

Optimization (Wrap-up), and Hyperplane based Classifiers (Perceptron and Support Vector Machines)

Piyush Rai

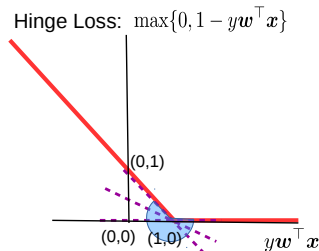
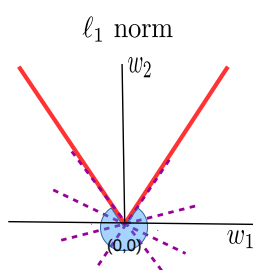
Introduction to Machine Learning (CS771A)

August 30, 2018



Recap: Subgradient Descent

- Use subgradient at non-differentiable points, use gradient elsewhere



- Left: each entry of (sub)gradient vector for $\|\mathbf{w}\|_1$, Right: (sub)gradient vector for hinge loss

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

Recap: Constrained Optimization via Lagrangian

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

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

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

Same as $f(\mathbf{w})$ when constraint satisfied

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

Lagrangian:

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha g(\mathbf{w})$$



Recap: Constrained Optimization via Lagrangian

- We minimize the Lagrangian $\mathcal{L}(\mathbf{w}, \alpha)$ w.r.t. \mathbf{w} and maximize w.r.t. α

$$\mathcal{L}(\mathbf{w}, \alpha) = f(\mathbf{w}) + \alpha g(\mathbf{w})$$

- For certain problems, the order of maximization and minimization does not matter
- Approach 1: Can first maximize w.r.t. α and then minimize w.r.t. \mathbf{w}

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\{ \max_{\alpha} \mathcal{L}(\mathbf{w}, \alpha) \right\}$$

- Approach 2: Can first minimize w.r.t. \mathbf{w} and then maximize w.r.t. α

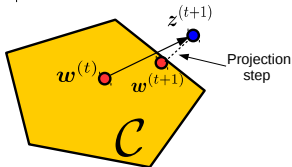
$$\hat{\alpha} = \arg \max_{\alpha} \left\{ \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \alpha) \right\}$$

- Approach 2 is known as **optimizing via the dual** (popular in SVM solvers; will see today!)
- **KKT condition**: At the optimal solution $\hat{\alpha} g(\hat{\mathbf{w}}) = 0$
- Multiple constraints (inequality/equality) can also be handled likewise

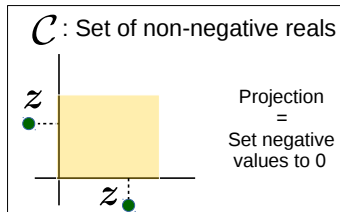
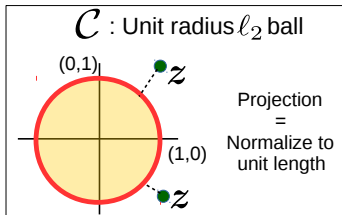


Recap: Projected Gradient Descent

- Same as GD + extra projection step we step out of the constraint set



- In some cases, the projection step is very easy



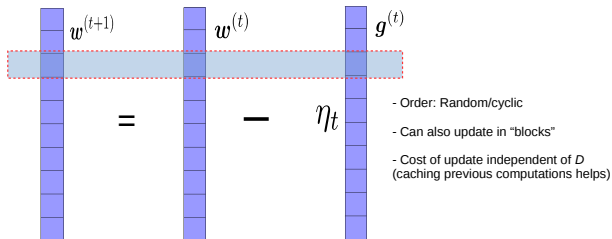
Co-ordinate Descent (CD)

- Standard GD update for $\mathbf{w} \in \mathbb{R}^D$ at each step

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

- CD: Each step updates one component (co-ordinate) at a time, keeping all others fixed

$$w_d^{(t+1)} = w_d^{(t)} - \eta_t g_d^{(t)}$$



Alternating Optimization

- Consider an optimization problems with several variables, say 2 variables \mathbf{w}_1 and \mathbf{w}_2

$$\{\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2\} = \arg \min_{\mathbf{w}_1, \mathbf{w}_2} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2)$$

