

# Iterative Optimization Algorithms

# The Big Picture

- Most of the optimization problems in this course don't have closed form solutions. Ridge regression is an exception.
- In these cases we resort to iterative optimization algorithms. Even when a closed form solution is available, an iterative solver can be more computationally efficient.
- Today we'll overview several iterative solvers

# Outline

- Gradient descent
- Stochastic gradient descent
- The subgradient method
- Coordinate descent

# ERM

- (*Regularized*) *empirical risk minimization* learns  $f$  by solving

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda \Omega(f),$$

- Different choices of  $L, \mathcal{F}, \Omega$  give rise to different methods.
- Example: ridge regression without offset

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

# Computational Complexity

- What is the computational complexity of computing

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda n \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

where

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$(n \times d)$

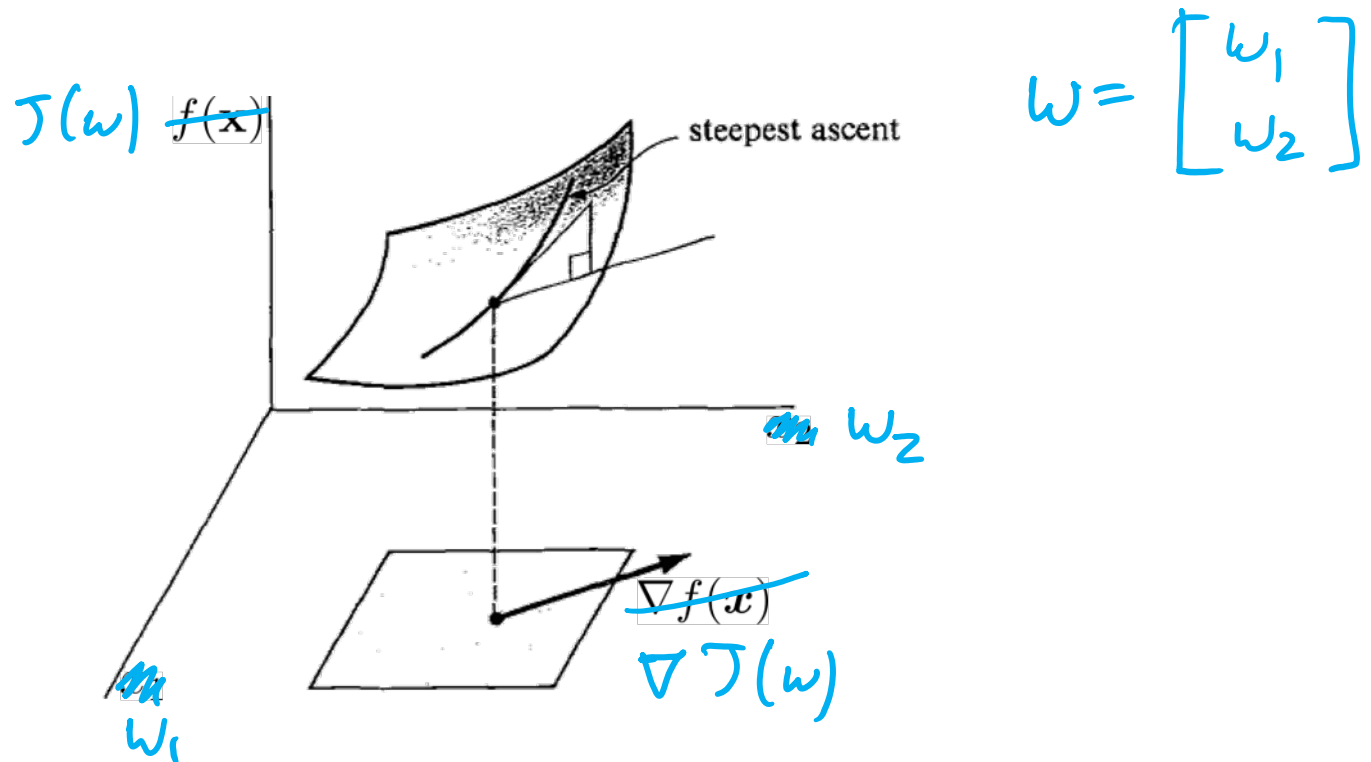
$(d \times d)$   $\mathbf{X}^T \mathbf{X} :$   $O(nd^2)$

$\mathbf{X}^T \mathbf{X} :$   $O(d^3)$

Overall:  $O(nd^2 + d^3)$

# Gradient

The gradient of a function is a vector that points in the direction of steepest ascent



