

Optimization: Stochastic Gradient Descent

Recap on Optimization

GD: simply follow the negative of the gradient

Recap on Optimization

GD: simply follow the negative of the gradient

AdaGrad — each dim has its own learning rate, adapted based on the cumulation of the past squared derivatives — help make progress along all axes.

Recap on Optimization

GD: simply follow the negative of the gradient

AdaGrad — each dim has its own learning rate, adapted based on the cumulation of the past squared derivatives — help make progress along all axes.

GD w/ momentum: think about gradient as “acceleration”, “velocity” is the exponential average of “acceleration” — help power through very flat region

Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

Recap on Gradient Descent

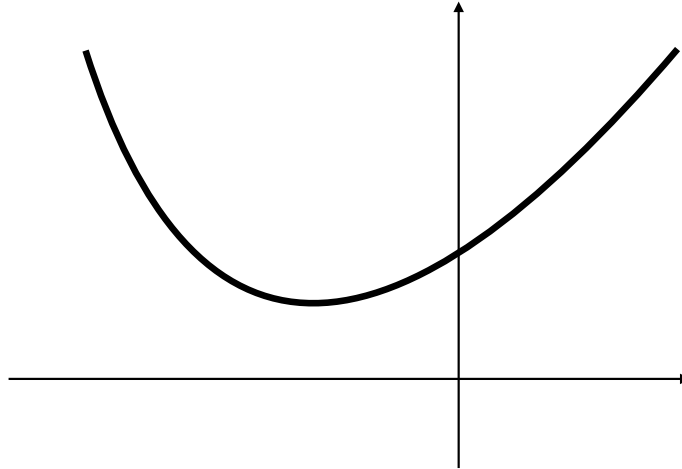
Gradient descent minimizes $\ell(w)$ iteratively:

$$w^{t+1} = w^t - \eta \nabla \ell(w) \big|_{w=w_t}$$

Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

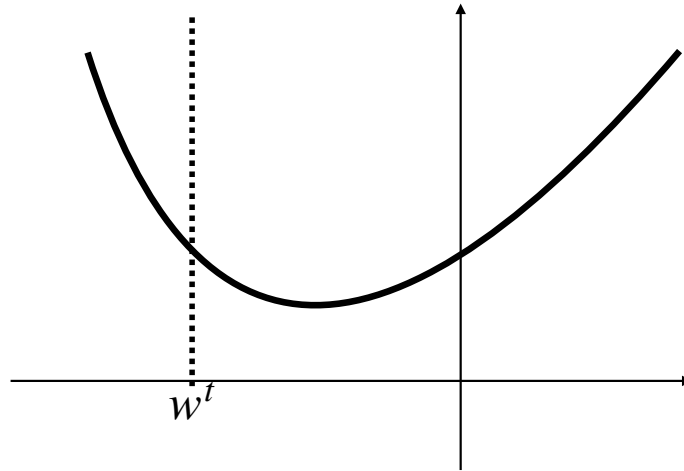
$$w^{t+1} = w^t - \eta \nabla \ell(w) \Big|_{w=w_t}$$



Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

$$w^{t+1} = w^t - \eta \nabla \ell(w) \Big|_{w=w^t}$$



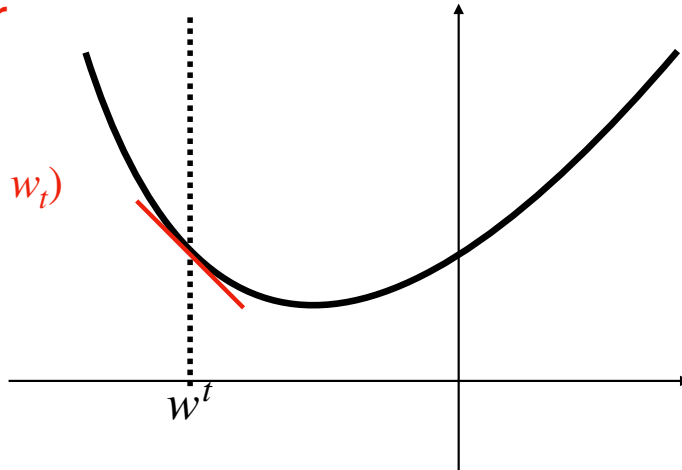
Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

