

# ECS289: Scalable Machine Learning

Cho-Jui Hsieh  
UC Davis

Sept 29, 2016

# Outline

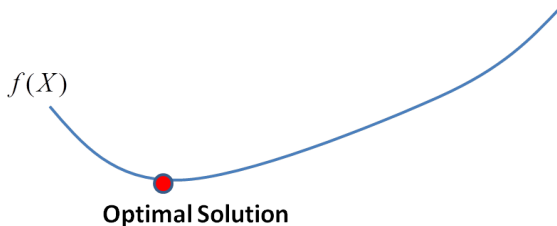
- Convex vs Nonconvex Functions
- Coordinate Descent
- Gradient Descent
- Newton's method
- Stochastic Gradient Descent

# Numerical Optimization

- Numerical Optimization:

$$\min_X f(X)$$

- Can be applied to computer science, economics, control engineering, operating research, ...
- **Machine Learning:** find a model that **minimizes** the prediction error.



# Properties of the Function

- Smooth function: functions with continuous derivative.

Example: ridge regression

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Non-smooth function: Lasso, primal SVM

$$\text{Lasso: } \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

$$\text{SVM: } \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

# Convex Functions

- A function is convex if:

$$\forall \mathbf{x}_1, \mathbf{x}_2, \forall t \in [0, 1], f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

- No local optimum (why?)

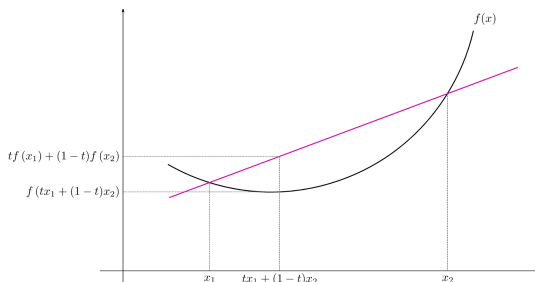


Figure from Wikipedia

# Convex Functions

- If  $f(\mathbf{x})$  is twice differentiable, then  
 $f$  is convex if and only if  $\nabla^2 f(\mathbf{x}) \succeq 0$

- Optimal solution may not be unique:  
has a set of optimal solutions  $\mathcal{S}$

- Gradient: capture the first order change of  $f$ :

$$f(\mathbf{x} + \alpha \mathbf{d}) = f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{d} + O(\alpha^2)$$

- If  $f$  is convex and differentiable, we have the following optimality condition:

$$\mathbf{x}^* \in \mathcal{S} \text{ if and only if } \nabla f(\mathbf{x}) = 0$$

# Strongly Convex Functions

- $f$  is strongly convex if there exists an  $m > 0$  such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

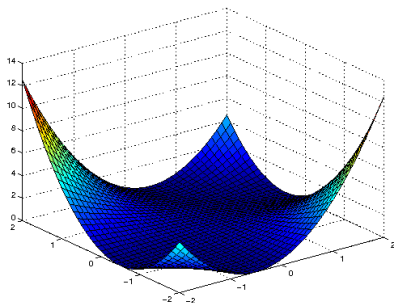
- A strongly convex function has a unique global optimum  $\mathbf{x}^*$  (why?)
- If  $f$  is twice differentiable, then

$f$  is strongly convex if and only if  $\nabla^2 f(\mathbf{x}) \succ mI$  for all  $\mathbf{x}$

- Gradient descent, coordinate descent will converge linearly (will see later)

# Nonconvex Functions

- If  $f$  is **nonconvex**, most algorithms can only converge to **stationary points**
- $\bar{x}$  is a stationary point if and only if  $\nabla f(\bar{x}) = 0$
- Three types of stationary points:
  - (1) Global optimum
  - (2) Local optimum
  - (3) Saddle point
- Example: matrix completion, neural network, ...
- Example:  $f(x, y) = \frac{1}{2}(xy - a)^2$





