiring to explicitly compute the mapping

Even if the expansion was finite, the complexity of the primal problem may be too large

On the other hand, the complexity of the dual problem only depends on the number of training points.

In practice, we are computing a linear separation surface in the expanded space, which corresponds to a non-linear separation surface in the original feature space

We have already seen an example of feature expansion that provides quadratic separation surface given by the mapping

$$x_i \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} x_1^2 x_2^2 x_3^2 \\ x_1 x_2 x_3 \\ x_1^2 x_2 x_3 \\ x_1 x_2^2 x_3 \\ x_1 x_2 x_3^2 \\ \vdots \\ x_1^d x_2^e x_3^f \end{bmatrix} \Phi(x) = \begin{bmatrix} \text{vec}(x x^T) \\ \sqrt{x} x \end{bmatrix} \text{ quadratic expansion}$$

$$\begin{aligned} \text{Note: } & \text{vec}(A)^T \text{vec}(B) \\ &= \text{Tr}(A^T \cdot B) \\ &= \langle A, B \rangle \end{aligned}$$

$$k(x_1, x_2) = (\text{vec}(x_1 x_2^T))^T (\text{vec}(x_2 x_1^T)) = \langle x_1 x_2^T, x_2 x_1^T \rangle + 2x_1^T x_2 + 1$$

$$\begin{aligned} & \langle x_1 x_2^T, x_2 x_1^T \rangle = \dots \text{dot product} \\ & = \text{Tr}((x_1 x_2^T)^T (x_2 x_1^T)) = / \\ & \rightarrow \text{Tr}(x_1 x_2^T x_2 x_1^T) = \frac{x_1 x_2^T \text{ after}}{x_1 x_2^T} \\ & = \text{Tr}((x_1^T x_2)(x_2 x_1^T)) = \\ & = (x_1^T x_2)^2 \end{aligned}$$

$$\Phi(x_1)^T \Phi(x_2) = (x_1^T x_2)^2 + 2x_1^T x_2 + 1 = (x_1^T x_2 + 1)^2 = k(x_1, x_2)$$

$\hookrightarrow O(d)$ complexity

So instead of doing this transformation explicitly we can simply use the Kernel function $k(x_1, x_2) = (x_1^T x_2 + 1)^2$ directly. This becomes especially valuable for higher-degree polynomials where the explicit feature space will grow exponentially with the degree.

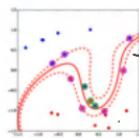
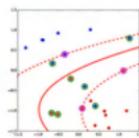
In general we can define the polynomial kernels of degree d as

$$k(x_1, x_2) = (x_1^T x_2 + 1)^d \quad O(d)$$

Note that the expanded space, although finite, grows very quickly

Poly — $d = 2, C = 1.0$

Poly — $d = 4, C = 10.0$



degree 4 polynomial hyperplane

If we know Φ , we can trivially define a corresponding kernel function $k(x_1, x_2) = \Phi(x_1)^T \Phi(x_2)$

However, we typically want to avoid defining explicit mappings

We would rather like to know which functions $k(x_1, x_2)$ correspond to a feature mapping Φ in a suitable space \mathcal{H}

Mercer's condition provides a sufficient condition for k to define a dot-product in an expanded space

Let $h(u, v)$ be a symmetric function in L^2 (the space of square integrable function)
 if for all function $g(u)$ such that $\int g(u)^2 du < \infty$, we have :

$$\int h(u, v) g(u) g(v) du dv > 0$$

[this are sufficient conditions]
ALSO IN A NON-SQUARE-INTEGRABLE CASE

then h defines a dot product in sum expanded space.

This mathematical condition essentially state that the Kernel function must be "positive definite" or more precisely "positive semi-definite". In practical term if we form a matrix H where $H_{ij} = h(x_i, x_j)$ for any set of point $\{x_1, \dots, x_n\}$ this matrix must have non-negative eigenvalues.

Mercer's condition is important because it guarantees that a function can be expressed as a dot product in some features space.

If a function satisfies Mercer's condition, there exist some feature mapping Φ such that :

$$h(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$$

Even if we never explicitly compute or define Φ .

However, Mercer's condition, doesn't provide a constructive way to build new kernels. Instead in practice

- we use well-established kernels that are known that satisfy the Mercer's condition.
- We apply operations that preserve the Kernel property to combine or modify existing Kernels (ex. Sum ($h_1 + h_2$) = $h_1(x_i, x_j) + h_2(x_i, x_j)$), Product $h_1 h_2$, scalar Multiplication, function application where f is function with a convergent Taylor series with non-negative coefficients.