$$w^{t+1} = w^t - \eta \nabla \ell(w) \big|_{w=w^t}$$

First-order Taylor
expansion at w_t :

$$\ell(w_t) + \nabla \ell(w_t)^\top (w - w_t)$$



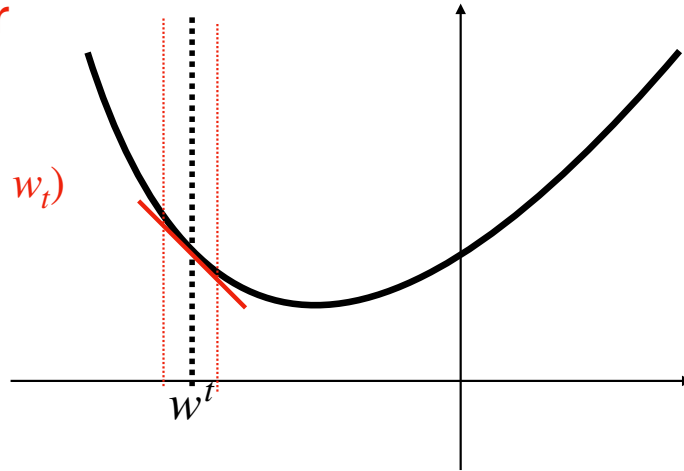
Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

$$w^{t+1} = w^t - \eta \nabla \ell(w) \big|_{w=w^t}$$

First-order Taylor
expansion at w_t :

$$\ell(w_t) + \nabla \ell(w_t)^\top (w - w_t)$$



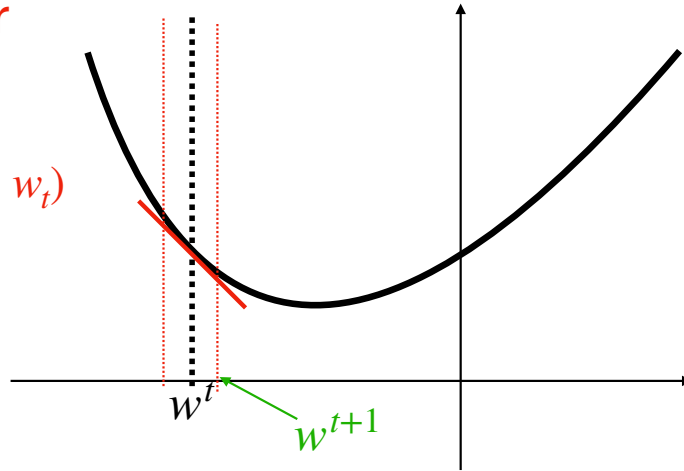
Recap on Gradient Descent

Gradient descent minimizes $\ell(w)$ iteratively:

$$w^{t+1} = w^t - \eta \nabla \ell(w) \big|_{w=w^t}$$

First-order Taylor
expansion at w_t :

$$\ell(w_t) + \nabla \ell(w_t)^\top (w - w_t)$$



Objective

Understand the Stochastic GD algorithm, its convergence, and its benefits over GD

Outline for Today

1. Stochastic Gradient Descent
2. Mini-Batch SGD

Loss minimization in ML

In ML, the loss we minimize typically has some special form, e.g., in LR:

$$\ell(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

← NLL of logistic regression

Loss minimization in ML

In ML, the loss we minimize typically has some special form, e.g., in LR:

$$\ell(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

Avg over n data points, i.e., $\sum_{i=1}^n \ell(x_i, y_i; w)/n$

Loss minimization in ML

In ML, the loss we minimize typically has some special form, e.g., in LR:

$$\ell(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

Avg over n data points, i.e., $\sum_{i=1}^n \ell(x_i, y_i; w) / n$

To compute the gradient $\nabla \ell(w)$, we need to enumerate all n training data points

Loss minimization in ML

In ML, the loss we minimize typically has some special form, e.g., in LR:

$$\ell(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

Avg over n data points, i.e., $\sum_{i=1}^n \ell(x_i, y_i; w)/n$

To compute the gradient $\nabla \ell(w)$, we need to enumerate all n training data points

Can be very slow!