# Gradient Descent

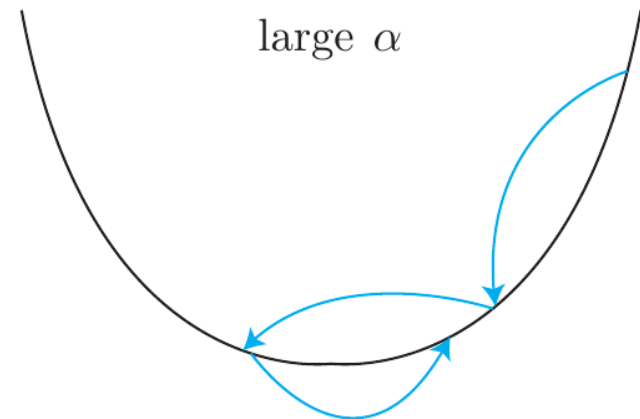
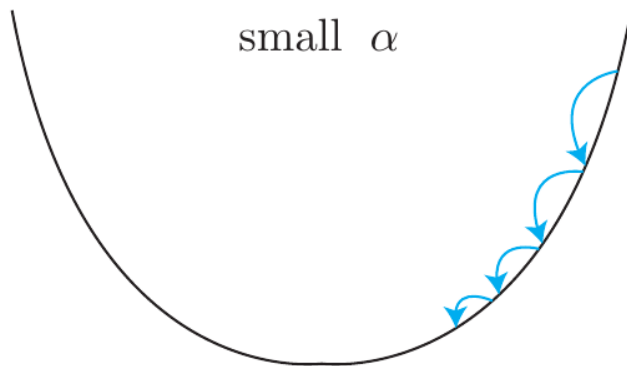
- Consider minimizing the generic objective function  $J(\theta)$
- Initial guess  $\theta_0$
- For  $t = 1, \dots, \text{max\_iter}$

$$\theta_t \leftarrow \theta_{t-1} - \alpha_t \nabla J(\theta_{t-1})$$

If convergence condition satisfied, exit

End

step size,  
learning rate



# Gradient Descent for Linear Regr.

The regularized least-squares (i.e., ridge regression) objective function can be written (~~after eliminating  $\hat{b}$~~ )

$$J(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{r}^T \mathbf{w} + c,$$

where  $\mathbf{A} = 2(\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})$ ,  $\mathbf{r} = -2\mathbf{X}^T \mathbf{y}$ , and  $c = \mathbf{y}^T \mathbf{y}$

1.  $\nabla J(\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{r} = 2(\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})\mathbf{w} - 2\mathbf{X}^T \mathbf{y}$
2. What is the computational complexity of gradient descent in terms of  $d$  and  $n$ ?

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \underbrace{\mathbf{X}^T}_{O(nd)} (\underbrace{\mathbf{X} \mathbf{w}}_{O(nd)})$$

conjugate  
gradient  
descent

Conclusion:

$O(nd)$  operations  
per iteration  
of GD.



# Stochastic Gradient Descent

- Suppose it is possible to write

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \mathcal{J}_i(\boldsymbol{\theta})$$

where  $J_i(\boldsymbol{\theta})$  depends on the training data only through  $(\mathbf{x}_i, y_i)$  ~~or, in the case of ridge regression,  $(\tilde{\mathbf{x}}_i, \tilde{y}_i)$ .~~

- *Stochastic gradient descent* is the following variation on gradient descent:

- Initialize  $\boldsymbol{\theta}_0$ , set  $t = 0$

- For  $j = 1, \dots, \text{max\_iter}$

For  $i = 1, \dots, n$  in random order

$$\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \alpha_j \nabla J_i(\boldsymbol{\theta}_{t-1})$$

$$t \leftarrow t + 1$$

End

If convergence condition satisfied, exit

End

$$\nabla J(\boldsymbol{\theta}) = \sum \nabla \mathcal{J}_i(\boldsymbol{\theta})$$

$$g = 0$$

For  $i = 1$  to  $n$

$$g = g + \nabla \mathcal{J}_i(\boldsymbol{\theta}_{t-1})$$

End

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \alpha g$$

GD

# SGD for Ridge Regression

$$\begin{aligned} J(w) &= \sum_{i=1}^n (y_i - w^T x_i)^2 + n\lambda \|w\|^2 \\ &= \sum_{i=1}^n J_i(w) \end{aligned}$$

where  $J_i(w) = (y_i - w^T x_i)^2 + \lambda \|w\|^2$

$$\nabla J_i(w) = -2(y_i - w^T x_i)x_i + 2\lambda w \in \mathbb{R}^d$$

$O(d)$  operations per update

# Poll

True or False: If the step-size  $\alpha_j$  is carefully chosen, then the ridge regression objective function decreases at every iteration  $j$  of gradient descent (unless you're already at a local min)

