CS-GY 6923: Lecture 6 Gradient Descent + Stochastic Gradient Descent

NYU Tandon School of Engineering, Prof. Christopher Musco

Goal: Minimize generic differentiable loss function:

$$L(\beta) = -\sum_{i=1}^{n} y_i \log(h(\beta^T \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\beta^T \mathbf{x}_i))$$

$$L(\beta) = \|\mathbf{X}\beta - \mathbf{y}\|_2^2 / L(\beta) = \|\mathbf{X}\beta - \mathbf{y}\|_1 + \lambda \|\beta\|_2^2 / L(\beta) = \|\mathbf{X}\beta - \mathbf{y}\|_1 + \lambda \|\beta\|_2^2 / L(\beta).$$

Gradient Descent: Most common iterative method for solving this problem.

Given a function *L* to minimize, assume we have routines for computing:

- Function oracle: Evaluate $L(\beta)$ for any β .
- **Gradient oracle**: Evaluate $\nabla \underline{L(\beta)}$ for any β

Gradient descent will use these routines in a black-box way to find the optimal β^* .

Basic Gradient descent algorithm:



 η is the <u>step-size</u> parameter or <u>learning rate</u>.

We came to an important observations: $L(\mathcal{G}^{(i+1)})$ 1. For small enough η , we always have that $L(\mathcal{G}^{(i+1)}) \leq \mathcal{W}^{(i-1)}$.

$$L(\beta + \underline{v}) - L(\beta \text{ Here}) \approx \langle \nabla L(\beta), v \rangle. = \langle \nabla L(\beta), -m \nabla L(\beta) \rangle$$

 $u = 1$ under $v = 0$ and $v = 0$
Conclusion: Gradient descent always converges to a local minimum or stationary point of *L*. Typically to a local minimum.

VISUALIZING IN 2D



6

STEEPEST DESCENT





CONVEX FUNCTION

In words: A function is convex if a line between any two points on the function lies above the function. Captures the notion that a function looks like a bowl.



This function **is not** convex.

CONVEX FUNCTION

In words: A function is convex if a line between any two points on the function lies above the function. Captures the notion that a function looks like a bowl.



This function is convex.

CONVEX FUNCTION

In words: A function is convex if a line between any two points on the function lies above the function. Captures the notion that a function looks like a bowl.



What functions are convex?

- Least squares loss for linear regression.
- ℓ_1 loss for linear regression.
- + Either of these with and ℓ_1 or ℓ_2 regularization penalty.
- Logistic regression! Logistic regression with regularization.
- Many other models in machine leaning.

CONVEXITY OF LEAST SQUARES REGRESSION LOSS

See notes from last week on proof that $L(\beta) = ||\mathbf{X}\beta - \mathbf{y}||_2^2$ is convex. *Electrow* just consider λ

Simpler problem: prove that $L(\beta) = \beta^2$ is convex.

$$((1-\Lambda)\mathcal{B}_{1}, + \Lambda\mathcal{B}_{2})^{2} \leq (1-\Lambda)\mathcal{B}_{1}^{2} + \Lambda\mathcal{B}_{2}^{2}$$

CONVERGENCE ANALYSIS FOR CONVEX FUNCTIONS

Assume:

- <u>L is conve</u>x.
- Lipschitz function: for all β , $(||\nabla L(\beta)||_2 \leq G)$
- Starting radius: $\|\beta^* \beta^{(0)}\|_2 \le R$.

Gradient descent:

- Choose number of steps T.
- Starting point $\beta^{(0)}$. E.g. $\beta^{(0)} = \mathbf{0}$.
- For $i = 0, \dots, T$: • $\beta^{(i+1)} = \beta^{(i)} - \eta \nabla L(\beta^{(i)})$
- Return $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}^{(i)}} L(\boldsymbol{\beta}).$



Proof is made tricky by the fact that $L(\beta^{(i)})$ does not improve monotonically. We can "overshoot" the minimum. This is why the step size needs to depend on 1/G.

