

Optimization Techniques for ML (wrap-up)

CS771: Introduction to Machine Learning

Piyush Rai

Today

- Some practical aspects for optimization for ML
- Constrained optimization
- Optimization of non-differentiable functions



Some Practical Aspects: Iterate Averaging for SGD ³

- SGD iterates $\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \mathbf{w}^{(3)}, \dots$ can be noisy (recall SGD computes gradients using randomly picked single training example, or a small minibatch)
- Polyak-Ruppert Averaging:** Average SGD iterates and use the average in the end

SGD/mini-batch SGD
update at iteration $t + 1$

Stochastic gradient on a
single/minibatch of examples

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

Running average weight
vector at iteration $t + 1$

Averaged weight vector
at previous iteration t

$$\bar{\mathbf{w}}^{(t+1)} = \frac{t}{t+1} \bar{\mathbf{w}}^{(t)} + \frac{1}{t+1} \mathbf{w}^{(t+1)}$$

This way of computing the average is the
same as doing $\bar{\mathbf{w}}^{(t+1)} = \sum_{i=1}^{t+1} \mathbf{w}^{(i)}$ but to
avoid storing the previous weights, so we
compute a running average

Sometimes, we don't start
averaging from iteration 1 but
after some warm-up iterations



- Averaging is quite popular for SGD. **Stochastic Weighted Averaging (SWA)** is another such recently proposed scheme (similar to Polyak-Ruppert Averaging) used for deep neural networks



Some Practical Aspects: Assessing Convergence

- Various ways to assess convergence, e.g. consider converged if
 - The objective's value (on train set) ceases to change much across iterations

$$L(\mathbf{w}^{(t+1)}) - L(\mathbf{w}^{(t)}) < \epsilon \quad (\text{for some small pre-defined } \epsilon)$$

- The parameter values cease to change much across iterations

$$\|\mathbf{w}^{(t+1)} - \mathbf{w}^{(t)}\| < \tau \quad (\text{for some small pre-defined } \tau)$$

- Above condition is also equivalent to saying that the gradients are close to zero

$$\|\mathbf{g}^{(t)}\| \rightarrow 0$$

Caution: May not yet be at the optima. Use at your own risk!

- The objective's value has become small enough that we are happy with 😊
- Use a validation set to assess if the model's performance is acceptable (early stopping)



Some Practical Aspects: Learning Rate (Step Size)

5

- Some guidelines to select good learning rate (a.k.a. step size) η_t
- For convex functions, setting η_t something like C/t or C/\sqrt{t} often works well
 - These step-sizes are actually theoretically optimal in some settings
 - In general, we want the learning rates to satisfy the following conditions
 - $\eta_t \rightarrow 0$ as t becomes very very large
 - $\sum \eta_t = \infty$ (needed to ensure that we can potentially reach anywhere in the parameter space)
 - Sometimes carefully chosen **constant learning rates** (usually small, or initially large and later small) also work well in practice
- Can also search for the “best” step-size by solving an opt. problem in each step

C is a hyperparameter

Also called
“line search”

$$\eta_t = \arg \min_{\eta \geq 0} f(\mathbf{w}^{(t)} - \eta \cdot \mathbf{g}^{(t)})$$

A one-dim optimization problem
(note that $\mathbf{w}^{(t)}$ and $\mathbf{g}^{(t)}$ are fixed)

- A faster alternative to line search is the **Armijo-Goldstein** rule
 - Starting with current (or some large) learning rate (from prev. iter), and try a few values in decreasing order until the objective's value has a sufficient reduction



Some Practical Aspects: Adaptive Gradient Methods⁶

- Can also use different learning rate in different dimensions

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

Vector of learning rates
along each dimension

Element-wise product of
two vectors

$$e_d^{(t)} = \frac{1}{\sqrt{\epsilon + \sum_{\tau=1}^t \left(g_d^{(\tau)}\right)^2}}$$

If some dimension had big updates recently (marked by large gradient values), slow down along those directions by using smaller learning rates - **AdaGrad** (Duchi et al, 2011)

- Can use a **momentum** term to stabilize gradients by reusing info from past grads
 - Move faster along directions that were previously good
 - Slow down along directions where gradient has changed abruptly

