

Optimization for ML (wrap-up)

CS771: Introduction to Machine Learning

Announcements

- Quiz 2 on Sept 7, 6:15-7:00pm in L20 and L19
 - Class will continue from 7:00pm-8:00pm in L20
- Quiz 1 grading will be completed soon (sorry for the delay)
- Mini-project 1 to be released later this week
- Mid-sem exam on Sept 15, 8:00am-10:00am



Plan today

- Optimization methods
 - **Constrained** optimization
 - Projected gradient descent
 - Lagrangian based method
 - Optimization of **non-differentiable** functions
 - **Subgradient** descent
 - Some examples of subgradient descent
- Large margin classification (Support Vector Machines (SVM))
 - Solving the SVM problem involves aspects of both optimization of non-differentiable functions as well as constrained optimization, and some other aspects of optimization too, so it is also a good “case study” 😊



Constrained Optimization



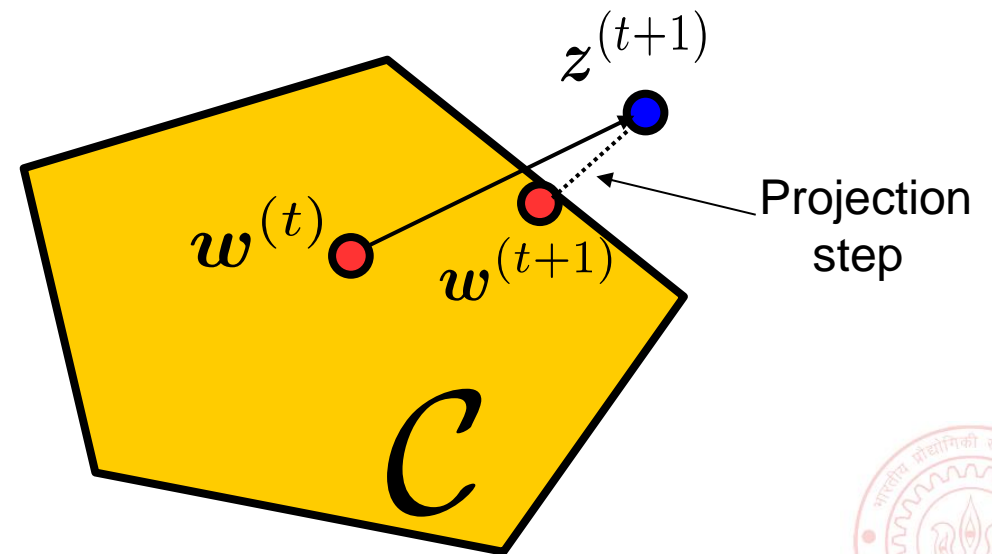
Projected Gradient Descent

- Consider an optimization problem of the form

$$\mathbf{w}_{opt} = \arg \min_{\mathbf{w} \in \mathcal{C}} L(\mathbf{w})$$

- Projected GD is very similar to GD with an extra **projection step**
- Each iteration t will be of the form

- Perform update: $\mathbf{z}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$
- Check if $\mathbf{z}^{(t+1)}$ satisfies constraints
 - If $\mathbf{z}^{(t+1)} \in \mathcal{C}$, set $\mathbf{w}^{(t+1)} = \mathbf{z}^{(t+1)}$
 - If $\mathbf{z}^{(t+1)} \notin \mathcal{C}$, project as $\mathbf{w}^{(t+1)} = \Pi_{\mathcal{C}}[\mathbf{z}^{(t+1)}]$



Projected GD: How to Project?

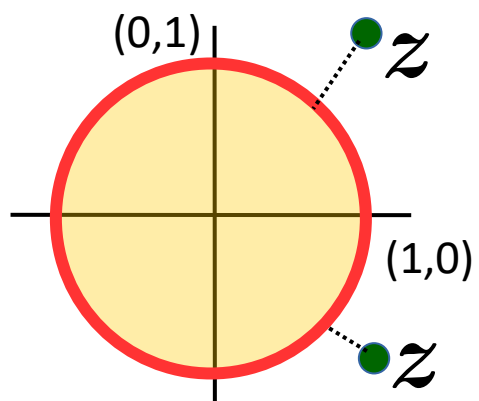
- Here projecting a point means finding the “closest” point from the constraint set

$$\Pi_{\mathcal{C}}[\mathbf{z}] = \arg \min_{\mathbf{w} \in \mathcal{C}} \|\mathbf{z} - \mathbf{w}\|^2$$

Another constrained optimization problem! But simpler to solve in some special cases! 😊

