Today   → Maximum Margin Classifier / Support vector

                                                Machines

Supervised learning
- → Max min formulation
- → Constrained formulation
- → Lagrangian functor (KKT conditions)
- → Solution through the Hinge loss

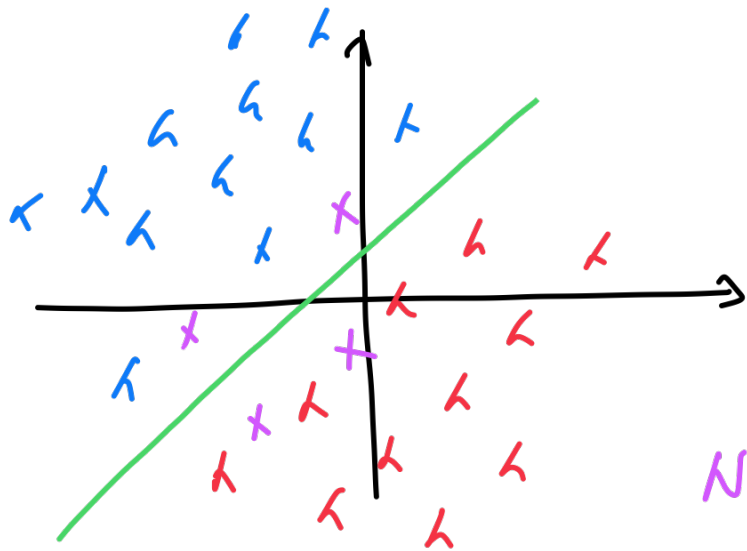Unsupervised learning
- → Clustering → Combinatorial approaches (including K means / K medoid)
- → Bump Hunting approaches (including A priori algorithm)

# Maximum Margin Classifier

Question: Given a binary classification problem
when should one position our classifier for it to be as
robust as possible toward new (unknown) examples?



One idea is to try to position the discriminant right in between the 2 classes.
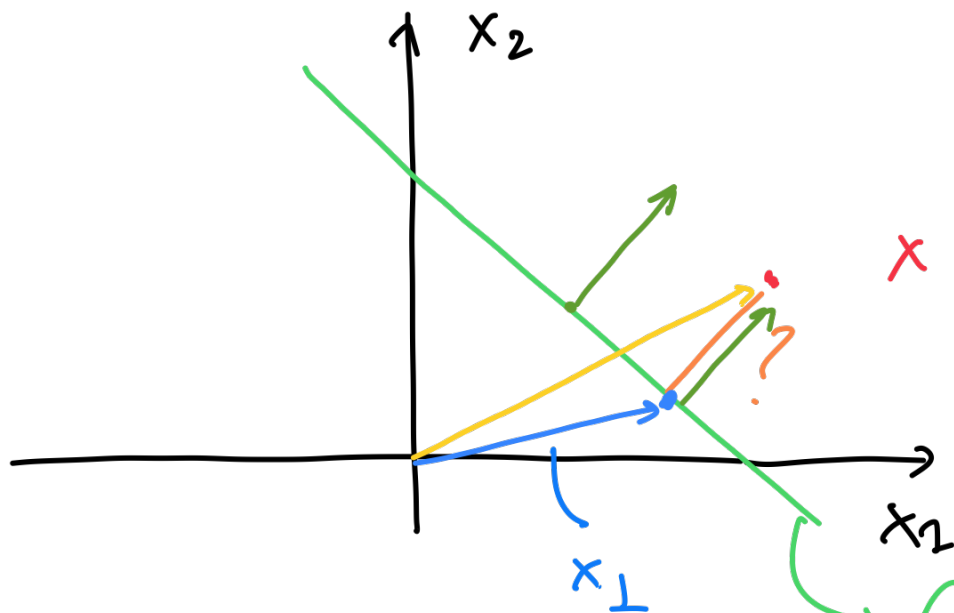
New examples are likely to be "close" to their corresponding classes and such a position is thus likely going to reduce the misclassification risk.

→ Putting the discriminant right in between the 2 classes
is equivalent to maximizing the distance between the
discriminant and the 2 classes

Which can be done by maximizing the distance between the
discriminant and the closest point from each class.

How can one express this distance ?

$$X = X_\perp + r \frac{\vec{\beta}}{\|\beta\|} \qquad (\ast)$$



For any plane of the
form $\beta_1 x_1 + \beta_2 x_2 + \beta_0 = 0$
the normal to the $y(x)$
plane is $[\beta_1, \beta_2]$

$$y(x) = \vec{\beta}^T x + \beta_0 \qquad \vec{\beta} = [\beta_1, \beta_2]$$

substituting (*)

$$y(x) = \vec{\beta}^T x + \beta_0$$

y(x) is prediction from the plane

$$= \vec{\beta}^T x_\perp + \vec{\beta}^T r \frac{\beta^T}{\|\beta\|} + \beta_0$$

$$= \boxed{\vec{\beta}^T x_\perp + \frac{\|\beta\|^2}{\|\beta\|} r + \beta_0}$$

$$= 0 \text{ as } x_\perp \text{ belongs to the plane } \beta_1 x_1 + \beta_2 x_2 + \beta_0 = 0$$