/ Different with different
def. of Kernels

The "Gaussian Radial Basis Function Kernel" is one of the most widely used Kernel in practices :

$$h(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

it's corresponds to a mapping (Φ) into an infinite-dimensional feature space

Aproximación

$$h(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

This kernel has a remarkable property: it corresponds to a mapping into an infinite-dimensional feature space. Let's understand why this is the case.

The RBF kernel can be expanded using the Taylor series of the exponential function:

$$e^{-\gamma \|x_i - x_j\|^2} = e^{-\gamma \|x_i\|^2} \cdot e^{-\gamma \|x_j\|^2} \cdot e^{2x_i^T x_j}$$

This can be rewritten as:

$$e^{-\gamma \|x_i - x_j\|^2} = e^{-\gamma \|x_i\|^2} \cdot e^{-\gamma \|x_j\|^2} \cdot e^{2x_i^T x_j}$$

The term $e^{2x_i^T x_j}$ can be expanded as:

$$e^{2x_i^T x_j} = \sum_{n=0}^{\infty} \frac{(2\gamma)^n (x_i^T x_j)^n}{n!}$$

Each term $(x_i^T x_j)^n$ corresponds to a feature mapping that includes all possible products of n features. As we sum from $n = 0$ to ∞ , we're effectively working in an infinite-dimensional feature space.

So a RBF implicitly maps the data to an infinite-dimensional space where linear separations becomes much more likely, allowing it to create highly flexible decision boundaries

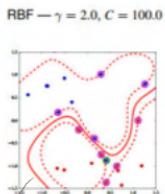
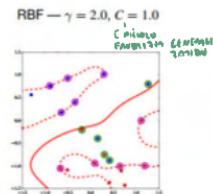
In the RBF the kernel value decrease with the distance between points. Point that are close in the input space have high kernel value (high similarity) while distant point have kernel value near 0.

The parameter γ controls the width of the kernel determining how far the influence of a single example reaches:

- Small γ : Wide Kernel, where each support vector influences point in a large region
- Large γ : Narrow Kernel, where the influence of support vectors is highly localized

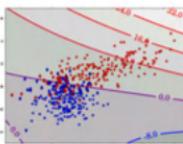
This mean that when γ is very large, the kernels becomes extremely narrow, and the SVM behaves similarly to a nearest neighbor classifier.

While when γ is very small, the kernel becomes wide and smooth, leading to nearly linear decision boundaries.

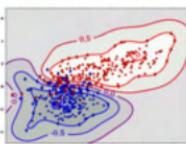


δ CONTROLS THE SMOOTHNESS

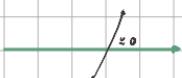
Poly — $d = 2, C = 1.0$



RBF — $\gamma = 0.5, C = 1.0$



$$\|x_i - x_k\| \leq \epsilon$$



$s(x) = \sum_i \alpha_i x_i^T x + b = 0$ this is the separation rule is control by the support vector
Accumulating more support vector means $x_i^T x + b$ is zero b

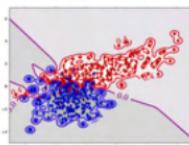
$\sum_i \alpha_i x_i^T x_j = 0$ if two point are close kernel is 2 if the are far is 0

1st point
Support Vector

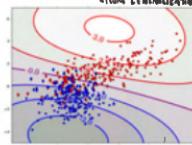


similarity
since the are closest
close to it
not support vector

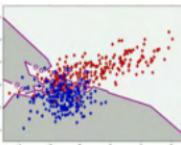
RBF (no bias) — $\gamma = 10, C = 1.0$



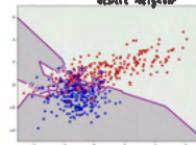
RBF (no bias) — $\gamma = 0.02, C = 1.0$
LARGE γ



RBF (no bias) — $\gamma = 200, C = 1.0$



1-NN
nearest neighbor



Practical Consideration

One of the most crucial aspect of applying SVMs effectively is selecting appropriate hyperparameters. The two main category of hyperparameters are:

① Kernel-specific parameters

Different kernel have different parameters that control their behavior:

- For polynomial kernels, the degree d and the constant c
- For RBF kernels: the width parameter γ (gamma)

The choice of γ in RBF kernels is particularly important:

- Too small γ : The kernel becomes too wide, creating overly smooth decision boundaries that may underfit the data
- Too large γ : The kernel becomes too narrow, potentially creating complex boundaries that overfit the training data

② Regularization parameter C

The parameter C controls the trade-off between maximizing the margin and minimizing classification errors:

• Small C : Prioritizes a wider margin, potentially at the cost of a smaller margin

• Large C : Focuses on reducing training errors, potentially at the cost of a smaller margin

Cross-validation is the recommended approach for selecting these hyperparameters. This typically involves:

① Splitting the training data into multiple folds

② Training models with different hyperparameter combinations

③ Selecting the combination that performs best on the validation sets

Despite the power of kernel SVMs, linear SVMs (using the linear kernel) are often sufficient and preferable in certain scenarios.

