CS-GY 9223 D: Lecture 5 Gradient Descent and Projected Gradient Descent

NYU Tandon School of Engineering, Prof. Christopher Musco



- Choose your partner and email me by **next Wednesday**, **10/14**.
- Topic and 1 page proposal due **11/04**.
- See project guidelines on course webpage for details.

What techniques did we learn?

Have some function $f : \mathbb{R}^d \to \mathbb{R}$. Want to find \mathbf{x}^* such that:

$$f(\mathbf{x}^*) = \min_{\mathbf{x}} f(\mathbf{x}).$$

Or at least $\hat{\mathbf{x}}$ which is close to a minimum. E.g. $f(\hat{\mathbf{x}}) \leq \min_{\mathbf{x}} f(\mathbf{x}) + \epsilon$

Often we have some additional constraints:

- **x** > 0.
- $\|\mathbf{x}\|_{2} \le R$, $\|\mathbf{x}\|_{1} \le R$.
- $\mathbf{a}^T \mathbf{x} > c$.

CONTINUOUS OPTIMIZATION



Dimension d = 2:



Continuouos optimization is the foundation of modern machine learning.

Supervised learning: Want to learn a model that maps inputs

- numerical data vectors
- images, video
- text documents

to predictions

- numerical value (probability stock price increases)
- label (is the image a cat? does the image contain a car?)
- decision (turn car left, rotate robotic arm)

Let M_x be a model with parameters $\mathbf{x} = \{x_1, \dots, x_k\}$, which takes as input a data vector **a** and outputs a prediction.

Example:

$$M_{\mathbf{x}}(\mathbf{a}) = \operatorname{sign}(\mathbf{a}^{\mathsf{T}}\mathbf{x})$$

MACHINE LEARNING MODEL

Example:



 $x \in \mathbb{R}^{(\text{\# of connections})}$ is the parameter vector containing all the network weights.

Classic approach in <u>supervised learning</u>: Find a model that works well on data that you already have the answer for (labels, values, classes, etc.).

- Model M_x parameterized by a vector of numbers x.
- Dataset $\mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(n)}$ with outputs $y^{(1)}, \ldots, y^{(n)}$.

Want to find $\hat{\mathbf{x}}$ so that $M_{\hat{\mathbf{x}}}(\mathbf{a}^{(i)}) \approx y^{(i)}$ for $i \in 1, ..., n$. How do we turn this into a function minimization problem? **Loss function** $L(M_x(\mathbf{a}), y)$: Some measure of distance between prediction $M_x(\mathbf{a})$ and target output y. Increases if they are further apart.

- Squared (ℓ_2) loss: $|M_x(\mathbf{a}) y|^2$
- Absolute deviation (ℓ_1) loss: $|M_x(a) y|$
- Hinge loss: $1 y \cdot M_x(a)$
- Cross-entropy loss (log loss).
- Etc.

Empirical risk minimization:

$$f(\mathbf{x}) = \sum_{i=1}^{n} L\left(M_{\mathbf{x}}(\mathbf{a}^{(i)}), y^{(i)}\right)$$

Solve the optimization problem $\min_{\mathbf{x}} f(\mathbf{x})$.

- $M_x(a) = x^T a$. x contains the regression coefficients.
- $L(z, y) = |z y|^2$.
- $f(\mathbf{x}) = \sum_{i=1}^{n} |\mathbf{x}^{T} \mathbf{a}^{(i)} y^{(i)}|^2$

$$f(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$$

where **A** is a matrix with $\mathbf{a}^{(i)}$ as its *i*th row and **y** is a vector with $y^{(i)}$ as its *i*th entry.

The choice of algorithm to minimize $f(\mathbf{x})$ will depend on:

- The form of $f(\mathbf{x})$ (is it linear, is it quadratic, does it have finite sum structure, etc.)
- If there are any additional constraints imposed on **x**. E.g. $\|\mathbf{x}\|_2 \leq c$.

