# Optimization Techniques for ML

CS771: Introduction to Machine Learning Piyush Rai



## Derivatives

Will sometimes use f'(x) to denote the derivative

• Magnitude of derivative at a point is the rate of change of the func at that point



- Derivative becomes zero at stationary points (optima or saddle points)
  - The function becomes "flat" ( $\Delta f(x) = 0$  if we change x by a very little at such points)
  - These are the points where the function has its maxima/minima (unless they are saddles)

# Saddle Points

Points where derivative is zero but are neither minima nor maxima



Saddle points are very common for loss functions of deep learning models

Need to be handled carefully during optimization

Second or higher derivative may help identify if a stationary point is a saddle

CS771: Intro to ML

# Rules of Derivatives

Some basic rules of taking derivatives

- Sum Rule: (f(x) + g(x))' = f'(x) + g'(x)
- Scaling Rule:  $(a \cdot f(x))' = a \cdot f'(x)$  if a is not a function of x
- Product Rule:  $(f(x) \cdot g(x))' = f'(x) \cdot g(x) + g'(x) \cdot f(x)$
- Quotient Rule:  $(f(x)/g(x))' = (f'(x) \cdot g(x) g'(x)f(x))/(g(x))^2$
- Chain Rule:  $(f(g(x)))' \stackrel{\text{\tiny def}}{=} (f \circ g)'(x) = f'(g(x)) \cdot g'(x)$



We already used some of these (sum, scaling and chain) when calculating the derivative for the linear regression model



## Derivatives

How the derivative itself changes tells us about the function's optima



- The second derivative f''(x) can provide this information
  - It is the rate of change of the first derivative



# Multivariate Functions

- Most functions that we see in ML are multivariate function
- Example: Loss fn L(w) in lin-reg was a multivar function of D-dim vector w $L(w): \mathbb{R}^D \to \mathbb{R}$
- Here is an illustration of a function of 2 variables (4 maxima and 5 minima)



Plot courtesy: http://benchmarkfcns.xyz/benchmarkfcns/griewankfcn.html

# Derivatives of Multivariate Functions

- Can define derivative for a multivariate functions as well via the gradient
- <u>Gradient</u> of a function  $f(\mathbf{x}): \mathbb{R}^D \to \mathbb{R}$  is a  $D \times 1$  <u>vector</u> of <u>partial</u> derivatives

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_D}\right) =$$

Each element in this gradient vector tells us how much f will change if we move a little along the corresponding (akin to one-dim case)

- Conditions for optima defined similar to one-dim case
  - Required properties for one-dim case must be satisfied for all components of x
- The second derivative in the multivariate case is known as the Hessian matrix



# The Hessian

• For a multivar scalar valued function  $f(\mathbf{x}): \mathbb{R}^D \to \mathbb{R}$ , Hessian is a  $D \times D$  matrix



- The Hessian matrix can be used to assess the optima/saddle points
  - $\nabla f(\mathbf{x}) = 0$  and  $\nabla^2 f(\mathbf{x})$  is a positive semi-definite (PSD) matrix then  $\mathbf{x}$  is a minima
  - $\nabla f(\mathbf{x}) = 0$ , and  $\nabla^2 f(\mathbf{x})$  is a negative semi-definite (NSD) matrix then  $\mathbf{x}$  is a maximal
  - The determinant of the Hessian is zero at saddle point (has at least one zero eig-val)

<sup>9</sup> 

# Convex and Non-Convex Functions

- A function being optimized can be either convex or non-convex
- Here are a couple of examples of convex functions



Here are a couple of examples of non-convex functions



Non-convex functions have multiple minima. Usually harder to optimize as compared to convex functions

Loss functions of most deep learning models are non-convex



#### Convex Sets

• A set S of points is a convex set, if for any two points  $x, y \in S$ , and  $0 \le \alpha \le 1$ 

z is also called a "convex combination" of two points 
$$z = \alpha x + (1 - \alpha)y \in S$$



