

Optimization for ML (2)

CS771: Introduction to Machine Learning

Piyush Rai

The Plan

- Some basic techniques for solving optimization problems
 - First-order optimality
 - Gradient descent
- Dealing with non-differentiable functions
 - Sub-gradients and sub-differential



Optimization Problems in ML

- The general form of an optimization problem in ML will usually be

Usually a sum of the
training error + regularizer

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

- Here $L(\mathbf{w})$ denotes the loss function to be optimized
- \mathcal{C} is the constraint set that the solution must belong to, e.g.,
 - Non-negativity constraint: All entries in \mathbf{w}_{opt} must be non-negative
 - Sparsity constraint: \mathbf{w}_{opt} is a sparse vector with at most K non-zeros
- If no \mathcal{C} is specified, it is an unconstrained optimization problem
- Constrained opt. probs can be converted into unconstrained opt. (will see later)
- For now, assume we have an unconstrained optimization problem

However, possible to have
linear/ridge regression
where solution has some
constraints (e.g., non-neg,
sparsity, or even both)

Linear and ridge regression
that we saw were
unconstrained (\mathbf{w}_{opt} was a
real-valued vector)

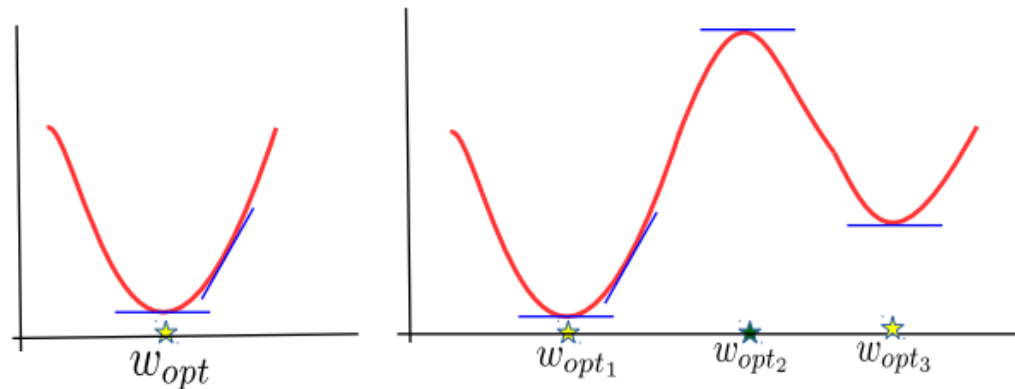


Methods for Solving Optimization Problems



Method 1: Using First-Order Optimality

- Very simple. Already used this approach for linear and ridge regression



Called “first order” since only gradient is used and gradient provides the first order info about the function being optimized



The approach works only for very simple problems where the objective is convex and there are no constraints on the values \mathbf{w} can take

- First order optimality: The gradient \mathbf{g} must be equal to zero at the optima

$$\mathbf{g} = \nabla_{\mathbf{w}}[L(\mathbf{w})] = \mathbf{0}$$

- Sometimes, setting $\mathbf{g} = \mathbf{0}$ and solving for \mathbf{w} gives a closed form solution
- If closed form solution is not available, the gradient vector \mathbf{g} can still be used in iterative optimization algos, like [gradient descent](#)



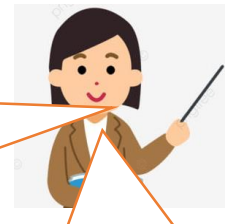
Method 2: Iterative Optimiz. via Gradient Descent⁶



Can I used this approach to solve **maximization** problems?

