



SOICT

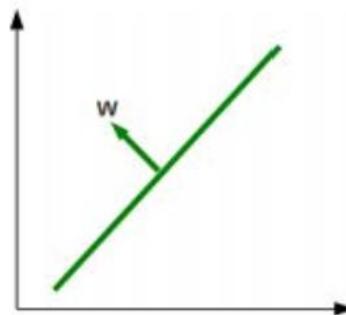
HUST
ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

Perceptrons and SVMs

Dam Quang Tuan

Hyperplane

- Separates a D -dimensional space into two half-spaces (positive and negative)
- Defined by a normal vector $\mathbf{w} \in \mathbb{R}^D$ (pointing towards positive half-space)



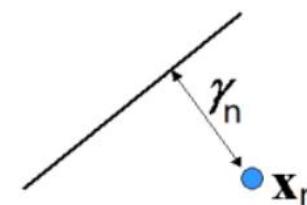
$b > 0$ means moving $\mathbf{w}^\top \mathbf{x} = 0$ along the direction of \mathbf{w} ; $b < 0$ means in opp. dir.

$$\mathbf{w}^\top \mathbf{x} + b = 0$$

- Equation of the hyperplane: $\mathbf{w}^\top \mathbf{x} = 0$
- Assumption: The hyperplane passes through origin. If not, add a bias term b
- Distance of a point \mathbf{x}_n from a hyperplane $\mathbf{w}^\top \mathbf{x} + b = 0$

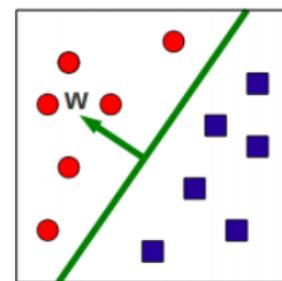
Can be positive or negative

$$\gamma_n = \frac{\mathbf{w}^\top \mathbf{x}_n + b}{\|\mathbf{w}\|}$$



Hyperplane based (binary) classification

- Basic idea: Learn to separate two classes by a hyperplane $\mathbf{w}^\top \mathbf{x} + b = 0$



Prediction Rule

$$y_* = \text{sign}(\mathbf{w}^\top \mathbf{x}_* + b)$$

For multi-class classification with hyperplanes, there will be multiple hyperplanes (e.g., one for each pair of classes); more on this later



- The hyperplane may be “implied” by the model, or learned directly
 - Implied: Prototype-based classification, nearest neighbors, generative classification, etc
 - Directly learned: Logistic regression, Perceptron, Support Vector Machine (SVM), etc
- The “direct” approach defines a model with params \mathbf{w} (and optionally a bias param b)
 - The parameters are learned by optimizing a classification loss function (will soon see examples)
 - These are also discriminative approaches – \mathbf{x} is not modeled but treated as fixed (given)
- The hyperplane need not be linear (e.g., can be made nonlinear using kernels; later)

Loss Functions for Classification

- In regression (assuming linear model $\hat{y} = \mathbf{w}^\top \mathbf{x}$), some common loss fn

$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

$$\ell(y, \hat{y}) = |y - \hat{y}|$$

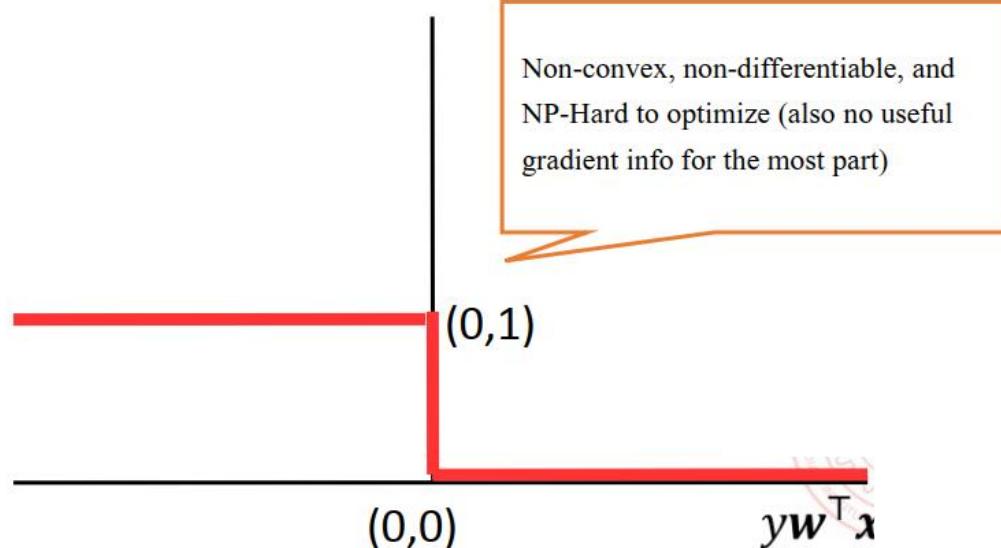
- These measure the difference between the true output and model's prediction
- What about loss functions for classification where $\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x})$?
- Perhaps the most natural classification loss function would be a "0-1 Loss"