What are some example algorithms for continuous optimization?

Gradient descent: A greedy algorithm for minimizing functions of multiple variables that often works amazingly well.



For i = 1, ..., d, let x_i be the ith entry of **x**. Let $e^{(i)}$ be the ith standard basis vector.

Partial derivative:

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{e}^{(i)}) - f(\mathbf{x})}{t}$$

Directional derivative:

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}$$

Gradient:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} rac{\partial f}{\partial x_1}(\mathbf{x}) \\ rac{\partial f}{\partial x_2}(\mathbf{x}) \\ \vdots \\ rac{\partial f}{\partial x_d}(\mathbf{x}) \end{bmatrix}$$

7

Directional derivative:

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \to 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^{\mathsf{T}}\mathbf{v}$$

Given a function *f* to minimize, assume we have:

- Function oracle: Evaluate $f(\mathbf{x})$ for any \mathbf{x} .
- Gradient oracle: Evaluate $\nabla f(\mathbf{x})$ for any \mathbf{x} .

We view the implementation of these oracles as black-boxes, but they can often require a fair bit of computation.

Linear least-squares regression:

- Given $\mathbf{a}^{(1)}, \dots \mathbf{a}^{(n)} \in \mathbb{R}^d$, $y^{(1)}, \dots y^{(n)} \in \mathbb{R}$.
- Want to minimize:

$$f(\mathbf{x}) = \sum_{i=1}^{n} \left(\mathbf{x}^{\mathsf{T}} \mathbf{a}^{(i)} - \mathbf{y}^{(i)} \right)^2 = \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2.$$

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^n 2\left(\mathbf{x}^T \mathbf{a}^{(i)} - y^{(i)}\right) \cdot a_j^{(i)} = (2\mathbf{A}\mathbf{x} - \mathbf{y})^T \boldsymbol{\alpha}^{(j)}$$

where $\alpha^{(j)}$ is the *i*th <u>column</u> of **A**.

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^T (\mathbf{A}\mathbf{x} - \mathbf{y})$$

What is the time complexity of a gradient oracle for $\nabla f(\mathbf{x})$?

Greedy approach: Given a starting point **x**, make a small adjustment that decreases $f(\mathbf{x})$. In particular, $\mathbf{x} \leftarrow \mathbf{x} + \eta \mathbf{v}$ and $f(\mathbf{x}) \leftarrow f(\mathbf{x} + \eta \mathbf{v})$.

What property do I want in **v**?

Leading question: When η is small, what's an approximation for $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x})$?

 $f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x}) \approx$

DIRECTIONAL DERIVATIVES

$$D_{\mathbf{v}}f(\mathbf{x}) = \lim_{t\to 0} \frac{f(\mathbf{x}+t\mathbf{v})-f(\mathbf{x})}{t} = \nabla f(\mathbf{x})^T \mathbf{v}.$$

So:

$$f(\mathbf{x} + \eta \mathbf{v}) - f(\mathbf{x}) \approx$$

How should we choose v so that $f(x + \eta v) < f(x)$?

Prototype algorithm:

- For i = 0, ..., T:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\mathbf{x}^{(T)}$.

 η is a step-size parameter, which is often adapted on the go. For now, assume it is fixed ahead of time.

1 dimensional example:

2 dimensional example:



For a convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations *T*, gradient descent will converge to a near global minimum:

 $f(\mathbf{x}^{(T)}) \leq f(\mathbf{x}^*) + \epsilon.$

Examples: least squares regression, logistic regression, kernel regression, SVMs.

For a non-convex function $f(\mathbf{x})$: For sufficiently small η and a sufficiently large number of iterations *T*, gradient descent will converge to a near stationary point:

 $\|\nabla f(\mathbf{x}^{(T)})\|_2 \leq \epsilon.$

Examples: neural networks, matrix completion problems, mixture models.