For max. problems we can use gradient **ascent**
 $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta_t \mathbf{g}^{(t)}$

Iterative since it requires several steps/iterations to find the optimal solution



Fact: Gradient gives the direction of **steepest change** in function's value

Will move in the direction of the gradient

For convex functions, GD will converge to the global minima

Good initialization needed for non-convex functions

Gradient Descent

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

Will see the justification shortly

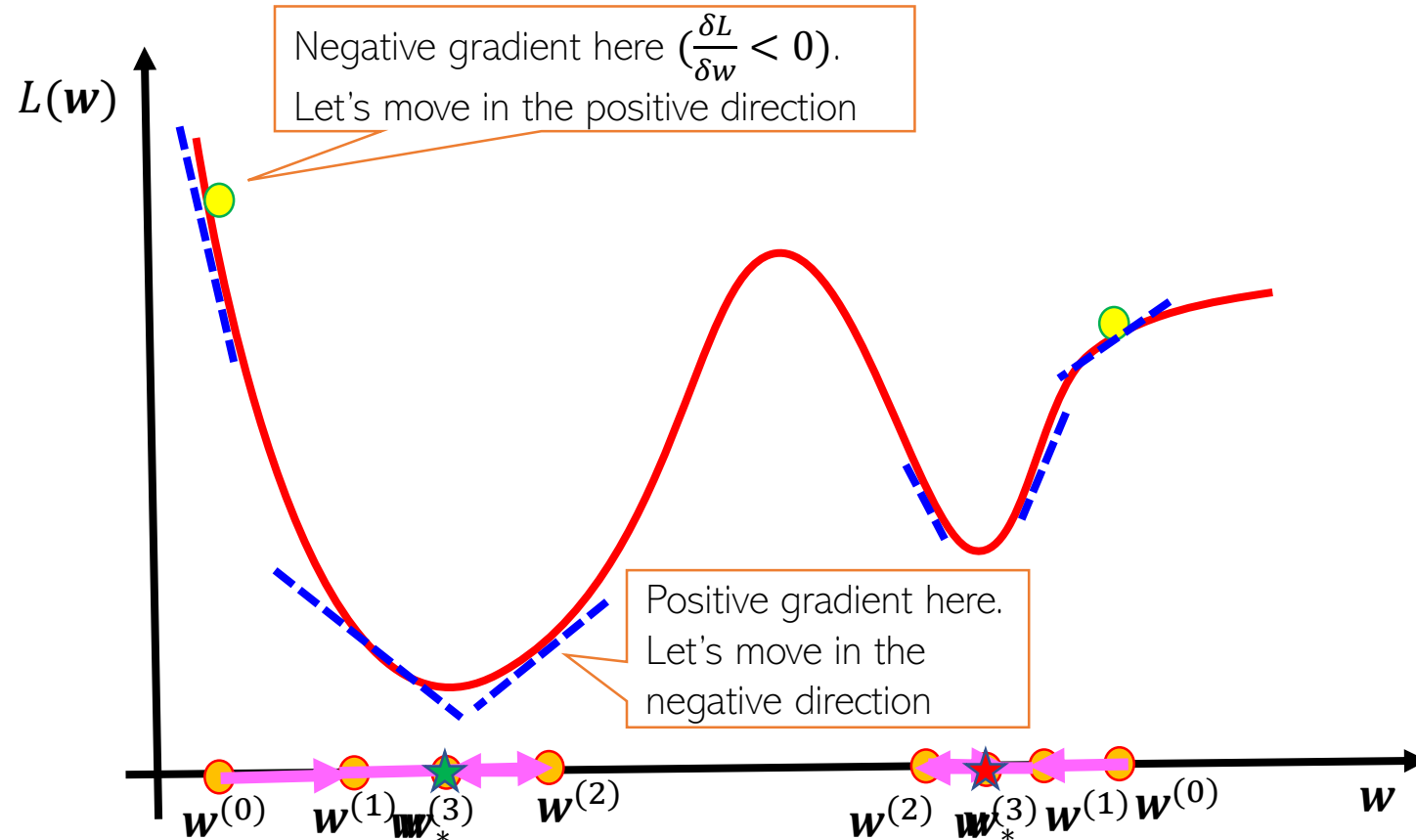
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)}$$

