

COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2023.

Lecture 22

Logistics

- Problem Set 4 core problems are due Friday. Challenge problems are due Monday.
- Quiz due Monday (last of the semester).
- Problem Set 5 will be released Friday or Saturday, and is optional. The core problems can be used to replace the lowest core problem grade on a previous problem set. It will contain three challenge problems as well.
- Final exam study material has been released on the course webpage/Moodle. I will announce additional office hours for final review shortly.
- Next Monday 12/4 at **3pm in CS 140** I will hold another linear algebra review session.
- Please fill out SRTIs (course reviews)!

Summary

Last Class Before Break: Fast computation of the SVD/eigendecomposition.

- Power method for approximating the top eigenvector of a matrix.
- Analysis of convergence rate – we didn't quite finish but we covered the most important ideas.
- Convergence rate depends on the gap between the largest and second largest eigenvalue.

Final Three Classes:

- General iterative algorithms for optimization, specifically **gradient descent** and its variants.
- What are these methods, when are they applied, and how do you analyze their performance?
- Small taste of what you can find in COMPSCI 5900P or 6900P.

Discrete vs. Continuous Optimization

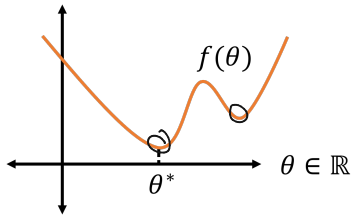
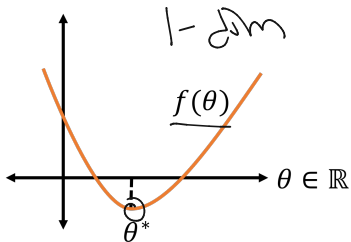
Discrete (Combinatorial) Optimization: (~~traditional CS algorithms~~)

- Graph Problems: min-cut, max flow, shortest path, matchings, maximum independent set, traveling salesman problem
- Problems with discrete constraints or outputs: bin-packing, scheduling, sequence alignment, submodular maximization
- Generally searching over a finite but exponentially large set of possible solutions. Many of these problems are NP-Hard.

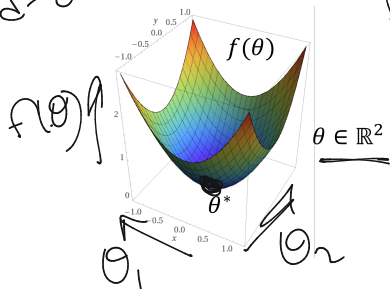
Continuous Optimization: (maybe seen in ML/advanced algorithms)

- Unconstrained convex and non-convex optimization.
- Linear programming, quadratic programming, semidefinite programming

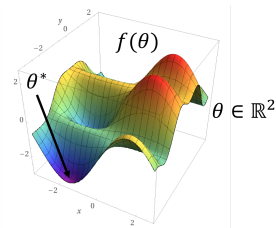
Continuous Optimization Examples



d-dim



$f: \mathbb{R}^2 \rightarrow \mathbb{R}$



Mathematical Setup

Given some function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, find $\vec{\theta}_*$ with:

$$\underline{f(\vec{\theta}_*) = \min_{\vec{\theta} \in \mathbb{R}^d} f(\vec{\theta})}$$

$$\vec{\theta}_* = \arg \min_{\theta \in \mathbb{R}^d} f(\theta)$$

Mathematical Setup

Given some function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, find $\vec{\theta}_*$ with:

$$f(\vec{\theta}_*) = \min_{\vec{\theta} \in \mathbb{R}^d} f(\vec{\theta}) + \epsilon$$

Typically up to some small approximation factor.

Mathematical Setup

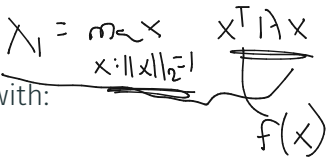
Given some function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, find $\vec{\theta}_*$ with:

$$f(\vec{\theta}_*) = \min_{\substack{\vec{\theta} \in \mathbb{R}^d \\ \|\vec{\theta}\|_2 \leq 1}} f(\vec{\theta}) + \epsilon$$

Typically up to some small approximation factor.

