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

NYU Tandon School of Engineering, Prof. Christopher Musco

Lots of great ideas. Two popular ones:

- Using a different "temperature" for sampling.
- Backing off to a simple model (e.g., short *n*-gram) some of the time to add randomness.

$$t:2$$

$$\{p_1 \ldots \ldots p_n\}$$

$$\sum_{i=1}^{n} p_i = 1 \quad v_3 \; v_2 \qquad v_n$$

$$\boxed{.9} \; 0 \; \boxed{.0} \; 0 \; 0 \; \boxed{.08} \; 0)^{t} \quad v_x \; \text{temprature} \quad t \in (0, \ldots, \infty)$$

$$\frac{.9^2}{.81} \quad \frac{.01^2}{.001} \quad \frac{.08^2}{.081}$$

$$p_i \;\rightarrow\; \frac{p_i^{t}}{\sum_{j=1}^{n} p_j^{t}}$$

2

$$\beta^* \colon \arg\min L(\beta)$$

**Goal**: Minimize the logistic loss:

$$\left( L(\beta) \right) = -\sum_{i=1}^{n} y_i \log(h(\beta^T x_i)) + (1 - y_i) \log(1 - h(\beta^T x_i))$$

I.e. find $\beta^* = \arg\min L(\beta)$. How should we do this?

Set all partial derivatives to 0! Recall that $\nabla L(\boldsymbol{\beta})$ is the length $d$ vector containing all partial derivatives evaluated at $\boldsymbol{\beta}$:

$$\nabla L(\boldsymbol{\beta}) = \begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix}$$

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

Let $\mathbf{X} \in \mathbb{R}^{d \times n}$ be our data matrix with $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathbb{R}^d$ as rows.

Let $\mathbf{y} = [y_1, \ldots, y_n]$. A calculation gives (see notes on webpage):

$$\nabla L(\boldsymbol{\beta}) = \mathbf{X}^T (h(\mathbf{X}\boldsymbol{\beta}) - \mathbf{y})$$

$$h(t) \; \frac{1}{1 + e^{-t}}$$

where $h(\mathbf{X}\boldsymbol{\beta}) = \frac{1}{1 + e^{-\mathbf{X}\boldsymbol{\beta}}}$. Here all operations are entrywise. I.e in Python you would compute:

```python
h = 1/(1 + np.exp(-X@beta))
grad = np.transpose(X)@(h - y)
```

To find $\boldsymbol{\beta}$ minimizing $L(\boldsymbol{\beta})$ we typically start by finding a $\boldsymbol{\beta}$ where:

$$\nabla L(\boldsymbol{\beta}) = \mathsf{X}^T \left( h(\mathsf{X}\boldsymbol{\beta}) - \mathsf{y} \right) = \mathbf{0}$$

- In contrast to what we saw when minimizing the squared loss for linear regression, there's no simple closed form expression for such a $\boldsymbol{\beta}$!
- This is the typical situation when minimizing loss in machine learning: linear regression was a lucky exception.
- **Main question:** How do we minimize a loss function $L(\boldsymbol{\beta})$ when we can't explicitly compute where it's gradient is $\mathbf{0}$?

$$\beta^{(0)} \longrightarrow \beta^{(1)} \longrightarrow \beta^{(2)} \cdots \beta^{(T)}$$

**Much better idea.** Use a <u>guided</u> search approach.

- Start with some $\beta^{(0)}$, and at each step try to change $\beta$ slightly to reduce $L(\beta)$.

- Hopefully find an approximate minimizer for $L(\beta)$ much more quickly than brute-force search.

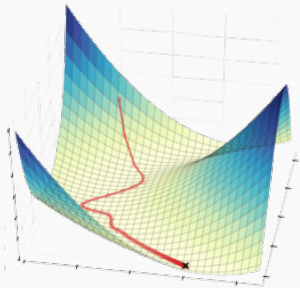- **Concrete goal:** Find $\widetilde{\beta}$ with

$$L(\widetilde{\beta}) < \min_{\beta} L(\beta) + \epsilon$$

for some small error term $\epsilon$.

$$\beta^{(i+1)} = \text{func}\left(\beta^{(i)}\right)$$

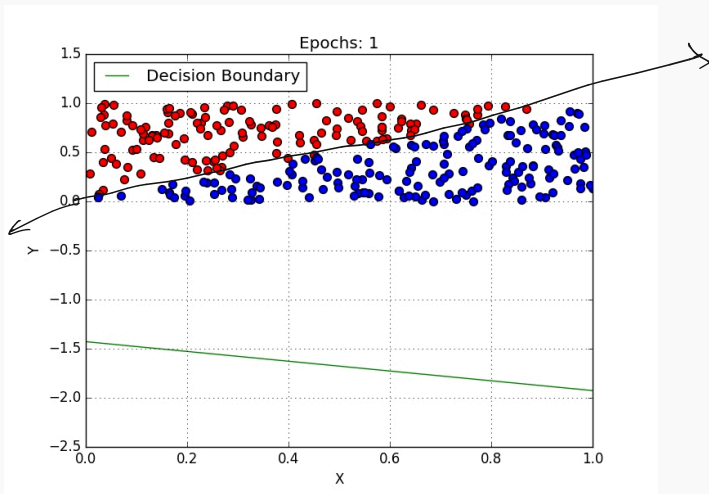**Gradient descent:** A greedy search algorithm for minimizing functions of multiple variables (including loss functions) that often works amazingly well. **What does greedy mean here?**.
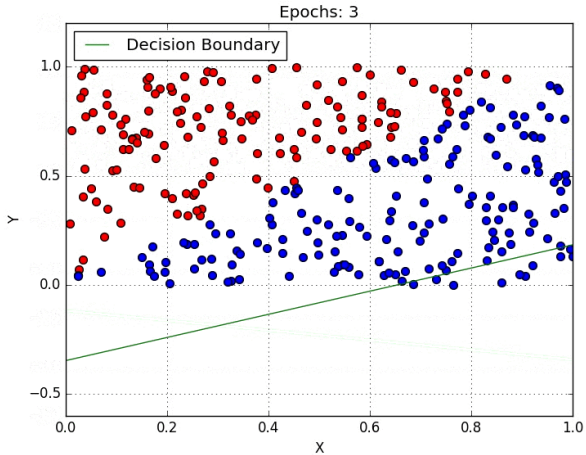