Stochastic GD to rescue

In ML, the loss we minimize typically has some special form, e.g., in LR:

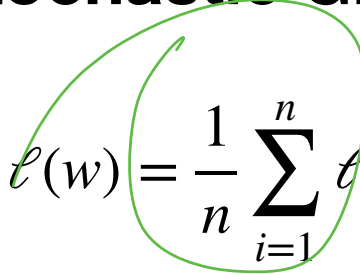
$$\ell(w) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i(w^\top x_i)))$$

Avg over n data points, i.e., $\sum_{i=1}^n \ell(x_i, y_i; w)/n$

Idea: randomly sample a data point (x, y) , use $\nabla \ell(x, y; w)$ to replace $\nabla \ell(w)$

Stochastic GD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$



Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

Stochastic GD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample a point (x_i, y_i) from the n data points

Stochastic GD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample a point (x_i, y_i) from the n data points
2. Compute noisy gradient $\tilde{g}^t = \nabla \ell(x_i, y_i; w) |_{w=w^t}$

GD: $\nabla \left(\frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w) \right)$

Stochastic GD

$$\text{Goal: minimize } \ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample a point (x_i, y_i) from the n data points
2. Compute noisy gradient $\tilde{g}^t = \nabla \ell(x_i, y_i; w) |_{w=w^t}$
3. Update (GD): $w^{t+1} = w^t - \eta \tilde{g}^t$

Intuition of why Stochastic GD can work

Claim: the random noisy gradient is an **unbiased** estimate of the true gradient

$$\nabla \ell(x_i, y_i; w)$$

$$\nabla \sum_{i=1}^n \ell(x_i, y_i; w) / n$$

Intuition of why Stochastic GD can work

Claim: the random noisy gradient is an **unbiased** estimate of the true gradient