Often under some constraints:

$$\left[\begin{array}{l} \cdot \|\vec{\theta}\|_2 \leq 1, \quad \|\vec{\theta}\|_1 \leq 1. \\ \cdot A\vec{\theta} \leq \vec{b}, \quad \vec{\theta}^T A \vec{\theta} \geq 0. \\ \cdot \sum_{i=1}^d \vec{\theta}(i) \leq c. \end{array} \right.$$

$$\lambda_1 = \max_{x: \|x\|_2=1} x^T A x$$


The diagram shows a handwritten parabola opening downwards, representing the function $f(x)$. A horizontal line is drawn across the peak of the parabola, labeled λ_1 . The x-axis is labeled $x: \|x\|_2=1$. The function is labeled $f(x)$ at the bottom right.

Why Continuous Optimization?

Modern machine learning centers around continuous optimization.

Typical Set Up: (supervised machine learning)

- Have a **model**, which is a function mapping inputs to predictions (neural network, linear function, low-degree polynomial etc).
- The model is parameterized by a **parameter vector** (weights in a neural network, coefficients in a linear function or polynomial)
- Want to **train** this model on input data, by picking a parameter vector such that the model does a good job mapping inputs to predictions on your training data.

This training step is typically formulated as a continuous optimization problem.

Optimization in ML

Example: Linear Regression

Optimization in ML

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=} \langle \vec{\theta}, \vec{x} \rangle$

$$\begin{bmatrix} 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=}} \langle \vec{\theta}, \vec{x} \rangle = \vec{\theta}(1) \cdot \vec{x}(1) + \dots + \vec{\theta}(d) \cdot \vec{x}(d)$.

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=}} \langle \vec{\theta}, \vec{x} \rangle = \vec{\theta}(1) \cdot \vec{x}(1) + \dots + \vec{\theta}(d) \cdot \vec{x}(d)$.

Parameter Vector: $\vec{\theta} \in \mathbb{R}^d$ (the regression coefficients)

Optimization in ML

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=}} \langle \vec{\theta}, \vec{x} \rangle = \vec{\theta}(1) \cdot \vec{x}(1) + \dots + \vec{\theta}(d) \cdot \vec{x}(d)$.

Parameter Vector: $\vec{\theta} \in \mathbb{R}^d$ (the regression coefficients)

Optimization Problem: Given data points (training points) $\vec{x}_1, \dots, \vec{x}_n$ (the rows of data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$) and labels $y_1, \dots, y_n \in \mathbb{R}$, find $\vec{\theta}_*$ minimizing the **loss function**:

$$L(\vec{\theta}, \mathbf{X}, \vec{y}) = \sum_{i=1}^n \ell(M_{\vec{\theta}}(\vec{x}_i), y_i)$$

Handwritten notes: "loss" above the sum, and $(m_{\theta}(x_i) - y_i)^2$ to the right of the sum.

where ℓ is some measurement of how far $M_{\vec{\theta}}(\vec{x}_i)$ is from y_i .

Optimization in ML

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=}} \langle \vec{\theta}, \vec{x} \rangle = \vec{\theta}(1) \cdot \vec{x}(1) + \dots + \vec{\theta}(d) \cdot \vec{x}(d)$.

Parameter Vector: $\vec{\theta} \in \mathbb{R}^d$ (the regression coefficients)

Optimization Problem: Given data points (training points) $\vec{x}_1, \dots, \vec{x}_n$ (the rows of data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$) and labels $y_1, \dots, y_n \in \mathbb{R}$, find $\vec{\theta}_*$ minimizing the **loss function**:

$$L(\vec{\theta}, \mathbf{X}, \vec{y}) = \sum_{i=1}^n \ell(M_{\vec{\theta}}(\vec{x}_i), y_i)$$

where ℓ is some measurement of how far $M_{\vec{\theta}}(\vec{x}_i)$ is from y_i .