- Loss = 1 if $\hat{y} \neq y$ and Loss = 0 if $\hat{y} = y$.
- Assuming labels as +1/-1, it means

$$\ell(y, \hat{y}) = \begin{cases} 1 & \text{if } y\mathbf{w}^\top \mathbf{x} < 0 \\ 0 & \text{if } y\mathbf{w}^\top \mathbf{x} \geq 0 \end{cases}$$

Same as $\mathbb{I}[y\mathbf{w}^\top \mathbf{x} < 0]$ or $\mathbb{I}[\text{sign}(\mathbf{w}^\top \mathbf{x}) \neq y]$

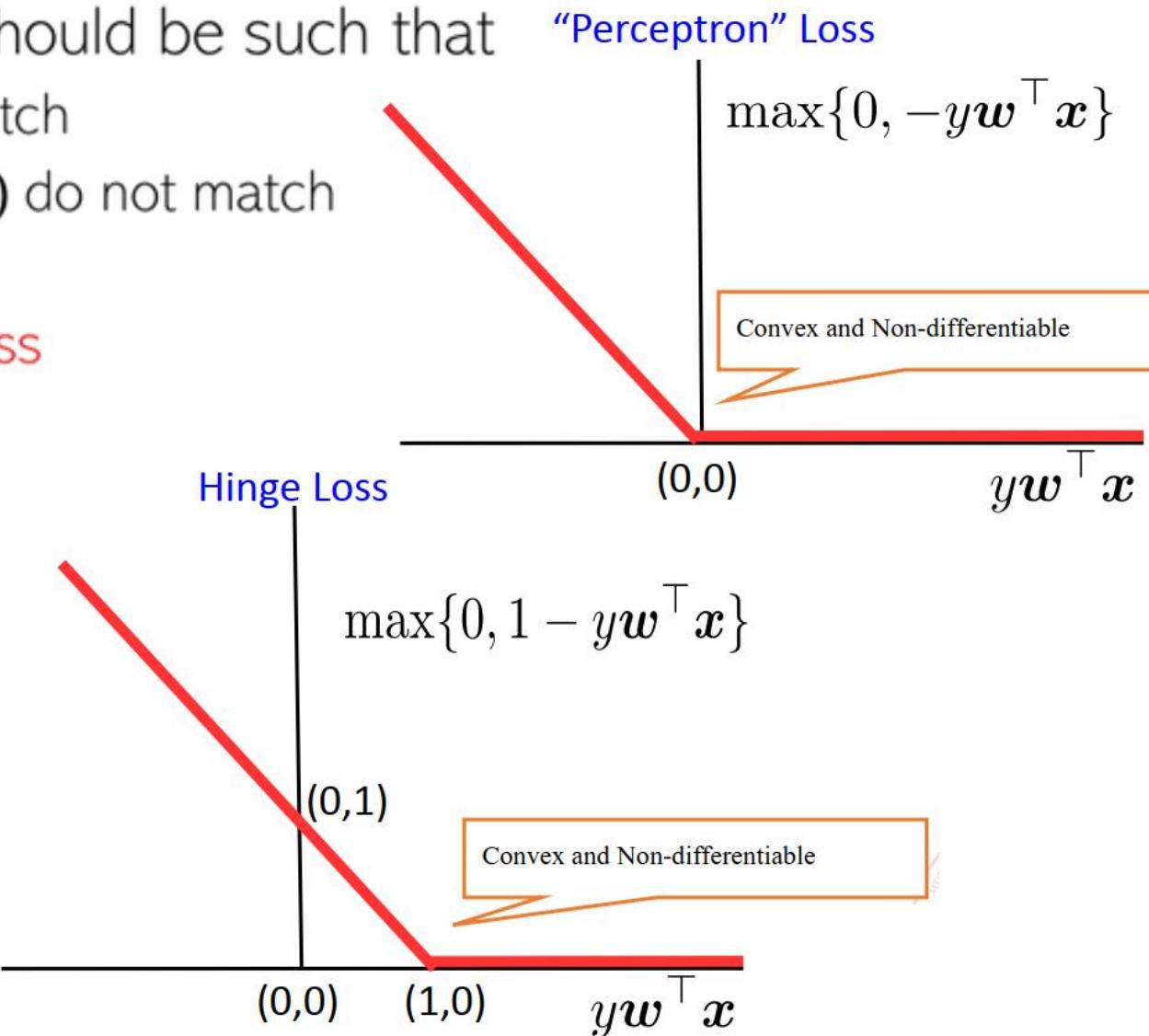
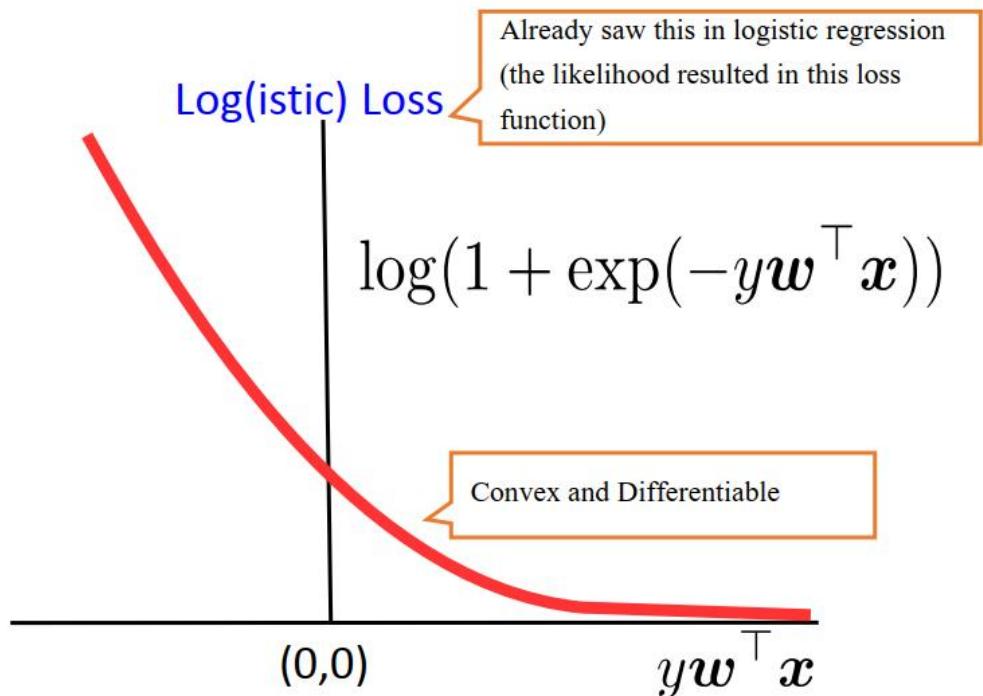
0-1 Loss