CONVEX VS. NON-CONVEX



One issue with non-convex functions is that they can have **local minima**. Even when they don't, convergence analysis requires different assumptions than convex functions.

We care about <u>how fast</u> gradient descent and related methods converge, not just that they do converge.

- Bounding iteration complexity requires placing some assumptions on *f*(**x**).
- Stronger assumptions lead to better bounds on the convergence.

Understanding these assumptions can help us design faster variants of gradient descent (there are many!).

Today, we will start with **convex functions** only.

CONVEXITY

Definition (Convex)

A function *f* is convex iff for any $\mathbf{x}, \mathbf{y}, \lambda \in [0, 1]$:

$$(1 - \lambda) \cdot f(\mathbf{x}) + \lambda \cdot f(\mathbf{y}) \ge f((1 - \lambda) \cdot \mathbf{x} + \lambda \cdot \mathbf{y})$$



GRADIENT DESCENT

Definition (Convex)

A function *f* is convex if and only if for any **x**, **y**:

 $f(\mathbf{x} + \mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^{\mathsf{T}} \mathbf{z}$

Equivalently:

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \nabla f(\mathbf{x})^{\mathsf{T}}(\mathbf{x} - \mathbf{y})$$



GRADIENT DESCENT ANALYSIS

Assume:

- *f* is convex.
- Lipschitz function: for all \mathbf{x} , $\|\nabla f(\mathbf{x})\|_2 \leq \mathbf{G}$.
- Starting radius: $\|\mathbf{x}^* \mathbf{x}^{(0)}\|_2 \le R$.

Gradient descent:

- Choose number of steps T.
- $\eta = \frac{R}{G\sqrt{T}}$
- For i = 0, ..., T:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg\min_i f(\mathbf{x}^{(i)})$.

GRADIENT DESCENT ANALYSIS

Claim (GD Convergence Bound) If $T \ge \frac{R^2G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Proof is made tricky by the fact that $f(\mathbf{x}^{(i)})$ does not improve monotonically. We can "overshoot" the minimum.

Claim (GD Convergence Bound) If $T \ge \frac{R^2G^2}{\epsilon^2}$ and $\eta = \frac{R}{G\sqrt{T}}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1: For all i = 0, ..., T,

$$f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \le \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Claim (GD Convergence Bound) If $T \ge \frac{R^2 G^2}{\epsilon^2}$ and $\eta = \frac{R}{G\sqrt{T}} = \frac{\epsilon}{G^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1: For all *i* = 0, ..., *T*,

$$f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \le \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Telescoping sum:

$$\sum_{i=0}^{T-1} \left[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \le \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2}$$
$$\frac{1}{T} \sum_{i=0}^{T-1} \left[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \le \frac{R^2}{2T\eta} + \frac{\eta G^2}{2}$$

GRADIENT DESCENT ANALYSIS

Claim (GD Convergence Bound) If $T \ge \frac{R^2G^2}{\epsilon^2}$ and $\eta = \frac{R}{G\sqrt{T}} = \frac{\epsilon}{G^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Final step:

$$\frac{1}{T} \sum_{i=0}^{T-1} \left[f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \right] \le \epsilon$$
$$\left[\frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)}) \right] - f(\mathbf{x}^*) \le \epsilon$$

We always have that $\min_i f(\mathbf{x}^{(i)}) \leq \frac{1}{T} \sum_{i=0}^{T-1} f(\mathbf{x}^{(i)})$, so this is what we return:

$$f(\hat{\mathbf{x}}) = \min_{i \in 1, \dots, T} f(\mathbf{x}^{(i)}) \le f(\mathbf{x}^*) + \epsilon.$$

Typical goal: Solve a <u>convex minimization problem</u> with additional <u>convex constraints</u>.

 $\min_{\mathbf{x}\in\mathcal{S}}f(\mathbf{x})$

where \mathcal{S} is a **convex set**.