- $\ell(M_{\vec{\theta}}(\vec{x}_i), y_i) = (M_{\vec{\theta}}(\vec{x}_i) - y_i)^2$ (least squares regression)
- $y_i \in \{-1, 1\}$ and $\ell(M_{\vec{\theta}}(\vec{x}_i), y_i) = \ln(1 + \exp(-y_i M_{\vec{\theta}}(\vec{x}_i)))$ (logistic regression)

Optimization in ML

Example: Linear Regression

Model: $M_{\vec{\theta}} : \mathbb{R}^d \rightarrow \mathbb{R}$ with $M_{\vec{\theta}}(\vec{x}) \stackrel{\text{def}}{=} \langle \vec{\theta}, \vec{x} \rangle = \vec{\theta}(1) \cdot \vec{x}(1) + \dots + \vec{\theta}(d) \cdot \vec{x}(d)$.

Parameter Vector: $\vec{\theta} \in \mathbb{R}^d$ (the regression coefficients)

Optimization Problem: Given data points (training points) $\vec{x}_1, \dots, \vec{x}_n$ (the rows of data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$) and labels $y_1, \dots, y_n \in \mathbb{R}$, find $\vec{\theta}_*$ minimizing the **loss function**:

$$L_{\mathbf{X}, \mathbf{y}}(\vec{\theta}) = L(\vec{\theta}, \mathbf{X}, \vec{y}) = \sum_{i=1}^n \ell(M_{\vec{\theta}}(\vec{x}_i), y_i)$$

where ℓ is some measurement of how far $M_{\vec{\theta}}(\vec{x}_i)$ is from y_i .

- $\ell(M_{\vec{\theta}}(\vec{x}_i), y_i) = (M_{\vec{\theta}}(\vec{x}_i) - y_i)^2$ (least squares regression)
- $y_i \in \{-1, 1\}$ and $\ell(M_{\vec{\theta}}(\vec{x}_i), y_i) = \ln(1 + \exp(-y_i M_{\vec{\theta}}(\vec{x}_i)))$ (logistic regression)

Optimization in ML

$$L_{\mathbf{X}, \mathbf{y}}(\vec{\theta}) = \sum_{i=1}^n \ell(M_{\vec{\theta}}(\vec{x}_i), y_i)$$

- Supervised means we have labels y_1, \dots, y_n for the training points.
- Solving the final optimization problem has many different names: likelihood maximization, empirical risk minimization, minimizing training loss, etc.
- Continuous optimization is also very common in unsupervised learning. (PCA, spectral clustering, etc.)
- Generalization tries to explain why minimizing the loss $L_{\mathbf{X}, \mathbf{y}}(\vec{\theta})$ on the *training points* minimizes the loss on future *test points*. I.e., makes us have good predictions on future inputs.

Optimization Algorithms

Choice of optimization algorithm for minimizing $f(\vec{\theta})$ will depend on many things:

- The form of f (in ML, depends on the model & loss function).
- Any constraints on $\vec{\theta}$ (e.g., $\|\vec{\theta}\| < c$).
- Computational constraints, such as memory constraints.

$$\underline{L_{\mathbf{x}, \mathbf{y}}(\vec{\theta})} = \sum_{i=1}^n \underline{\ell(M_{\vec{\theta}}(\vec{x}_i), y_i)}$$

Optimization Algorithms

Choice of optimization algorithm for minimizing $f(\vec{\theta})$ will depend on many things:

- The form of f (in ML, depends on the model & loss function).
- Any constraints on $\vec{\theta}$ (e.g., $\|\vec{\theta}\| < c$).
- Computational constraints, such as memory constraints.

$$f(x) = x^2$$

$\frac{df}{dx} = 2x = 0$ when $x=0$

$$L_{\mathbf{x}, \mathbf{y}}(\vec{\theta}) = \sum_{i=1}^n \ell(M_{\vec{\theta}}(\vec{x}_i), y_i)$$



What are some popular optimization algorithms?