- Often, this “joint” optimization is hard/impossible. We can consider an alternating scheme

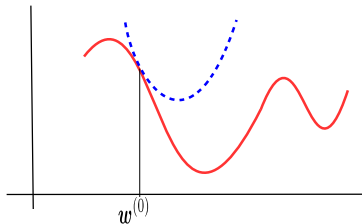
ALT-OPT

- 1 Initialize one of the variables, e.g., $\mathbf{w}_2 = \mathbf{w}_2^{(0)}, t = 0$
 - 2 Solve $\mathbf{w}_1^{(t+1)} = \arg \min_{\mathbf{w}_1} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2^{(t)})$ // \mathbf{w}_2 “fixed” at its most recent value $\mathbf{w}_2^{(t)}$
 - 3 Solve $\mathbf{w}_2^{(t+1)} = \arg \min_{\mathbf{w}_2} \mathcal{L}(\mathbf{w}_1^{(t+1)}, \mathbf{w}_2)$ // \mathbf{w}_1 “fixed” at its most recent value $\mathbf{w}_1^{(t+1)}$
 - 4 $t = t + 1$. Go to step 2 if not converged yet.
- Usually converges to a **local optima** of $\mathcal{L}(\mathbf{w}_1, \mathbf{w}_2)$. Also connections to **EM** (will see later)
 - VERY VERY useful!!!** Also extends to more than 2 variables. CD is somewhat like ALT-OPT.

Newton's Method

- **Newton's method** uses second-order information (second derivative a.k.a. Hessian)
- At each point $\mathbf{w}^{(t)}$, minimize the **quadratic** (second-order) approximation of the function

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}^{(t)}) + \nabla f(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 f(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)}) \right\}$$

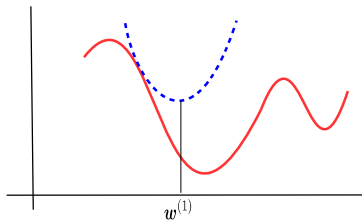


- **Exercise:** Verify that $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$. Also no learning rate needed!
- Converges much faster than GD. But also expensive due to Hessian computation/inversion.
- Many ways to approximate the Hessian (e.g., using previous gradients); also look at L-BFGS etc.

Newton's Method

- **Newton's method** uses second-order information (second derivative a.k.a. Hessian)
- At each point $\mathbf{w}^{(t)}$, minimize the **quadratic** (second-order) approximation of the function

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}^{(t)}) + \nabla f(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 f(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)}) \right\}$$

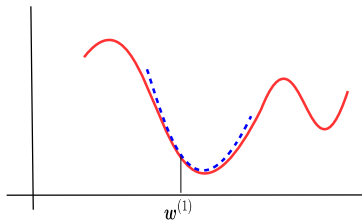


- **Exercise:** Verify that $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$. Also no learning rate needed!
- Converges much faster than GD. But also expensive due to Hessian computation/inversion.
- Many ways to approximate the Hessian (e.g., using previous gradients); also look at L-BFGS etc.

Newton's Method

- **Newton's method** uses second-order information (second derivative a.k.a. Hessian)
- At each point $\mathbf{w}^{(t)}$, minimize the **quadratic** (second-order) approximation of the function

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}^{(t)}) + \nabla f(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 f(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)}) \right\}$$

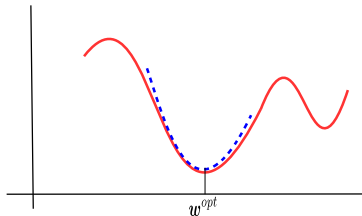


- **Exercise:** Verify that $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$. Also no learning rate needed!
- Converges much faster than GD. But also expensive due to Hessian computation/inversion.
- Many ways to approximate the Hessian (e.g., using previous gradients); also look at L-BFGS etc.