① High-dimensional data: When working with very high-dimensional feature space, the data often becomes nearly linearly separable. This phenomenon is related to the "curse of dimensionality" becoming a blessing for linear separability.

② Computational efficiency: Linear SVMs are significantly faster to train and evaluate. The scoring function is a simple dot product between the weight vector and the input features, making predictions extremely efficient. (mush)

SVMs are not invariant to affine transformation of the data, meaning that scaling and transformations of features can significantly impact performance. Proper preprocessing typically includes:

Centering: Subtracting the mean of each feature to center the data around the origin

Scaling: Normalizing features to have comparable ranges or standardizing them to have zero mean and unit variance

Whitening: Transforming the data to have uncorrelated features with unit variance

These preprocessing steps are particularly important for distance-based kernels like the RBF kernel, which are sensitive to the scale of the input features.

Unlike logistic regression, SVM scores do not have a direct probabilistic interpretation. The raw output of an SVM is simply the signed distance from the decision boundary, not a probability.

When probabilistic outputs are needed (for instance, to make risk-weighted decisions or combine SVM predictions with other probabilistic models), post-processing methods can be applied.

Platt scaling: Fitting a logistic regression model to transform SVM scores into probabilities

Isotonic regression: A more flexible non-parametric approach to calibrate scores to probabilities

These methods typically require a held-out validation set to avoid overfitting.

When working with imbalanced datasets (where one class is much more frequent than the other), standard SVMs may produce biased results favoring the majority class. There are several approaches to address this:

- ② **Class-weighted sum:** Assigning different C values to different classes to balance their importance.
- ③ **Sample-weighted sum:** Giving individual samples different weights based on their importance.

To re-balance the classes we can use a different value of C for the different classes (or even for different samples):

$$\arg \min_{w,b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n C_i [1 - z_i(w^T x_i + b)]_+ = \max_{w,b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n C_i \max[0, 1 - z_i(w^T x_i + b)]$$

The corresponding dual problem is⁵

$$\begin{aligned} & \arg \max_{\alpha} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T H \alpha \\ \text{s.t. } & \sum_{i=1}^n \alpha_i z_i = 0, \quad 0 \leq \alpha_i \leq C_i, \quad i = 1, \dots, n \end{aligned}$$

⁵In the laboratory we will omit the constraint related to the bias, since we will not explicitly model the bias term

For class balancing, we can set $C_i = C_T$ for samples of class \mathcal{H}_T and $C_i = C_F$ for samples of class \mathcal{H}_F .

We can select $C_T = \frac{\pi_T}{\pi_T^{emp}}$ and $C_F = \frac{\pi_F}{\pi_F^{emp}}$.

π_T^{emp} and π_F^{emp} are the empirical priors (i.e. sample proportions) for the two classes computed over the training set.

In practice, we are simulating the application class proportions at training time.

These are the target (desired) class probabilities

This approach effectively simulates the target class distribution during training, helping the model focus appropriately on both classes regardless of their representation in the training data.

MNIST — Comparison of average pairwise EER for different kernels (raw features)

Kernel	$C = 10$	$C = 1$	$C = 0.1$	LogReg ($\lambda = 1e^{-3}$)
Linear	1.6%	1.3% <i>wt. just right</i>	1.1% <i>(0.9%)</i>	1.2% <i>(0.9%)</i>
Poly ($d = 2$)	0.3%	0.3%	0.3%	—
Poly ($d = 3$)	0.4%	0.4%	0.4%	—
RBF ($\gamma = 0.2$)	0.7%	0.7%	0.9%	—
RBF ($\gamma = 1.0$)	0.5%	0.5%	0.6%	—
RBF ($\gamma = 2.0$)	0.1%	0.1%	0.1%	—

we need few to reduce dimension to 50

⁶Quadratic expansion of PCA-reduced features

Kernel Choice Matters: Nonlinear kernels significantly outperform linear kernels on this image recognition task. This suggests that the relationship between pixel values and digit classes is inherently nonlinear.

RBF Kernel with $\gamma = 2.0$ Achieves Best Performance: This configuration achieves an impressive 0.1% error rate across different C values, showing that with the right kernel parameters, SVMs can be extremely effective for image classification.

