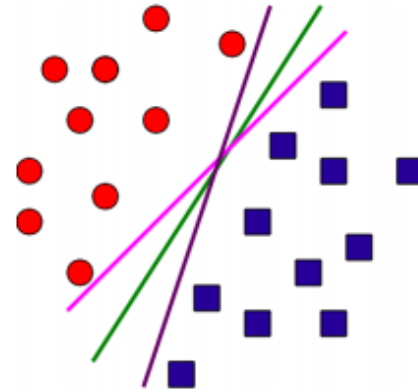


Large-Margin Classification and Support Vector Machine

CS771: Introduction to Machine Learning

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

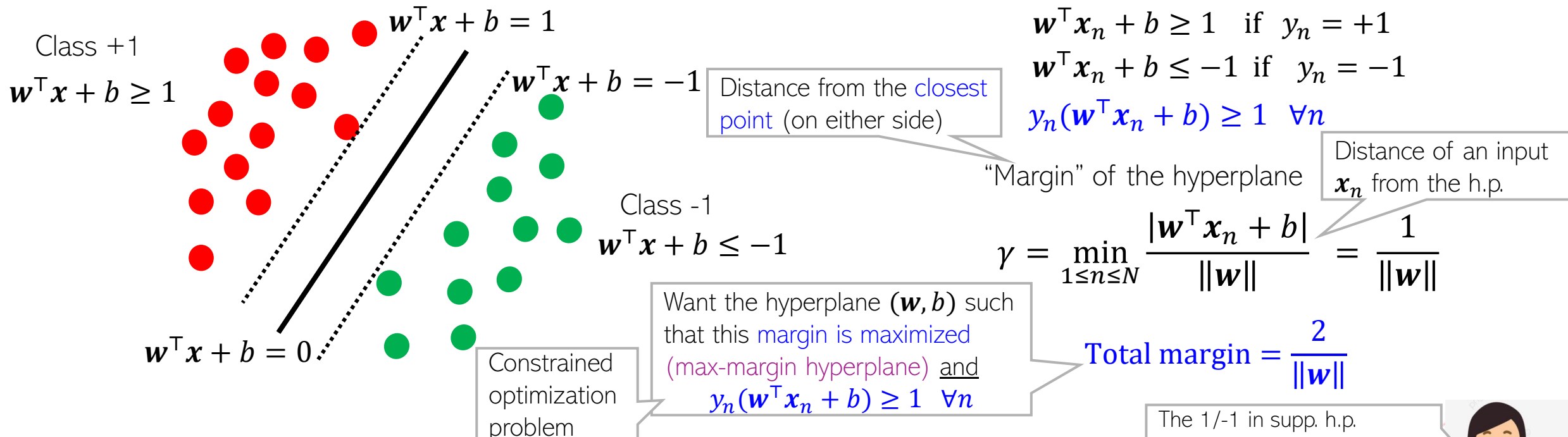


Support Vector Machine (SVM)

3

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 $\mathbf{w}^T \mathbf{x} + b = 0$ (nonlinear ext. later)