Which of these is convex?

CONSTRAINED CONVEX OPTIMIZATION



Definition (Convex set)

A set S is convex if for any $\mathbf{x}, \mathbf{y} \in S, \lambda \in [0, 1]$:

 $(1-\lambda)\mathbf{x} + \lambda \mathbf{y} \in \mathcal{S}.$

Gradient descent:

- For i = 0, ..., T:
 - $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} \eta \nabla f(\mathbf{x}^{(i)})$
- Return $\hat{\mathbf{x}} = \arg\min_i f(\mathbf{x}^{(i)})$.

Even if we start with $\mathbf{x}^{(0)} \in S$, there is no guarantee that $\mathbf{x}^{(0)} - \eta \nabla f(\mathbf{x}^{(0)})$ will remain in our set.

Extremely simple modification: Force $\mathbf{x}^{(i)}$ to be in S by **projecting** onto the set.

Given a function f to minimize and a convex constraint set S, assume we have:

- Function oracle: Evaluate $f(\mathbf{x})$ for any \mathbf{x} .
- Gradient oracle: Evaluate $\nabla f(\mathbf{x})$ for any \mathbf{x} .
- **Projection oracle**: Evaluate $P_{\mathcal{S}}(\mathbf{x})$ for any \mathbf{x} .

 $P_{\mathcal{S}}(\mathbf{x}) = \underset{\mathbf{y} \in \mathcal{S}}{\arg\min \|\mathbf{x} - \mathbf{y}\|_2}$

PROJECTION ORACLES

- How would you implement $P_{\mathcal{S}}$ for $\mathcal{S} = \{\mathbf{y} : \|\mathbf{y}\|_2 \leq 1\}$.
- How would you implement $P_{\mathcal{S}}$ for $\mathcal{S} = \{ \mathbf{y} : \mathbf{y} = \mathbf{Q}\mathbf{z} \}$.



Given function $f(\mathbf{x})$ and set S, such that $\|\nabla f(\mathbf{x})\|_2 \leq G$ for all $\mathbf{x} \in S$ and starting point $\mathbf{x}^{(0)}$ with $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq R$.

Projected gradient descent:

- Select starting point $\mathbf{x}^{(0)}$, $\eta = \frac{R}{G\sqrt{T}}$.
- For i = 0, ..., T:

•
$$\mathbf{z} = \mathbf{x}^{(i)} - \eta \nabla f(\mathbf{x}^{(i)})$$

$$\cdot \mathbf{x}^{(i+1)} = P_{\mathcal{S}}(\mathbf{z})$$

• Return $\hat{\mathbf{x}} = \arg\min_i f(\mathbf{x}^{(i)})$.

Claim (PGD Convergence Bound)

If f, S are convex and $T \ge \frac{R^2G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Analysis is almost identical to standard gradient descent! We just need one additional claim:

Claim (Contraction Property of Convex Projection)

If S is convex, then for <u>any</u> $\mathbf{y} \in S$,

 $\|\mathbf{y} - P_{\mathcal{S}}(\mathbf{x})\|_2 \leq \|\mathbf{y} - \mathbf{x}\|_2.$



GRADIENT DESCENT ANALYSIS

Claim (PGD Convergence Bound)

If f, S are convex and $T \ge \frac{R^2 G^2}{\epsilon^2}$, then $f(\hat{\mathbf{x}}) \le f(\mathbf{x}^*) + \epsilon$.

Claim 1: For all i = 0, ..., T, $f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*) \le \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{z} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$ $\le \frac{\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 - \|\mathbf{x}^{(i+1)} - \mathbf{x}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$

Same telescoping sum argument:

$$\left[\frac{1}{T}\sum_{i=0}^{T-1}f(\mathbf{x}^{(i)})\right]-f(\mathbf{x}^*)\leq \frac{R^2}{2T\eta}+\frac{\eta G^2}{2}.$$