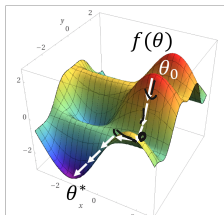
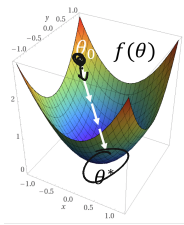
grid search
random search
gradient descent

ADAM
ADAGRAD
LBFGS

Gradient Descent

Next few classes: Gradient descent (and some important variants)

- An extremely simple greedy iterative method, that can be applied to almost any continuous function we care about optimizing.
- Often not the 'best' choice for any given function, but it is the approach of choice in ML since it is simple, general, and often works very well.
- At each step, tries to move towards the lowest nearby point in the function that is can – in the opposite direction of the gradient.



Multivariate Calculus Review

Let $\vec{e}_i \in \mathbb{R}^d$ denote the i^{th} standard basis vector,

$$\vec{e}_i = \underbrace{[0, 0, 1, 0, 0, \dots, 0]}_{1 \text{ at position } i}.$$

Multivariate Calculus Review

Let $\vec{e}_i \in \mathbb{R}^d$ denote the i^{th} standard basis vector,

$$\vec{e}_i = \underbrace{[0, 0, 1, 0, 0, \dots, 0]}_{1 \text{ at position } i}.$$

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

Partial Derivative:

$$\frac{\partial f}{\partial \theta(i)} = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta} + \epsilon \cdot \vec{e}_i) - f(\vec{\theta})}{\epsilon}.$$

$\frac{\partial f}{\partial \theta(i)}$

$$\begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 2 + \epsilon \\ 1 \\ 3 \end{bmatrix}$$

Multivariate Calculus Review

Let $\vec{e}_i \in \mathbb{R}^d$ denote the i^{th} standard basis vector,

$$\vec{e}_i = \underbrace{[0, 0, 1, 0, 0, \dots, 0]}_{\text{1 at position } i}.$$

Partial Derivative:

$$\frac{\partial f}{\partial \vec{\theta}(i)} = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta} + \epsilon \cdot \vec{e}_i) - f(\vec{\theta})}{\epsilon}.$$

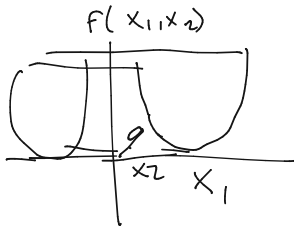
Directional Derivative:

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \epsilon \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad \frac{D_{\vec{v}} f(\vec{\theta})}{\epsilon} = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta} + \epsilon \vec{v}) - f(\vec{\theta})}{\epsilon}.$$

Multivariate Calculus Review

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \quad \nabla f \in \mathbb{R}^d$$

Gradient: Just a 'list' of the partial derivatives.



$$\vec{\nabla} f(\vec{\theta}) = \begin{bmatrix} \frac{\partial f}{\partial \theta(1)} \\ \frac{\partial f}{\partial \theta(2)} \\ \vdots \\ \frac{\partial f}{\partial \theta(d)} \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} 0 \\ x_2 \end{bmatrix}$$

Multivariate Calculus Review

Gradient: Just a 'list' of the partial derivatives.

$$\vec{\nabla} f(\vec{\theta}) = \begin{bmatrix} \frac{\partial f}{\partial \theta(1)} \\ \frac{\partial f}{\partial \theta(2)} \\ \vdots \\ \frac{\partial f}{\partial \theta(d)} \end{bmatrix}$$

Directional Derivative in Terms of the Gradient:

$$D_{\vec{v}} f(\vec{\theta}) = \langle \vec{v}, \vec{\nabla} f(\vec{\theta}) \rangle.$$

$$v(1) \cdot \frac{\partial f}{\partial \theta(1)} + \frac{v(2) \partial f}{\partial \theta(2)}$$

Function Access

Often the functions we are trying to optimize are very complex (e.g., a neural network). We will assume access to:

Function Evaluation: Can compute $f(\vec{\theta})$ for any $\vec{\theta}$.