The single most important computational tool in machine learning. And it's remarkable simple + easy to implement.
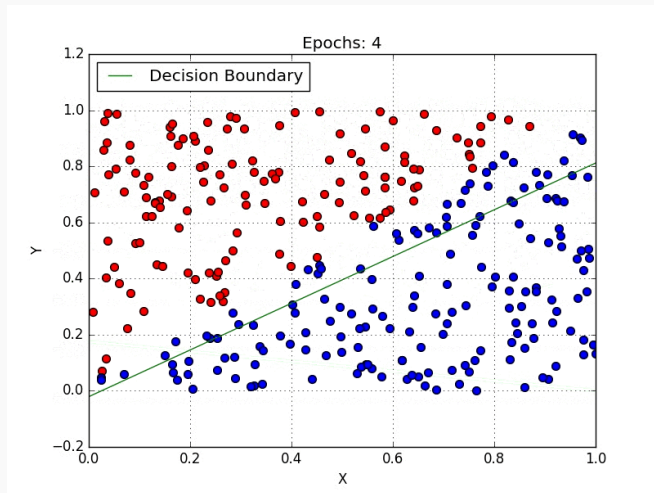
$p$ = iteration

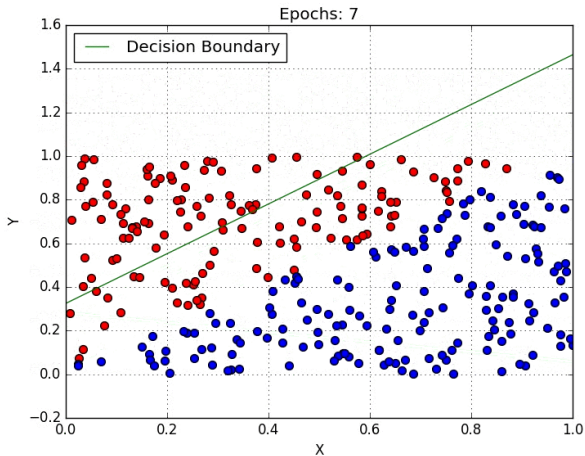Epochs: 3400

First order oracle model: Given a function $L$ to minimize, assume we can:

$\rightarrow$ returns scalar

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

$\rightarrow$ returns vector in $\mathbb{R}^d$

These are very general assumptions. Gradient descent will not use <u>any other information</u> about the loss function $L$ when trying to find a $\beta$ which minimizes $L$.

Basic Gradient descent algorithm:

- Choose starting point $\beta^{(0)}$.
- For $i = 0, \ldots, T$:
    - $\beta^{(i+1)} = \beta^{(i)} - \eta \nabla L(\beta^{(i)})$
- Return $\beta^{(T)}$.

$\eta > 0$ is a step-size parameter. Also called the learning rate.

$\beta^{(0)} = 0$

$L(\beta)$ $\qquad \forall L(\beta) \in \mathbb{R}^d$

$\downarrow$

$\mathbb{R}^d$

scalar

## Why does this method work?

**First observation:** if we actually reach the minimizer $\beta^*$ then we stop.
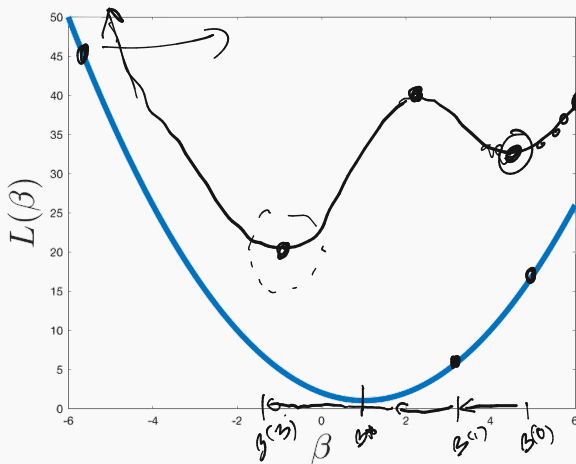
Consider a 1-dimensional loss function. I.e. where $\beta$ is just a single value. Our update step is $\beta^{(i+1)} = \beta^{(i)} - \eta \underline{L'(\beta^{(i)})}$

$$\beta^{(1)} = \beta^{(0)} - \eta$$

$\downarrow$ positive #



12

Mathematical way of thinking about it: $L'(\beta) \approx L(\beta + \Delta) - L(\beta)$

$\underbrace{\phantom{L(\beta+\Delta)-L(\beta)}}_{\Delta}$

By definition, $L'(\beta) = \lim_{\Delta \to 0} \frac{L(\beta + \Delta) - L(\beta)}{\Delta}$. So for small values of $\Delta$, we expect that:

$$L(\beta + \Delta) - L(\beta) \approx \Delta \cdot L'(\beta).$$

We want $L(\beta + \Delta)$ to be smaller than $L(\beta)$, so we want $\Delta \cdot L'(\beta)$ to be negative.

This can be achieved by choosing $\Delta = -L'(\beta)$ or really $\Delta = -\eta \cdot L'(\beta)$ for positive step size $\eta$.