- For some sets \mathcal{C} , the projection step is easy

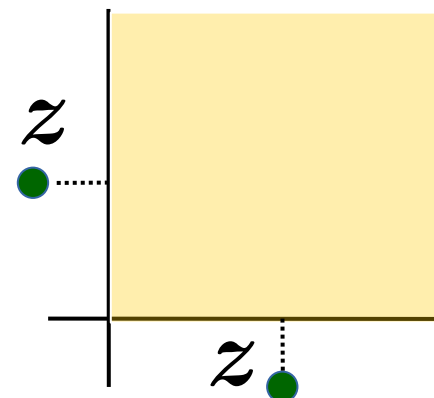
\mathcal{C} : Unit radius ℓ_2 ball



Projection = Normalize to unit Euclidean length vector

$$\hat{\mathbf{z}} = \begin{cases} \mathbf{z} & \text{if } \|\mathbf{z}\|_2 \leq 1 \\ \frac{\mathbf{z}}{\|\mathbf{z}\|_2} & \text{otherwise} \end{cases}$$

\mathcal{C} : Set of non-negative reals



Projection = Set each negative entry in \mathbf{z} to be zero

$$\hat{z}_i = \begin{cases} z_i & \text{if } z_i \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



Constrained Opt. via Lagrangian

- Consider the following constrained minimization problem (using f instead of L)

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} f(\mathbf{w}), \quad \text{s.t.} \quad g(\mathbf{w}) \leq 0$$

- Note: If constraints of the form $g(\mathbf{w}) \geq 0$, use $-g(\mathbf{w}) \leq 0$
- Can handle multiple inequality and equality constraints too (will see later)
- Can transform the above into the following equivalent unconstrained problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} f(\mathbf{w}) + c(\mathbf{w})$$

$$c(\mathbf{w}) = \max_{\alpha \geq 0} \alpha g(\mathbf{w}) = \begin{cases} \infty, & \text{if } g(\mathbf{w}) > 0 \quad (\text{constraint violated}) \\ 0 & \text{if } g(\mathbf{w}) \leq 0 \quad (\text{constraint satisfied}) \end{cases}$$

- Our problem can now be written as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \max_{\alpha \geq 0} \alpha g(\mathbf{w}) \right\}$$



Constrained Opt. via Lagrangian

- Therefore, we can write our original problem as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \max_{\alpha \geq 0} \alpha g(\mathbf{w}) \right\} = \arg \min_{\mathbf{w}} \left\{ \max_{\alpha \geq 0} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\}$$

The Lagrangian: $\mathcal{L}(\mathbf{w}, \alpha)$

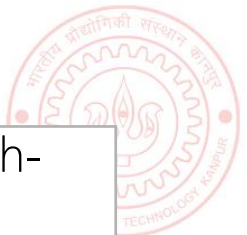
- The Lagrangian is now optimized w.r.t. \mathbf{w} and α (Lagrange multiplier)
- We can define **Primal** and **Dual** problem as

$$\begin{aligned} \hat{\mathbf{w}}_P &= \arg \min_{\mathbf{w}} \left\{ \max_{\alpha \geq 0} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\} && \text{(Primal Problem)} \\ \hat{\mathbf{w}}_D &= \arg \max_{\alpha \geq 0} \left\{ \min_{\mathbf{w}} \{ f(\mathbf{w}) + \alpha g(\mathbf{w}) \} \right\} && \text{(Dual Problem)} \end{aligned}$$

Both equal if $f(\mathbf{w})$ and the set $g(\mathbf{w}) \leq 0$ are convex

$$\alpha_D g(\hat{\mathbf{w}}_D) = 0$$

complimentary slackness/Karush-Kuhn-Tucker (KKT) condition



Constrained Opt. with Multiple Constraints

- We can also have multiple inequality and equality constraints

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} f(\mathbf{w}) \\ \text{s.t.} \quad &g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, K \\ &h_j(\mathbf{w}) = 0, \quad j = 1, \dots, L \end{aligned}$$

Note that, unlike α , Lagrange variables β of the equality constraints are unconstrained

Exercise: Justify why!



- Introduce Lagrange multipliers $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_K]$ and $\beta = [\beta_1, \beta_2, \dots, \beta_L]$
- The Lagrangian based primal and dual problems will be

$$\begin{aligned} \hat{\mathbf{w}}_P &= \arg \min_{\mathbf{w}} \left\{ \max_{\alpha \geq 0, \beta} \left\{ f(\mathbf{w}) + \sum_{i=1}^K \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^L \beta_j h_j(\mathbf{w}) \right\} \right\} \\ \hat{\mathbf{w}}_D &= \arg \max_{\alpha \geq 0, \beta} \left\{ \min_{\mathbf{w}} \left\{ f(\mathbf{w}) + \sum_{i=1}^K \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^L \beta_j h_j(\mathbf{w}) \right\} \right\} \end{aligned}$$