Newton's Method

- **Newton's method** uses second-order information (second derivative a.k.a. Hessian)
- At each point $\mathbf{w}^{(t)}$, minimize the **quadratic** (second-order) approximation of the function

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w}} \left\{ f(\mathbf{w}^{(t)}) + \nabla f(\mathbf{w}^{(t)})^\top (\mathbf{w} - \mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^\top \nabla^2 f(\mathbf{w}^{(t)}) (\mathbf{w} - \mathbf{w}^{(t)}) \right\}$$



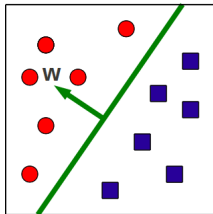
- **Exercise:** Verify that $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - (\nabla^2 f(\mathbf{w}^{(t)}))^{-1} \nabla f(\mathbf{w}^{(t)})$. Also no learning rate needed!
- Converges much faster than GD. But also expensive due to Hessian computation/inversion.
- Many ways to approximate the Hessian (e.g., using previous gradients); also look at L-BFGS etc.

Summary

- Gradient methods are simple to understand and implement
- More sophisticated optimization methods often use gradient methods
 - **Backpropagation algorithm** used in deep neural nets is **GD + chain rule** of differentiation
- Use **subgradient** methods if function not differentiable
- Constrained optimization require methods such as **Lagrangian** or **projected gradient**
- **Second order methods** such as Newton's method are much faster but computationally expensive
- But computing all this gradient related stuff looks scary to me. Any help?
 - Don't worry. **Automatic Differentiation (AD)** methods available now
 - AD only requires **specifying the loss function** (useful for complex models like deep neural nets)
 - Many packages such as Tensorflow, PyTorch, etc. provide AD support
 - But having a good understanding of optimization is still helpful



Hyperplane based Classification



All **linear models for classification** are basically about learning hyperplanes!

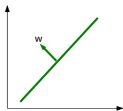
Already saw logistic regression (probabilistic linear classifier).

Will look at some more today - Perceptron, SVM
(also how some of the optimization methods we saw can be applied in these cases)



Hyperplanes

- Separates a D -dimensional space into two **half-spaces** (positive and negative)

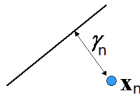


- Defined by normal vector $\mathbf{w} \in \mathbb{R}^D$ (pointing towards positive half-space)
- Equation of the hyperplane: $\mathbf{w}^\top \mathbf{x} = 0$
- Assumption: The hyperplane passes through origin. If not, add a **bias** term $b \in \mathbb{R}$

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

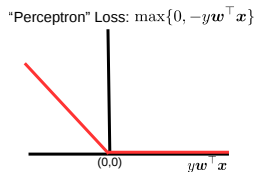
- $b > 0$ means moving it parallelly in the direction of \mathbf{w} ($b < 0$ means moving in opposite direction)
- **Distance** of a point \mathbf{x}_n from a hyperplane (can be +ve/-ve)

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



A Mistake-Driven Method for Learning Hyperplanes

- Let's ignore the bias term b for now. So the hyperplane is simply $\mathbf{w}^\top \mathbf{x} = 0$
- Consider SGD to learn a hyperplane based model with loss: $\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^\top \mathbf{x}_n\}$



- Loss not differentiable at $y_n \mathbf{w}^\top \mathbf{x}_n = 0$, so we will use subgradients there. The (sub)gradient will be

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

- If we use $k = 0$ then $\mathbf{g}_n = 0$ for $y_n \mathbf{w}^\top \mathbf{x}_n \geq 0$, and $\mathbf{g}_n = -y_n \mathbf{x}_n$ if $y_n \mathbf{w}^\top \mathbf{x}_n < 0$
- Thus \mathbf{g}_n nonzero only when $y_n \mathbf{w}^\top \mathbf{x}_n < 0$ (mistake). SGD will update \mathbf{w} only in these cases!

Mistake-Driven Learning of Hyperplanes

- The complete SGD algorithm for a model with this loss function will be

Stochastic SubGD

- 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., $y_n \mathbf{w}^{(t)\top} \mathbf{x}_n < 0$
$$\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + y_n \mathbf{x}_n \\ t &= t + 1\end{aligned}$$
- 4 If not converged, go to step 2.