$\Delta \cdot L'(\beta)$

$= -\left(L'(\beta)\right)^2$

$\underbrace{\phantom{xx}}_{}$

$\underbrace{\phantom{xxxx}}_{\text{negative.}}$

$$\beta^{(i+1)} = \beta^{(i)} - \eta L'(\beta^{(i)})$$

13

For high dimensional functions ($\boldsymbol{\beta} \in \mathbb{R}^d$), our update involves a vector $\mathbf{v} \in \mathbb{R}^d$. At each step:

$$\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} + \mathbf{v}.$$

$$-\mu \, \nabla L(\beta)$$

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

**Question:** When $\mathbf{v}$ is small, what's an approximation for $L(\boldsymbol{\beta} + \mathbf{v}) - L(\boldsymbol{\beta})$?

$$\mu \cdot e_1$$

$$\underbrace{L(\boldsymbol{\beta} + \mathbf{v}) - L(\boldsymbol{\beta})} \approx$$

$$\begin{bmatrix} \partial L / \partial \beta_1 \\ \partial L / \partial \beta_2 \\ \vdots \\ \partial L / \partial \beta_d \end{bmatrix}$$

$$\frac{\partial L}{\partial \beta_1} = \lim_{\mu \to 0} \frac{L(\beta + \mu e_1) - L(\beta)}{\mu}$$

If $\mu$ is small $\qquad L(\beta + \mu e_1) - L(\beta) \approx \mu \cdot \dfrac{\partial L}{\partial \beta_1} = \langle \mu e_1, \nabla L(\beta) \rangle$

14

We have $v = \begin{bmatrix} \vdots \\ v_i \end{bmatrix} = v_1 e_1 + v_2 e_n + \ldots v_d e_d$

$L(\beta + v_1 e_1) - L(\beta)$

$$L(\boldsymbol{\beta} + \mathbf{v}) - L(\boldsymbol{\beta}) \approx \frac{\partial L}{\partial \beta_1} v_1 + \frac{\partial L}{\partial \beta_2} v_2 + \ldots + \frac{\partial L}{\partial \beta_d} v_d$$

$L(\beta + v_1 e_1 + v_2 e_2)$
$- L(\beta + v_1 e_1)$

$$= \langle \nabla L(\boldsymbol{\beta}), \mathbf{v} \rangle.$$

**How should we choose v so that $L(\boldsymbol{\beta} + \mathbf{v}) < L(\boldsymbol{\beta})$?**

$$L(\beta + v) - L(\beta) \le 0 \qquad v = -\nabla L(\beta)$$

$$\approx \langle \nabla L(\beta), v \rangle = \langle \nabla L(\beta), -\nabla L(\beta) \rangle$$

$$= -\langle \nabla L(\beta), \nabla L(\beta) \rangle$$

$$= -\|\nabla L(\beta)\|^2$$

[0] Formally, you might remember that we can define the **directional derivative** of a multivariate function: $D_{\mathbf{v}} L(\boldsymbol{\beta}) = \lim_{\Delta \to 0} \frac{L(\boldsymbol{\beta} + \Delta \mathbf{v}) - L(\boldsymbol{\beta})}{\Delta}$. We have that $D_{\mathbf{v}} L(\boldsymbol{\beta}) = \langle \nabla L(\boldsymbol{\beta}), \mathbf{v} \rangle$.

## Claim (Gradient descent = Steepest descent[1])

$$\frac{-\nabla L(\beta)}{\|\nabla L(\beta)\|_2} = \arg\min_{\mathbf{v} : \|\mathbf{v}\|_2 = 1} \langle \nabla L(\beta), \mathbf{v} \rangle$$

$\|\nabla L(\beta)\|_2 \cdot \|\mathbf{v}\|_2 \cdot \cos(\theta)$

**Recall:** For two vectors $\mathbf{a}, \mathbf{b}$,

$$\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|_2 \|\mathbf{b}\|_2 \cdot \cos(\theta)$$



[1]We could have restricted $\mathbf{v}$ using a different norm. E.g. $\|\mathbf{v}\|_1 \leq 1$ or $\|\mathbf{v}\|_\infty = 1$. These choices lead to variants of generalized steepest descent.

Level sets of $L(\beta)$ $\rightarrow [\beta_1, \beta_2]$

$L(\beta) = 10$
$L(\beta) = 9$
$L(\beta) = 8$
$\vdots$
$L(\beta) = 4$

$\beta^*$

$\beta_2$

$\beta_1$
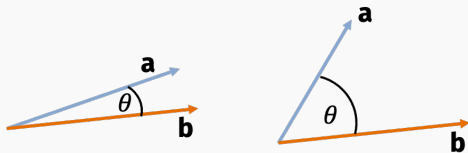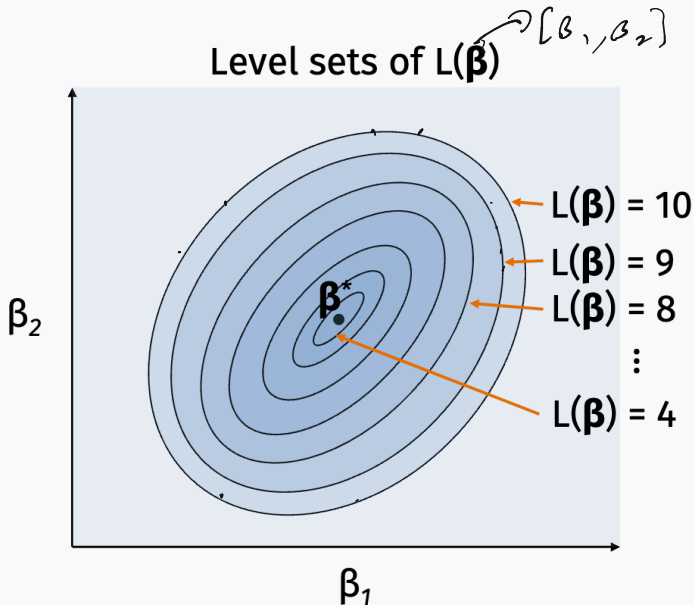
Claim (Gradient descent = Steepest descent)

$$\frac{-\nabla L(\beta)}{\|\nabla L(\beta)\|_2} = \arg\min_{\mathbf{v}, \|\mathbf{v}\|_2=1} \langle \nabla L(\beta), \mathbf{v} \rangle$$

**Level sets of L($\beta$)**



$L(\beta + v) - L(\beta)$

$\approx \langle \nabla L(\beta), v \rangle$

for small $v$.