(A) True ✓

(B) False

line search  
backtracking

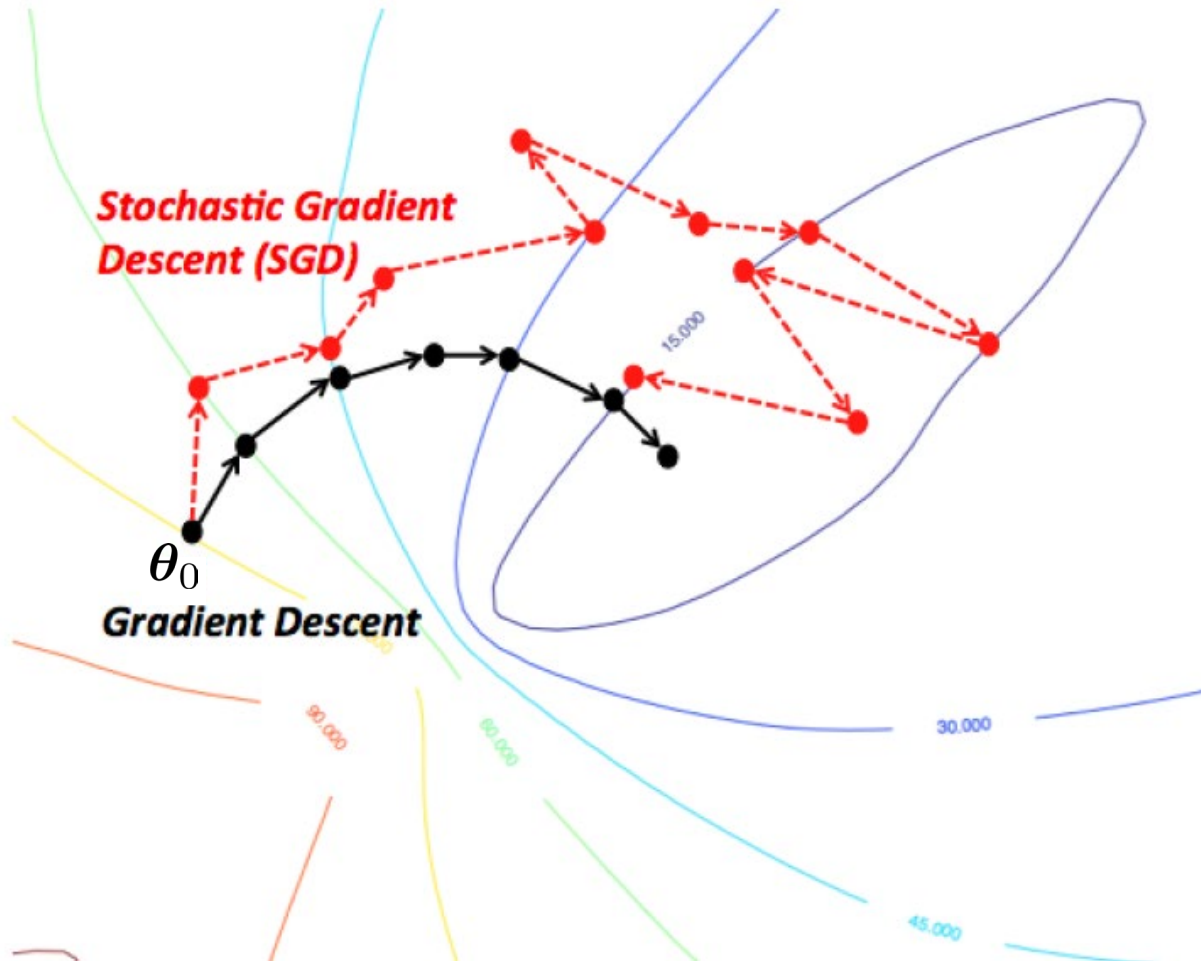
# Poll

True or False: If the step-size  $\alpha_j$  is carefully chosen, then the ridge regression objective function decreases at every iteration of stochastic gradient descent