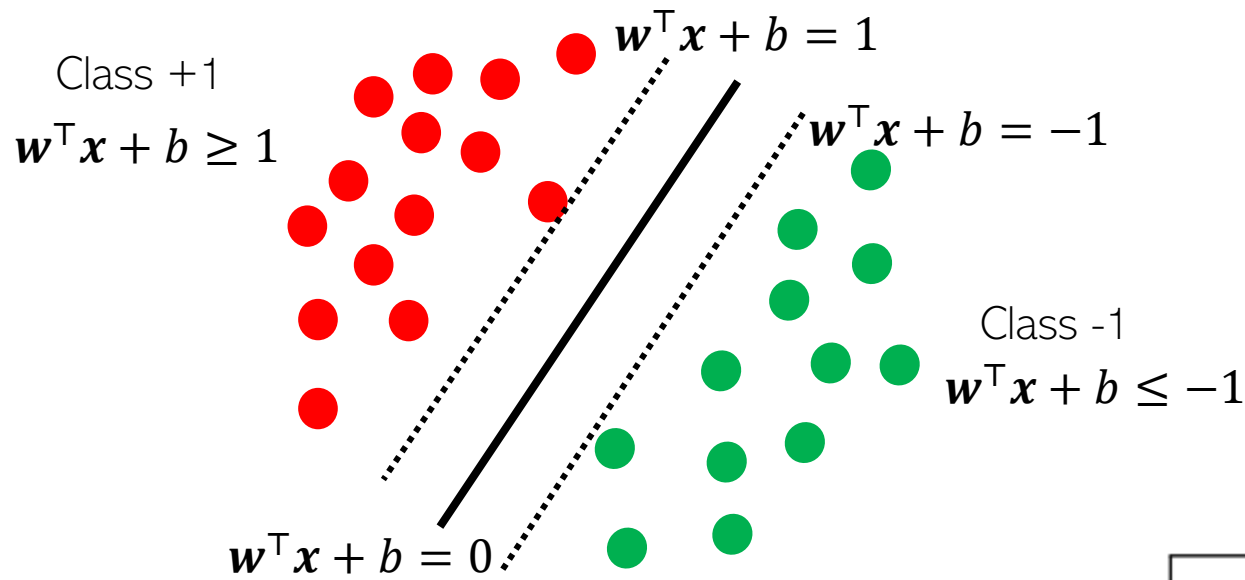


- Two other "supporting" hyperplanes defining a "no man's land"
 - Ensure that zero training examples fall in this region (will relax later)
 - The SVM idea: Position the hyperplane s.t. this region is as "wide" as possible



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



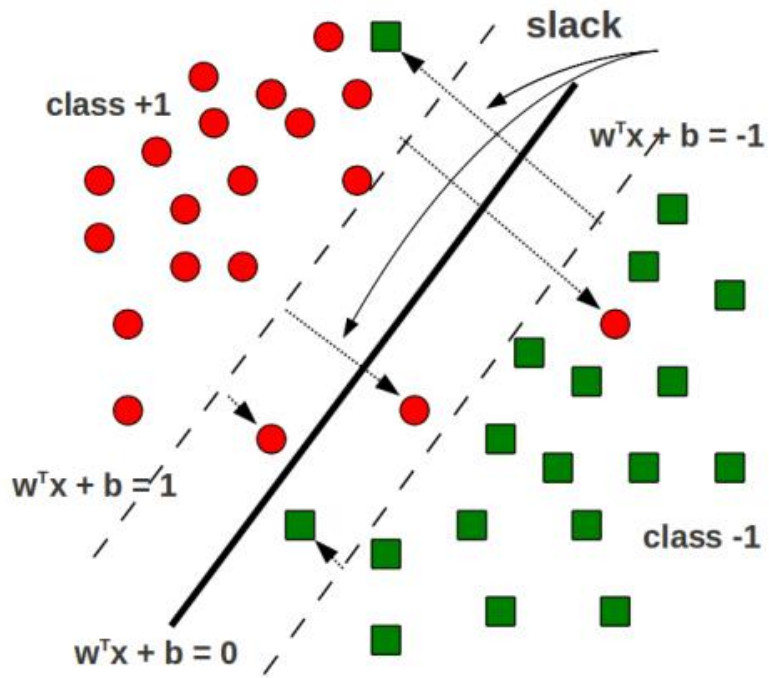
- Also want to maximize margin $2\gamma = \frac{2}{\|\mathbf{w}\|}$
- Equivalent to minimizing $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$
- The objective func. for hard-margin SVM

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

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

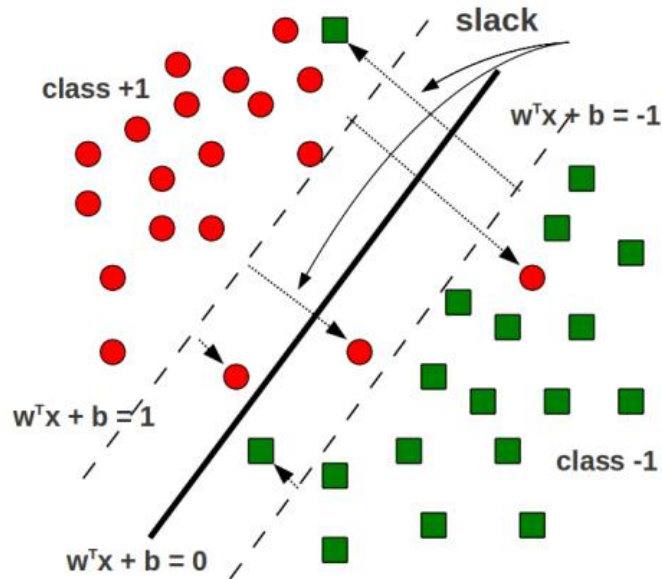
$$w^\top x_n + b \leq -1 + \xi_n \quad \text{if } y_n = -1$$

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



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)



- 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 (pointing to $\frac{\|\mathbf{w}\|^2}{2}$)
- Trade-off hyperparam (pointing to C)
- training error (pointing to $\sum_{n=1}^N \xi_n$)
- Sum of slacks is like the training error (pointing to $\sum_{n=1}^N \xi_n$)
- Lagrange based optimization can be used to solve it (pointing to the constraints)
- Constrained optimization problem with $2N$ inequality constraints. Objective and constraints both are convex (pointing to the entire problem)

- 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)



Solving the SVM Problem



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

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

- $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_N]$ 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

9

- The dual problem (min then max) is

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

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 co-ordinate ascent can't easily be applied

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

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y_n = 0$$

α_n tells us how important training example (\mathbf{x}_n, y_n) is

- The solution \mathbf{w} is simply a weighted sum of all the training inputs
- Substituting $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ in the Lagrangian, we get the dual problem as (verify)

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

IMPORTANT: inputs appear only as pairwise dot products. This will be useful later on when we make SVM nonlinear using kernel methods



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

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

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

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.