Gradient Evaluation: Can compute $\vec{\nabla}f(\vec{\theta})$ for any $\vec{\theta}$.

Function Access

Often the functions we are trying to optimize are very complex (e.g., a neural network). We will assume access to:

Function Evaluation: Can compute $f(\vec{\theta})$ for any $\vec{\theta}$.

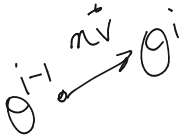
Gradient Evaluation: Can compute $\vec{\nabla}f(\vec{\theta})$ for any $\vec{\theta}$.

In neural networks:

- Function evaluation is called a **forward pass** (propagate an input through the network).
- Gradient evaluation is called a **backward pass** (compute the gradient via chain rule, using backpropagation).

Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\underline{\theta^{(0)}}$, in each iteration let $\underline{\theta^{(i)}} = \underline{\theta^{(i-1)}} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\underline{\theta^{(i-1)}} + \eta \vec{v})$.



Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta} + \epsilon \vec{v}) - f(\vec{\theta})}{\epsilon}.$$

Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}^{(i-1)}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta}^{(i-1)} + \epsilon \vec{v}) - f(\vec{\theta}^{(i-1)})}{\epsilon} \Bigg|$$

Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}^{(i-1)}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta}^{(i-1)} + \epsilon \vec{v}) - f(\vec{\theta}^{(i-1)})}{\epsilon}.$$

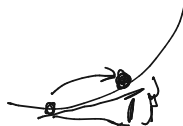
So for small η :

$$\underline{f(\vec{\theta}^{(i)})} - \underline{f(\vec{\theta}^{(i-1)})} = \underline{f(\vec{\theta}^{(i-1)} + \eta \vec{v})} - \underline{f(\vec{\theta}^{(i-1)})}$$

Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}^{(i-1)}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta}^{(i-1)} + \epsilon \vec{v}) - f(\vec{\theta}^{(i-1)})}{\epsilon}.$$



So for small η :

$$f(\vec{\theta}^{(i)}) - f(\vec{\theta}^{(i-1)}) = f(\vec{\theta}^{(i-1)} + \eta \vec{v}) - f(\vec{\theta}^{(i-1)}) \approx \eta \cdot \underline{D_{\vec{v}} f(\vec{\theta}^{(i-1)})}$$

Gradient Descent Greedy Approach

Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}^{(i-1)}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta}^{(i-1)} + \epsilon \vec{v}) - f(\vec{\theta}^{(i-1)})}{\epsilon}.$$

So for small η :

$$\begin{aligned} \underline{f(\vec{\theta}^{(i)}) - f(\vec{\theta}^{(i-1)})} &= f(\vec{\theta}^{(i-1)} + \eta \vec{v}) - f(\vec{\theta}^{(i-1)}) \approx \eta \cdot \underline{D_{\vec{v}} f(\vec{\theta}^{(i-1)})} \\ &= \eta \cdot \underline{\langle \vec{v}, \vec{\nabla} f(\vec{\theta}^{(i-1)}) \rangle}. \end{aligned}$$

Gradient Descent Greedy Approach

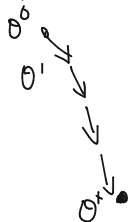
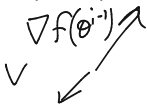
Gradient descent is a **greedy** iterative optimization algorithm:
Starting at $\vec{\theta}^{(0)}$, in each iteration let $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} + \eta \vec{v}$, where η is a (small) 'step size' and \vec{v} is a direction chosen to minimize $f(\vec{\theta}^{(i-1)} + \eta \vec{v})$.

$$D_{\vec{v}} f(\vec{\theta}^{(i-1)}) = \lim_{\epsilon \rightarrow 0} \frac{f(\vec{\theta}^{(i-1)} + \epsilon \vec{v}) - f(\vec{\theta}^{(i-1)})}{\epsilon}.$$