GRADIENT DESCENT



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

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



Claim 1(a): For all i = 0, ..., T,

$$\nabla L(\boldsymbol{\beta}^{(i)})^{\mathsf{T}}(\boldsymbol{\beta}^{(i)} - \boldsymbol{\beta}^{*}) \leq \frac{\|\boldsymbol{\beta}^{(i)} - \boldsymbol{\beta}^{*}\|_{2}^{2} - \|\boldsymbol{\beta}^{(i+1)} - \boldsymbol{\beta}^{*}\|_{2}^{2}}{2\eta} + \frac{\eta G^{2}}{2\eta}$$

Claim 1 follows from Claim 1(a) by our new definition of convexity. $\lfloor (\mathcal{G}^{(1)}) - \lfloor (\mathcal{G}^{*}) \leq \nabla \lfloor (\mathcal{G}^{(1)})^{\mathsf{T}} (\mathcal{G}^{(1)} - \mathcal{B}^{*})$

Claim (GD Convergence Bound)

$$If T \ge \frac{R^2 G^2}{\epsilon^2} \text{ and } \eta = \frac{R}{G\sqrt{T}}, \text{ then } L(\hat{\beta}) \le L(\beta^*) + \epsilon.$$

Claim 1(a): For all i = 0, ..., T, ¹

$$\begin{split} & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}}(\beta^{(i)} - \beta^{*}) \leq \frac{\|\beta^{(i)} - \beta^{*}\|_{2}^{2} - \|\beta^{(i+1)} - \beta^{*}\|_{2}^{2}}{(2\eta)^{\mathsf{T}}} + \frac{\eta G^{2}}{(2\eta)^{\mathsf{T}}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} - \mathcal{L}^{*} \right\|_{2}^{\mathsf{T}} = \| \mathcal{L}^{(i)} - \beta^{*} \|_{2}^{\mathsf{T}} - \| \mathcal{L}^{(i)} \|_{2}^{\mathsf{T}} + \frac{\eta G^{2}}{(2\eta)^{\mathsf{T}}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} - \mathcal{L}^{*} \right\|_{2}^{\mathsf{T}} = \| \mathcal{L}^{(i)} - \mathcal{L}^{*} \|_{2}^{\mathsf{T}} - 2m \nabla \mathcal{L}(\beta^{(i)})^{\mathsf{T}} (\beta^{(i)} - \beta^{*}) + \| \underline{m} \nabla \mathcal{L}(\beta^{(i)}) \|_{2}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} + \| \mathcal{L}^{(i)} - \beta^{*} \|_{2}^{\mathsf{T}} - 2m \nabla \mathcal{L}(\beta^{(i)})^{\mathsf{T}} (\beta^{(i)} - \beta^{*}) + \| \underline{m} \nabla \mathcal{L}(\beta^{(i)}) \|_{2}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} + (\beta^{(i)} - \beta^{*}) \right\|_{2}^{\mathsf{T}} - 2m \nabla \mathcal{L}(\beta^{(i)})^{\mathsf{T}} (\beta^{(i)} - \beta^{*}) + \| \underline{m} \nabla \mathcal{L}(\beta^{(i)}) \|_{2}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} + \beta^{*} - \beta^{*} \right\|_{2}^{\mathsf{T}} - \| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} \|_{2}^{\mathsf{T}} + m^{2} \mathcal{L}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} + \beta^{*} \right\|_{2}^{\mathsf{T}} - \| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} \|_{2}^{\mathsf{T}} + m^{2} \mathcal{L}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} + \beta^{*} \right\|_{2}^{\mathsf{T}} + m^{2} \mathcal{L}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} \|_{2}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} \|_{2}^{\mathsf{T}} \|_{2}^{\mathsf{T}} \\ & \left\| \mathcal{L}(\beta^{(i)})^{\mathsf{T}} \|_{2}^{\mathsf{T}} \|_$$

18

Claim (GD Convergence Bound)

If T
$$\geq rac{R^2G^2}{\epsilon^2}$$
 and $\eta=rac{R}{G\sqrt{T}}$, then L $(\hat{m{eta}})\leq$ L $(m{eta}^*)+\epsilon$.