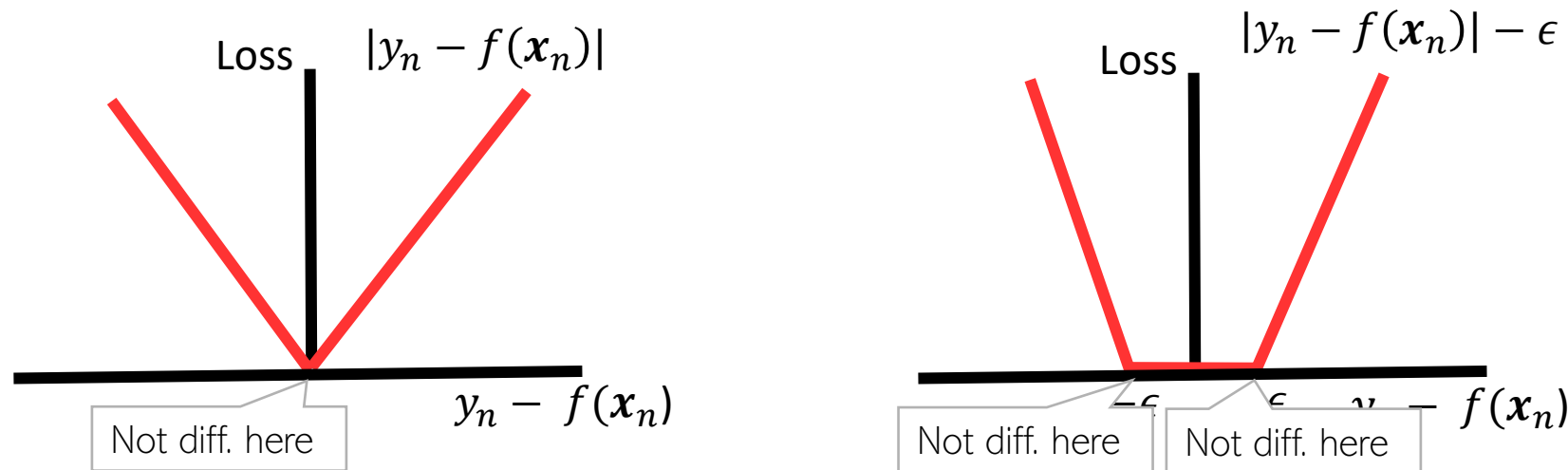


Optimization of Non-differentiable Functions



Dealing with Non-differentiable Functions

- In many ML problems, the objective function will be non-differentiable
- Some examples that we have already seen: Linear regression with absolute loss, or ϵ -insensitive loss; even ℓ_1 norm regularizer is non-diff

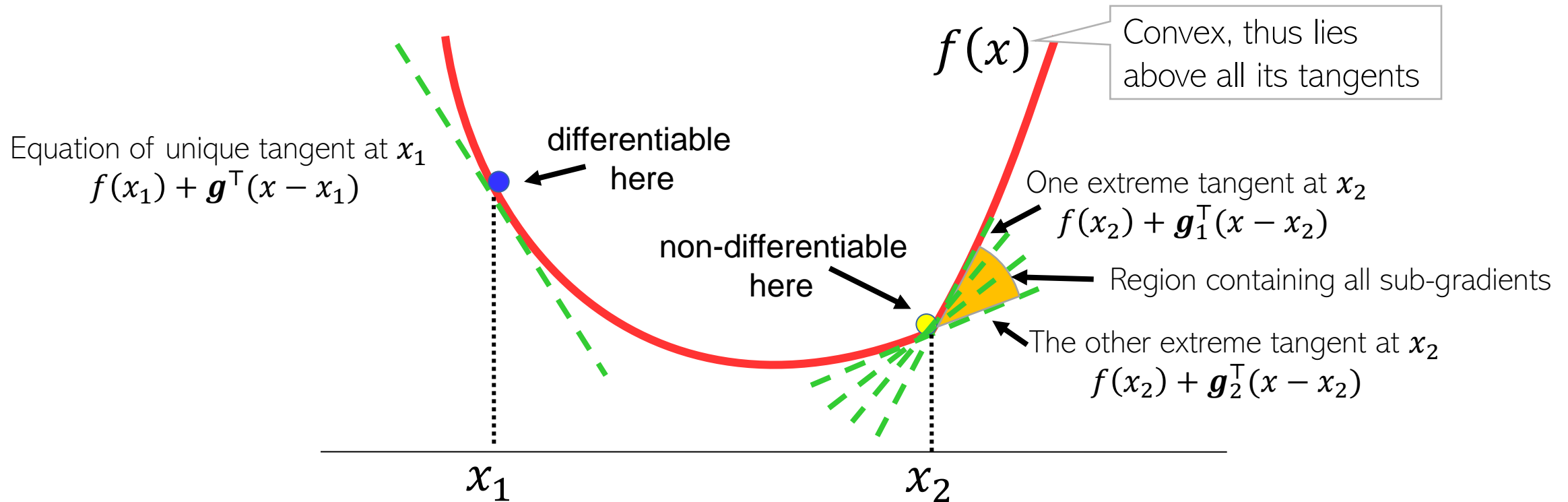


- Basically, any function in which there are points with kink is non-diff
 - At such points, the function is non-differentiable and thus **gradients not defined**
 - Reason: Can't define a unique tangent at such points