- This is the **Perceptron algorithm**. An example of an **online learning** algorithm
- Note: Assuming $\mathbf{w}^{(0)} = 0$, easy to see the final \mathbf{w} has the form $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$
 - .. where 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 will also lead to $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ (for some suitable α_n 's)



Perceptron: Corrective Updates and Convergence

- Suppose true $y_n = +1$ (positive example) and the model mispredicts, i.e., $\mathbf{w}^{(t)\top} \mathbf{x}_n < 0$
- After the update $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_n \mathbf{x}_n = \mathbf{w}^{(t)} + \mathbf{x}_n$

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

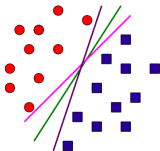
.. which is **less negative** than $\mathbf{w}^{(t)\top} \mathbf{x}_n$ (so the model has improved)

- **Exercise:** Verify that the model also improves after updating on a mistake on negative example
- Note: If training data is **linearly separable**, Perceptron converges in finite iterations
 - Proof: **Block & Novikoff theorem** (will provide the proof in a separate note)
 - What this means: It will eventually classify every training example correctly
 - Speed of convergence depends on the margin of separation (and on nothing else, such as N , D)
 - Note: In practice, we might want to stop sooner (to avoid overfitting)



Perceptron and (Lack of) Margins

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



- The one learned will depend on the initial \mathbf{w}
- Standard Perceptron doesn't guarantee any “margin” around the hyperplane
- Note: Possible to “artificially” introduce a margin in the Perceptron
 - Simply change the Perceptron mistake condition to

$$y_n \mathbf{w}^T \mathbf{x}_n < \gamma$$

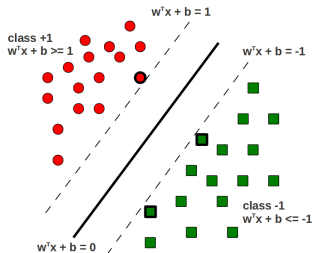
where $\gamma > 0$ is a **pre-specified margin**. For standard Perceptron, $\gamma = 0$

- **Support Vector Machine (SVM)** does this directly by learning the **maximum margin hyperplane**



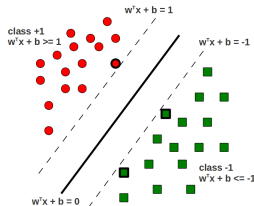
Support Vector Machine (SVM)

- SVM is a hyperplane based (linear) classifier that ensures a **large margin** around the hyperplane
- Note: We will assume the hyperplane to be of the form $\mathbf{w}^T \mathbf{x} + b = 0$ (will keep the bias term b)



- Note: SVMs can also learn **nonlinear decision boundaries** using **kernel methods** (will see later)
- Reason behind the name “Support Vector Machine”?
 - SVM optimization discovers the most important examples (called “**support vectors**”) in training data
 - These examples act as “balancing” the margin boundaries (hence called “support”)

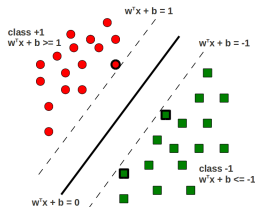
Learning a Maximum Margin Hyperplane



- Suppose we want a hyperplane $\mathbf{w}^T \mathbf{x} + b = 0$ such that
 - $\mathbf{w}^T \mathbf{x}_n + b \geq 1$ for $y_n = +1$
 - $\mathbf{w}^T \mathbf{x}_n + b \leq -1$ for $y_n = -1$
 - Equivalently, $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n$
 - Define the margin on each side: $\gamma = \min_{1 \leq n \leq N} \frac{|\mathbf{w}^T \mathbf{x}_n + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$
 - Total margin $= 2\gamma = \frac{2}{\|\mathbf{w}\|}$
- Want the hyperplane (\mathbf{w}, b) that gives the largest possible margin
- **Note:** Can replace $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$ by $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq m$ for some $m > 0$. It won't change the solution for \mathbf{w} , will just scale it by a constant, without changing the direction of \mathbf{w} (exercise)

Hard-Margin SVM

- Hard-Margin: **Every** training example has to fulfil the margin condition $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$