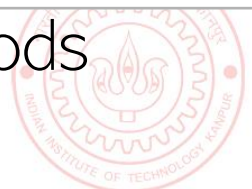
β usually set
as 0.9

The “momentum” term.
Set to 0 at initialization

$$\begin{aligned} \mathbf{m}^{(t)} &= \beta \mathbf{m}^{(t-1)} + \eta_t \mathbf{g}^{(t)} \\ \mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \mathbf{m}^{(t)} \end{aligned}$$

In an even faster version of this, $\mathbf{g}^{(t)}$ is replaced by the gradient computed at the next step if previous direction were used, i.e., $\nabla L(\mathbf{w}^{(t)} - \beta \mathbf{m}^{(t-1)})$.
Called Nesterov’s Accelerated Gradient (NAG) method

- Also exist several more advanced methods that combine the above methods
 - RMS-Prop: AdaGrad + Momentum, Adam: NAG + RMS-Prop
 - These methods are part of packages such as PyTorch, Tensorflow, etc



Constrained Optimization



Projected Gradient Descent

- Consider an optimization problem of the form

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

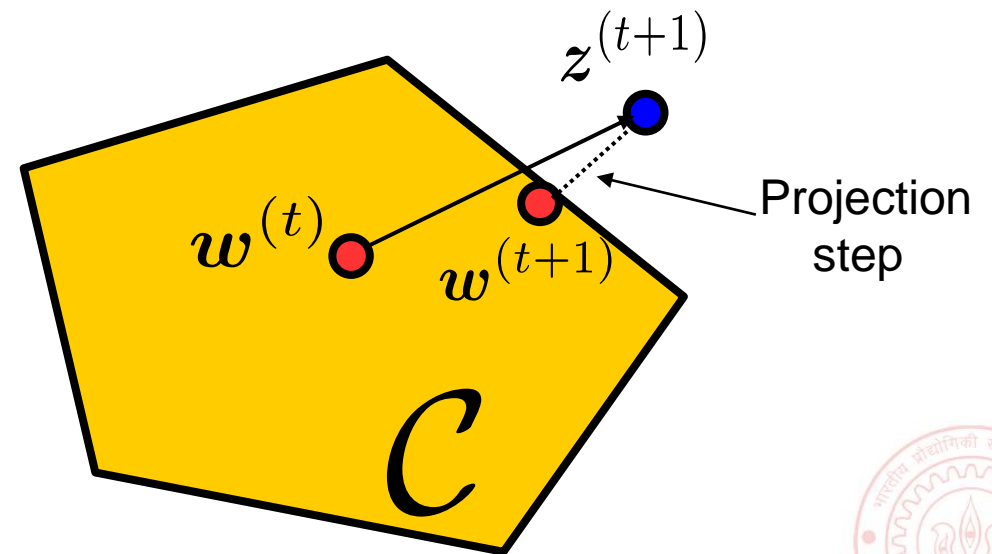
- Projected GD is very similar to GD with an extra **projection step**
- Each iteration t will be of the form

- 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 $\mathbf{w}^{(t+1)} = \Pi_{\mathcal{C}}[\mathbf{z}^{(t+1)}]$

Projection
operator



Projected GD: How to Project?

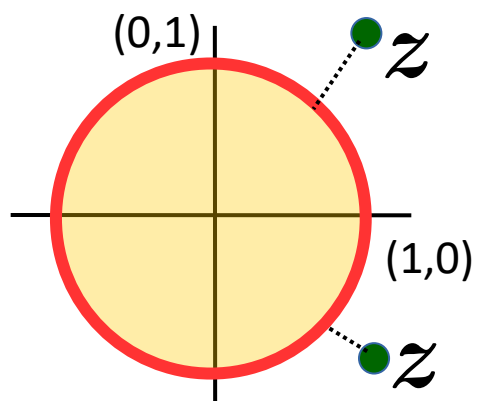
- Here projecting a point means finding the “closest” point from the constraint set