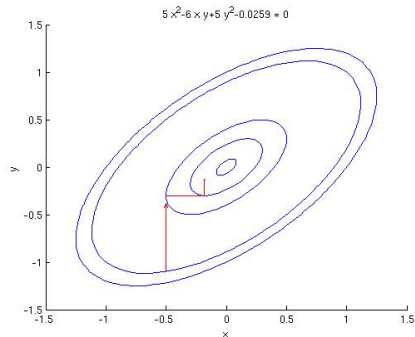
# Coordinate Descent

# Coordinate Descent

- Update one variable at a time
- Coordinate Descent: repeatedly perform the following loop
  - Step 1: pick an index  $i$
  - Step 2: compute a step size  $\delta^*$  by (approximately) minimizing

$$\operatorname{argmin}_{\delta} f(\mathbf{x} + \delta \mathbf{e}_i)$$

Step 3:  $x_i \leftarrow x_i + \delta^*$



# Coordinate Descent (update sequence)

- Three types of updating order:
- Cyclic: update sequence

$$\underbrace{x_1, x_2, \dots, x_n}_{\text{1st outer iteration}}, \underbrace{x_1, x_2, \dots, x_n}_{\text{2nd outer iteration}}, \dots$$

- Randomly permute the sequence for each outer iteration (faster convergence in practice)

Some interesting theoretical analysis for this recently

C.-P. Lee and S. J. Wright, Random Permutations Fix a Worst Case for Cyclic Coordinate Descent. 2016

# Coordinate Descent (update sequence)

- Three types of updating order:
- Cyclic: update sequence

$$\underbrace{x_1, x_2, \dots, x_n}_{\text{1st outer iteration}}, \underbrace{x_1, x_2, \dots, x_n}_{\text{2nd outer iteration}}, \dots$$

- Randomly permute the sequence for each outer iteration (faster convergence in practice)

Some interesting theoretical analysis for this recently

C.-P. Lee and S. J. Wright, Random Permutations Fix a Worst Case for Cyclic Coordinate Descent. 2016

- Random: each time pick a random coordinate to update
  - Typical way: sample from uniform distribution
  - Sample from uniform distribution vs sample from biased distribution

P. Zhao and T. Zhang, Stochastic Optimization with Importance Sampling for Regularized Loss Minimization. In ICML 2015

D. Csiba, Z. Qu and P. Richtarik, Stochastic Dual Coordinate Ascent with Adaptive Probabilities. In ICML 2015

# Greedy Coordinate Descent

- Greedy: choose the most “important” coordinate to update
- How to measure the importance?
  - By first derivative:  $|\nabla_i f(\mathbf{x})|$
  - By first and second derivative:  $|\nabla_i f(\mathbf{x})/\nabla_{ii}^2 f(\mathbf{x})|$
  - By maximum reduction of objective function

$$i^* = \operatorname{argmax}_{i=1,\dots,n} \left( f(\mathbf{x}) - \min_{\delta} f(\mathbf{x} + \delta \mathbf{e}_i) \right)$$

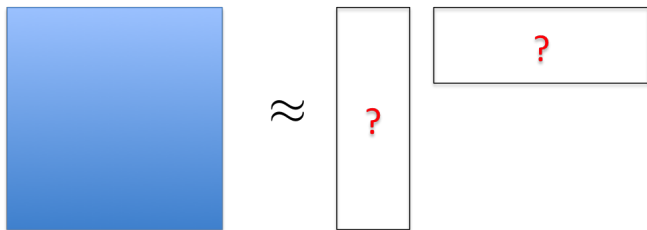
- Need to consider the time complexity for variable selection
- Useful for kernel SVM (see lecture 6)

## Extension: block coordinate descent

- Variables are divided into blocks  $\{\mathcal{X}_1, \dots, \mathcal{X}_p\}$ , where each  $\mathcal{X}_i$  is a subset of variables and

$$\mathcal{X}_1 \cup \mathcal{X}_2, \dots, \mathcal{X}_p = \{1, \dots, n\}, \quad \mathcal{X}_i \cap \mathcal{X}_j = \varnothing, \quad \forall i, j$$