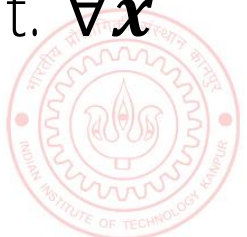
Sub-gradients

- For convex non-diff fn, can define **sub-gradients** at point(s) of non-differentiability



- For a convex, non-diff function $f(\mathbf{x})$, sub-gradient at \mathbf{x}_* is any vector \mathbf{g} s.t. $\forall \mathbf{x}$

$$f(\mathbf{x}) \geq f(\mathbf{x}_*) + \mathbf{g}^T(\mathbf{x} - \mathbf{x}_*)$$



Sub-gradients, Sub-differential, and Some Rules

- Set of all sub-gradient at a non-diff point \mathbf{x}_* is called the **sub-differential**

$$\partial f(\mathbf{x}_*) \triangleq \{\mathbf{g} : f(\mathbf{x}) \geq f(\mathbf{x}_*) + \mathbf{g}^\top (\mathbf{x} - \mathbf{x}_*) \quad \forall \mathbf{x}\}$$

- Some basic rules of sub-diff calculus to keep in mind

- Scaling rule: $\partial(c \cdot f(\mathbf{x})) = c \cdot \partial f(\mathbf{x}) = \{c \cdot \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$

- Sum rule: $\partial(f(\mathbf{x}) + g(\mathbf{x})) = \partial f(\mathbf{x}) + \partial g(\mathbf{x}) = \{\mathbf{u} + \mathbf{v} : \mathbf{u} \in \partial f(\mathbf{x}), \mathbf{v} \in \partial g(\mathbf{x})\}$

- Affine trans: $\partial f(\mathbf{a}^\top \mathbf{x} + b) = \mathbf{a} \cdot \partial f(t) = \{\mathbf{a} \cdot c : c \in \partial f(t)\}$, where $t = \mathbf{a}^\top \mathbf{x} + b$

- Max rule: If $h(\mathbf{x}) = \max\{f(\mathbf{x}), g(\mathbf{x})\}$ then we calculate $\partial h(\mathbf{x})$ at \mathbf{x}_* as

- If $f(\mathbf{x}_*) > g(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \partial f(\mathbf{x}_*)$, If $g(\mathbf{x}_*) > f(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \partial g(\mathbf{x}_*)$

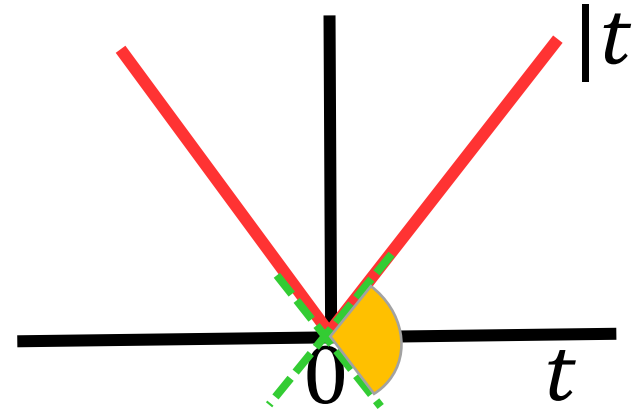
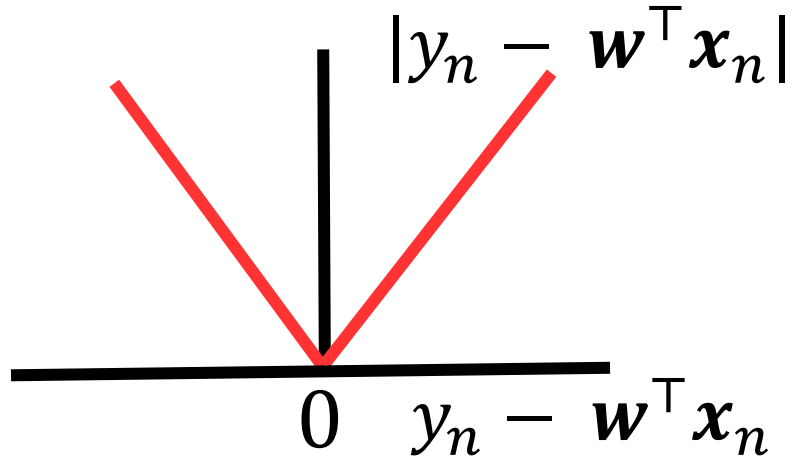
- If $f(\mathbf{x}_*) = g(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \{\alpha \mathbf{a} + (1 - \alpha) \mathbf{b} : \mathbf{a} \in \partial f(\mathbf{x}_*), \mathbf{b} \in \partial g(\mathbf{x}_*), \alpha \in [0, 1]\}$

- \mathbf{x}_* is a stationary point for a non-diff function $f(\mathbf{x})$ if the zero vector belongs to the sub-differential at \mathbf{x}_* , i.e., $\mathbf{0} \in \partial f(\mathbf{x}_*)$