Loss Functions for Classification

- An ideal loss function for classification should be such that

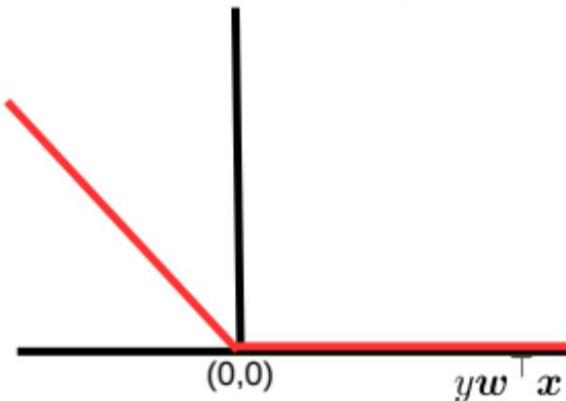
- Loss is small/zero if y and $\text{sign}(\mathbf{w}^\top \mathbf{x})$ match
- Loss is large/non-zero if y and $\text{sign}(\mathbf{w}^\top \mathbf{x})$ do not match
- Large positive $y\mathbf{w}^\top \mathbf{x} \Rightarrow$ small/zero loss**
- Large negative $y\mathbf{w}^\top \mathbf{x} \Rightarrow$ large/non-zero loss**



Learning by Optimizing Perceptron Loss

- Let's ignore the bias term b for now. So the hyperplane is simply $\mathbf{w}^\top \mathbf{x} = 0$
- The Perceptron loss function: $L(w) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^\top \mathbf{x}_n\}$. Let's do SGD

"Perceptron" Loss: $\max\{0, -y \mathbf{w}^\top \mathbf{x}\}$



Subgradients w.r.t. \mathbf{w}

$$\mathbf{g}_n = \begin{cases} 0, & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n > 0 \\ -y_n \mathbf{x}_n, & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n < 0 \\ \mathbf{k} y_n \mathbf{x}_n & \text{for } y_n \mathbf{w}^\top \mathbf{x}_n = 0 \quad (\text{where } k \in [-1, 0]) \end{cases}$$

One randomly chosen example in each iteration

- If we use $k = 0$ then $\mathbf{g}_n = \mathbf{0}$ for $y_n \mathbf{w}^\top \mathbf{x}_n \geq 0$, and $\mathbf{g}_n = -y_n \mathbf{x}_n$ for $y_n \mathbf{w}^\top \mathbf{x}_n < 0$
- Non-zero gradients only when the model makes a mistake on current example (\mathbf{x}_n, y_n)
- Thus SGD will update \mathbf{w} only when there is a mistake (mistake-driven learning)