- Each time update a  $\mathcal{X}_i$  by (approximately) solving the subproblem within the block
- Example: alternating minimization for matrix completion (2 blocks). (See lecture 7)



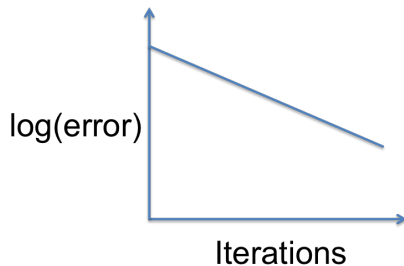
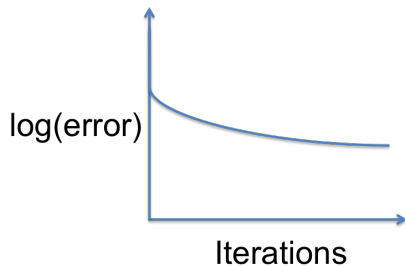
# Coordinate Descent (convergence)

- Converge to optimal if  $f(\cdot)$  is convex and smooth
- Has a linear convergence rate if  $f(\cdot)$  is strongly convex
- Linear convergence: error  $f(\mathbf{x}^t) - f(\mathbf{x}^*)$  decays as

$$\beta, \beta^2, \beta^3, \dots$$

for some  $\beta < 1$ .

- Local linear convergence: an algorithm converges linearly after  $\|\mathbf{x} - \mathbf{x}^*\| \leq K$  for some  $K > 0$



## Coordinate Descent: other names

- Alternating minimization (matrix completion)
- Iterative scaling (for log-linear models)
- Decomposition method (for kernel SVM)
- Gauss Seidel (for linear system when the matrix is positive definite)
- ...



# Gradient Descent

# Gradient Descent

- Gradient descent algorithm: repeatedly conduct the following update:

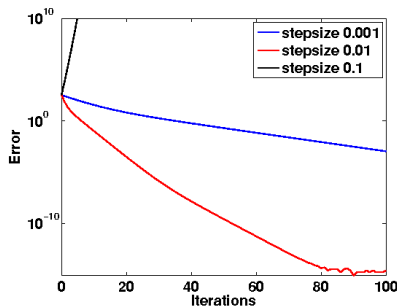
$$\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t - \alpha \nabla f(\mathbf{x}^t)$$

where  $\alpha > 0$  is the **step size**

- It is a fixed point iteration method:

$$\mathbf{x} - \alpha \nabla f(\mathbf{x}) = \mathbf{x} \text{ if and only if } \mathbf{x} \text{ is an optimal solution}$$

- Step size too large  $\Rightarrow$  diverge; too small  $\Rightarrow$  slow convergence



# Gradient Descent: successive approximation

- At each iteration, form an approximation of  $f(\cdot)$ :

$$\begin{aligned} f(\mathbf{x}^t + \mathbf{d}) &\approx \tilde{f}_{\mathbf{x}^t}(\mathbf{d}) := f(\mathbf{x}^t) + \nabla f(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \left( \frac{1}{\alpha} \mathbf{I} \right) \mathbf{d} \\ &= f(\mathbf{x}^t) + \nabla f(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2\alpha} \mathbf{d}^T \mathbf{d} \end{aligned}$$

- Update solution by  $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$
- $\mathbf{d}^* = -\alpha \nabla f(\mathbf{x}^t)$  is the minimizer of  $\operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$
- $\mathbf{d}^*$  will decrease the original objective function  $f$   
if  $\alpha$  (step size) is small enough

# Gradient Descent: successive approximation

- However, the function value will decrease if

Condition 1:  $\tilde{f}_{\mathbf{x}}(\mathbf{d}) \geq f(\mathbf{x} + \mathbf{d})$  for all  $\mathbf{d}$