Regularization Impact Varies by Kernel: For the linear kernel, the choice of C makes a noticeable difference (error rates ranging from 1.1% to 1.6%). However, for nonlinear kernels, performance is much more stable across different C values. This suggests that once the right kernel is selected, the exact regularization strength becomes less critical.

Comparison with Logistic Regression: SVMs with nonlinear kernels consistently outperform logistic regression on this task. Even with quadratic feature expansion, logistic regression achieves only 0.9% error rate, compared to 0.1% for the best SVM configuration.

Extending SVMs to Multi-class Problem

① One-versus-All (OVA) Approach

This is the most common strategy:

- Train N different binary SVMs, each distinguishing one class from all others.
- For a new sample, evaluate all N classifier.
- Assign the sample to the class whose classifier produce the highest score.

While simple to implement, this approach may be affected by imbalance issues (since each classifier sees many negative examples but few positive ones) and may produce ambiguous regions where multiple classifiers give conflicting decisions.

② One-versus-One (OVO) Approach

It handle multiclass classification by breaking it down into multiple binary classification problems - specifically one binary classifier for each pair of classes.

- Training Phase

- For a problem with N classes, we train $N(N-1)/2$ binary SVM classifiers.
- Each classifier is trained to distinguish between just two classes, ignoring all other classes.

For example, with 10 classes (like digits 0-9) we would train 45 different binary classifiers (one for each pair: 0-vs-1, 0-vs-2, ... 8-vs-9)

- Production Phase

For a new sample, we run it through all $N(N-1)/2$ classifiers, each classifier votes for one of the two classes it was trained on, then we count the votes received by each class.

The sample is assigned to the class that received the most votes.

Another possible way to operate is the Direct Multiclass Formulation

The direct multiclass formulation extends the SVM concept to handle all classes simultaneously in a single optimization problem, rather than breaking it down into multiple binary problems.

- Introduce an objective function that maximizes the margin between each class and all the others. For example, (we ignore the bias terms)

$$\min_{\mathbf{W}} \frac{1}{2} \|\mathbf{W}\|_2^2 + C \sum_{i=1}^n \max_{y'} [(\mathbf{w}_{y'} - \mathbf{w}_{z_i})^T \mathbf{x}_i + \delta_{y', z_i}]$$

Compare with multiclass logistic regression

$$\min_{\mathbf{W}} \frac{\lambda}{2} \|\mathbf{W}\|_2^2 + \frac{1}{n} \sum_i \log \left(\sum_{y'} e^{(\mathbf{w}_{y'} - \mathbf{w}_{z_i})^T \mathbf{x}_i} \right)$$

- $(\mathbf{w}_{y'})^T \mathbf{x}_i$ is the score for some class y' .
- $(\mathbf{w}_{z_i})^T \mathbf{x}_i$ is the score for the true class z_i .
- $(\mathbf{w}_{y'} - \mathbf{w}_{z_i})^T \mathbf{x}_i$ is the difference between these two scores.
- δ_{y', z_i} equals 1 when $y' = z_i$ (same class) and 0 otherwise.

The term $\max_y [(\mathbf{w}_y - \mathbf{w}_{z_i})^T \mathbf{x}_i + \delta_{y', z_i}]$ is calculating the loss based on the most problematic class.

Note:-

"The most competing/problematic class, is the incorrect class that receives the highest score for a given input example/training point."

In multi-class classification, each class has its own weight vector (decision function). When we classify an example, we compute a score for each possible class and typically assign the example to the class with the highest score.

For an input example \mathbf{x}_i with true class z_i , we would calculate scores for all possible classes:

- Score for class 1: $\mathbf{w}_1^T \mathbf{x}_i$
- Score for class 2: $\mathbf{w}_2^T \mathbf{x}_i$
- Score for class 3: $\mathbf{w}_3^T \mathbf{x}_i$
- And so on...
I

Ideally, we want the score for the true class ($\mathbf{w}_{z_i}^T \mathbf{x}_i$) to be higher than the scores for all other classes. But among all the incorrect classes, one of them will have the highest score. This is what I'm calling "the most competing class" - it's the incorrect class that is most likely to be confused with the true class for that particular example.

The max of finds the class that maximizes this expression; in other words, it finds the class that most violates our desire to have the true class score higher than all other class scores by some margin.

The SVM is specially focused on increasing the margin between the true class and the most competing class rather than warning about classes that are already easily distinguished.

the incorrect class with the higher score

Kernel methods are not restricted to SVMs

- Logistic Regression can be extended to incorporate kernels
- In general, we can apply the kernel tricks to problems which depend only on dot-products
 - Kernel PCA
 - Gaussian Processes
 - ...