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

Telescoping sum:

$$\sum_{i=0}^{T-1} \left[L(\beta^{(i)}) - L(\beta^{*}) \right] \leq \frac{|\beta^{(0)} - \beta^{*}||_{2}^{2}}{2\eta} + \frac{|\beta^{(0)} - \beta^{*}||$$

Claim (GD Convergence Bound)

If
$$T \geq rac{R^2G^2}{\epsilon^2}$$
 and $\eta = rac{R}{G\sqrt{7}}$, then L($oldsymbol{\hat{eta}}$) \leq L($oldsymbol{eta}^*$) + ϵ .



20

Claim (GD Convergence Bound) If $T \ge \frac{R^2 G^2}{\epsilon^2}$ and $\eta = \frac{R}{G\sqrt{\pi}}$, then $L(\hat{\beta}) \le L(\beta^*) + \epsilon$. $\frac{1}{T} \frac{\tilde{z}_{i}}{\tilde{z}_{i,0}} \frac{L(\mathcal{B}^{(i)})}{\tilde{z}_{i,0}} \approx \frac{1}{T} \frac{\tilde{z}_{i,0}}{\tilde{z}_{i,0}} \frac{\operatorname{Wiy} L(\mathcal{B}^{(i)})}{\tilde{z}_{i,0}}$ $= \operatorname{Wiy} L(\mathcal{B}^{(i)})$ Final step: $\frac{1}{T}\sum_{i=1}^{t-1} \left[L(\boldsymbol{\beta}^{(i)}) - L(\boldsymbol{\beta}^*) \right] \leq \epsilon$ I/F $\left|\frac{1}{T}\sum_{i=1}^{l-1}L(\boldsymbol{\beta}^{(i)})\right| - L(\boldsymbol{\beta}^*) \leq \epsilon$ We always have that $\underline{\min_{i} L(\beta^{(i)})} \leq (\frac{1}{T} \sum_{i=0}^{T-1} L(\beta^{(i)}))$, so this is what we return:

$$L(\hat{\boldsymbol{\beta}}) = \min_{i \in 1,...,T} L(\boldsymbol{\beta}^{(i)}) \le L(\boldsymbol{\beta}^*) + \epsilon.$$

Gradient descent algorithm for minimizing $L(\beta)$:

- Choose arbitrary starting point $\beta^{(0)}$.
- For i = 1, ..., T:

•
$$\boldsymbol{\beta}^{(i+1)} = \boldsymbol{\beta}^{(i)} - \eta \nabla L(\boldsymbol{\beta}^{(i)})$$

• Return $\beta^{(t)}$.

In practice we don't set the <u>step-size/learning rate</u> parameter $\eta = \frac{R}{G\sqrt{T}}$, since we typically don't know these parameters. The above analysis can also be loose for many functions.

 η needs to be chosen sufficiently small for gradient descent to converge, but too small will slow down the algorithm.

LEARNING RATE

Precision in choosing the learning rate η is not super important, but we do need to get it to the right order of magnitude.



LEARNING RATE

"Overshooting" can be a problem if you choose the step-size too high.



Often a good idea to plot the <u>entire optimization</u> curve for diagnosing what's going on.

We will have a lab on gradient descent optimization after the midterm we're you'll get practice doing this.

Just as in regularization, search over a grid of possible parameters: $\frac{1}{2000}$

$$\eta = [\underbrace{2^{-5}, 2^{-4}, 2^{-3}, \dots, 2^9, 2^{10}}_{-1}].$$

Or tune by hand based on the optimization curve.