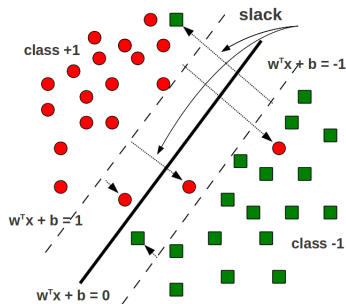
- Also want to maximize the margin $\gamma \propto \frac{1}{\|\mathbf{w}\|}$. Equivalent to **minimizing** $\|\mathbf{w}\|^2$ or $\frac{\|\mathbf{w}\|^2}{2}$
- The objective for hard-margin SVM

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

- Constrained optimization with N inequality constraints (note: function and constraints are convex)

Soft-Margin SVM (More Commonly Used)

- Allow some training examples to fall **within the margin region**, or be even **misclassified** (i.e., fall on the wrong side). Preferable if **training data is noisy**

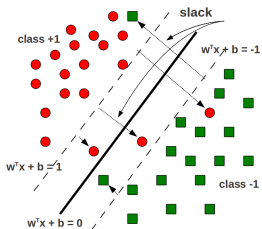


- Each training example (\mathbf{x}_n, y_n) given a “slack” $\xi_n \geq 0$ (distance by which it “violates” the margin). If $\xi_n > 1$ then \mathbf{x}_n is totally on the wrong side
 - Basically, we want a **soft-margin condition**: $y_n(w^T \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0$



Soft-Margin SVM (More Commonly Used)

- Goal: Maximize the margin, while also minimizing the **sum of slacks** (don't want too many training examples violating the margin condition)



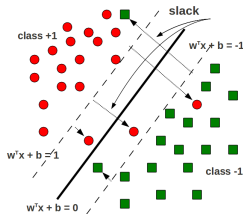
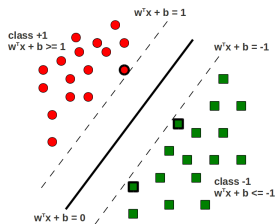
- The primal objective for soft-margin SVM can thus be written as

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

- Constrained optimization with $2N$ inequality constraints
- Parameter C controls the trade-off between large margin vs small training error



Summary: Hard-Margin SVM vs Soft-Margin SVM



- Objective for the hard-margin SVM (unknowns are \mathbf{w} and b)

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

- Objective for the soft-margin SVM (unknowns are \mathbf{w} , b , and $\{\xi_n\}_{n=1}^N$)

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

- In either case, we have to solve a constrained, convex optimization problem



Solving Hard-Margin SVM



Solving Hard-Margin SVM

- The hard-margin SVM optimization problem is:

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

- A constrained optimization problem. Can solve using Lagrange's method
- Introduce **Lagrange Multipliers** α_n ($n = \{1, \dots, N\}$), one for each constraint, and solve

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

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



Solving Hard-Margin SVM

- The dual problem (min then max) is

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

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

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

- Important: Note the form of the solution \mathbf{w} - it is simply a **weighted sum of all the training inputs** $\mathbf{x}_1, \dots, \mathbf{x}_N$ (and α_n is like the “importance” of \mathbf{x}_n)
- Substituting $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ in Lagrangian, we get the dual problem as (verify)

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



Solving Hard-Margin SVM

- Can write the objective more compactly in vector/matrix form as

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

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

- **Good news:** This is **maximizing a concave function** (or minimizing a convex function - verify that the Hessian is \mathbf{G} , which is p.s.d.). Note that our original SVM objective was also convex
- **Important:** Inputs \mathbf{x} 's only appear as **inner products** (helps to “kernelize”; more when we see kernel methods)
- Can solve[†] the above objective function for α using various methods, e.g.,
 - Treating the objective as a **Quadratic Program** (QP) and running some off-the-shelf QP solver such as quadprog (MATLAB), CVXOPT, CPLEX, etc.
 - Using **(projected) gradient methods** (projection needed because the α 's are constrained). Gradient methods will usually be much faster than QP methods.