The Perceptron Algorithm

- Stochastic Sub-grad desc on Perceptron loss is also known as the Perceptron algorithm

Stochastic SubGD

Note: An example may get chosen several times during the entire run

- ① Initialize $\mathbf{w} = \mathbf{w}^{(0)}$, $t = 0$, set $\eta_t = 1, \forall t$
- ② Pick some (\mathbf{x}_n, y_n) randomly.
- ③ If current \mathbf{w} makes a **mistake** on (\mathbf{x}_n, y_n) , i.e., $y_n \mathbf{w}^{(t)^\top} \mathbf{x}_n < 0$

Mistake condition

$$\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + y_n \mathbf{x}_n \\ t &= t + 1\end{aligned}$$
- ④ If not converged, go to step 2.

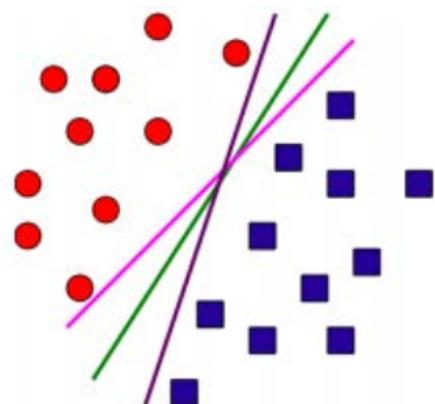
Updates are “corrective”: If $y_n = +1$ and $\mathbf{w}^\top \mathbf{x}_n < 0$, after the update $\mathbf{w}^\top \mathbf{x}_n$ will be less negative. Likewise, if $y_n = -1$ and $\mathbf{w}^\top \mathbf{x}_n > 0$, after the update $\mathbf{w}^\top \mathbf{x}_n$ will be less positive

If training data is linearly separable, the Perceptron algo will converge in a finite number of iterations
(Block & Novikoff theorem)

- An example of an **online learning** algorithm (processes one training ex. at a time)
 - Assuming $\mathbf{w}^{(0)} = 0$, easy to see that the final \mathbf{w} has the form $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$
 - α_n is total number of mistakes made by the algorithm on example (\mathbf{x}_n, y_n)
 - As we'll see, many other models also have weights \mathbf{w} in the form $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$
- Meaning of α_n may be different

Perceptron and (lack of) Margins

- Perceptron would learn a hyperplane (of many possible) that separates the classes



Basically, it will learn the hyperplane which corresponds to the \mathbf{w} that minimizes the Perceptron loss

Kind of an “unsafe” situation to have – ideally would like it to be reasonably away from closest training examples from either class

- Doesn't guarantee any “margin” around the hyperplane

- The hyperplane can get arbitrarily close to some training example(s) on either side
- This may not be good for generalization performance

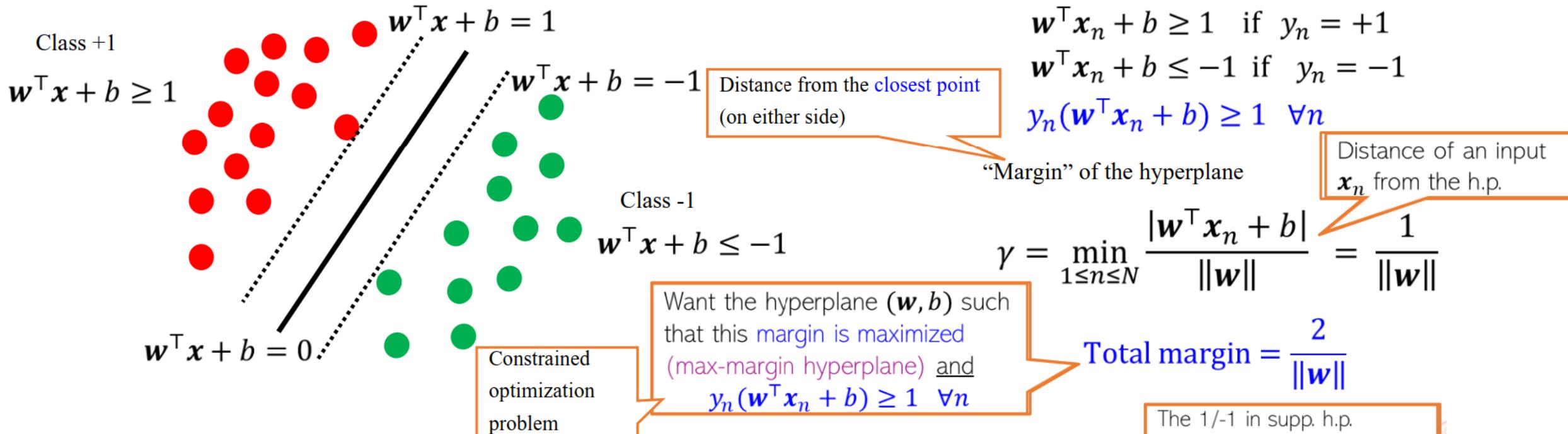
$\gamma > 0$ is some pre-specified margin