18

Basic Gradient descent (GD) algorithm:

- Choose starting point $\boldsymbol{\beta}^{(0)}$.
- For $i = 0, \ldots, T$:
    - $\boldsymbol{\beta}^{(i+1)} = \boldsymbol{\beta}^{(i)} - \underline{\eta \nabla L(\boldsymbol{\beta}^{(i)})}$
- Return $\boldsymbol{\beta}^{(T)}$.



- **Theoretical questions:** Does gradient descent always converge to the minimum of the loss function $L$? Can you prove how quickly?
- **Practical questions:** How to choose $\eta$? Any other modifications needed for good practical performance?

19

- For sufficiently small $\eta$, every step of GD either
    1. Decreases the function value.
    2. Gets stuck because the gradient term equals 0

## Claim

*For sufficiently <u>small</u> $\eta$ and a sufficiently large number of iterations T, gradient descent will converge to a <u>local minimum</u> or stationary point of the loss function $\tilde{\boldsymbol{\beta}}^*$. I.e. with*

$$\nabla L(\tilde{\boldsymbol{\beta}}^*) = \mathbf{0}.$$

You can have stationary points that are not minima (local maxima, saddle points). In practice, always converge to local minimum.



Very unlikely to land precisely on another stationary point and get stuck. Non-minimal stationary points are "unstable".

For a broad class of functions, GD converges to global minima.

**Definition (Convex)**

A function $L$ is convex iff for any $\beta_1, \beta_2, \lambda \in [0, 1]$:

$$(1 - \lambda) \cdot L(\beta_1) + \lambda \cdot L(\beta_2) \geq L((1 - \lambda) \cdot \beta_1 + \lambda \cdot \beta_2)$$



$$L((1-\lambda)\beta_1 + \lambda\beta_2)$$

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.

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.

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.

What functions are convex?

$$V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix}$$

$$v_{d-1} \frac{\partial L}{\partial \beta_{d-1}}$$

$$v_d \frac{\partial L}{\partial \beta_d}$$

- Least squares loss for linear regression.
- $\ell_1$ loss for linear regression.
- Either of these with and $\ell_1$ or $\ell_2$ regularization penalty.
- Logistic regression! Logistic regression with regularization.
- Many other models in machine leaning.

$$L(2 + \Delta v_{d-2})$$
$$- L(2) / \Delta$$

$$\lim_{\Delta \to 0} \frac{L(\beta + \Delta V) - L(\beta)}{\Delta}$$

$$= \lim_{\Delta \to 0} L(\beta + \Delta v_1 e_1 + \ldots + \Delta v_d e_d) - \underbrace{L(\beta + \Delta v_1 e_1 + \ldots \Delta v_{d-1} e_{d-1})}_{2}$$

$$+ L(\beta + \Delta v_1 e_1 + \ldots + \Delta v_{d-1} e_{i-1}) - L(\beta + \Delta v_1 e_1 + \ldots + \Delta v_{d-2} e_{d-2}) / \Delta$$

$$+ \ldots$$

$$+ L(\beta + \Delta v_1 e_1) - L(\beta) / \Delta$$

26

What functions in machine learning are not convex? Loss functions involving neural networks, matrix completion problems, mixture models, many more.

Vary in how "bad" the non-convexity is. For example, some matrix factorization problems are non-convex but still only have global minima.
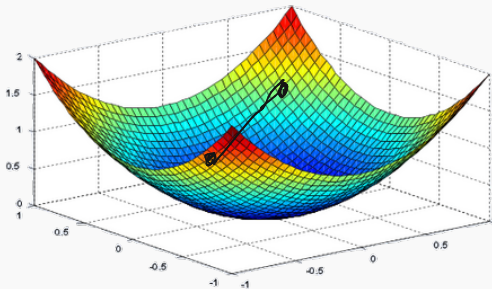
## CONVEXITY WARM UP

Prove that $L(\beta) = \beta^2$ is convex.

**To show:** For any $\beta_1, \beta_2, \lambda \in [0, 1]$,
$\lambda L(\beta_1) + (1 - \lambda)L(\beta_2) \geq L(\lambda \cdot \beta_1 + (1 - \lambda) \cdot \beta_2)$



$\lambda = \frac{1}{2}$

$1 - \lambda = \frac{1}{2}$

**AM-GM** Inequality: $c, d$ $\boxed{\sqrt{cd} \leq \frac{c+d}{2}}$

$(\sqrt{c} - \sqrt{d})^2 \geq 0$

LHS: $\frac{1}{2}\beta_1^2 + \frac{1}{2}\beta_2^2$

LHS $-$ BHS $\geq 0$

$c - 2\sqrt{c}\sqrt{d} + d \geq 0$

BHS: $\left(\frac{1}{2}\beta_1 + \frac{1}{2}\beta_2\right)^2$

$\frac{c+d}{2} \geq \sqrt{cd}$

$= \frac{1}{4}\beta_1^2 + \frac{1}{2}\beta_1\beta_2 + \frac{1}{4}\beta_2^2$

LHS $-$ BHS $= \frac{1}{4}\beta_1^2 + \frac{1}{4}\beta_2^2 - \frac{1}{2}\beta_1\beta_2$

To show: $\beta_1^2 + \beta_2^2 - 2\beta_1\beta_2 \geq 0$ $\longrightarrow$ Apply AM·GM to $\beta_1^2$ and $\beta_2^2$

28

Prove that $L(\beta) = \beta^2$ is convex.