(A) True

(B) False

# GD vs SGD



# The Lasso

- How can we solve

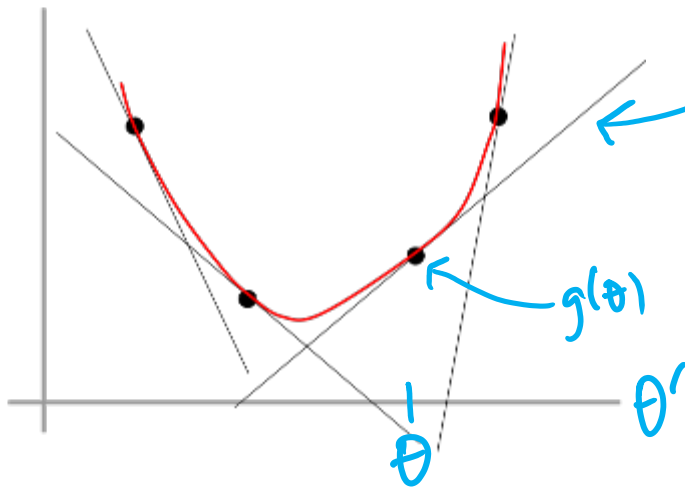
$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|_1?$$

# Subgradient Methods

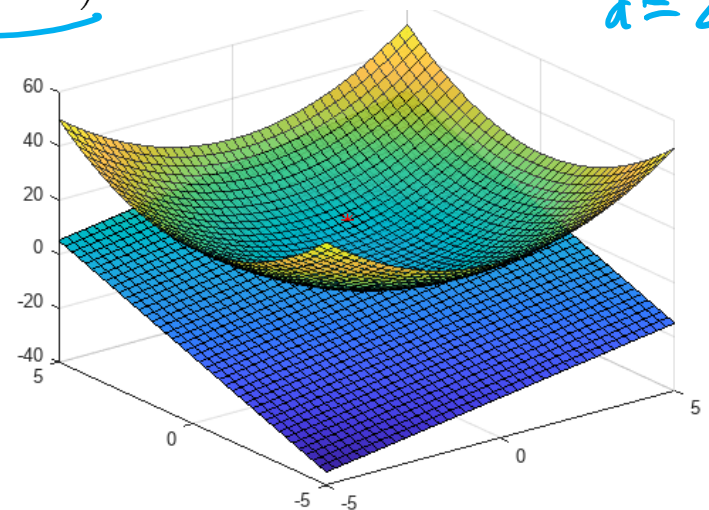
- The *subgradient method* is a generalization of gradient descent that applies to *nondifferentiable*, *convex* objective functions, like the lasso or ERM with hinge loss
- Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex, and let  $\theta \in \mathbb{R}^d$ . If  $g$  is differentiable, then  $u = \nabla g(\theta)$  is the only vector such that

$$g(\theta') \geq g(\theta) + u^T(\theta' - \theta) \quad \forall \theta'$$

$d=1$



$d=2$

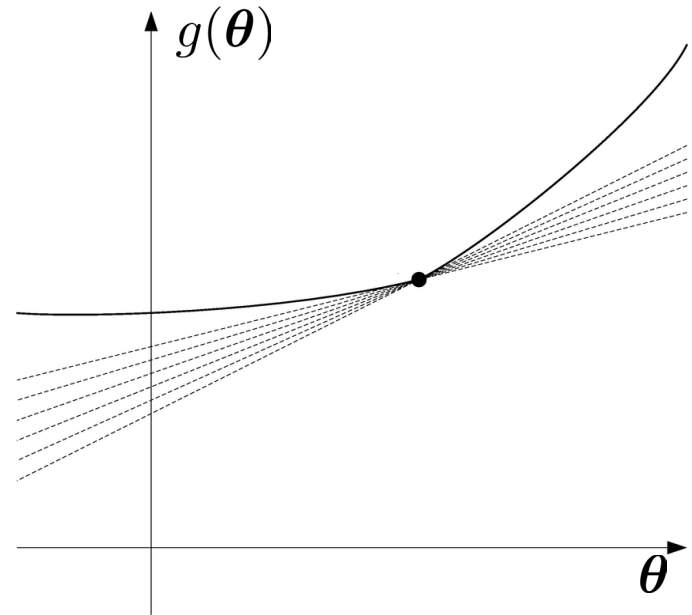


<https://mathoverflow.net/questions/327940/existence-of-a-strictly-convex-function-interpolating-given-gradients-and-values>