The affine transform rule is a special case of the more general chain rule



Sub-Gradient For Absolute Loss Regression



Using max rule of sub-differentials and using $|t| = \max\{t, -t\}$

$$\partial|t| = \begin{cases} 1 & \text{if } t > 0 \\ -1 & \text{if } t < 0 \\ [-1, +1] & \text{if } t = 0 \end{cases}$$

- The loss function for linear reg. with absolute loss: $L(\mathbf{w}) = |y_n - \mathbf{w}^T \mathbf{x}_n|$
- Non-differentiable at $y_n - \mathbf{w}^T \mathbf{x}_n = 0$
- Can use the affine transform and max rule of sub-diff calculus
- Assume $t = y_n - \mathbf{w}^T \mathbf{x}_n$. Then $\partial L(\mathbf{w}) = -\mathbf{x}_n \partial|t|$
 - $\partial L(\mathbf{w}) = -\mathbf{x}_n \times 1 = -\mathbf{x}_n$ if $t > 0$
 - $\partial L(\mathbf{w}) = -\mathbf{x}_n \times -1 = \mathbf{x}_n$ if $t < 0$
 - $\partial L(\mathbf{w}) = -\mathbf{x}_n \times c = -c\mathbf{x}_n$ where $c \in [-1, +1]$ if $t = 0$



Sub-Gradient Descent

- Suppose we have a non-differentiable function $L(\mathbf{w})$
- Sub-gradient descent is almost identical to GD except we use subgradients

Sub-Gradient Descent

- Initialize \mathbf{w} as $\mathbf{w}^{(0)}$
- For iteration $t = 0, 1, 2, \dots$ (or until convergence)
 - Calculate the sub-gradient $\mathbf{g}^{(t)} \in \partial L(\mathbf{w}^{(t)})$
 - Set the learning rate η_t
 - Move in the opposite direction of subgradient

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$$

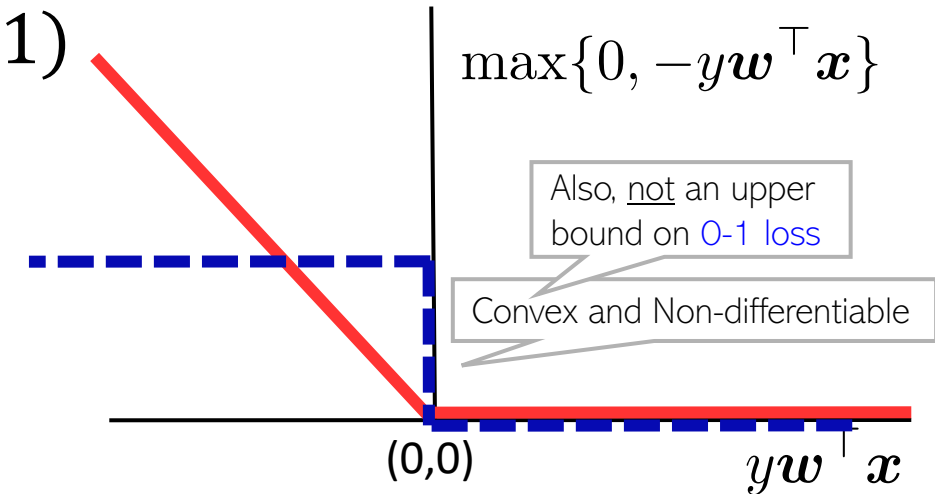


Detour: Other Loss Functions for Binary Classification

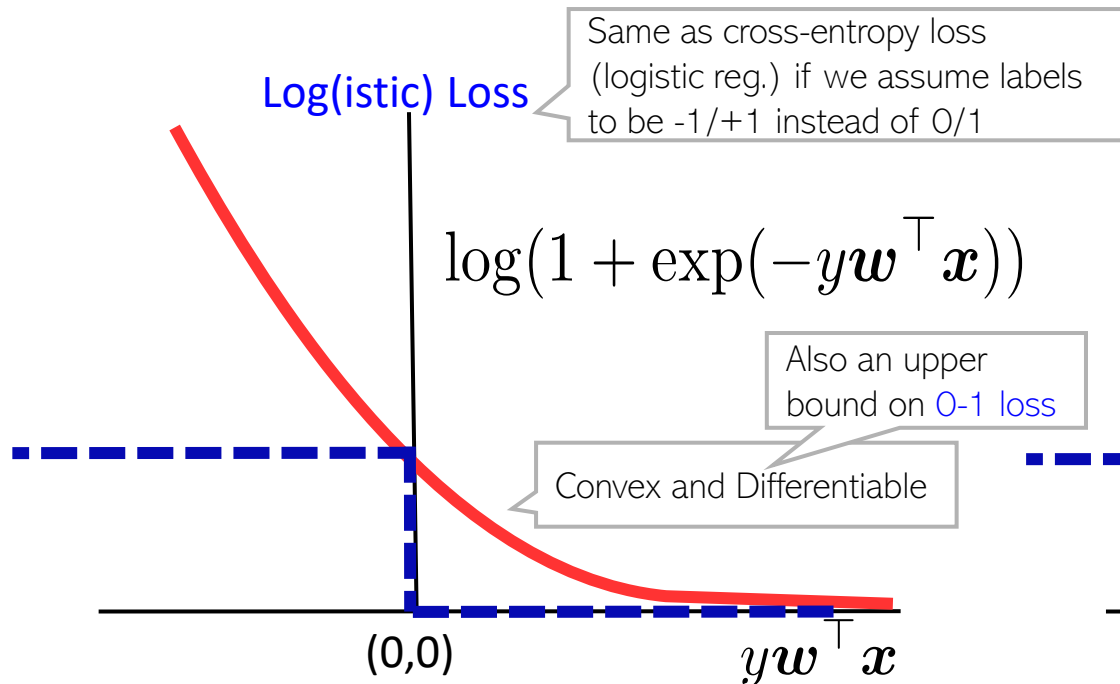
- For an ideal loss function, assuming $y_n \in (-1, +1)$