**To show:** For any $\beta_1, \beta_2, \lambda \in [0, 1]$,
$\lambda L(\beta_1) + (1 - \lambda)L(\beta_2) \geq L(\lambda \cdot \beta_1 + (1 - \lambda) \cdot \beta_2)$

**AM-GM** Inequality:

Trick for differentiable <u>single variable</u> functions: $L(\beta)$ is convex if and only if $\underline{\underline{L''(\beta) \geq 0}}$ for all $\beta$.

$L(\beta) = \beta^2 \quad L'(\beta) = 2\beta$

$\underline{\underline{L''(\beta) = 2}}$

$\fbox{3:49 break}$

$H$ is positive semidefinite

if $\underline{x^T H x \geq 0}$ for $\underline{\underline{all}}$

vectors $x \in \mathbb{R}^d$.

$\Longleftrightarrow$ all of $H$'s eigenvalues $\geq 0$.

Analog for higher dimensional functions is clunky. Need to prove that the Hessian matrix, $H \in \mathbb{R}^{d \times d}$, is <u>positive semi-definite</u>.

$$H_{ij} = \frac{\partial^2 L}{\partial \beta_i \partial \beta_j}$$

30

Prove that $L(\boldsymbol{\beta}) = \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2$ is convex. I.e. that:

$$\|\mathbf{X}(\lambda\boldsymbol{\beta}_1 + (1-\lambda)\boldsymbol{\beta}_1) - \mathbf{y}\|_2^2 \leq \lambda\|\mathbf{X}\boldsymbol{\beta}_1 - \mathbf{y}\|_2^2 + (1-\lambda)\|\mathbf{X}\boldsymbol{\beta}_2 - \mathbf{y}\|_2^2$$

Left hand side:

$$\|\mathbf{X}(\lambda\boldsymbol{\beta}_1 + (1-\lambda)\boldsymbol{\beta}_1) - \mathbf{y}\|_2^2 = \lambda^2\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + 2\lambda(1-\lambda)\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2 + (1-\lambda)^2\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$
$$+ \mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T(\lambda\mathbf{X}\boldsymbol{\beta}_1 + (1-\lambda)\lambda\mathbf{X}\boldsymbol{\beta}_2)$$

Right hand side:

$$\lambda\|\mathbf{X}\boldsymbol{\beta}_1 - \mathbf{y}\|_2^2 + (1-\lambda)\|\mathbf{X}\boldsymbol{\beta}_2 - \mathbf{y}\|_2^2 = \lambda\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + \lambda\mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T(\lambda\mathbf{X}\boldsymbol{\beta}_1) + (1-\lambda)\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$
$$+ (1-\lambda)\mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T((1-\lambda)\mathbf{X}\boldsymbol{\beta}_2)$$

Need to show:

$$\lambda^2\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + 2\lambda(1-\lambda)\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2 + (1-\lambda)^2\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2 \leq \lambda\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + (1-\lambda)\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$

Vector version of AM-GM:

$$\|\mathbf{a} - \mathbf{b}\|_2^2 = \mathbf{a}^T\mathbf{a} - 2\mathbf{a}^T\mathbf{b} + \mathbf{b}^T\mathbf{b} \geq 0$$
$$2\mathbf{a}^T\mathbf{b} \leq \mathbf{a}^T\mathbf{a} + \mathbf{b}^T\mathbf{b}$$

$$\lambda^2\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + 2\lambda(1-\lambda)\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2 + (1-\lambda)^2\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$
$$\leq \lambda^2\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + \lambda(1-\lambda)(\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2) + (1-\lambda)^2\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$
$$= \lambda\boldsymbol{\beta}_1^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_1 + (1-\lambda)\boldsymbol{\beta}_2^T\mathbf{X}^T\mathbf{X}\boldsymbol{\beta}_2$$

**Good exercise:** Prove that $L(\boldsymbol{\beta}) = \alpha\|\boldsymbol{\beta}\|_2^2$ is convex.

**Claim:** For any convex function $L(\boldsymbol{\beta})$, gradient descent with sufficiently small step size $\eta$ converges to the global minimum $\boldsymbol{\beta}^*$ of $L$.

- Choose starting point $\boldsymbol{\beta}^{(0)}$.
- For $i = 1, \ldots, T$:
  - $\boldsymbol{\beta}^{(i+1)} = \boldsymbol{\beta}^{(i)} - \eta \nabla L(\boldsymbol{\beta}^{(i)})$
- Return $\boldsymbol{\beta}^{(T)}$.

$\mathcal{E}$

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

- Bounding iteration complexity requires placing some assumptions on $L(\beta)$.
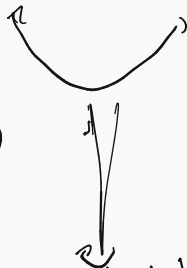- Stronger assumptions lead to better bounds on the convergence.

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

**Assume:**

$$L(\hat{\beta}) \leq L(\beta^{\phi}) + \varepsilon$$

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

**Gradient descent:**

- Choose number of steps $T$.

$T = \#$ of iterations

- Starting point $\beta^{(0)}$. E.g. $\beta^{(0)} = 0$.
- $\eta = \frac{R}{G\sqrt{T}}$

$\beta^{(T)}$

- For $i = 0, \ldots, T$:
  - $\beta^{(i+1)} = \beta^{(i)} - \eta \nabla L(\beta^{(i)})$
- Return $\hat{\beta} = \arg\min_{\beta^{(i)}} L(\beta)$.

$\beta^{(1)} \beta^{(2)} \ldots \beta^{(T)}$

35

Claim (GD Convergence Bound)

If $T \geq \frac{R^2 G^2}{\epsilon^2}$, then $L(\hat{\boldsymbol{\beta}}) \leq L(\boldsymbol{\beta}^*) + \epsilon$.