- Can artificially introduce margin by changing the mistake condition to $y_n \mathbf{w}^\top \mathbf{x}_n < \gamma$
- Support Vector Machine (SVM) does it directly by learning the max. margin hyperplane

Support Vector Machine (SVM)

SVM originally proposed by Vapnik and colleagues in early 90s

- Hyperplane based classifier. Ensures a large margin around the hyperplane
- Will assume a linear hyperplane to be of the form $\mathbf{w}^\top \mathbf{x} + b = 0$ (nonlinear ext. later)

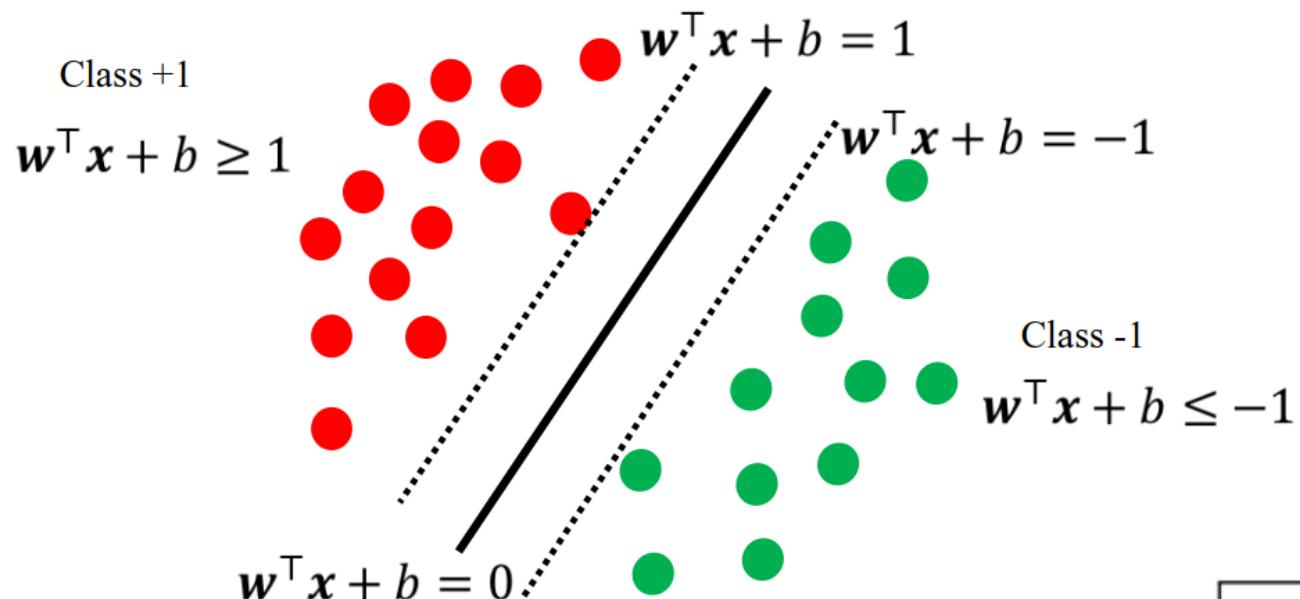


- Two other “supporting” hyperplanes defining a “no man’s land”
 - Ensure that zero training examples fall in this region (will relax later)
 - The SVM idea: Position the hyperplane s.t. this region is as “wide” as possible



Hard-Margin SVM

- Hard-Margin: Every training example must fulfil margin condition $y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$
- Meaning: Must not have any example in the no-man's land



Constrained optimization problem with N inequality constraints. Objective and constraints both are convex

- Also want to maximize margin $2\gamma = \frac{2}{\|\mathbf{w}\|}$
- Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$
- The objective func. for hard-margin SVM

$$\begin{aligned} & \min_{\mathbf{w}, b} f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ & \text{subject to } y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N \end{aligned}$$

Solving Hard-Margin SVM

- The hard-margin SVM optimization problem is

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & f(\mathbf{w}, b) = \frac{\|\mathbf{w}\|^2}{2} \\ \text{subject to} \quad & 1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) \leq 0, \quad n = 1, \dots, N \end{aligned}$$

- A constrained optimization problem. One option is to solve using Lagrange's method
- Introduce Lagrange multipliers α_n ($n = 1, \dots, N$), one for each constraint, and solve

$$\min_{\mathbf{w}, b} \max_{\boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{\|\mathbf{w}\|^2}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\}$$