<https://www.mathworks.com/help/matlab/math/calculate-tangent-plane-to-surface.html>

# Subgradients

- If  $g$  is convex but *not* differentiable, then for some  $\theta$ , there may be many  $\mathbf{u}$  satisfying the previous inequality.
- We define the *subdifferential* of  $g$  at  $\theta$ , denoted  $\partial g(\theta)$ , to be the set of all  $\mathbf{u}$  satisfying the inequality.
- A *subgradient* is any element of the subdifferential.
- In the figure, the subdifferential is the interval  $[g'_-(\theta), g'_+(\theta)]$  where  $g'_-(\theta), g'_+(\theta)$  denote the left and right derivatives.





# Subgradient Method

- In the *subgradient method*, we update the parameter just as in gradient descent, but where the gradient is replaced by *any* subgradient.
- Pseudo-code for minimizing  $g(\boldsymbol{\theta})$ :
  - initialize  $\boldsymbol{\theta}_0$
  - $t \leftarrow 0$
  - Repeat
    - \* select  $\mathbf{u}_t \in \partial g(\boldsymbol{\theta}_t)$
    - \*  $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \alpha_t \mathbf{u}_t$
    - \*  $t \leftarrow t + 1$
  - Until stopping criterion satisfied
- If it is possible to write

$$g(\boldsymbol{\theta}) = \sum_{i=1}^n g_i(\boldsymbol{\theta})$$

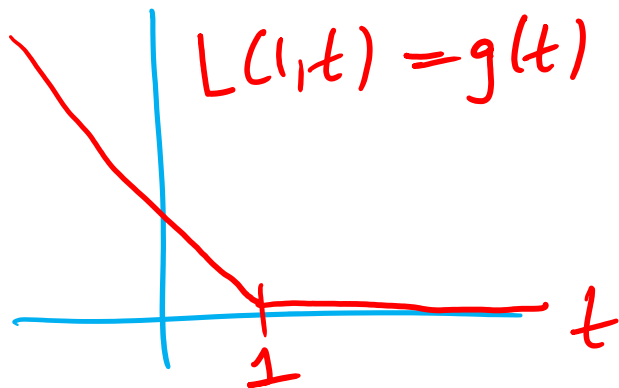
then we can also have a *stochastic subgradient method*, analogous to SGD.

# Subgradients

- What is the subdifferential of the hinge loss?

$$L(y, t) = \max(0, 1 - yt)$$

Take  $y=1$  for concreteness

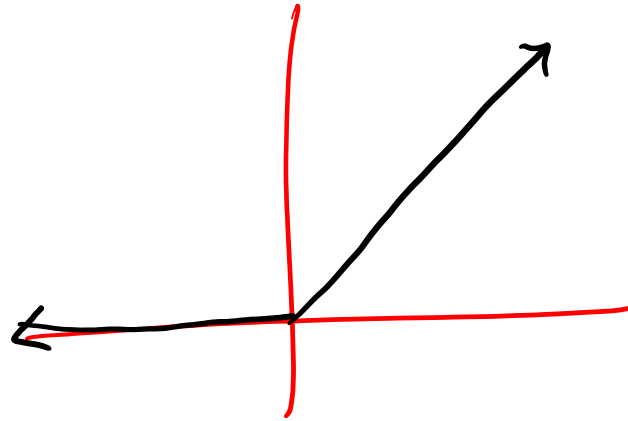


$$\partial g(t) = \begin{cases} -1 & t < 1 \\ [-1, 0] & t = 1 \\ 0 & t > 1 \end{cases}$$

$$\begin{aligned} t &< 1 \\ t &= 1 \\ t &> 1 \end{aligned}$$

# Poll

- The ReLU activation function is defined by  $\sigma(t) = \max(0, t)$ . The subdifferential of the ReLU function at  $t = 0$  is
  - (A)  $[-1, 0]$
  - (B)  $[0, 1]$  ✓
  - (C)  $[-1, 1]$
  - (D) None of the above



# Subgradients

- What is the subdifferential of  $\|w\|_1$ ?