[†] If interested in more details of the solver, see: “Support Vector Machine Solvers” by Bottou and Lin



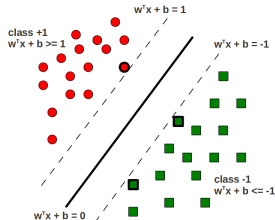
Hard-Margin SVM: The Solution

- Once we have the α_n 's, \mathbf{w} and b can be computed as:

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

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

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



- Reason: **Karush-Kuhn-Tucker (KKT) conditions**

- For the optimal α_n 's

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

- α_n is **non-zero** only if \mathbf{x}_n lies on one of the two **margin boundaries**, i.e., for which $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
- These examples are called **support vectors**
- Recall the support vectors “support” the margin boundaries



Solving Soft-Margin SVM



Solving Soft-Margin SVM

- Recall the soft-margin SVM optimization problem:

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

- Note: $\boldsymbol{\xi} = [\xi_1, \dots, \xi_N]$ is the vector of slack variables
- Introduce **Lagrange Multipliers** α_n, β_n ($n = \{1, \dots, N\}$), for constraints, and solve the Lagrangian:

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

- Note: The terms in red above were not present in the hard-margin SVM
- Two sets of dual variables $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]$ and $\boldsymbol{\beta} = [\beta_1, \dots, \beta_N]$. We'll eliminate the primal variables $\mathbf{w}, b, \boldsymbol{\xi}$ to get dual problem containing the dual variables (just like in the hard margin case)

Solving Soft-Margin SVM

- The Lagrangian problem to solve

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

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

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

- Note: Solution of \mathbf{w} again has the same form as in the hard-margin case (weighted sum of all inputs with α_n being the importance of input \mathbf{x}_n)
- Note: Using $C - \alpha_n - \beta_n = 0$ and $\beta_n \geq 0 \Rightarrow \alpha_n \leq C$ (recall that, for the hard-margin case, $\alpha \geq 0$)
- Substituting these in the Lagrangian \mathcal{L} gives the **Dual** problem

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



Solving Soft-Margin SVM

- Interestingly, the dual variables β don't appear in the objective!
- Just like the hard-margin case, we can write the dual more compactly as

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

where \mathbf{G} is an $N \times N$ matrix with $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$, and $\mathbf{1}$ is a vector of 1s

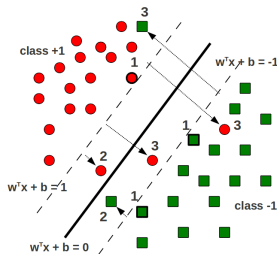
- Like hard-margin case, solving the dual requires concave maximization (or convex minimization)
- Can be solved[†] the same way as hard-margin SVM (except that $\alpha \leq C$)
 - Can solve for α using QP solvers or (projected) gradient methods
- Given α , the solution for \mathbf{w}, b has the same form as hard-margin case
- **Note:** α is again sparse. Nonzero α_n 's correspond to the support vectors

[†] If interested in more details of the solver, see: "Support Vector Machine Solvers" by Bottou and Lin



Support Vectors in Soft-Margin SVM

- The hard-margin SVM solution had only one type of support vectors
 - .. ones that lie on the margin boundaries $\mathbf{w}^T \mathbf{x} + b = -1$ and $\mathbf{w}^T \mathbf{x} + b = +1$
- The soft-margin SVM solution has **three types of support vectors**



- ① Lying on the margin boundaries $\mathbf{w}^T \mathbf{x} + b = -1$ and $\mathbf{w}^T \mathbf{x} + b = +1$ ($\xi_n = 0$)
- ② Lying within the margin region ($0 < \xi_n < 1$) but still on the correct side
- ③ Lying on the wrong side of the hyperplane ($\xi_n \geq 1$)



SVMs via Dual Formulation: Some Comments

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

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

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

- The dual formulation is nice due to two primary reasons:
 - Allows conveniently handling the margin based constraint (via Lagrangians)
 - Important:** Allows learning nonlinear separators by replacing inner products (e.g., $G_{mn} = y_m y_n \mathbf{x}_m^\top \mathbf{x}_n$) by kernelized similarities (kernelized SVMs)
- However, the dual formulation can be expensive if N is large. Have to solve for N variables $\alpha = [\alpha_1, \dots, \alpha_N]$, and also need to store an $N \times N$ matrix \mathbf{G}
- A lot of work[†] on speeding up SVM in these settings (e.g., can use co-ord. descent for α)

[†] See: "Support Vector Machine Solvers" by Bottou and Lin



SVM: Some Notes

- A hugely (perhaps the most!) 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