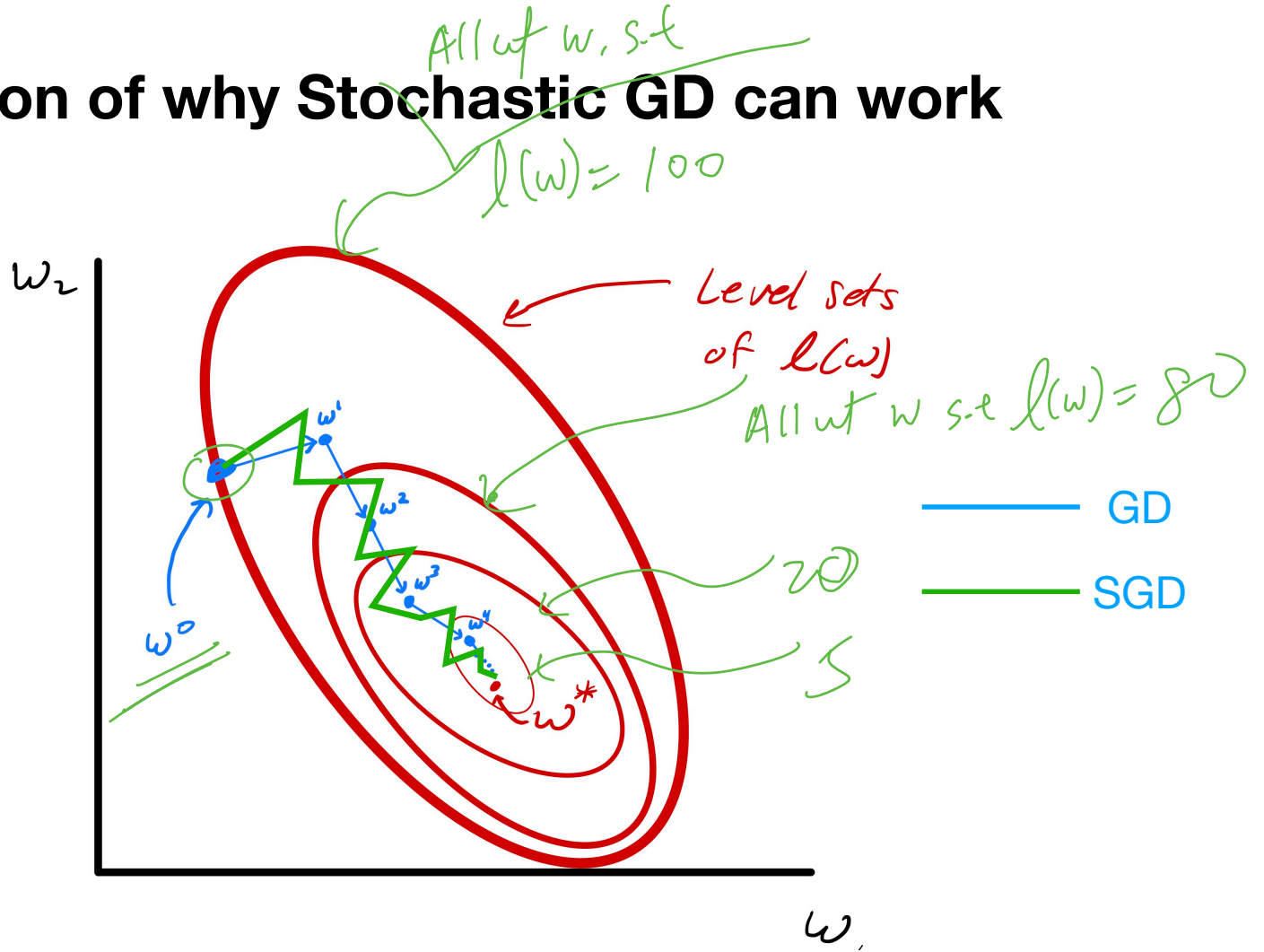
Note the point (x_i, y_i) is uniformly random sampled from n data points, we have:

$$\begin{aligned} & \mathbb{E}_{i \sim \text{Uniform}(1, 2, \dots, n)} \left[\nabla \ell(x_i, y_i; w) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \nabla \ell(x_i, y_i; w) = \nabla \left[\underbrace{\frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)}_{\ell(w)} \right] = \nabla \ell(w) \end{aligned}$$

Handwritten notes:

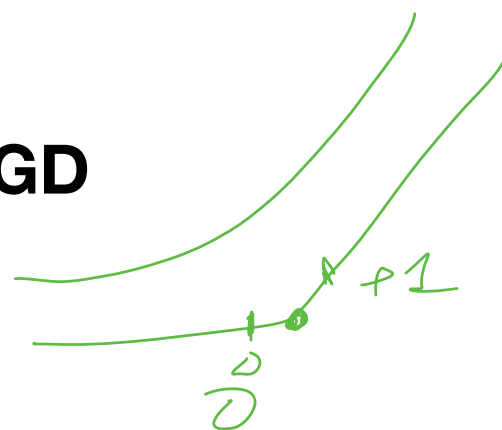
- A green arrow points from the expectation operator \mathbb{E} to the expression $i \sim \text{Uniform}(1, 2, \dots, n)$.
- A green arrow points from the expression $\frac{1}{n} \sum_{i=1}^n$ to the fraction $\frac{1}{n}$ in the equation.
- Green circles highlight $\mathbb{E} \nabla \ell(x_i, y_i; w)$ and the entire right-hand side of the equation.
- A green arrow points from the underbraced term $\frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$ to the label $\ell(w)$.

Intuition of why Stochastic GD can work



Theoretical Guarantee of SGD

(Informal theorem and no proof)



Consider a function w/ β -Lipschitz gradient, i.e., $\|\nabla \ell(w) - \nabla \ell(w')\|_2 \leq \beta \|w - w'\|_2$.