We are then left with

$$y(x) = \|\beta\| r$$

SIGNED DISTANCE

$$r = \frac{y(x)}{\|\beta\|} \quad (*)$$

( y(x) can be both positive or negative depending on whether x lies above or below the plane.

if we consider a classification problem where the points above the plane are given a +1 target and the points below are given a -1 target, then the "unsigned" distance of a point $x$ to the plane is given by multiplying $r$ by the target $t^{(i)}$ of $x$

$$\text{dist}(x, \text{plane}) = \underbrace{r}_{(*)} t(x) = \frac{y(x)}{\|\beta\|} t(x)$$

Given this expression for the distance, in order to learn a robust classifier, we can just find the points that are the closest to the plane and then "push" the plane as far as possible from these points.

Correctly

$$\beta^*, \beta_0^* = \max_{\beta, \beta_0} \min_x \frac{y(x)\, t(x)}{\|\beta\|}$$

finding the distance to the closest point

$$\mathbb{1} \text{ here } \vec{\beta} = (\beta_1, \beta_2)$$

(not including $\beta_0$)

push the place as far as possible away from that point

MAXIMUM MARGIN CLASSIFIER

$$\beta^*, \beta_0^* = \underset{\vec{\beta}, \beta_0}{\arg\max} \; \min_x \frac{(\beta_1 x_1 + \beta_2 x_2 + \beta_0)\, t(x)}{\|(\beta_1, \beta_2)\|}$$

$(**)$

Note that for a 2D space $(x_1, x_2)$
the plane is the set of points $(x_1, x_2)$ satisfying $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$
which can also read as $x_2 = -\dfrac{\beta_1 x_1 - \beta_0}{\beta_2}$

(**) is not easy to solve as (1) it involves a ratio of
functions in the optimization
variables

(2) it involves competing
maximization and minimization
processes

In order to simplify this formulation we can however use the fact that the ratio that appears in the objective in (**) is invariant under any rescaling of $\vec{\beta}$

letting $\beta = \alpha \beta$, we get

$$\frac{(\alpha \beta_1 x_2 + \alpha \beta_2 x_2 + \alpha \beta_0) \, t(x)}{\| (\alpha \beta_1 , \alpha \beta_2) \|}$$

$$= \frac{\alpha (\beta_1 x_1 + \beta_2 x_2 + \beta_0) \, t(x)}{\alpha \| (\beta_1, \beta_2) \|}$$

in particular this means we can fix the scaling $\alpha$ and optimize over all vectors $\beta$ satisfying that fixed scaling

One possible choice is to choose the scaling such that for the closest point to the plane, we get

$$(\beta_0 + \beta_1 x_1 + \beta_2 x_2) \, t(x) = 1$$

In other words, we fix the Scaling ambiguity by requiring that for the closest point to the plane, the distance to the plane should be $\frac{1}{\|\beta\|}$ (or equivalently the numerator should be 1)

Following from this choice, all the training examples must then satisfy

$$(\beta_0 + \beta_1 x_1 + \beta_2 x_2) \, t(x) \geq 1$$

Our original formulation now reduces to a constrained formulation (a.k.a Canonical formulation)

$$(*) \quad \begin{array}{c} \text{argmax} \\ \vec{\beta}, \beta_0 \end{array} \quad \boxed{\frac{1}{\|\beta\|}}$$

DISTANCE OF CLOSEST POINT TO THE PLANE

$$\text{s.t} \quad (\beta_0 + \beta_1 x_1 + \beta_2 x_2) \, t(x) \geq 1$$

LINEAR CONSTRAINTS ( all the training examples must be at a distance larger then $\frac{1}{\|\beta\|}$ (closest distance)

CANONICAL FORMULATION

$(*)$ can be further simplified by squaring the $\|\beta\|$ and replacing the maximization with a minimization

$$\beta^*, \beta_0^* = \begin{array}{c} \text{argmin} \\ \vec{\beta}, \beta_0 \end{array} \quad \|\beta\|^2$$

$(*)$

$$\text{s.t} \quad (\beta_0 + \beta_1 x_1 + \beta_2 x_2) \, t(x) \geq 1$$

for all $x$ in training

Given (∗) we can analyze the solution(s) by relying on the Lagrangian function and the associated Karush-Kuhn-Tucker conditions
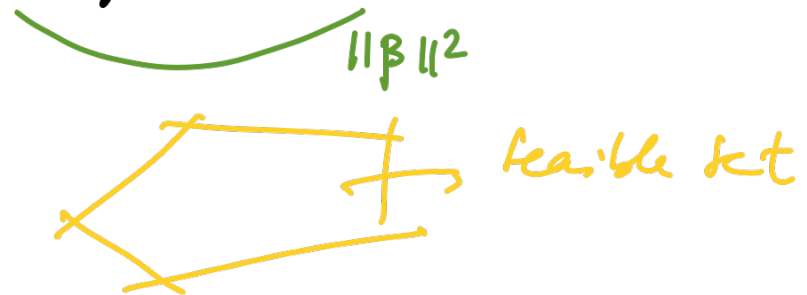
Given the constrained problem (∗) the lagrangian is defined by introducing a set of (lagrange) multipliers (one for each constraint) $\lambda_1, \lambda_2, \ldots, \lambda_N \geq 0$ as (L)

$$\mathcal{L}(\bar{\beta}, \beta_0, \lambda) = \|\beta\|^2 - \sum_{i=2}^{N} \lambda_i \left( \left( \beta_0 + \beta_1 x_2^{(i)} + \beta_2 x_2^{(i)} \right) t(x) - 1 \right)$$