The learning rate very imp. Should be set carefully (fixed or chosen adaptively). Will discuss some strategies later

Sometimes may be tricky to to assess convergence? Will see some methods later

Gradient Descent: An Illustration

7



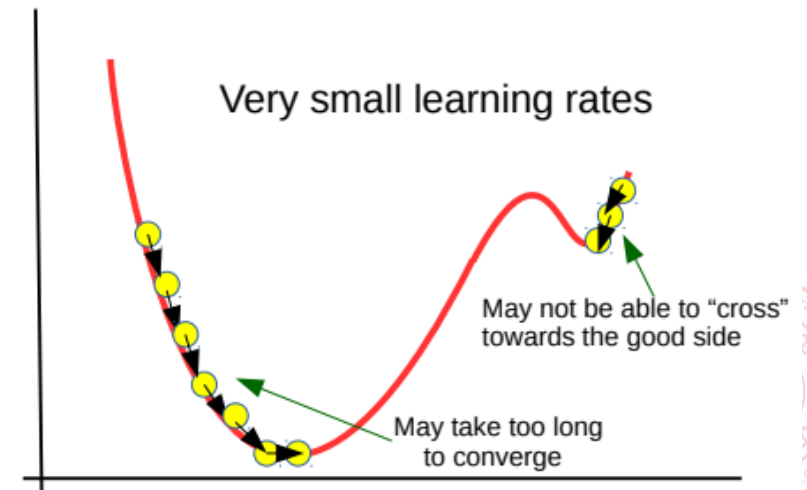
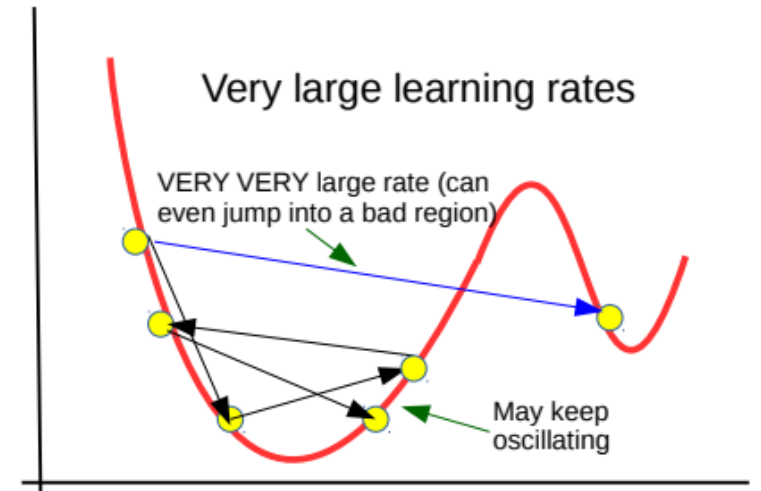
Woohoo! 😊 Global minima found!!!

GD thanks you for the good initialization 😊

Stuck at a local minima 😞

Good initialization is very important

Learning rate is very important



GD: An Example

- Let's apply GD for least squares linear regression

$$\mathbf{w}_{ridge} = \arg \min_{\mathbf{w}} L_{reg}(\mathbf{w}) = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- The gradient: $\mathbf{g} = -\sum_{n=1}^N 2(y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n$

- Each GD update will be of the form

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta_t \sum_{n=1}^N 2 \left(y_n - \mathbf{w}^{(t)\top} \mathbf{x}_n \right) \mathbf{x}_n$$

Prediction error of current model $\mathbf{w}^{(t)}$ on the n^{th} training example