So for small η :

$$\begin{aligned} \underline{f(\vec{\theta}^{(i)}) - f(\vec{\theta}^{(i-1)})} &= \underline{f(\vec{\theta}^{(i-1)} + \eta \vec{v}) - f(\vec{\theta}^{(i-1)})} \approx \eta \cdot D_{\vec{v}} f(\vec{\theta}^{(i-1)}) \\ &= \eta \cdot \langle \vec{v}, \underline{\vec{\nabla} f(\vec{\theta}^{(i-1)})} \rangle. \end{aligned}$$

We want to choose \vec{v} **minimizing** $\langle \vec{v}, \vec{\nabla} f(\vec{\theta}^{(i-1)}) \rangle$ – i.e., pointing in the direction of $\vec{\nabla} f(\vec{\theta}^{(i-1)})$ but with the opposite sign.



Gradient Descent Psuedocode

Gradient Descent

- Choose some initialization $\vec{\theta}^{(0)}$.
- For $i = 1, \dots, t$
 - $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} - \eta \nabla f(\vec{\theta}^{(i-1)})$
- Return $\vec{\theta}^{(t)}$, as an approximate minimizer of $f(\vec{\theta})$.

Step size η is chosen ahead of time or adapted during the algorithm (details to come.)

Gradient Descent Psuedocode

Gradient Descent

- Choose some initialization $\vec{\theta}^{(0)}$.
- For $i = 1, \dots, t$
 - $\vec{\theta}^{(i)} = \vec{\theta}^{(i-1)} - \eta \nabla f(\vec{\theta}^{(i-1)})$
- Return $\vec{\theta}^{(t)}$, as an approximate minimizer of $f(\vec{\theta})$.

Step size η is chosen ahead of time or adapted during the algorithm (details to come.)

- For now assume η stays the same in each iteration.

When Does Gradient Descent Work?

$$\frac{\theta_0 - \theta^*}{\nabla f(\theta_0)} = \mu^*$$

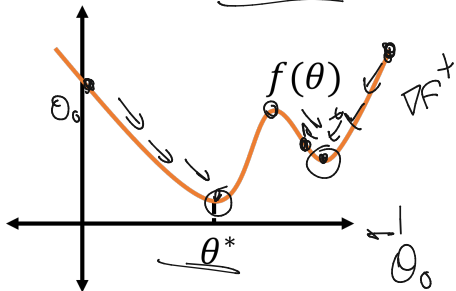
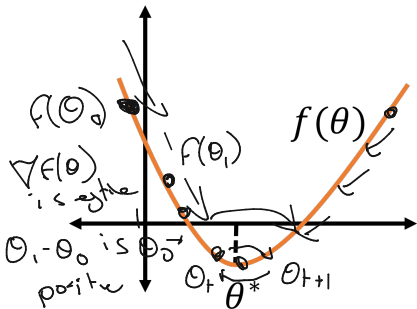
convex

$$\theta \in \mathbb{R}$$

$$\frac{df}{d\theta} = f'(\theta)$$

$$\nabla f(\theta) \in \mathbb{R}$$

non-convex



Gradient Descent Update: $\vec{\theta}_{i+1} = \vec{\theta}_i - \eta \nabla f(\vec{\theta}_i)$

$$\underline{\mu} \cdot \nabla f(\theta^{i-1}) = \underline{\mu} \cdot f'(\theta^{i-1})$$



Convexity

Definition – Convex Function: A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$\checkmark \Rightarrow \overset{.5}{(1-\lambda)} \cdot f(\vec{\theta}_1) + \overset{.5}{\lambda} \cdot f(\vec{\theta}_2) \geq f\left(\underbrace{(1-\lambda) \cdot \vec{\theta}_1 + \lambda \cdot \vec{\theta}_2}\right) \checkmark$$