- $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N]$ denotes the vector of Lagrange multipliers
- It is easier (and helpful; we will soon see why) to solve the dual: min and then max

Solving Hard-Margin SVM

- The dual problem (min then max) is

$$\max_{\alpha \geq 0} \min_{w, b} \mathcal{L}(w, b, \alpha) = \frac{w^T w}{2} + \sum_{n=1}^N \alpha_n \{1 - y_n(w^T x_n + b)\}$$

Note: if we ignore the bias term b then we don't need to handle the constraint $\sum_{n=1}^N \alpha_n y_n = 0$
(problem becomes a bit more easy to solve)



Otherwise, the α_n 's are coupled and some opt. techniques such as coordinate ascent can't easily be applied

- Take (partial) derivatives of \mathcal{L} w.r.t. w and b and setting them to zero gives (verify)

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_{n=1}^N \alpha_n y_n x_n$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0$$

α_n tells us how important training example (x_n, y_n) is

- The solution w is simply a weighted sum of all the training inputs

- Substituting $w = \sum_{n=1}^N \alpha_n y_n x_n$ in the Lagrangian, we get the dual problem as (verify)

This is also a "quadratic program" (QP) – a quadratic function of the variables α

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n y_m y_n (x_m^T x_n)$$

Maximizing a concave function
(or minimizing a convex function)
s.t. $\alpha \geq 0$ and $\sum_{n=1}^N \alpha_n y_n = 0$.
Many methods to solve it.

$$\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T G \alpha$$

G is an $N \times N$ p.s.d. matrix, also called the **Gram Matrix**, $G_{nm} = y_n y_m x_n^T x_m$, and $\mathbf{1}$ is a vector of all 1s

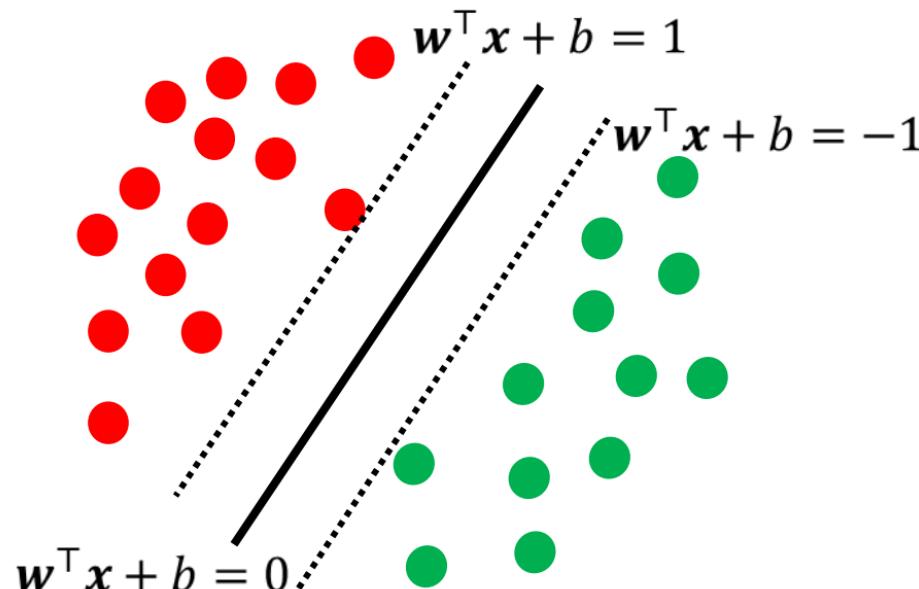
Solving Hard-Margin SVM

- Once we have the α_n 's by solving the dual, we can get \mathbf{w} and b as

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad (\text{we already saw this})$$