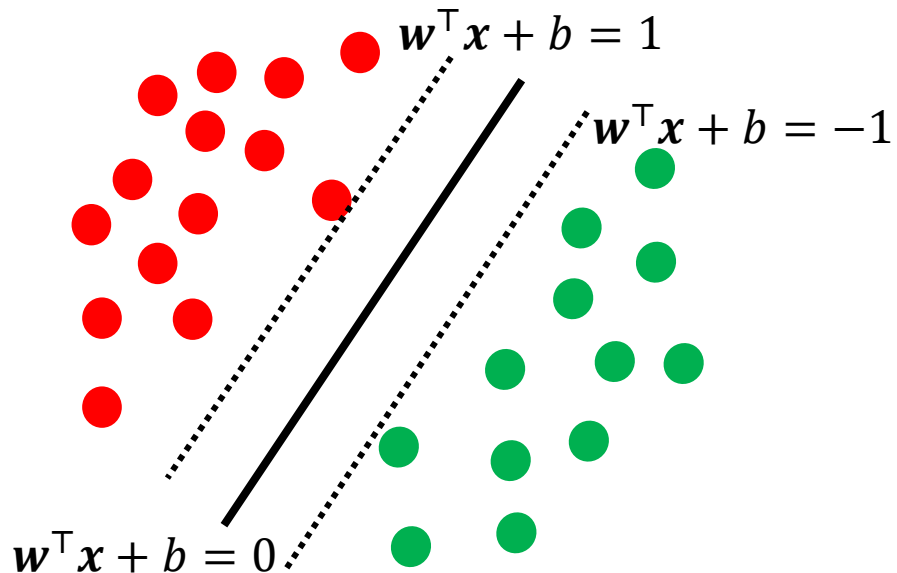
Solving Hard-Margin SVM

- 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

$$\alpha_n \{1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)\} = 0$$

- 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
- These examples are called support vectors



Solving Soft-Margin SVM

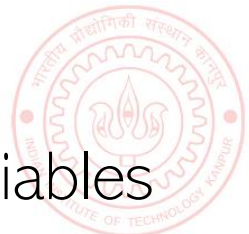
- Recall the soft-margin SVM optimization problem

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

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

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

- 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 $\mathbf{w}, b, \boldsymbol{\xi}$ to get dual problem containing the dual variables



Solving Soft-Margin SVM



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)

- The Lagrangian problem to solve

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

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

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

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

Weighted sum of training inputs

- 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

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

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

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

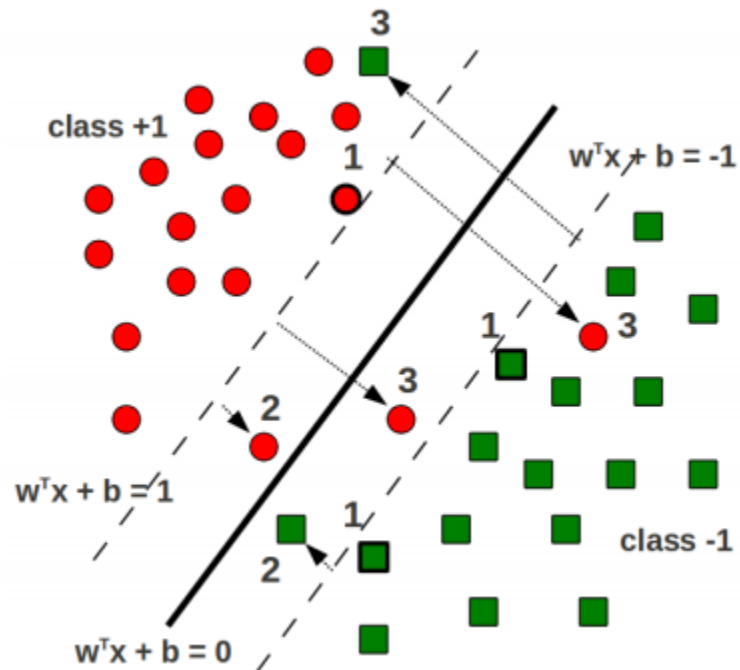
Maximizing a concave function (or minimizing a convex function) s.t. $\alpha \leq C$ and $\sum_{n=1}^N \alpha_n y_n = 0$. Many methods to solve it.