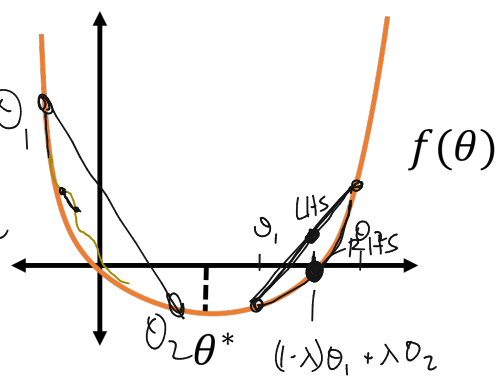
$$f''(\theta) \geq 0 \text{ For all } \theta$$

$$f(x) = 3x^2$$

$$((1-\lambda)3\theta_1 + \lambda 3\theta_2)$$

$$= 3((1-\lambda)\theta_1 + \lambda\theta_2)$$

$$= (1-\lambda)3\theta_1 + \lambda 3\theta_2$$



$\lambda = .5$

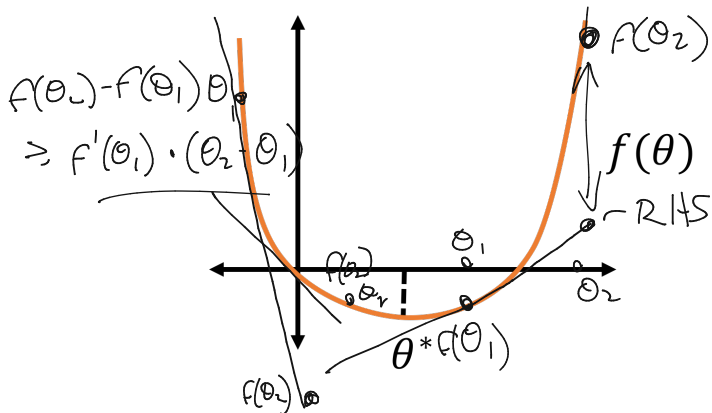
Convexity

Corollary – Convex Function: A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if and only if, for any $\vec{\theta}_1, \vec{\theta}_2 \in \mathbb{R}^d$ and $\lambda \in [0, 1]$:

$$f(\vec{\theta}_2) - f(\vec{\theta}_1) \geq \nabla f(\vec{\theta}_1)^T (\vec{\theta}_2 - \vec{\theta}_1)$$

$$f(x) = x^2$$

$$\leftarrow \nabla f(\theta_1, \theta_2 - \theta_1)$$



Conditions for Gradient Descent Convergence

Convex Functions: After sufficient iterations, if the step size η is chosen appropriately, gradient descent will converge to an **approximate minimizer** $\hat{\theta}$ with:

$$\underbrace{f(\hat{\theta})} \leq \underbrace{f(\vec{\theta}_*)} + \epsilon = \min_{\vec{\theta}} f(\vec{\theta}) + \epsilon.$$

Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

Conditions for Gradient Descent Convergence

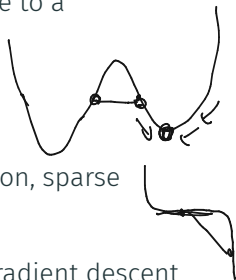
Convex Functions: After sufficient iterations, if the step size η is chosen appropriately, gradient descent will converge to an **approximate minimizer** $\hat{\theta}$ with:

$$\underline{f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta}} f(\vec{\theta}) + \epsilon.}$$

Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

Non-Convex Functions: After sufficient iterations, gradient descent will converge to an **approximate stationary point** $\hat{\theta}$ with:

$$\underline{\|\nabla f(\hat{\theta})\|_2 \leq \epsilon.}$$



Conditions for Gradient Descent Convergence

Convex Functions: After sufficient iterations, if the step size η is chosen appropriately, gradient descent will converge to a **approximate minimizer** $\hat{\theta}$ with:

$$f(\hat{\theta}) \leq f(\vec{\theta}_*) + \epsilon = \min_{\vec{\theta}} f(\vec{\theta}) + \epsilon.$$

Examples: least squares regression, logistic regression, sparse regression (lasso), regularized regression, SVMs,...

Non-Convex Functions: After sufficient iterations, gradient descent will converge to a **approximate stationary point** $\hat{\theta}$ with:

$$\|\nabla f(\hat{\theta})\|_2 \leq \epsilon.$$

Examples: neural networks, clustering, mixture models.