Assume for all iteration t , \tilde{g}^t is unbiased, and $\mathbb{E} \|\tilde{g}^t\|_2^2 \leq \sigma^2$,

Theoretical Guarantee of SGD

(Informal theorem and no proof)

Consider a function w/ β -Lipschitz gradient, i.e., $\|\nabla \ell(w) - \nabla \ell(w')\|_2 \leq \beta \|w - w'\|_2$.

Assume for all iteration t , \tilde{g}^t is unbiased, and $\mathbb{E}\|\tilde{g}^t\|_2^2 \leq \sigma^2$,

then with $\eta = \sqrt{\frac{1}{\beta\sigma^2 T}}$, SGD satisfies:

Theoretical Guarantee of SGD

(Informal theorem and no proof)

Consider a function w/ β -Lipschitz gradient, i.e., $\|\nabla \ell(w) - \nabla \ell(w')\|_2 \leq \beta \|w - w'\|_2$.

Assume for all iteration t , \tilde{g}^t is unbiased, and $\mathbb{E} \|\tilde{g}^t\|_2^2 \leq \sigma^2$,

then with $\eta = \sqrt{\frac{1}{\beta\sigma^2 T}}$, SGD satisfies:

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \|\nabla \ell(w^t)\|_2 \right] \leq 2\sqrt{\frac{\beta\sigma^2}{T}}$$

lim
T \rightarrow ∞
T: # of Iterations

Theoretical Guarantee of SGD

(Informal theorem and no proof)

Consider a function w/ β -Lipschitz gradient, i.e., $\|\nabla \ell(w) - \nabla \ell(w')\|_2 \leq \beta \|w - w'\|_2$.

Assume for all iteration t , \tilde{g}^t is unbiased, and $\mathbb{E} \|\tilde{g}^t\|_2^2 \leq \sigma^2$,

then with $\eta = \sqrt{\frac{1}{\beta\sigma^2 T}}$, SGD satisfies:

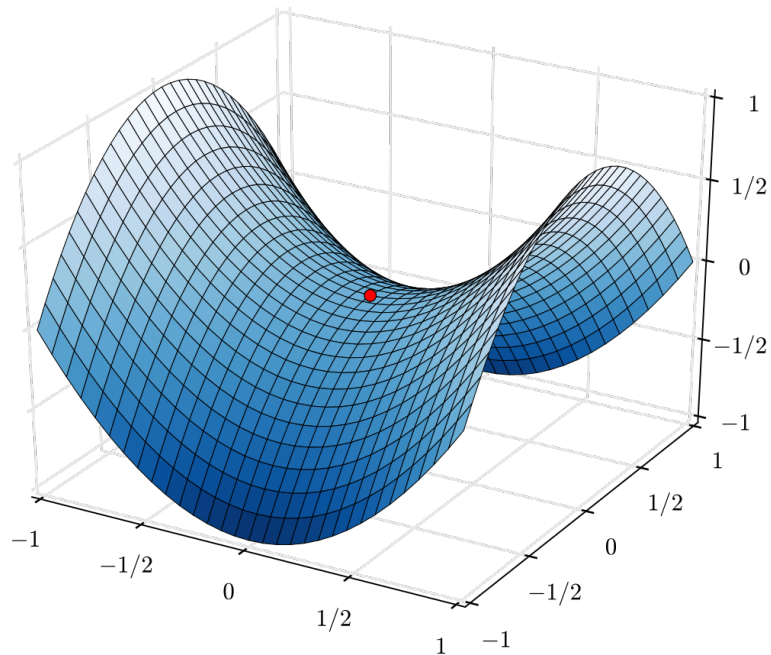
Larger variance can
make SGD slower

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \|\nabla \ell(w^t)\|_2 \right] \leq 2 \sqrt{\frac{\beta\sigma^2}{T}}$$