Condition 2:  $\tilde{f}_{\mathbf{x}}(\mathbf{0}) = f(\mathbf{x})$

- Why?

$$\begin{aligned} f(\mathbf{x}^t + \mathbf{d}^*) &\leq \tilde{f}_{\mathbf{x}^t}(\mathbf{d}^*) \\ &\leq \tilde{f}_{\mathbf{x}^t}(\mathbf{0}) \\ &= f(\mathbf{x}^t) \end{aligned}$$

- Condition 2 is satisfied by construction of  $\tilde{f}_{\mathbf{x}^t}$
- Condition 1 is satisfied if  $\frac{1}{\alpha}I \succeq \nabla^2 f(\mathbf{x})$  for all  $\mathbf{x}$

# Gradient Descent: step size

- A twice differentiable function has  $L$ -Lipchitz continuous gradient if and only if

$$\nabla^2 f(\mathbf{x}) \leq LI \quad \forall \mathbf{x}$$

- In this case, Condition 2 is satisfied if  $\alpha < \frac{1}{L}$
- **Theorem:** gradient descent converges if  $\alpha < \frac{1}{L}$
- **Theorem:** gradient descent converges linearly with  $\alpha < \frac{1}{L}$  if  $f$  is strongly convex

# Gradient Descent

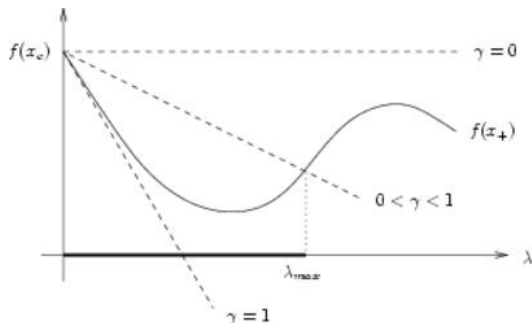
- In practice, we do not know  $L$  . . . . .
- Step size  $\alpha$  too large: the algorithm **diverges**
- Step size  $\alpha$  too small: the algorithm **converges very slowly**

# Gradient Descent: line search

- $\mathbf{d}^*$  is a “descent direction” if and only if  $(\mathbf{d}^*)^T \nabla f(\mathbf{x}) < 0$
- Armijo rule backtracking line search:  
Try  $\alpha = 1, \frac{1}{2}, \frac{1}{4}, \dots$  until it satisfies

$$f(\mathbf{x} + \alpha \mathbf{d}^*) \leq f(\mathbf{x}) + \gamma \alpha (\mathbf{d}^*)^T \nabla f(\mathbf{x})$$

where  $0 < \gamma < 1$



# Gradient Descent: line search

- Gradient descent with line search:
  - Converges to optimal solutions if  $f$  is smooth
  - Converges linearly if  $f$  is strongly convex
- However, each iteration requires evaluating  $f$  several times
- Several other step-size selection approaches  
(an ongoing research topic, especially for stochastic gradient descent)



# Gradient Descent: applying to ridge regression

**Input:**  $X \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^N$ , initial  $\mathbf{w}^{(0)}$

**Output:** Solution  $\mathbf{w}^* := \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$

1:  $t = 0$

2: **while** not converged **do**

3:   Compute the gradient

$$\mathbf{g} = X^T(X\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}$$

4:   Choose step size  $\alpha^t$

5:   Update  $\mathbf{w} \leftarrow \mathbf{w} - \alpha^t \mathbf{g}$

6:    $t \leftarrow t + 1$

7: **end while**

Time complexity:  $O(\operatorname{nnz}(X))$  per iteration

# Proximal Gradient Descent

- How can we apply gradient descent to solve the Lasso problem?