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

Another constrained optimization problem! But simpler to solve! 😊

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

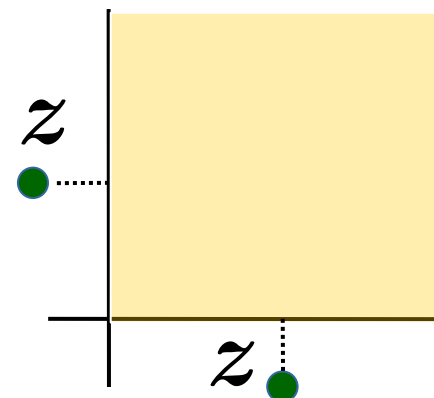
\mathcal{C} : Unit radius ℓ_2 ball



Projection = Normalize to unit Euclidean length vector

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

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



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

$$\hat{x}_i = \begin{cases} x_i & \text{if } x_i \geq 0 \\ 0 & \text{if } x_i < 0 \end{cases}$$

Projected GD commonly used only when the projection step is simple and efficient to compute



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

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

- 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 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

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

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

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

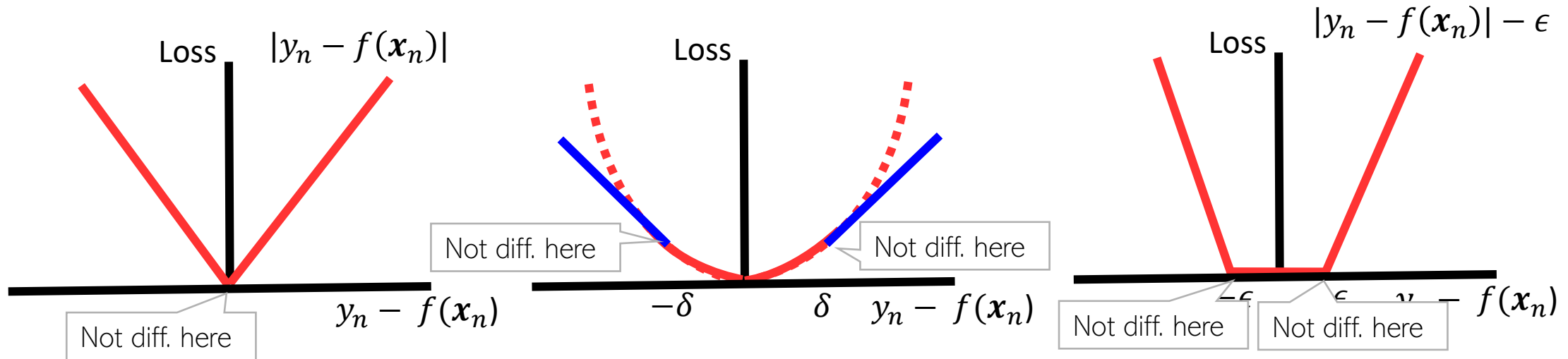


Optimization of Non-differentiable Functions



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 Huber loss, or ϵ -insensitive loss; even ℓ_1 norm regularizer is non-diff