$\sigma \rightarrow \infty$

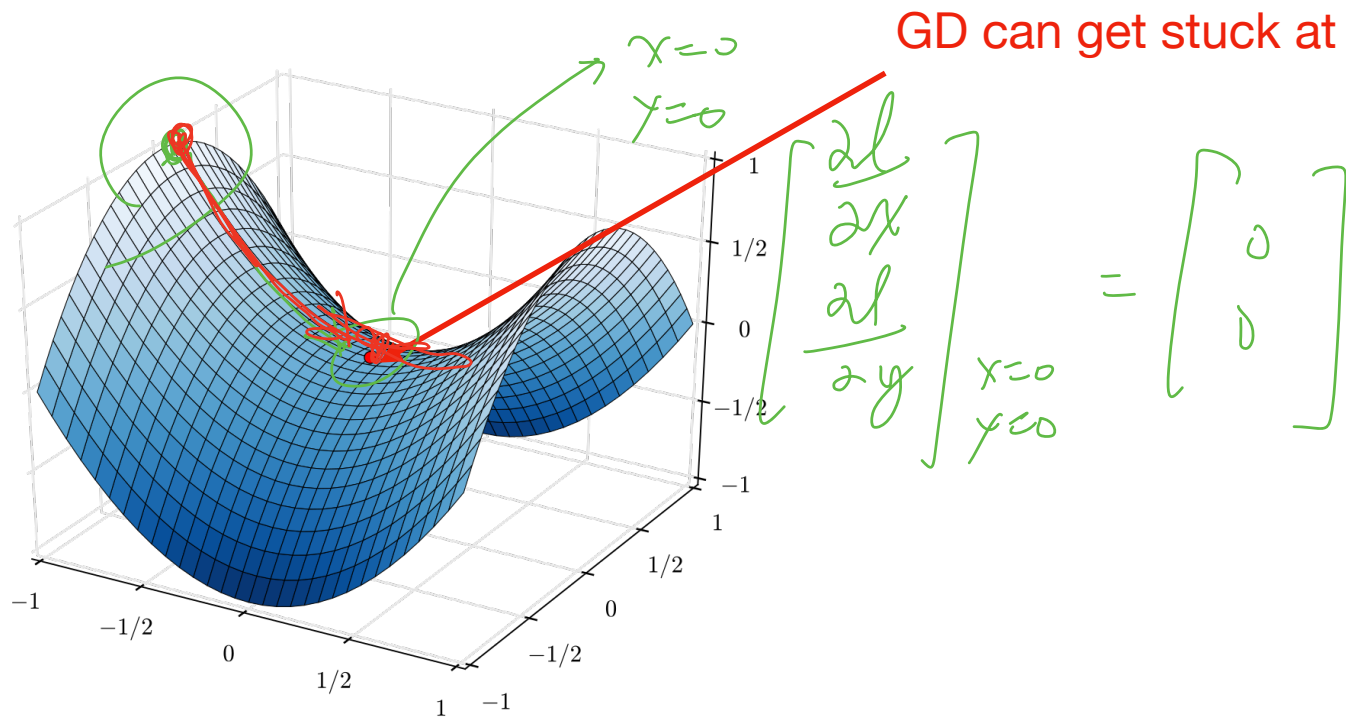
Empirical Benefit of SGD on Non-Convex Optimization

e.g., saddle point $\ell(x, y) = x^2 - y^2$



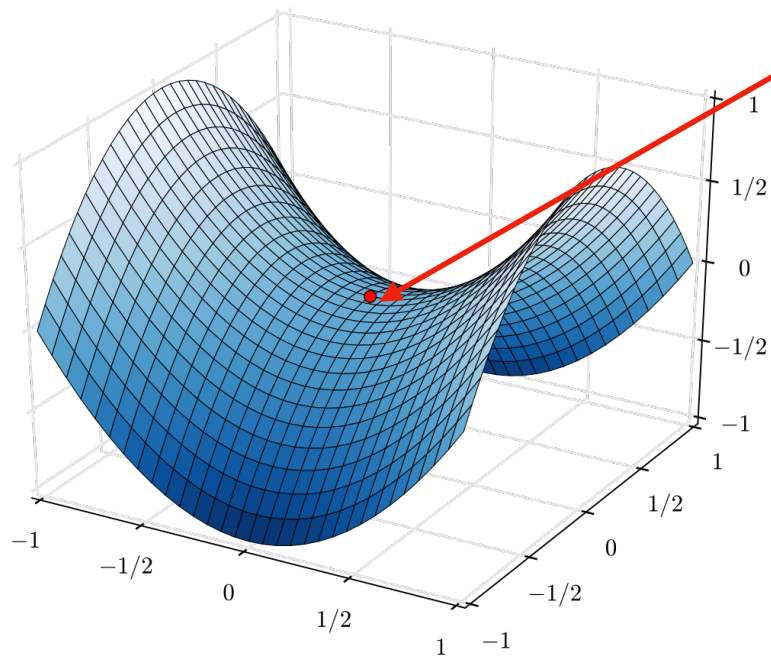
Empirical Benefit of SGD on Non-Convex Optimization