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$\|w\|_1 = \sum_{j=1}^d |w_j|$$

$$\partial_{w_j} \|w\|_1 = \partial |w_j| = \begin{cases} -1 & w_j < 0 \\ [-1, 1] & w_j = 0 \\ 1 & w_j > 0 \end{cases}$$

$$\begin{aligned} w_j &< 0 \\ w_j &= 0 \\ w_j &> 0 \end{aligned}$$

$$\partial_w \|w\|_1 = \begin{bmatrix} \text{sign}(w_1) \\ \vdots \\ \text{sign}(w_d) \end{bmatrix}$$

$$\text{where } \text{sign}(0) = [-1, 1]$$

# Coordinate Descent

$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

- The subgradient method can converge slowly for the lasso.
- Coordinate descent cycles through the coordinates of  $\theta$ , updating one coordinate while leaving the others fixed.
- If  $J$  is the objective function, and  $\theta^{(0)}$  is the initial iterate, then  $\theta^{(1)}$  is obtained by

$$\theta_1^{(1)} = \arg \min_{\phi \in \mathbb{R}} J(\phi, \theta_2^{(0)}, \theta_3^{(0)}, \dots, \theta_p^{(0)})$$

$$\theta_2^{(1)} = \arg \min_{\phi \in \mathbb{R}} J(\theta_1^{(1)}, \phi, \theta_3^{(0)}, \dots, \theta_p^{(0)})$$

$\vdots$

# Coordinate Descent

- In general, if  $\theta^{(t)}$  is the  $t^{th}$  iterate, then

$$\theta_j^{(t)} = \underset{\theta_j}{\operatorname{argmin}} \quad J(\theta_1^{(t)}, \dots, \theta_{j-1}^{(t)}, \theta_j, \theta_{j+1}^{(t-1)}, \dots, \theta_p^{(t-1)})$$

Apply to lasso w/  $\theta = \begin{bmatrix} b \\ w \end{bmatrix} \in \mathbb{R}^{d+1}$

# Coordinate Descent for the Lasso

- Let's apply CD to minimize

$$\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|_1$$

- $b$  update:  $b^{(t)} = \bar{y} - (\mathbf{w}^{(t-1)})^T \bar{\mathbf{x}}$

- To update  $w_j^{(t)}$  we need to solve

$$\min_{w_j} g(w_j) := \frac{1}{n} \sum_{i=1}^n (y_i - [w_1^{(t)} \dots w_{j-1}^{(t)} w_j w_{j+1}^{(t-1)} \dots w_d^{(j-1)}] \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix} - b^{(t)})^2 + \lambda |w_j|$$

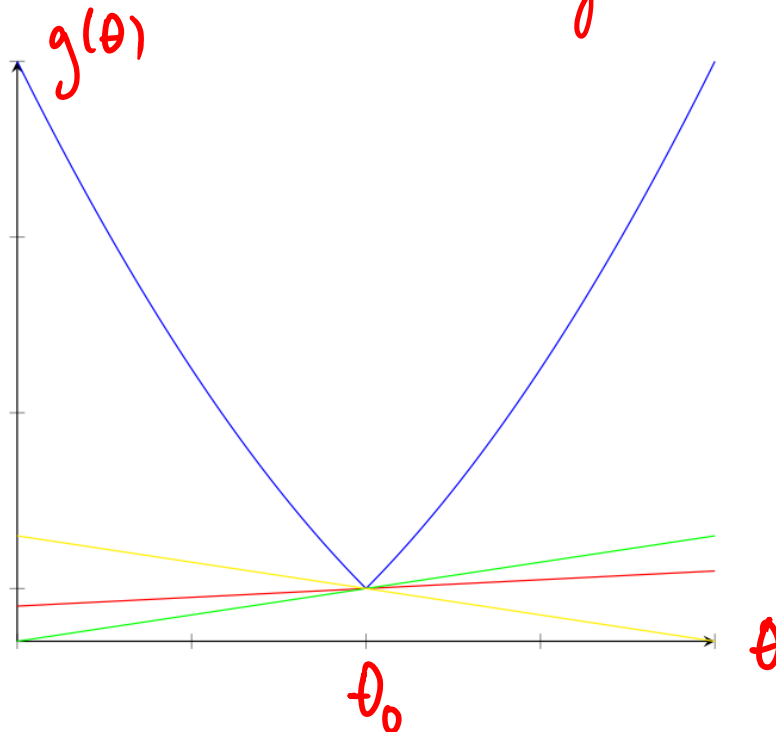
- How to solve?