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \underbrace{\|\mathbf{w}\|_1}_{\text{non-differentiable}}$$

- General composite function minimization:

$$\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) := \{g(\mathbf{x}) + h(\mathbf{x})\}$$

where  $g$  is smooth and convex,  $h$  is convex but may be non-differentiable

- Usually assume  $h$  is simple (for computational efficiency)

# Proximal Gradient Descent: successive approximation

- At each iteration, form an approximation of  $f(\cdot)$ :

$$\begin{aligned} f(\mathbf{x}^t + \mathbf{d}) &\approx \tilde{f}_{\mathbf{x}^t}(\mathbf{d}) := g(\mathbf{x}^t) + \nabla g(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \left( \frac{1}{\alpha} I \right) \mathbf{d} + h(\mathbf{x}^t + \mathbf{d}) \\ &= g(\mathbf{x}^t) + \nabla g(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2\alpha} \mathbf{d}^T \mathbf{d} + h(\mathbf{x}^t + \mathbf{d}) \end{aligned}$$

- Update solution by  $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$
- This is called “proximal” gradient descent
- Usually  $\mathbf{d}^* = \operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$  has a closed form solution

# Proximal Gradient Descent: $\ell_1$ -regularization

- The subproblem:

$$\begin{aligned}\mathbf{x}^{t+1} &= \mathbf{x}^t + \underset{\mathbf{d}}{\operatorname{argmin}} \nabla g(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2\alpha} \mathbf{d}^T \mathbf{d} + \lambda \|\mathbf{x}^t + \mathbf{d}\|_1 \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{u} - (\mathbf{x}^t - \alpha \nabla g(\mathbf{x}^t))\|^2 + \lambda \alpha \|\mathbf{u}\|_1 \\ &= \mathcal{S}(\mathbf{x}^t - \alpha \nabla g(\mathbf{x}^t), \alpha \lambda),\end{aligned}$$

where  $\mathcal{S}$  is the soft-thresholding operator defined by

$$\mathcal{S}(a, z) = \begin{cases} a - z & \text{if } a > z \\ a + z & \text{if } a < -z \\ 0 & \text{if } a \in [-z, z] \end{cases}$$

# Proximal Gradient: soft-thresholding

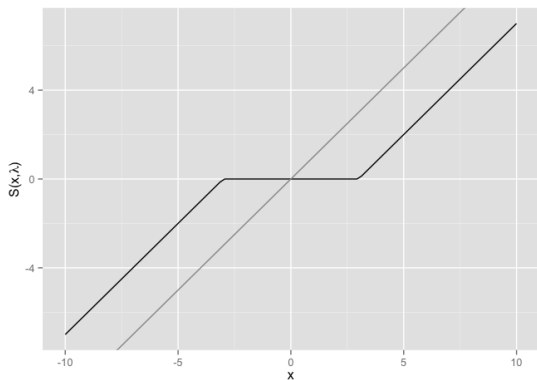


Figure from <http://jocelynchi.com/soft-thresholding-operator-and-the-lasso-solution/>

# Proximal Gradient Descent for Lasso

**Input:**  $X \in \mathbb{R}^{N \times d}$ ,  $\mathbf{y} \in \mathbb{R}^N$ , initial  $\mathbf{w}^{(0)}$

**Output:** Solution  $\mathbf{w}^* := \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$

1:  $t = 0$

2: **while** not converged **do**

3:   Compute the gradient

$$\mathbf{g} = X^T(X\mathbf{w} - \mathbf{y})$$

4:   Choose step size  $\alpha^t$

5:   Update  $\mathbf{w} \leftarrow \mathcal{S}(\mathbf{w} - \alpha^t \mathbf{g}, \alpha^t \lambda)$

6:    $t \leftarrow t + 1$

7: **end while**

Time complexity:  $O(\operatorname{nnz}(X))$  per iteration

# Newton's Method

# Newton's Method

- Iteratively conduct the following updates:

$$\mathbf{x} \leftarrow \mathbf{x} - \alpha \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

where  $\alpha$  is the step size

- If  $\alpha = 1$ : converges quadratically when  $\mathbf{x}^t$  is close enough to  $\mathbf{x}^*$ :

