Subgradient Descent: Some Examples

Large-Margin Classification (SVM)

CS771: Introduction to Machine Learning Pivush Rai

Sub-gradients

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



Sub-gradients, Sub-differential, and Some Rules

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

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

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

Scaling rule: $\partial (c \cdot f(\mathbf{x})) = c \cdot \partial f(\mathbf{x}) = \{c \cdot \mathbf{v} : \mathbf{v} \in \partial f(\mathbf{x})\}$ is a special case of the more general chain rule

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

- Affine trans: $\partial f(\mathbf{a}^{\mathsf{T}}\mathbf{x} + b) = \mathbf{a} \cdot \partial f(t) = \{\mathbf{a} \cdot c : c \in \partial f(t)\}$, where $t = \mathbf{a}^{\mathsf{T}}\mathbf{x} + b$
- Max rule: If $h(x) = \max\{f(x), g(x)\}$ then we calculate $\partial h(x)$ at 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 $f(\mathbf{x}_*) = g(\mathbf{x}_*), \partial h(\mathbf{x}_*) = \{\alpha \mathbf{a} + (1 \alpha)\mathbf{b} : \mathbf{a} \in \partial f(\mathbf{x}_*), \mathbf{b} \in \partial g(\mathbf{x}_*), \alpha \in [0, 1]\}$
- x_* is a stationary point for a non-diff function f(x) if the zero vector belongs to the sub-differential at x_* , i.e., $0 \in \partial f(x_*)$

The affine transform rule

Subgradient For Regression with Absolute Loss

• Absolute Loss for regression: $L(w) = |y_n - w^T x_n|$



- Can use chain and max rule of sub-diff calculus
- Assume $t = y_n w^T x_n$. Then $\boldsymbol{g}_n = \partial L(\boldsymbol{w}) = -x_n \partial |t|$
 - $\partial L(w) = -x_n \times 1 = -x_n$ if t > 0
 - $\partial L(w) = -x_n \times -1 = x_n$ if t < 0
 - $\partial L(w) = -x_n \times c = -cx_n$ where $c \in [-1, +1]$ if t = 0
 - If we pick c = 0 then $g_n = -\operatorname{sign}(t)x_n$

Can now use subgradient descent or <u>stochastic</u> subgradient descent using these (sub)gradients





Subgradient for Classification with Perceptron Loss⁵

• Perceptron loss for binary classification: $L(w) = \sum_{n=1}^{N} \max\{0, -y_n w^{\mathsf{T}} x_n\}$



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

if
$$y_n w^{(t)} x_n \leq 0$$

Weaning: The current weight
vector does not correctly
predict the label of x_n
 $w^{(t+1)} = w^{(t)} + \eta_t y_n x_n$
Weaning: The current weight
weights make a mistake)

Perceptron Algorithm for Binary Classification

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



An example of an online learning algorithm (processes one training ex. at a time)

• Assuming $w^{(0)} = 0$, easy to see that the final w has the form $w = \sum_{n=1}^{N} \alpha_n y_n x_n$

- α_n is total number of mistakes made by the algorithm on example (x_n, y_n)
- Important: **w** in many classification/regression models can be written as $\mathbf{w} = \sum_{n=1}^{N} \alpha_n y_n \mathbf{x}_n$

Meaning of α_n may be different depending on the problem

Thus **w** is In the span

CS771: Intro to ML

of training inputs

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 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

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

- 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 $y_n w^T x_n \leq \gamma$
- Methods like logistic regression also do not guarantee large margins
- Support Vector Machine (SVM) does it directly by learning the max. margin hyperplane

Learning Large-Margin Hyperplanes



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



Hard-Margin SVM

- Hard-Margin: Every training example must fulfil margin condition $y_n(w^T x_n + b) \ge 1$
- Meaning: Must not have any example in the no-man's land



Soft-Margin SVM (More Commonly Used)



for each training example is just the hinge loss



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

Allow some training examples to fall within the Helps in getting a wider no-man's land (margin region) 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 $(\boldsymbol{x}_n, \boldsymbol{y}_n)$ is known as slack $\xi_n \geq 0$
- $\xi_n > 1$ means totally on the wrong side

 $w'x_n + b \ge 1 - \xi_n$ if $y_n = +1$ $w^{\top} x_n + b \leq -1 + \xi_n$ if $y_n = -1$ $y_n(\mathbf{w}^{\mathsf{T}}\mathbf{x}_n+b) \geq 1 - \xi_n \quad \forall n$ Soft-margin constraint CS771: Intro to ML

Soft-Margin SVM (Contd)