- Large positive $y_n \mathbf{w}^\top \mathbf{x}_n \Rightarrow$ small/zero loss
- Large negative $y_n \mathbf{w}^\top \mathbf{x}_n \Rightarrow$ large/non-zero loss
- Small (large) loss if predicted probability of the true label is large (small)

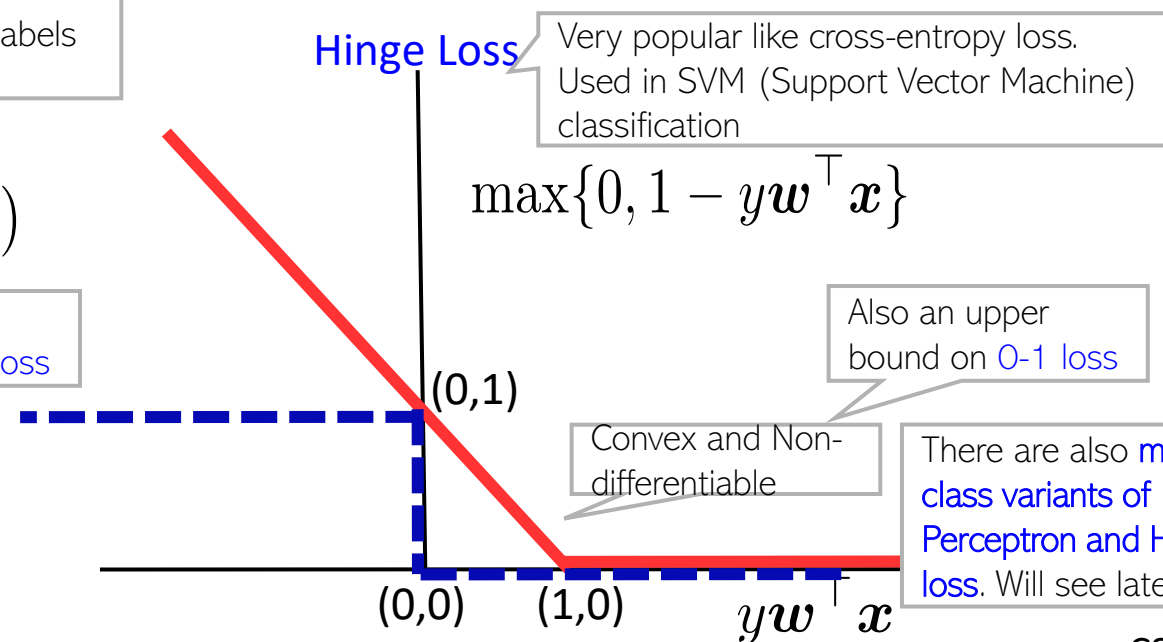
"Perceptron" Loss



Log(istic) Loss



Hinge Loss

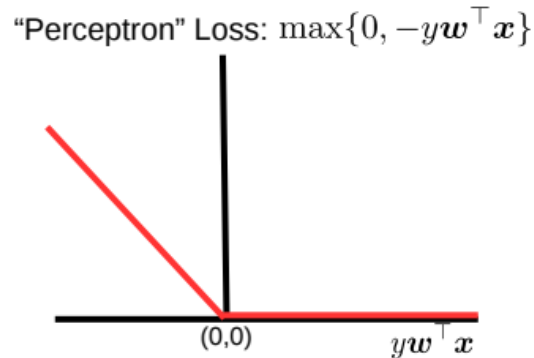


There are also multi-class variants of Perceptron and Hinge loss. Will see later



Subgradient Descent for Perceptron Loss

- Perceptron loss for binary classification: $L(\mathbf{w}) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^\top \mathbf{x}_n\}$



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

Using max rule of sub-differentials

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

- If we pick $c = -1$ then $\mathbf{g}_n = -\mathbb{I}[y_n \mathbf{w}^\top \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$
- Stochastic subgradient descent for Perceptron loss will have updates of the form

if $y_n \mathbf{w}^{(t)\top} \mathbf{x}_n \leq 0$

Meaning: The current weight vector does not correctly predict the label of \mathbf{x}_n