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\| \leq K \|\mathbf{x}^t - \mathbf{x}^*\|^2$$

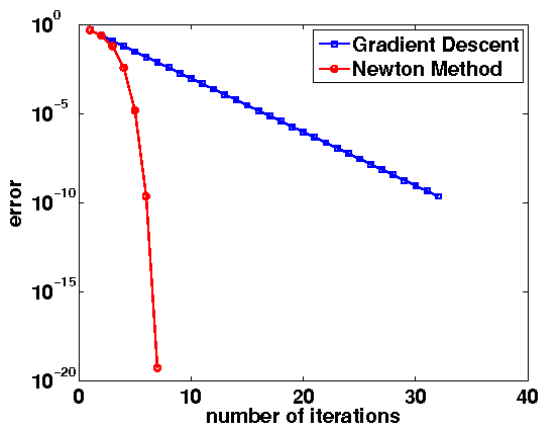
for some constant  $K$ . This means the error  $f(\mathbf{x}^t) - f(\mathbf{x}^*)$  decays quadratically:

$$\beta, \beta^2, \beta^4, \beta^8, \beta^{16}, \dots$$

- Only need few iterations to converge in this “quadratic convergence region”



# Newton's Method



However, Newton's update rule is more expensive than gradient descent/coordinate descent

# Newton's Method

- Need to compute  $\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$
- Closed form solution:  $O(d^3)$  for solving a  $d$  dimensional linear system
- Usually solved by another iterative solver:
  - gradient descent
  - coordinate descent
  - conjugate gradient method
  - ...
- Examples: primal L2-SVM/logistic regression,  $\ell_1$ -regularized logistic regression, ...

# Newton's Method

- At each iteration, form an approximation of  $f(\cdot)$ :

$$f(\mathbf{x}^t + \mathbf{d}) \approx \tilde{f}_{\mathbf{x}^t}(\mathbf{d}) := f(\mathbf{x}^t) + \nabla f(\mathbf{x}^t)^T \mathbf{d} + \frac{1}{2\alpha} \mathbf{d}^T \nabla^2 f(\mathbf{x}) \mathbf{d}$$

- Update solution by  $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$
- When  $\mathbf{x}$  is far away from  $\mathbf{x}^*$ , needs line search to guarantee convergence
- Assume  $L\mathbf{I} \succeq \nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}$  for all  $\mathbf{x}$ , then  $\alpha \leq \frac{m}{L}$  guarantee the objective function value decrease because

$$\frac{L}{m} \nabla^2 f(\mathbf{x}) \succeq \nabla^2 f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y}$$

- In practice, we often just use line search.

# Proximal Newton

- What if  $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$  and  $h(\mathbf{x})$  is non-smooth ( $h(\mathbf{x}) = \|\mathbf{x}\|_1$ )?
- At each iteration, form an approximation of  $f(\cdot)$ :

$$f(\mathbf{x}^t + \mathbf{d}) \approx \tilde{f}_{\mathbf{x}^t}(\mathbf{d}) := g(\mathbf{x}^t) + \nabla g(\mathbf{x}^t)^T \mathbf{d} + \frac{\alpha}{2} \mathbf{d}^T \nabla^2 g(\mathbf{x}) \mathbf{d} + h(\mathbf{x} + \mathbf{d})$$

- Update solution by  $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t + \operatorname{argmin}_{\mathbf{d}} \tilde{f}_{\mathbf{x}^t}(\mathbf{d})$
- Need another iterative solver for solving the subproblem

# Stochastic Gradient Method

# Stochastic Gradient Method: Motivation

- Widely used for machine learning problems (with large number of samples)
- Given training samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , we usually want to solve the following empirical risk minimization (ERM) problem:

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \ell_i(\mathbf{x}_i),$$

where each  $\ell_i(\cdot)$  is the loss function

- Minimize the summation of individual loss on each sample

# Stochastic Gradient Method

- Assume the objective function can be written as

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- Stochastic gradient method:

Iterative conducts the following updates

- Choose an index  $i$  (uniform) randomly
  - $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t - \eta^t \nabla f_i(\mathbf{x}^t)$
- $\eta^t > 0$  is the step size

# Stochastic Gradient Method

- Assume the objective function can be written as

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- Stochastic gradient method:

Iterative conducts the following updates

- Choose an index  $i$  (uniform) randomly
  - $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t - \eta^t \nabla f_i(\mathbf{x}^t)$
- $\eta^t > 0$  is the step size
  - Why does SG work?

$$E_i[\nabla f_i(\mathbf{x})] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}) = \nabla f(\mathbf{x})$$



# Stochastic Gradient Method

- Assume the objective function can be written as

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

- Stochastic gradient method:

Iterative conducts the following updates

- Choose an index  $i$  (uniform) randomly
  - $\mathbf{x}^{t+1} \leftarrow \mathbf{x}^t - \eta^t \nabla f_i(\mathbf{x}^t)$
- $\eta^t > 0$  is the step size
  - Why does SG work?

$$E_i[\nabla f_i(\mathbf{x})] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}) = \nabla f(\mathbf{x})$$

- Is it a fixed point method? No if  $\eta > 0$  because  $\mathbf{x}^* - \eta \nabla_i f(\mathbf{x}^*) \neq \mathbf{x}^*$
- Is it a descent method? No, because  $f(\mathbf{x}^{t+1}) \not\leq f(\mathbf{x}^t)$

# Stochastic Gradient

- Step size  $\eta$  has to decay to 0  
(e.g.,  $\eta^t = Ct^{-a}$  for some constant  $a, C$ )  
SGD converges sub-linearly
- Many variants proposed recently
  - SVRG, SAGA (2013, 2014): variance reduction
  - AdaGrad (2011): adaptive learning rate
  - RMSProp (2012): estimate learning rate by a running average of gradient.
  - Adam (2015): adaptive moment estimation
- Widely used in machine learning

# SGD in Deep Learning

- AdaGrad: adaptive step size for each parameter

Update rule at the  $t$ -th iteration:

- Compute  $\mathbf{g} = \nabla f(\mathbf{x})$
  - Estimate the second moment:  $G_{ii} = G_{ii} + g_i^2$  for all  $i$
  - Parameter update:  $x_i \leftarrow x_i - \frac{\eta}{\sqrt{G_{ii}}} g_i$  for all  $i$
- Proposed for convex optimization in:
    - “Adaptive subgradient methods for online learning and stochastic optimization” (JMLR 2011)
    - “Adaptive Bound Optimization for Online Convex Optimization” (COLT 2010)

# SGD in Deep Learning

- AdaGrad: adaptive step size for each parameter

Update rule at the  $t$ -th iteration:

- Compute  $\mathbf{g} = \nabla f(\mathbf{x})$
  - Estimate the second moment:  $G_{ii} = G_{ii} + g_i^2$  for all  $i$
  - Parameter update:  $x_i \leftarrow x_i - \frac{\eta}{\sqrt{G_{ii}}} g_i$  for all  $i$
- Adam (Kingma and Ba, 2015): maintain both first moment and second moment

Update rule at the  $t$ -th iteration:

- Compute  $\mathbf{g} = \nabla f(\mathbf{x})$
- $\mathbf{m} = \beta_1 \mathbf{m} + (1 - \beta_1) \mathbf{g}$
- $\hat{\mathbf{m}} = \mathbf{m}_t / (1 - \beta_1^t)$  (estimate of first moment)
- $\mathbf{v} = \beta_2 \mathbf{v} + (1 - \beta_2) \mathbf{g}^2$
- $\hat{\mathbf{v}} = \mathbf{v}_t / (1 - \beta_2^t)$  (estimate of second moment)
- $\mathbf{x} \leftarrow \mathbf{x} - \eta \hat{\mathbf{m}} / (\hat{\mathbf{v}} + \epsilon)$

All the operations are element-wise

# Stochastic Gradient: Variance Reduction

- SGD is not a fixed point method: even if  $\mathbf{x}^*$  is optimal,

$$\mathbf{x}^* - \eta \nabla_i f(\mathbf{x}^*) \neq \mathbf{x}^*$$

- Reason:  $E[\nabla_i f(\mathbf{x}^*)] = \nabla f(\mathbf{x}^*) = 0$  but variance can be large
- Variance reduction: reduce the variance so that variance  $\rightarrow 0$  when  $t \rightarrow \infty$ .