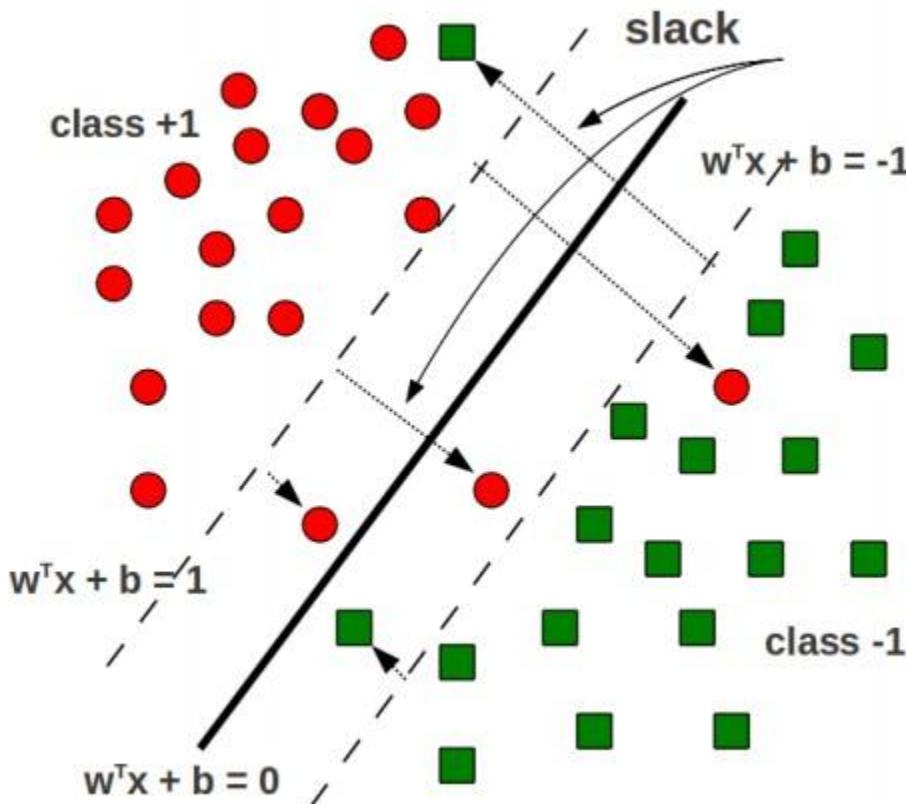
$$b = -\frac{1}{2} (\min_{n:y_n=+1} \mathbf{w}^T \mathbf{x}_n + \max_{n:y_n=-1} \mathbf{w}^T \mathbf{x}_n) \quad (\text{exercise})$$

- A nice property: Most α_n 's in the solution will be zero (sparse solution)



- Reason: KKT conditions
 - For the optimal α_n 's, we must have
 - Thus α_n nonzero only if $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$, i.e., the training example lies on the boundary
- $$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$
- These examples are called support vectors

Soft-Margin SVM (More Commonly Used)



Soft-margin constraint:

- Allow some training examples to fall within the no-man's land (margin region)
- Even okay for some training examples to fall totally on the wrong side of h.p.
- Extent of “violation” by a training input (\mathbf{x}_n, y_n) is known as **slack** $\xi_n \geq 0$
- $\xi_n > 1$ means totally on the wrong side

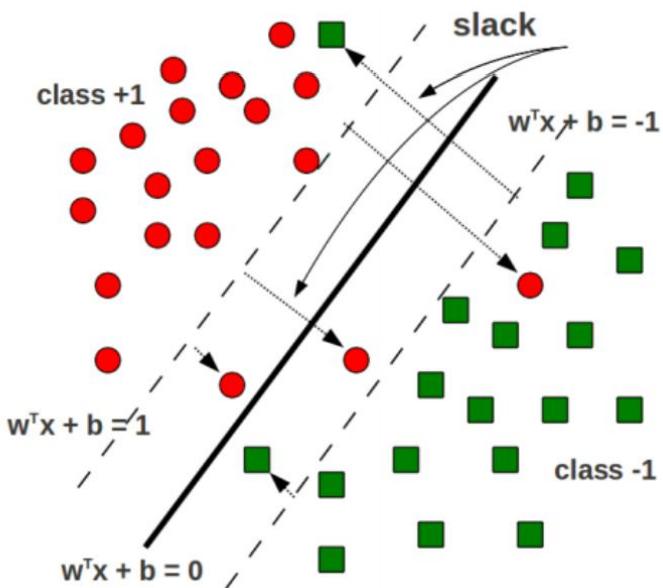
$$\mathbf{w}^\top \mathbf{x}_n + b \geq 1 - \xi_n \quad \text{if } y_n = +1$$

$$\mathbf{w}^\top \mathbf{x}_n + b \leq -1 + \xi_n \quad \text{if } y_n = -1$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n \quad \forall n$$

Soft-Margin SVM (Contd)

- Goal: Still want to maximize the margin such that
 - Soft-margin constraints $y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n$ are satisfied for all training ex.
 - Do not have too many margin violations (sum of slacks $\sum_{n=1}^N \xi_n$ should be small)