BACKTRACKING LINE SEARCH/ARMIJO RULE





Approximation holds true for small η . If it holds, error monotonically decreases.

BACKTRACKING LINE SEARCH/ARMIJO RULE

Gradient descent with backtracking line search:

- Choose arbitrary starting point $oldsymbol{eta}$.
- Choose starting step size η .
- Choose $\tau, \underline{c} < 1$ (typically both $\underline{c} = 1/2$ and $\tau = 1/2$)

• For
$$\underline{i = 1, ..., T}$$
:
• $\underline{\beta}^{(new)} = \beta - \eta \nabla L(\beta)$
• If $\underline{L}(\underline{\beta}^{(new)}) \leq L(\beta) - c\eta \underbrace{\nabla L(\beta)}_{\bullet} \qquad c = 0$
• $\eta \leftarrow \tau^{-1}\underline{\eta} \qquad \mathcal{M} \leftarrow 2\mathcal{M}$
• Else

$$\cdot \underline{\eta} \leftarrow \tau \underline{\eta}$$

Always decreases objective value, works very well in practice.

BACKTRACKING LINE SEARCH/ARMIJO RULE



Always decreases objective value, works very well in practice.

COMPLEXITY OF GRADIENT DESCENT

Complexity of computing the gradient will depend on you loss function. 、んらい内マ

Example 1: Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a data matrix. $2 \underbrace{\mathbf{X}^{\mathsf{T}} \mathbf{X}}_{\mathcal{S}} - \mathbf{X}^{\mathsf{T}} \underbrace{\mathbf{Y}}_{\mathcal{S}}$ $L(\boldsymbol{\beta}) = \underbrace{\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_{2}^{2}} \qquad \nabla L(\boldsymbol{\beta}) = 2 \underbrace{\mathbf{X}^{\mathsf{T}} \left(\underbrace{\mathbf{X}}_{\mathcal{\beta}} - \mathbf{y} \right)}_{\mathcal{S}} = \mathcal{O}$



- Runtime of closed form solution $\beta^* = (X^T X)^{-1} X^T y$: $\mathcal{O}(n d^2)$
- Runtime of one GD step: $O(d_n)$

Complexity of computing the gradient will depend on you loss function.

Example 1: Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a data matrix.

$$L(\boldsymbol{\beta}) = -\sum_{i=1}^{n} y_i \log(h(\boldsymbol{\beta}^{\mathsf{T}} \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\boldsymbol{\beta}^{\mathsf{T}} \mathbf{x}_i))$$
$$\nabla L(\boldsymbol{\beta}) = \mathbf{X}^{\mathsf{T}}(\underline{h}(\mathbf{X}\boldsymbol{\beta}) - \mathbf{y})$$

- · No closed form solution.
- Runtime of one GD step: $O(\gamma_{M} \downarrow)$

Frequently the complexity is <u>O(nd)</u> if you have n data-points and d parameters in your model.

Not bad, but the dependence on <u>n</u> can be a lot! n might be on the order of thousands, or millions.

Stochastic Gradient Descent (SGD).

• Powerful randomized variant of gradient descent used to train machine learning models when *n* is large and thus computing a full gradient is expensive.



STOCHASTIC GRADIENT DESCENT

Let
$$L_{j}(\beta)$$
 denote $\ell(\beta, \mathbf{x}_{j}, \mathbf{y}_{j})$.
Claim: If $j \in 1, ..., n$ is chosen uniformly at random. Then:

$$\mathbb{E}\left[n \cdot \nabla L_{j}(\beta)\right] = \nabla L(\beta).$$

$$= \sum_{j=1}^{\infty} fr\left((1_{1} = 0 > e_{-j})\right) \cdot n \nabla L_{j}(\beta)$$

$$= \sum_{j=1}^{\infty} \frac{1}{2} \cdot n \nabla L_{j}(\beta) = \sum_{j=1}^{\infty} \nabla L_{j}(\beta) = \nabla L(\beta).$$

$$\nabla L_{j}(\beta) \text{ is called a stochastic gradient.}$$

$$\sum_{j=1}^{\infty} \sqrt{1} \int_{0}^{\infty} \sqrt{1}$$

STOCHASTIC GRADIENT DESCENT

 $\mathcal{B}^{(f)} - \mathcal{M}, \qquad \begin{pmatrix} 1 \\ 1 \\ -1 \\ 2 \end{pmatrix}$