# More Sugradients

- Let  $g(\theta)$  be a convex function. Then a point  $\theta_0$  is a global minimizer of  $g$  if and only if

$$0 \in \partial g(\theta).$$

*This follows directly from the definitions of global min and subdifferential.*





# Lasso Subproblem Solution

- Introduce the following notation:

$$\mathbf{w}_{-j}^{(t)} = \begin{bmatrix} w_1^{(t)} \\ \vdots \\ w_{j-1}^{(t)} \\ w_{j+1}^{(t-1)} \\ \vdots \\ w_d^{(t-1)} \end{bmatrix} \quad \mathbf{x}_{i,-j} = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,j-1} \\ x_{i,j+1} \\ \vdots \\ x_{i,d} \end{bmatrix}.$$

- You will show on HW4 that the subdifferential of  $g(w_j)$  is

$$\partial g(w_j) = \begin{cases} a_j^{(t)} w_j - c_j^{(t)} - \lambda, & w_j < 0 \\ [a_j^{(t)} w_j - c_j^{(t)} - \lambda, a_j^{(t)} w_j - c_j^{(t)} + \lambda], & w_j = 0 \\ a_j^{(t)} w_j - c_j^{(t)} + \lambda, & w_j > 0 \end{cases}.$$

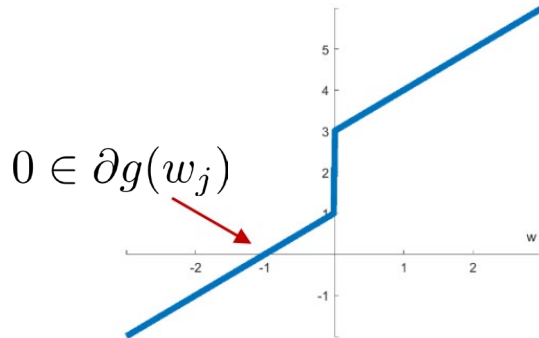
where

$$a_j^{(t)} = \frac{2}{n} \sum_i x_{ij}^2, \quad c_j^{(t)} = \frac{2}{n} \sum_i x_{ij} (y_i - (\mathbf{w}_{-j}^{(t)})^T \mathbf{x}_{i,-j} - b^{(t)}).$$

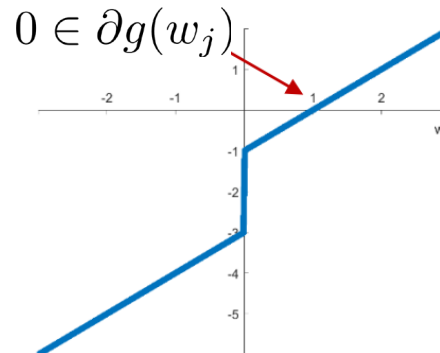
# Lasso Subproblem Solution

- On HW4, you will show that there is a unique value of  $w_j$  such that  $0 \in \partial g(w_j)$ .
- There are three cases to consider, shown below.