e.g., saddle point $\ell(x, y) = x^2 - y^2$



Empirical Benefit of SGD on Non-Convex Optimization

e.g., saddle point $\ell(x, y) = x^2 - y^2$



GD can get stuck at the saddle point

Using a noisy gradient, we can escape this saddle point

Outline for Today

1. Stochastic Gradient Descent

2. Mini-Batch SGD

Reduce the variance via mini-batch

Reduce the variance via mini-batch

SGD's convergence typically depend on the second moment of \tilde{g} , i.e., $\mathbb{E}\|\tilde{g}\|_2^2$

Larger variance implies slower convergence

Reduce the variance via mini-batch

SGD's convergence typically depend on the second moment of \tilde{g} , i.e., $\mathbb{E}\|\tilde{g}\|_2^2$

Larger variance implies slower convergence

Solution: we can reduce the variance using a **mini-batch**

Reduce the variance via mini-batch

$$m \geq 1 \quad m \in \mathbb{N}^+$$

Randomly sample m data points from the dataset, denoted as \mathcal{B}

Reduce the variance via mini-batch

Randomly sample m data points from the dataset, denoted as \mathcal{B}

$$\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x, y)$$



Reduce the variance via mini-batch

Randomly sample m data points from the dataset, denoted as \mathcal{B}

$$\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x, y)$$

Averaging over m points reduce the variance

Reduce the variance via mini-batch

Randomly sample m data points from the dataset, denoted as \mathcal{B}

$$\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x, y)$$

Averaging over m points reduce the variance

Claim: \tilde{g} is still unbiased, and variance of \tilde{g} decreases as m increases

Mini-batch SGD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:



Mini-batch SGD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample m points, denoted as mini-batch \mathcal{B}

Mini-batch SGD

$$\text{Goal: minimize } \ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample m points, denoted as mini-batch \mathcal{B}

2. Compute gradient $\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x_i, y_i) |_{w=w^t}$

Mini-batch SGD

$$\text{Goal: minimize } \ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

1. Randomly sample m points, denoted as mini-batch \mathcal{B}

2. Compute gradient $\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x_i, y_i) |_{w=w^t}$

3. Update (GD): $w^{t+1} = w^t - \eta \tilde{g}^t$

Mini-batch SGD

Goal: minimize $\ell(w) = \frac{1}{n} \sum_{i=1}^n \ell(x_i, y_i; w)$

Initialize $w^0 \in \mathbb{R}^d$ randomly

Iterate until convergence:

Batch size m & learning rate η are very important hyper-parameters!

1. Randomly sample m points, denoted as mini-batch \mathcal{B}

2. Compute gradient $\tilde{g} = \frac{1}{m} \sum_{(x,y) \in \mathcal{B}} \nabla \ell(w; x_i, y_i) |_{w=w^t}$

3. Update (GD): $w^{t+1} = w^t - \eta \tilde{g}^t$

$$\left(\frac{1}{n} \sum_{i=1}^n f(x_i) \right)' = \frac{1}{n} \sum_{i=1}^n f'(x_i)$$

Closing Remarks on Optimization

1. Min-batch SGD is the foundation of today's deep learning
2. Can use Stochastic gradients together w/ AdaGrad, GD w/ Momentum, and Adam