• Above means that all points on the line-segment between x and y lie within S



The domain of a convex function needs to be a convex set



## **Convex Functions**

• Informally, f(x) is convex if all of its chords lie above the function everywhere



Note: "Chord lies above function" more formally means If *f* is convex then given  $\alpha_1, \ldots, \alpha_n$  s.t  $\sum_{i=1}^n \alpha_i = 1$  $f\left(\sum_{i=1}^n \alpha_i x_i\right) \le \sum_{i=1}^n \alpha_i f(x_i)$ Jensen's Inequality

Formally, (assuming differentiable function), some tests for convexity:

• First-order convexity (graph of f must be above all the tangents)

$$f(y) \qquad f(y) = f(x) + \nabla f(x)^{T}(y - x)$$

$$f(x) + \nabla f(x)^{T}(y - x)$$

Evercical Chave that

#### Some Basic Rules for Convex Functions

- Some basic rules to check if f(x) is convex or not
  - All linear and affine functions (e.g., ax + b) are convex
  - $\exp(ax)$  is convex for  $x \in \mathbb{R}$ , for any  $a \in \mathbb{R}$
  - log(x) is concave (not convex) for x > 0
  - $x^a$  is convex for x > 0, for any  $a \ge 1$  and a < 0, concave for  $0 \le a \le 1$
  - $|x|^a$  is convex for  $x \in \mathbb{R}$ , for any  $a \ge 1$
  - All norms in  $\mathbb{R}^D$  are convex
  - Non-negative weighted sum of convex functions is also a convex function
  - Affine transformation preserves convexity: if f(x) is convex then f(x) = f(ax + b) is also convex
  - Some rules to check whether composition f(x) = h(g(x)) of two functions h and g is convex

f is convex if h is convex and nondecreasing, and g is convex, f is convex if h is convex and nonincreasing, and g is concave, f is concave if h is concave and nondecreasing, and g is concave, f is concave if h is concave and nonincreasing, and g is convex.

# Optimization Problems in ML

- The general form of an optimization problem in ML will usually be Usually a sum of the training error + regularizer  $w_{opt} = \arg \min_{w \in C} L(w)$
- Here L(w) denotes the loss function to be optimized
- C is the constraint set that the solution must belong to, e.g.,
  - Non-negativity constraint: All entries in  $w_{opt}$  must be non-negative
  - Sparsity constraint:  $w_{opt}$  is a sparse vector with atmost K non-zeros
- If no  $\boldsymbol{C}$  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)

14

Linear and ridge regression that we saw were unconstrained ( $w_{opt}$  was a real-valued vector)

# Methods for Solving Optimization Problems



15

# 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  $\boldsymbol{w}$  can take

ullet First order optimality: The gradient  $oldsymbol{g}$  must be equal to zero at the optima

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

-Sometimes, setting g = 0 and solving for w gives a closed form solution

- If closed form solution is not available, the gradient vector  $\boldsymbol{g}$  can still be used in iterative optimization algos, like gradient descent

16



# Gradient Descent: An Illustration



18

# GD: An Example

Let's apply GD for least squares linear regression

$$\mathbf{w}_{ridge}$$
 = arg min<sub>w</sub>  $L_{reg}(\mathbf{w})$  = arg min<sub>w</sub>  $\sum_{n=1}^{N} (y_n - \mathbf{w}^{\mathsf{T}} \mathbf{x}_n)^2$ 

• The gradient: 
$$\boldsymbol{g} = -\sum_{n=1}^{N} 2(y_n - \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n) \boldsymbol{x}_n$$

Each GD update will be of the form

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

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

$$w^{(t+1)} = w^{(t)} + \eta_t \sum_{n=1}^N 2(y_n - w^{(t)^{\mathsf{T}}} x_n) x_n$$

• Exercise: Assume N = 1, and show that GD update improves prediction on the training input  $(x_n, y_n)$ , i.e.,  $y_n$  is closer to  $w^{(t+1)^{\mathsf{T}}} x_n$  than to  $w^{(t)^{\mathsf{T}}} x_n$ 

 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)

# Next class

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