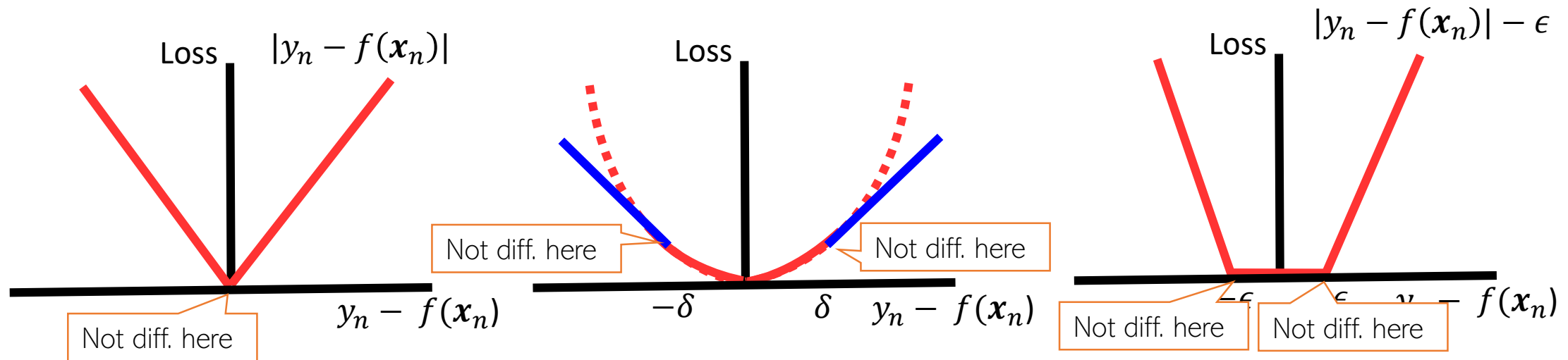
Training examples on which the current model's error is large contribute more to the update

- Exercise: Assume $N = 1$, and show that GD update improves prediction on the training input (\mathbf{x}_n, y_n) , i.e, y_n is closer to $\mathbf{w}^{(t+1)\top} \mathbf{x}_n$ than to $\mathbf{w}^{(t)\top} \mathbf{x}_n$
 - This is sort of a proof that GD updates are “corrective” in nature (and it actually is true not just for linear regression but can also be shown for various other ML models)