Proof is made tricky by the fact that $L(\boldsymbol{\beta}^{(i)})$ does not improve monotonically. We can "overshoot" the minimum. This is why the step size needs to depend on $1/G$.
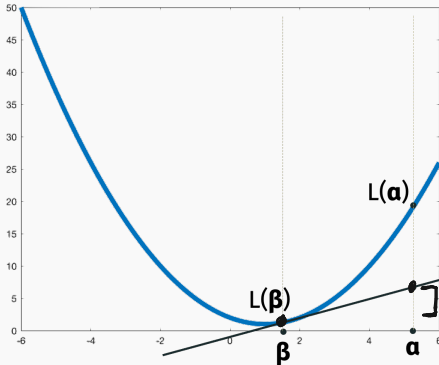
36

## Definition (Alternative Convexity Definition)

"First derivative
definition"

A function $L$ is convex if and only if for any $\underline{\beta}, \alpha$: $\in \mathbb{R}^d$

$$L(\alpha) - L(\beta) \geq \nabla L(\beta)^T (\alpha - \beta)$$

$L(\beta) - L(\alpha) \leq \nabla L(\beta)^T (\beta - \alpha)$



L(α)

L(β)

$\}$ $L(\alpha) - L(\beta)$

β        α

## Claim (GD Convergence Bound)

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

**Claim 1:** For all $i = 0, \ldots, T$,

want large positive

$$L(\beta^{(i)}) - L(\beta^*) \leq \frac{\|\beta^{(i)} - \beta^*\|_2^2 - \|\beta^{(i+1)} - \beta^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

"If you are far away, you make progress towards the optimum".

**Claim 1(a):** For all $i = 0, \ldots, T$,

$$\leq \nabla L(\beta^{(i)})^T(\beta^{(i)} - \beta^*) \leq \frac{\|\beta^{(i)} - \beta^*\|_2^2 - \|\beta^{(i+1)} - \beta^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Claim 1 follows from Claim 1(a) by our new definition of convexity.

## Claim (GD Convergence Bound)

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

RHS: right hand side
LHS: left hand side

**Claim 1(a):** For all $i = 0, \dots, T$, [2]

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

$$\nabla L(\beta^{(i)})^T(\beta^{(i)} - \beta^*) \leq \frac{\|\beta^{(i)} - \beta^*\|_2^2 - \|\beta^{(i+1)} - \beta^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

Let $w = \beta^{(i)} - \beta^*$, $z = \eta \nabla L(\beta^{(i)})$. So, RHS: $\frac{\|w\|_2^2 - \|w-z\|_2^2}{2\eta} + \frac{\eta G^2}{2}$.

$\|w-z\|_2 = \|w\|_2^2 - 2w^T z + \|z\|_2^2$.

LHS: $\frac{z^T w}{\eta}$ =

So, RHS: $\frac{\|w\|_2^2 - (\|w\|_2^2 - 2z^T w + \|z\|_2^2)}{2\eta} + \frac{\eta G^2}{2} = \frac{z^T w}{\eta} - \frac{\|z\|_2^2}{2\eta} + \frac{\eta G^2}{2}$

$\geq \frac{z^T w}{\eta} - \frac{\eta^2 G^2}{2\eta} + \frac{\eta G^2}{2}$

$= \text{LHS}$.  Concluded: RHS $\geq$ LHS

---

[2] Recall that $\|x - y\|_2^2 = \|x\|_2^2 - 2x^T y + \|y\|_2^2$.

$\rightarrow$ Here used that $\|z\|_2^2 = \eta^2 \|\nabla L(\beta^{(i)})\|_2^2 \leq \eta^2 G^2$ by assumption.

## Claim (GD Convergence Bound)

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

$i : 0 \cdots T - 1$

**Claim 1:** For all $i = 0, \ldots, T$,

$$L(\boldsymbol{\beta}^{(i)}) - L(\boldsymbol{\beta}^*) \leq \frac{\|\boldsymbol{\beta}^{(i)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(i+1)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

**Telescoping sum:**

$$\sum_{i=0}^{T-1} \left[ L(\boldsymbol{\beta}^{(i)}) - L(\boldsymbol{\beta}^*) \right] \leq \frac{\|\boldsymbol{\beta}^{(0)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

$$+ \frac{\|\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(2)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

$$+ \frac{\|\boldsymbol{\beta}^{(2)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(3)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

$$\vdots$$

$$+ \frac{\|\boldsymbol{\beta}^{(T-1)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(T)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{\eta G^2}{2}$$

40

### Claim (GD Convergence Bound)

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

$$\leq \frac{\|\beta^{(0)} - \beta^*\|_v^2}{2\eta} \leq \frac{R^2}{2\eta}$$

Telescoping sum:

$$\sum_{i=0}^{T-1} \left[ L(\boldsymbol{\beta}^{(i)}) - L(\boldsymbol{\beta}^*) \right] \leq \frac{\|\boldsymbol{\beta}^{(0)} - \boldsymbol{\beta}^*\|_2^2 - \|\boldsymbol{\beta}^{(T)} - \boldsymbol{\beta}^*\|_2^2}{2\eta} + \frac{T\eta G^2}{2}$$

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

$$\eta = \frac{R}{G\sqrt{T}} \qquad T \geq \frac{R^2 G^2}{\epsilon^2}$$

$$= \frac{RG}{2\sqrt{T}} + \frac{RG}{2\sqrt{T}} = \frac{RG}{\sqrt{T}} \leq \epsilon$$

41

Claim (GD Convergence Bound)                         $O\left(\frac{1}{\epsilon^2}\right)$