"mistake-driven" update
(update only when the current weights make a mistake)

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta_t y_n \mathbf{x}_n$$



Perceptron Algorithm for Binary Classification

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

- Initialize \mathbf{w} as $\mathbf{w}^{(0)}$, set $\eta_t = 1$
- For iteration $t = 0, 1, 2, \dots$ (or until convergence)
 - Pick a training example (\mathbf{x}_n, y_n) randomly
 - If $y_n \mathbf{w}^{(t)\top} \mathbf{x}_n \leq 0$ (i.e., mistake) then update

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_n \mathbf{x}_n$$

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

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 $\mathbf{w}^{(0)} = \mathbf{0}$, 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)
 - Important: \mathbf{w} in many classification/regression models can be written as $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$

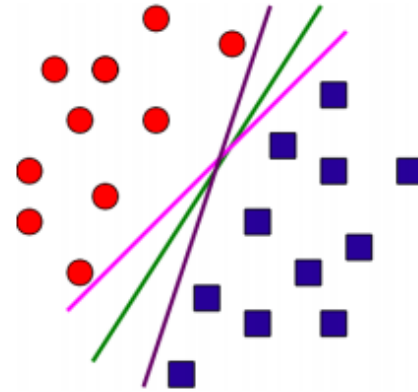
Thus \mathbf{w} is in the **span** of training inputs

Meaning of α_n may be different depending on the problem



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
- Can artificially introduce margin by changing the mistake condition to $\mathbf{y}_n \mathbf{w}^T \mathbf{x}_n \leq \gamma$
- Methods like logistic regression also do not guarantee large margins
- Support Vector Machine (SVM) does it directly by learning the **max. margin hyperplane**

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



Learning Large-Margin Hyperplanes

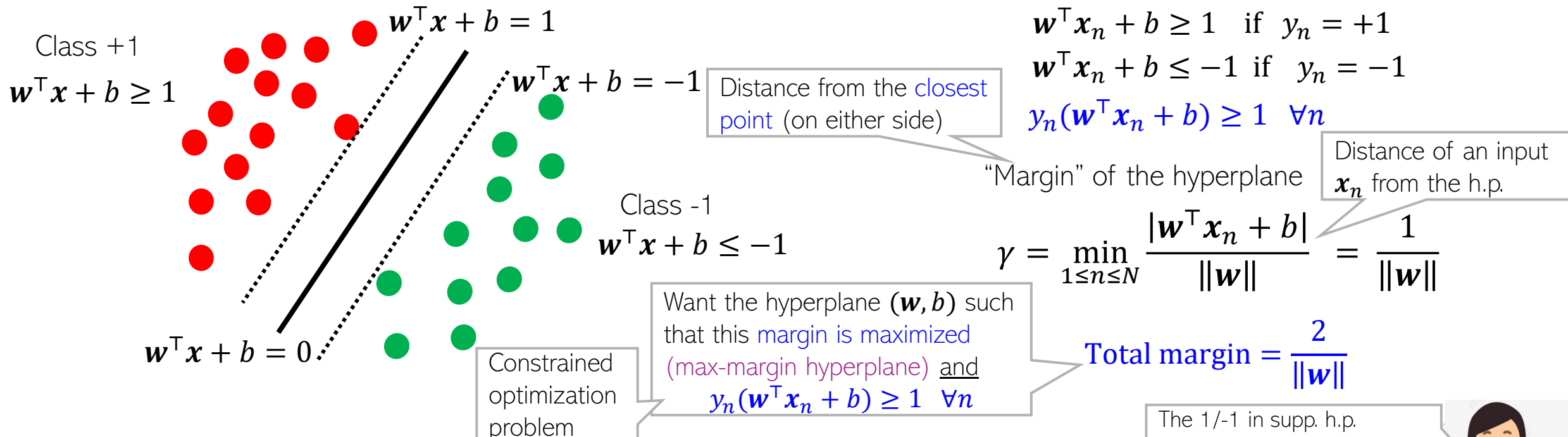


Support Vector Machine (SVM)

21

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)