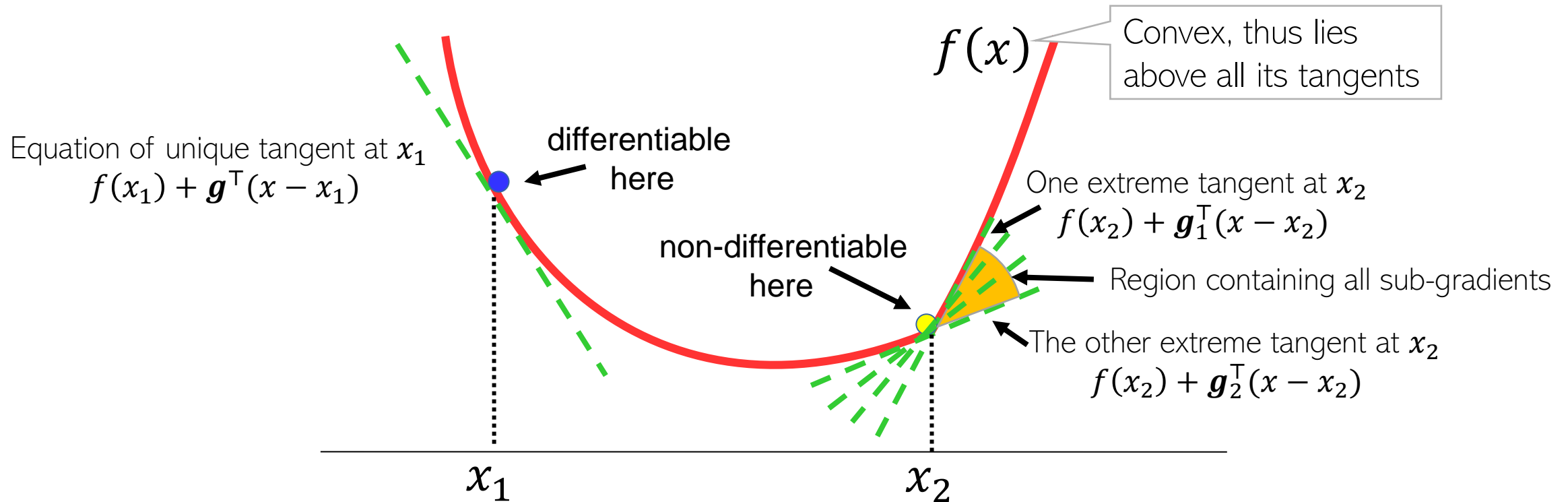


- Basically, any function in which there are points with kink is non-diff
 - At such points, the function is non-differentiable and thus **gradients not defined**
 - Reason: Can't define a unique tangent at such points



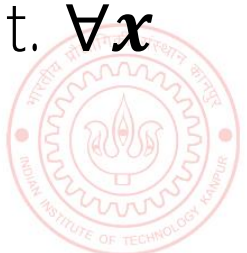
Sub-gradients

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

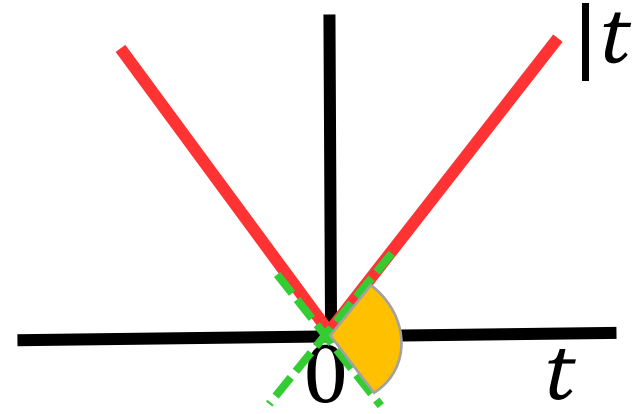
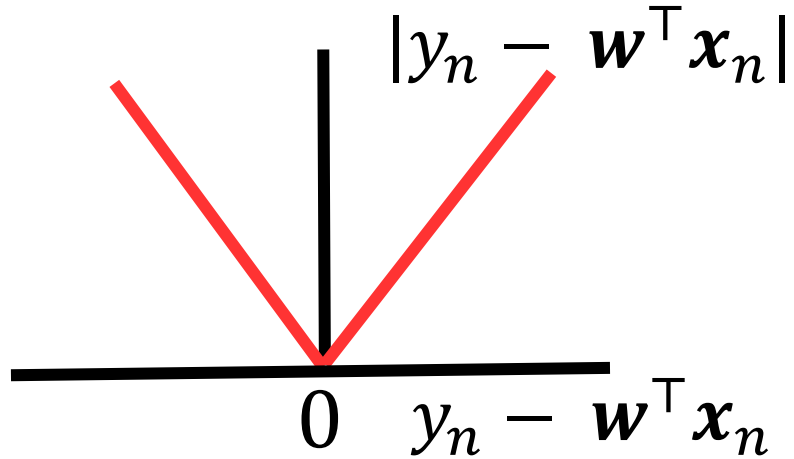
- If $f(\mathbf{x}_*) = g(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \{\alpha \mathbf{a} + (1 - \alpha) \mathbf{b} : \mathbf{a} \in \partial f(\mathbf{x}_*), \mathbf{b} \in \partial g(\mathbf{x}_*), \alpha \in [0, 1]\}$

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

- \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., $\mathbf{0} \in \partial f(\mathbf{x}_*)$



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

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



Optimization for ML: Some Final Comments

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