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

Final step:

$$\frac{1}{T}\sum_{i=0}^{T-1} L(\beta^{(i)}) - \frac{1}{T}\sum_{i=0}^{T-1} L(\beta^*)$$

$$= L(\beta^*)$$

$$\frac{1}{T}\sum_{i=0}^{T-1}\left[L(\boldsymbol{\beta}^{(i)}) - L(\boldsymbol{\beta}^*)\right] \leq \epsilon$$

$$\left[\frac{1}{T}\sum_{i=0}^{T-1} L(\boldsymbol{\beta}^{(i)})\right] - L(\boldsymbol{\beta}^*) \leq \epsilon$$

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

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

42

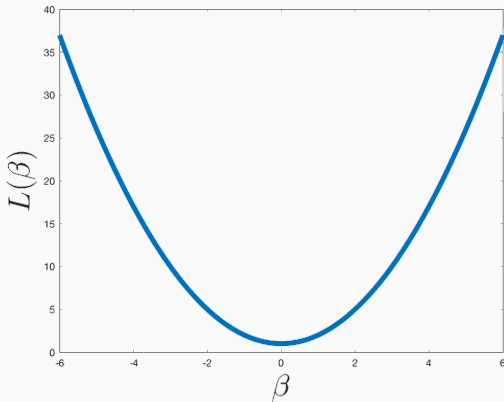Gradient descent algorithm for minimizing $L(\beta)$:

- Choose arbitrary starting point $\beta^{(0)}$.
- For $i = 1, \ldots, T$:
    - $\beta^{(i+1)} = \beta^{(i)} - \eta\nabla L(\beta^{(i)})$
- Return $\beta^{(T)}$.

In practice we don't set the step-size/learning rate 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.

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

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

Assume:

- $L$ is convex.
- Lipschitz function: for all $\boldsymbol{\beta}$, $\|\nabla L(\boldsymbol{\beta})\|_2 \leq G$.
- Starting radius: $\|\boldsymbol{\beta}^* - \boldsymbol{\beta}^{(0)}\|_2 \leq R$.

Gradient descent:

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

This result tells us exactly how to set the learning rate $\eta$.

45

But...

- We don't usually know $R$ or $G$ in advance. We might not even know $T$.
- Even if we did, setting $\eta = \frac{R}{G\sqrt{T}}$ tends to be a very conservative in practice. The choice 100% leads to convergence, but usually to fairly slow convergence.
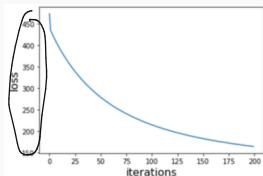- What if $L$ is not convex?

Just as in regularization, search over a grid of possible parameters:

$$\left( \eta = [2^{-5}, 2^{-4}, 2^{-3}, \ldots, 2^{9}, 2^{10}]. \right)$$

Can manually check if we are converging too slow or undershooting by plotting the optimization curve.
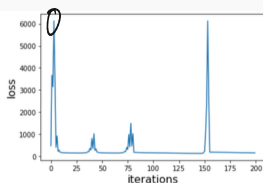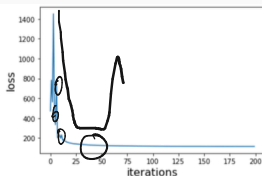
$\mathcal{M}$

Plot's of loss vs. number of iterations for three difference
choices of step size.



smallest                                                          largest

Recall: If we set $\beta^{(i+1)} \leftarrow \beta^{(i)} - \eta \nabla L(\beta^{(i)})$ then:

$$L(\beta^{(i+1)}) \approx L(\beta^{(i)}) - \eta \left\langle \nabla L(\beta^{(i)}), \nabla L(\beta^{(i)}) \right\rangle$$

$$= L(\beta^{(i)}) - \eta \|\nabla L(\beta^{(i)})\|_2^2.$$

Approximation holds for small $\eta$. If it holds, maybe we could get away with a larger $\eta$. If it doesn't, we should probably reduce $\eta$.

$$L(\beta^{(i+1)}) - L(\beta^{(i)}) \leq \eta \|\nabla L(\beta^{(i)})\|_2^2$$
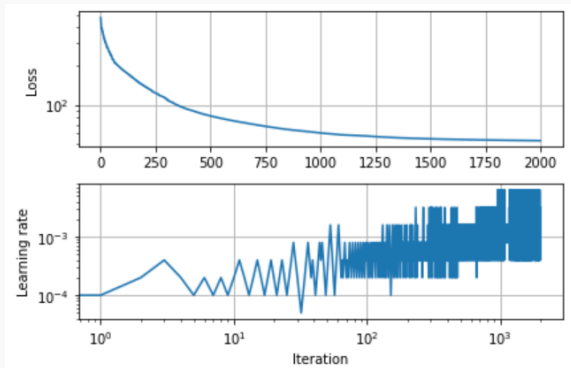
Gradient descent with backtracking line search:

- Choose arbitrary starting point $\beta$.
- Choose starting step size $\eta$
- Choose $c < 1$ (typically both $c = 1/2$)
- For $i = 1, \ldots, T$:
  - $\beta^{(new)} = \beta - \eta \nabla L(\beta)$ $\rightarrow$ standard GD step
  - If $L(\beta^{(new)}) \leq L(\beta) - c \cdot \eta \|\nabla L(\beta)\|_2^2$
    - $\beta \leftarrow \beta^{(new)}$
    - $\eta \leftarrow 2\eta$ $\qquad \eta \leftarrow \eta/c$
  - Else
    - $\eta \leftarrow \eta/2$ $\qquad \leftarrow \eta \cdot c$

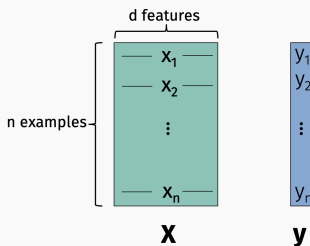Always decreases objective value, works very well in practice.