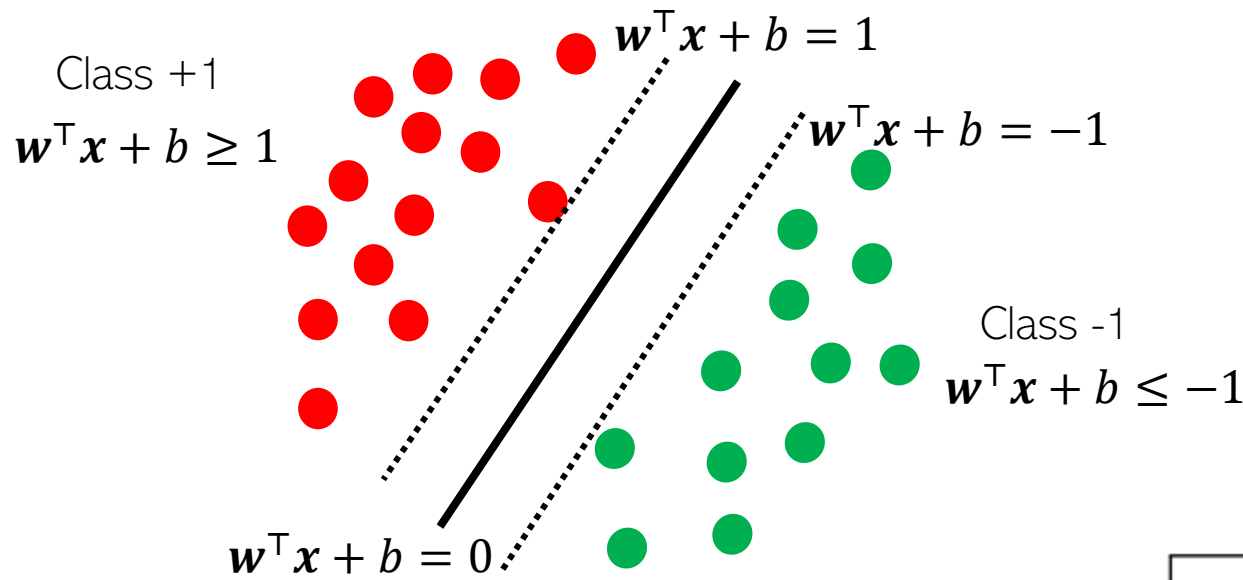
- Two other “supporting” hyperplanes defining a “no man’s land”
 - Ensure that zero training examples fall in this region (will relax later)
 - The SVM idea: Position the hyperplane s.t. this region is as “wide” as possible

The 1/-1 in supp. h.p. equations is arbitrary; can replace by any scalar m/-m and solution won't change, except a simple scaling of \mathbf{w}



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



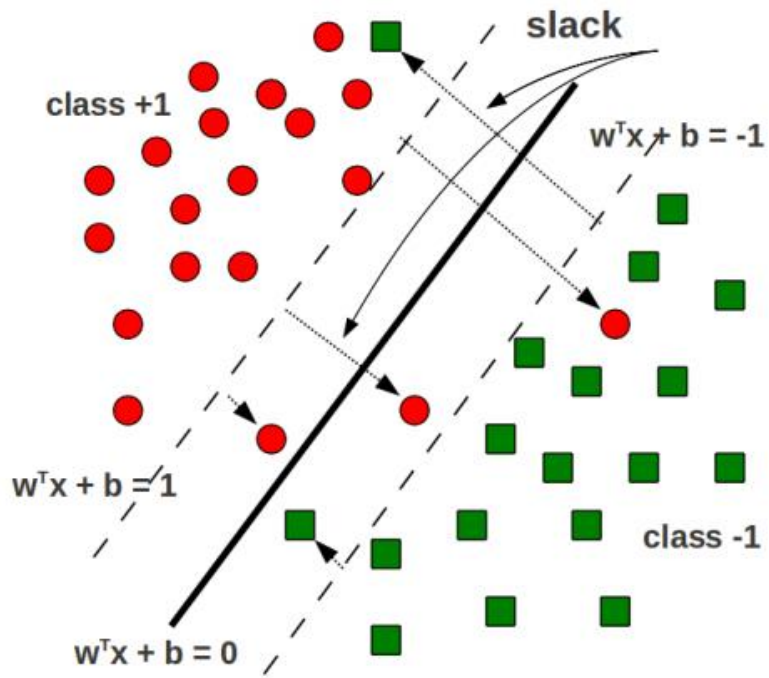
- Also want to maximize margin $2\gamma = \frac{2}{\|\mathbf{w}\|}$
- Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$
- The objective func. for hard-margin SVM

Lagrange based optimization can be used to solve it

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

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

Soft-Margin SVM (More Commonly Used)



Note/verify that the slack for each training example is just the hinge loss



$$\xi_n = \max\{0, 1 - y_n(w^\top x_n + b)\}$$

Soft-margin constraint

- Allow some training examples to fall within the no-man's land (margin region) Helps in getting a wider margin (and better generalization)
- 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

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

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

$$y_n(w^\top x_n + b) \geq 1 - \xi_n \quad \forall n$$



Soft-Margin SVM (Contd)

- Goal: Still want to maximize the margin such that
 - Soft-margin constraints $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n$ 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)

Sum of slacks is like the training error

- The objective func. for soft-margin SVM

$$\min_{\mathbf{w}, b, \xi} f(\mathbf{w}, b, \xi) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n$$

subject to $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0 \quad n = 1, \dots, N$

Annotations:

- Inversely prop. to margin: $\frac{\|\mathbf{w}\|^2}{2}$
- Trade-off hyperparam: C
- training error: $\sum_{n=1}^N \xi_n$
- Lagrange based optimization can be used to solve it
- Constrained optimization problem with $2N$ inequality constraints. Objective and constraints both are convex

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