- Goal: Still want to maximize the margin such that
 - Soft-margin constraints $y_n(w^T x_n + b) \ge 1 \xi_n$ 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)



- Lagrange based optimization The objective func. for soft-margin SVM can be used to solve it Trade-off hyperparam Inversely prop. training Constrained optimization to margin error **|w**||² problem with 2N inequality $f(w, b, \xi) =$ min w,b,ξ constraints. Objective and constraints both are convex subject to $y_n(\boldsymbol{w}^T\boldsymbol{x}_n+b) \geq 1-\boldsymbol{\xi}_n, \quad \boldsymbol{\xi}_n \geq 0$ $n=1,\ldots,N$
- 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)

CS771: Intro to ML

Sum of slacks is like

the training error

Solving the SVM Problem



Solving Hard-Margin SVM

The hard-margin SVM optimization problem is

$$\min_{\boldsymbol{w},b} f(\boldsymbol{w},b) = \frac{||\boldsymbol{w}||^2}{2}$$

subject to $1 - y_n(\boldsymbol{w}^T \boldsymbol{x}_n + b) \leq 0, \qquad n = 1, \dots, N$

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

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

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

 $\alpha > 0$

The dual problem (min then max) is

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)

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

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. \boldsymbol{w} and \boldsymbol{b} and setting them to zero gives (verify)

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

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

- The solution \boldsymbol{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) **IMPORTANT:** inputs appear only as pairwise This is also a "quadratic $\max_{\boldsymbol{\alpha} \geq 0} \mathcal{L}_D(\boldsymbol{\alpha}) = \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \alpha_m \alpha_n y_m y_n(\boldsymbol{x}_m^T \boldsymbol{x}_n)$ dot products. This will be useful later on when program" (QP) – a quadratic we make SVM nonlinear using kernel methods function of the variables α $\max_{\alpha \geq 0} \mathcal{L}_D(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha$ G is an $N \times N$ p.s.d. matrix, also called the Gram Maximizing a concave function Matrix, $G_{nm} = y_n y_m x_n^{\mathsf{T}} x_m$, and **1** is a vector of all 1s (or minimizing a convex function) s.t. $\boldsymbol{\alpha} \geq \mathbf{0}$ and $\sum_{n=1}^{N} \alpha_n y_n = 0$. CS771: Intro to ML Many methods to solve it. (Note: For various SVM solvers, can see "Support Vector Machine Solvers" by Bottou and Lin)

Solving Hard-Margin SVM

One we have the α_n 's by solving the dual, we can get w and b as

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

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

• 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 $y_n(w^T x_n + b) = 1$, i.e., the training example lies on the boundary

 $\boldsymbol{\alpha}_n\{1-\boldsymbol{y}_n(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}_n+b)\}=0$

These examples are called support vectors



Solving Soft-Margin SVM

Recall the soft-margin SVM optimization problem

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

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

- 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

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

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

Solving Soft-Margin SVM
The Lagrangian problem to solve

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

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

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

The dual variables β don't

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

Given
$$\boldsymbol{\alpha}$$
, \boldsymbol{w} and \boldsymbol{b} can be
found just like the hard-margin
SVM case

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

Maximizing a convex function)
(or minimizing a convex f

Support Vectors in Soft-Margin SVM

- The hard-margin SVM solution had only one type of support vectors
 - All lied on the supporting hyperplanes $w^T x_n + b = 1$ and $w^T x_n + b = -1$
- The soft-margin SVM solution has <u>three</u> 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:
$$\max_{\boldsymbol{\alpha} \geq 0} \mathcal{L}_D(\boldsymbol{\alpha}) = \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{G} \boldsymbol{\alpha}$$

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 $G_{nm} = y_n y_m x_n^T x_m$ 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 $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, ..., \alpha_N]$
 - Need to pre-compute and store $N \times N$ gram matrix **G**
- Lot of work on speeding up SVM in these settings (e.g., can use co-ord. descent for lpha)

SVM: At Test Time

Prediction for a test point



• 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 *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

$$\arg\min_{\boldsymbol{w},b,\boldsymbol{\xi}} \frac{||\boldsymbol{w}||^2}{2} + C \sum_{n=1}^{N} \xi_n$$

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

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

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

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

A Co-ordinate Ascent Algorithm for SVM

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

λŢ

$$\operatorname{argmax}_{0 \leq \alpha \leq C} \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{m,n=1}^{N} \alpha_m \alpha_n y_m y_n x_m^{\mathsf{T}} x_n$$

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

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

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

$$\operatorname{Can compute these in the beginning itself}$$

$$\operatorname{Can compute these in th$$

- 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 ③) 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