Under appropriate conditions our Original Constrained formulation (∗) is equivalent to

$$\arg\min_{\bar{\beta}, \beta_0} \max_{\lambda} \mathcal{L}(\bar{\beta}, \beta_0, \lambda)$$

$\|\beta\|^2$

feasible set

**(A)** For any $\beta$ which violate at least one of the constraints we have $((\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) t(x) - 1) < 0$

and hence $-\lambda_i ((\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) t(x) - 1) \geq 0$

and can be increased to $+\infty$ by growing $\lambda_i$

in this **(A)** setting we thus have $\max_\lambda L(\vec{\beta}, \beta_0, \lambda) = +\infty$

**(B)** On the contrary, for a $\beta$ that satisfies all the constraints (i.e which lies inside the feasible set) we have

$$((\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) t(x) - 1) > 0$$

hence $-\lambda_i ((\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) t(x) - 1) < 0$

and the best the max over $\lambda_i$ can do is to bring that quantity to zero

in this Ⓑ setting, we thus have

$$\max_{\lambda} \mathcal{L}(\vec{\beta}, \beta_0, \lambda) = \|\beta\|^2$$

Grouping Ⓐ + Ⓑ we see that

$$\arg\min_{\vec{\beta}, \beta_0} \max_{\lambda_i} \mathcal{L}(\vec{\beta}, \beta_0, \lambda) = \arg\min_{\vec{\beta}, \beta_0} \begin{cases} +\infty & \text{if } \vec{\beta}, \beta_0 \text{ are not feasible} \\ \|\beta\|^2 & \text{if } (\beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)}) \, t(x^{(i)}) \geq 1 \\ & \forall \ x^{(i)} \end{cases}$$

which is exactly equivalent to
our original (i.e constrained) formulation

Now that we have the lagrangian function we can turn to the following theorem (known as the Karush–Kuhn–Tucker conditions) which connect the solutions of the constrained problem with the lagrangian function $L$.

**Theorem** (KKT Conditions for linearly constrained problems)

Consider a constrained problem

$$(*) \begin{cases} \min \ f(\beta) \\[2mm] \text{s.t} \quad a_i^T \beta \le b_i \quad i = 1, \dots, N \end{cases}$$

Assume that $f(\beta)$ is convex and continuously differentiable over $\mathbb{R}^{\delta+1}$. Let $\beta^*$ be a feasible point for $(*)$ (i.e $\beta^*$ satisfies $a_i^T \beta^* \le b_i \ \forall i$)

then $\beta^*$ is an optimal solution of the constrained problem (*) if and only if there exists a set of multipliers $\lambda_1^*, ..., \lambda_N^* \geq 0$ such that

$$\nabla f(\beta^*) + \sum_{i=1}^{N} \lambda_i^* \vec{a}_i = 0 \quad \text{and} \quad \lambda_i^*(a_i^T \beta - b_i) = 0$$
$$\forall \; i \in \{1, .., N\}$$

Using the KKT Conditions above, and applying those conditions to (*) we find that our classifier has to satisfy the equation

$$\frac{\partial L}{\partial \beta_1} = 2\beta_1 - \sum_{i=1}^{N} \lambda_i^* t(x^{(i)}) x_1^{(i)} = 0 \quad *$$

$$\frac{\partial L}{\partial \beta_2} = 2\beta_2 - \sum_{i=1}^{N} \lambda_i^* t(x^{(i)}) x_2^{(i)} = 0 \quad **$$

Compactly we thus have
$$\nabla_{(\beta_1, \beta_2)} L = 2\beta - \sum_{i=1}^{N} \lambda_i^* \bar{x}^{(i)} t(x^{(i)})$$

Doing the same for $\beta_0$ ( which does not appear in $\|\beta\|^2$ )

we get
$$\frac{\partial L}{\partial \beta_0} = -\sum_{i=1}^{N} \lambda_i^* t(x^{(i)}) \cdot 1 = 0 \qquad \text{***}$$

Solving (*), (**) & (***) for $\beta_1, \beta_2$ and substituting

in the lagrangian Ⓛ we get a formulation ( maximization

problem ) that only depends on the $\lambda_i$'s known as the

dual ( as opposed to the primal formulation (*))

this formulation is a quadratic optimization problem

Both the primal and dual formulations can be solved efficiently (e.g. through CVX see recitation).