Gradient descent with backtracking line search:



Always decreases objective value, works very well in practice. We will see this in a lab.

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

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

$$L(\boldsymbol{\beta}) = \|X\boldsymbol{\beta} - y\|_2^2 \qquad \nabla L(\boldsymbol{\beta}) = 2X^T(X\boldsymbol{\beta} - y)$$



- Runtime of closed form solution $\boldsymbol{\beta}^* = (X^T X)^{-1} X^T y$:
- Runtime of one GD step:

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

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

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

$$\nabla L(\boldsymbol{\beta}) = X^T (h(X\boldsymbol{\beta}) - \mathbf{y})$$

· No closed form solution.
· Runtime of one GD step:

Frequently the complexity is $O(nd)$ if you have $n$ data-points and $d$ parameters in your model. This will also be the case for neural networks.

Not bad, but the dependence on $n$ can be a lot! $n$ might be on the order of thousands, or millions, or trillions.

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.

Applies to any loss with <u>finite sum</u> structure:

$$L(\boldsymbol{\beta}) = \sum_{j=1}^{n} \ell(\boldsymbol{\beta}, \mathbf{x}_j, y_j)$$

Let $L_j(\boldsymbol{\beta})$ denote $\ell(\boldsymbol{\beta}, \mathbf{x}_j, y_j)$.

**Claim:** If $j \in 1, \ldots, n$ is chosen uniformly at random. Then:

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

$\nabla L_j(\boldsymbol{\beta})$ is called a stochastic gradient.

SGD iteration:

- Initialize $\boldsymbol{\beta}^{(0)}$.
- For $i = 0, \ldots, T-1$:
  - Choose $j$ uniformly at random from $\{1, 2, \ldots, n\}$.
  - Compute stochastic gradient $\mathbf{g} = \nabla L_j(\boldsymbol{\beta}^{(i)})$.
  - Update $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \eta \cdot n\mathbf{g}$

Move in direction of steepest descent <u>in expectation.</u>

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

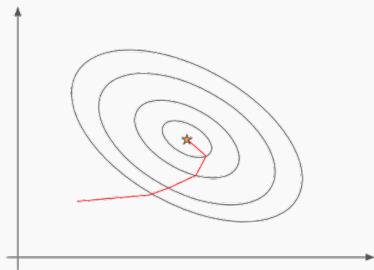**Example:** Let $X \in \mathbb{R}^{n \times d}$ be a data matrix.

$$L(\boldsymbol{\beta}) = \|X\boldsymbol{\beta} - y\|_2^2 = \sum_{j=1}^{n}(y_j - \boldsymbol{\beta}^T x_j)^2$$

- Runtime of one SGD step:

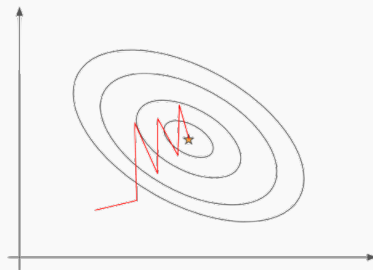**Gradient descent:** Fewer iterations to converge, higher cost per iteration.

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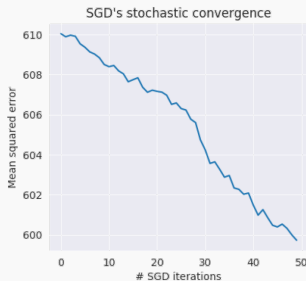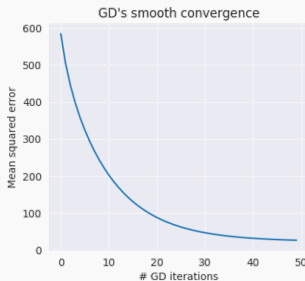


Gradient Descent                         Stochastic Gradient Descent

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 $j = 1, \ldots, n$. 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: $j = 1, \ldots, j = n$.
- Convergence rates for training ML models are often discussed in terms of epochs instead of iterations.

Practical Modification: **Mini-batch Gradient Descent.**
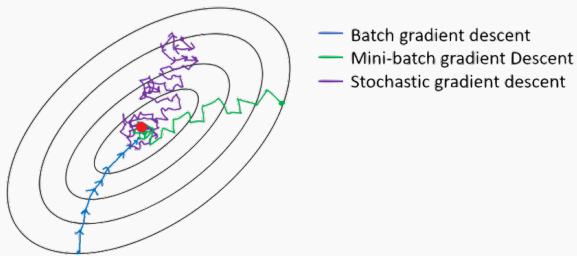
Observe that for any <u>batch size</u> $s$,

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

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?

- Overall faster convergence (fewer iterations needed).

Practical Mod. 2: **Per-parameter adaptive learning rate.**

Let $\mathbf{g} = \begin{bmatrix} g_1 \\ \vdots \\ g_p \end{bmatrix}$ be a stochastic or batch stochastic gradient. Our

typical parameter update looks like:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \eta \mathbf{g}.$$

We've already seen a simple method for adaptively choosing
the learning rate/step size $\eta$.

Practical Mod. 2: **Per-parameter adaptive learning rate.**

In practice, ML lost functions can often be optimized much faster by using "adaptive gradient methods" like Adagrad, Adadelta, RMSProp, and ADAM. These methods make updates of the form:

$$\boldsymbol{\beta}_{t+1} = \boldsymbol{\beta}_t - \begin{bmatrix} \eta_1 \cdot g_1 \\ \vdots \\ \eta_d \cdot g_d \end{bmatrix}$$

So we have a separate learning rate for each entry in the gradient (e.g. parameter in the model). And each $\eta_1, \ldots, \eta_p$ is chosen adaptively.

- 1.5 hours long, but should take 1 hour. Here in the classroom.
- Will have a short lecture after exam/break.
- 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.
- Covers everything through last class. Don't need to know gradient descent or optimization.