- SGD iteration:
 - Initialize $\beta^{(0)}$.
 - For i = 0, ..., T 1:
 - Choose j uniformly at random. $f_{ac} \rightarrow \dots \rightarrow \dots$
 - Compute stochastic gradient $\mathbf{g} = \nabla L_i(\boldsymbol{\beta}^{(i)})$.

• Update
$$\beta^{(t+1)} = \beta^{(t)} - \eta \cdot n\mathbf{g}$$

Move in direction of steepest descent in expectation.

Cost of computing \underline{g} is <u>independent</u> of n!

$$\nabla [b] = \tilde{z} (\nabla [j (G))$$

COMPLEXITY OF STOCHASTIC GRADIENT DESCENT

Example: Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a data matrix.

· Runtime of one SGD step:

$$L(\beta) = ||X\beta - y||_{2}^{2} = \sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}$$

$$L(\beta) = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$U(\beta) = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

$$V = \frac{||X\beta - y||_{2}^{2}}{\sum_{j=1}^{n} (y_{j} - \beta^{T} x_{j})^{2}}$$

>(&)

STOCHASTIC GRADIENT DESCENT

Gradient descent: Fewer iterations to converge, higher cost per iteration.

Stochastic Gradient descent: More iterations to converge, lower cost per iteration.



Gradient Descent

Stochastic Gradient Descent

Gradient descent: Fewer iterations to converge, higher cost per iteration.

Stochastic Gradient descent: More iterations to converge, lower cost per iteration.



Typical implementation: Shuffled Gradient Descent

Instead of choosing *j* independently at random for each iteration, randomly permute (shuffle) data and set j = 1, ..., n. After every *n* iterations, reshuffle data and repeat.

- Relatively similar convergence behavior to standard SGD.
- Important term: one epoch denotes one pass over all training examples: j = 1, ..., j = n.
- Convergence rates for training ML models are often discussed in terms of epochs instead of iterations.

STOCHASTIC GRADIENT DESCENT IN PRACTICE

 $\nabla L_i(G)$

$$L(G) = \sum_{j=1}^{2} h_j(G) = (G^T x_j - y_j)^2$$
Practical Modification: Mini-batch Gradient Descent.
Observe that for any batch size s, S << y, O(d)

$$\mathbb{E}\left[\frac{n}{s}\sum_{i=1}^{s}\nabla L_{j_i}(\beta)\right] = \underbrace{\nabla L(\beta)}_{\nabla L_{\mu}(\mathcal{C})}.$$

if j_1, \ldots, j_s are chosen independently and uniformly at random from $1, \ldots, n$.

Instead of computing a full stochastic gradient, compute the average gradient of a small random set (a <u>mini-batch</u>) of training data examples.

Question: Why might we want to do this?

MINI-BATCH GRADIENT DESCENT



- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

• Overall faster convergence (fewer iterations needed).

MIDTERM

$$\left[\left| \frac{1}{2} \right| \right] = \frac{2 \chi^{T} (\chi_{G-2})}{2 \chi^{T} (\chi_{G-2})}$$

- 1 hour long, here in the classroom. We will have lecture after.
- You can bring in a single, 2-sided cheat sheet with terms, definitions, etc.
- Mix of short answer questions (true/false, matching, etc.) and questions similar to the homework but easier.
- Might need to write some easy pseudocode.
- Covers everything through today. Don't need to know gradient descent proof of convergence.