- The objective func. for soft-margin SVM

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

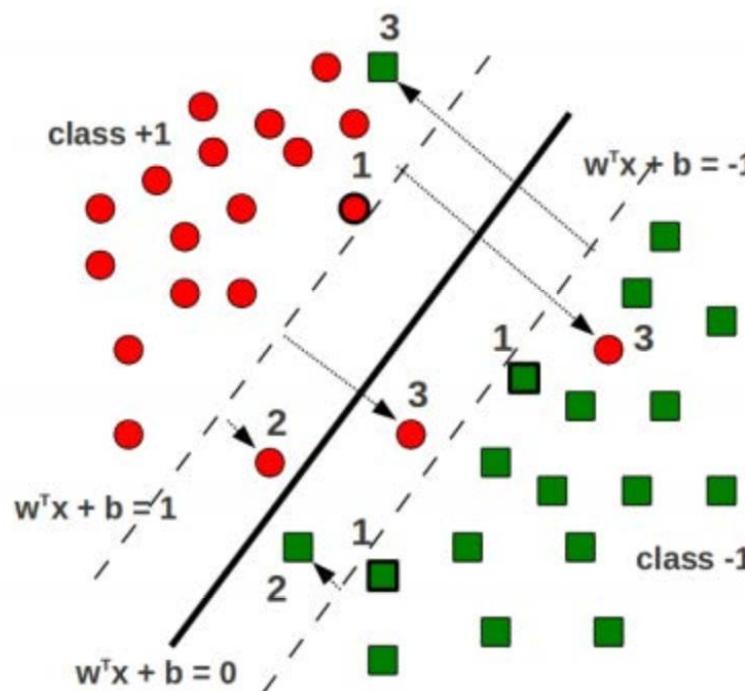
Annotations for the equation:

- Inversely prop. to margin: Points to the term $\frac{\|\mathbf{w}\|^2}{2}$.
- Trade-off hyperparam: Points to the term C .
- training error: Points to the term $\sum_{n=1}^N \xi_n$.
- Constrained optimization problem with $2N$ inequality constraints. Objective and constraints both are convex: Points to the subject to part of the equation.

- Hyperparameter C controls the trade off between large margin and small training error (need to tune)
 - Large C : small training error but also small margin (bad)
 - Small C : large margin but large training error (bad)

Support Vectors in Soft-Margin SVM

- The hard-margin SVM solution had only one type of support vectors
 - All lied on the supporting hyperplanes $\mathbf{w}^T \mathbf{x}_n + b = 1$ and $\mathbf{w}^T \mathbf{x}_n + b = -1$
- The soft-margin SVM solution has three types of support vectors (with nonzero α_n)



1. Lying on the supporting hyperplanes
2. Lying within the margin region but still on the correct side of the hyperplane
3. Lying on the wrong side of the hyperplane (misclassified training examples)

Solving Soft-Margin SVM

- Recall the soft-margin SVM optimization problem

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

- Here $\xi = [\xi_1, \xi_2, \dots, \xi_N]$ is the vector of slack variables
- Introduce Lagrange multipliers α_n, β_n for each constraint and solve Lagrangian

$$\min_{\mathbf{w}, b, \xi} \max_{\alpha \geq 0, \beta \geq 0} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) - \xi_n\} - \sum_{n=1}^N \beta_n \xi_n$$

- The terms in red color above were not present in the hard-margin SVM
- Two set of dual variables $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N]$ and $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_N]$
- Will eliminate the primal var \mathbf{w}, b, ξ to get dual problem containing the dual variables



Solving Soft-Margin SVM

- The Lagrangian problem to solve

$$\min_{\mathbf{w}, b, \xi} \max_{\alpha \geq 0, \beta \geq 0} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n \{1 - y_n(\mathbf{w}^\top \mathbf{x}_n + b) - \xi_n\} - \sum_{n=1}^N \beta_n \xi_n$$

Note: if we ignore the bias term b then we don't need to handle the constraint $\sum_{n=1}^N \alpha_n y_n = 0$ (problem becomes a bit more easy to solve)



Otherwise, the α_n 's are coupled and some opt. techniques such as co-ordinate aspect can't easily applied

- Take (partial) derivatives of \mathcal{L} w.r.t. \mathbf{w} , b , and ξ_n and setting to zero gives

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n, \quad \text{Weighted sum of training inputs}$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0, \quad \frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \Rightarrow C - \alpha_n - \beta_n = 0$$

- Using $C - \alpha_n - \beta_n = 0$ and $\beta_n \geq 0$, we have $\alpha_n \leq C$ (for hard-margin, $\alpha_n \geq 0$)
- Substituting these in the Lagrangian \mathcal{L} gives the Dual problem

Given α , \mathbf{w} and b can be found just like the hard-margin SVM case

$$\max_{\alpha \leq C, \beta \geq 0} \mathcal{L}_D(\alpha, \beta) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n y_m y_n (\mathbf{x}_m^\top \mathbf{x}_n) \quad \text{s.t.} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

The dual variables β don't appear in the dual problem!

Maximizing a concave function (or minimizing a convex function) s.t. $\alpha \leq C$ and $\sum_{n=1}^N \alpha_n y_n = 0$. Many methods to solve it.

$$\max_{\alpha \leq C} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha$$

In the solution, α will still be sparse just like the hard-margin SVM case. Nonzero α_n correspond to the support vectors

References

CS771: Intro to Machine Learning (Fall 2021), Nisheeth Srivastava, IIT Kanpur