- SVRG (Johnson and Zhang, 2013)
- For  $t = 1, 2, \dots$ 
  - 1  $\tilde{\mathbf{x}} = \mathbf{x}$
  - 2  $\tilde{\mathbf{g}} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{\mathbf{x}})$
  - 3 For  $s = 1, 2, \dots$ 
    - 1 Randomly pick  $i$
    - 2  $\mathbf{x} \leftarrow \mathbf{x} - \eta \underbrace{(\nabla f_i(\mathbf{x}) - \nabla f_i(\tilde{\mathbf{x}}) + \tilde{\mathbf{g}})}_{\mathbf{a}_i}$
- Can easily show
  - $E[\mathbf{a}_i] = 0$  (unbiased)
  - $\text{Var}[\mathbf{a}_i] = 0$  when  $\mathbf{x} = \tilde{\mathbf{x}} = \mathbf{x}^*$
- Have linear convergence rate

# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- How to write as  $\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ ?
- How to decompose into  $n$  components?

# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- How to write as  $\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ ?
- First approach:  $f_i(\mathbf{w}) = (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$
- Update rule:

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i - 2\eta^t \lambda \mathbf{w} \\ &= (1 - 2\eta^t \lambda) \mathbf{w} - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i \end{aligned}$$



# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- How to write as  $\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ ?
- First approach:  $f_i(\mathbf{w}) = (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$
- Update rule:

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i - 2\eta^t \lambda \mathbf{w} \\ &= (1 - 2\eta^t \lambda) \mathbf{w} - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i \end{aligned}$$

- Need  $O(d)$  complexity per iteration even if data is sparse

# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- How to write as  $\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ ?
- First approach:  $f_i(\mathbf{w}) = (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$
- Update rule:

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i - 2\eta^t \lambda \mathbf{w} \\ &= (1 - 2\eta^t \lambda) \mathbf{w} - 2\eta^t (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i \end{aligned}$$

- Need  $O(d)$  complexity per iteration even if data is sparse
- Solution: store  $\mathbf{w} = s\mathbf{v}$  where  $s$  is a scalar

# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- Second approach:

define  $\Omega_i = \{j \mid X_{ij} \neq 0\}$  for  $i = 1, \dots, n$

define  $n_j = |\{i \mid X_{ij} \neq 0\}|$  for  $j = 1, \dots, d$

define  $f_i(\mathbf{w}) = (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \sum_{j \in \Omega_i} \frac{\lambda n}{n_j} w_j^2$

# Stochastic Gradient: applying to ridge regression

- Objective function:

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

- Second approach:

define  $\Omega_i = \{j \mid X_{ij} \neq 0\}$  for  $i = 1, \dots, n$

define  $n_j = |\{i \mid X_{ij} \neq 0\}|$  for  $j = 1, \dots, d$

define  $f_i(\mathbf{w}) = (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \sum_{j \in \Omega_i} \frac{\lambda n}{n_j} w_j^2$

- Update rule when selecting index  $i$ :

$$w_j^{t+1} \leftarrow w_j^t - 2\eta^t (\mathbf{x}_i^T \mathbf{w}^t - y_i) X_{ij} - \frac{2\eta^t \lambda n}{n_j} w_j^t, \quad \forall j \in \Omega_i$$

- Solution: update can be done in  $O(|\Omega_i|)$  operations

## Coming up

- Next class: Parallel Optimization

Questions?