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

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

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)



SVMs via Dual Formulation: Some Comments

- Recall the final dual objectives for hard-margin and soft-margin SVM

Hard-Margin SVM: $\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha$

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

Note: Both these ignore the bias term b otherwise will need another constraint $\sum_{n=1}^N \alpha_n y_n = 0$

- The dual formulation is nice due to two primary reasons
 - Allows conveniently handling the margin based constraint (via Lagrangians)
 - Allows learning nonlinear separators by replacing inner products in $G_{nm} = y_n y_m \mathbf{x}_n^\top \mathbf{x}_m$ by general kernel-based similarities (more on this when we talk about kernels)
- However, dual formulation can be expensive if N is large (esp. compared to D)
 - Need to solve for N variables $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_N]$
 - Need to pre-compute and store $N \times N$ gram matrix \mathbf{G}
- Lot of work on speeding up SVM in these settings (e.g., can use co-ord. descent for α)



SVM: At Test Time

- Prediction for a test point

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

Dot product similarity of the test input \mathbf{x}_* with the training input \mathbf{x}_n

(Approach 1)

$$= \text{sign} \left(\sum_{n=1}^N \alpha_n y_n \mathbf{x}_n^\top \mathbf{x}_* + b \right)$$

(Approach 2)

- For linear SVMs, we usually prefer approach 1 since it is faster (just one dot product)
- The second approach's cost scales in the number of support vectors found by SVM (i.e., training examples with nonzero α_n). Also need to store them at test time
- The second approach is useful (and has to be used) for **nonlinear SVMs** where **\mathbf{w} cannot** usually be expressed as a finite dimensional vector (more when we talk about kernel methods)



Solving for SVM in the Primal

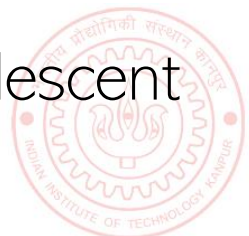
- Maximizing margin subject to constraints led to the soft-margin formulation of SVM

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

- Note that slack ξ_n is the same as $\max\{0, 1 - y_n(\mathbf{w}^\top \mathbf{x}_n + b)\}$, i.e., hinge loss for (\mathbf{x}_n, y_n)
- Thus the above is equivalent to minimizing the ℓ_2 regularized hinge loss

$$\mathcal{L}(\mathbf{w}, b) = \sum_{n=1}^N \max\{0, 1 - y_n(\mathbf{w}^\top \mathbf{x}_n + b)\} + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Sum of slacks is like sum of hinge losses, C and λ play similar roles
- Can learn (\mathbf{w}, b) directly by minimizing $\mathcal{L}(\mathbf{w}, b)$ using (stochastic) (sub)grad. descent
 - Hinge-loss version preferred for linear SVMs, or with other regularizers on \mathbf{w} (e.g., ℓ_1)



A Co-ordinate Ascent Algorithm for SVM

- Recall the dual objective of soft-margin SVM (assuming no bias b)

$$\operatorname{argmax}_{\mathbf{0} \leq \boldsymbol{\alpha} \leq C} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$$

Note that $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$

- Focusing on just one of the components of $\boldsymbol{\alpha}$ (say α_n), the objective becomes

$$\operatorname{argmax}_{\mathbf{0} \leq \alpha_n \leq C} \alpha_n - \frac{1}{2} \alpha_n^2 \|\mathbf{x}_n\|^2 - \frac{1}{2} \alpha_n y_n \sum_{m \neq n} \alpha_m y_m \mathbf{x}_m^\top \mathbf{x}_n$$

Can compute these in the beginning itself

Can efficiently compute it if we also store \mathbf{w} . It is equal to $\mathbf{w}^\top \mathbf{x}_n - \alpha_n y_n \|\mathbf{x}_n\|^2$

- The above is a simple quadratic maximization of a concave function: Global maxima
- If constraint violated, project α_n in $[0, C]$: If $\alpha_n < 0$, set it to 0, if $\alpha_n > C$, set it to C
- Can cycle through each coordinate α_n in a random or cyclic fashion



SVM: Summary

- A hugely (perhaps the most, before deep learning became fashionable 😊) popular classification algorithm
- Reasonably mature, highly optimized SVM softwares freely available (perhaps the reason why it is more popular than various other competing algorithms)
- Some popular ones: libSVM, LIBLINEAR, sklearn also provides SVM
- Lots of work on scaling up SVMs* (both large N and large D)
- Extensions beyond binary classification (e.g., multiclass, structured outputs)
- Can even be used for regression problems (Support Vector Regression)
- Nonlinear extensions possible via kernels



* See: “Support Vector Machine Solvers” by Bottou and Lin