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

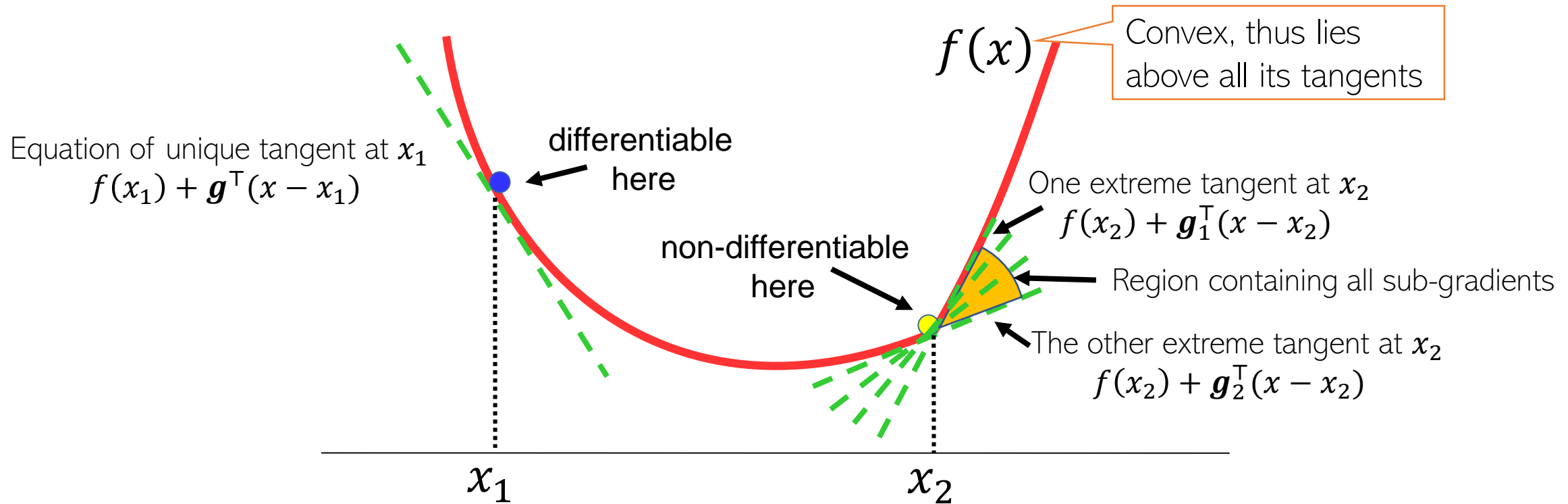


- Basically, any function in which there are points with kink is non-diff
 - At such points, the function is non-differentiable and thus **gradients not defined**
 - Reason: Can't define a unique tangent at such points



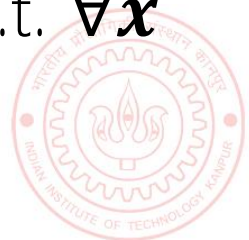
Sub-gradients

- 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: $\partial f(\mathbf{a}^\top \mathbf{x} + b) = \partial f(t) \cdot \mathbf{a} = \{c \cdot \mathbf{a} : c \in \partial f(t)\}$, 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 $g(\mathbf{x}_*) > f(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \partial g(\mathbf{x}_*)$

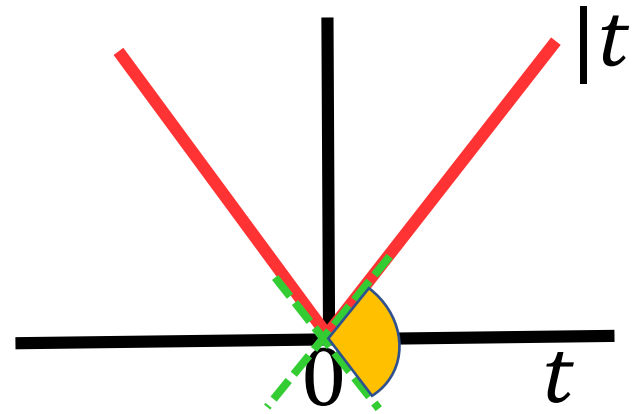
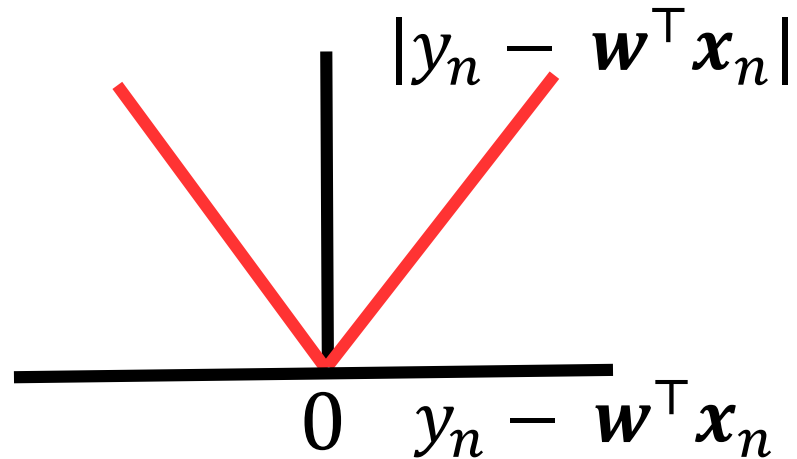
- If $f(\mathbf{x}_*) = g(\mathbf{x}_*)$, $\partial h(\mathbf{x}_*) = \{\alpha \mathbf{a} + (1 - \alpha) \mathbf{b} : \mathbf{a} \in \partial f(\mathbf{x}_*), \mathbf{b} \in \partial g(\mathbf{x}_*), \alpha \in [0, 1]\}$

- \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



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

- 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 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

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



Coming up next

- Making GD faster: Stochastic gradient descent
- Constrained optimization
- Co-ordinate descent
- Alternating optimization
- Practical issue in optimization for ML