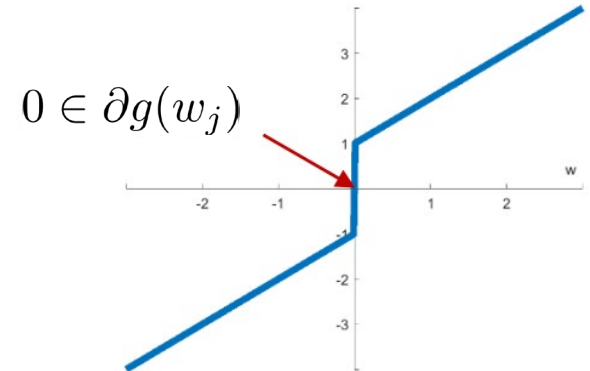
$$c_j^{(t)} < -\lambda$$



$$c_j^{(t)} > \lambda$$



$$c_j^{(t)} \in [-\lambda, \lambda]$$



# CD for Lasso: Final Algorithm

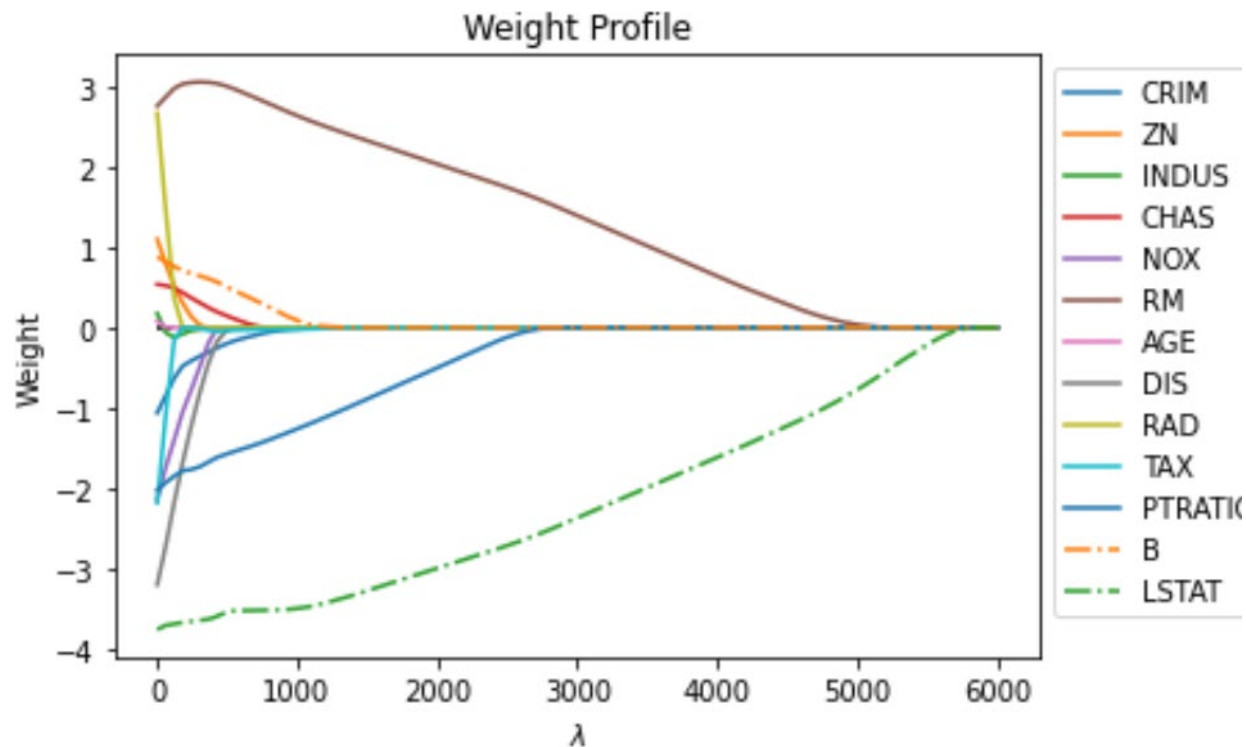
- Initialize  $\mathbf{w}^{(0)}, b^{(0)}$
- For  $t = 1, 2, \dots$ 
  - $b^{(t)} = \bar{y} - (\mathbf{w}^{(t-1)})^T \bar{\mathbf{x}}$
  - For  $j = 1, \dots, d$ 
    - \* Compute  $c_j^{(t)}, a_j^{(t)}$
    - \*  $w_j^{(t)} = \text{soft}(c_j^{(t)} / a_j^{(t)}, \lambda / a_j^{(t)})$
  - If stopping criterion met, break

Soft thresholding:

$$\text{soft}(\alpha, \beta) := \begin{cases} \alpha - \beta & \text{if } \alpha > \beta \\ 0 & \text{if } \alpha \in [-\beta, \beta] \\ \alpha + \beta & \text{if } \alpha < -\beta \end{cases}$$

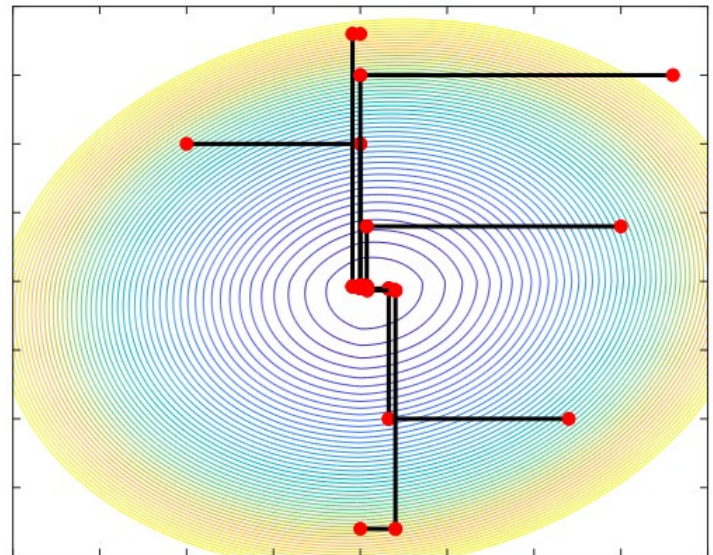
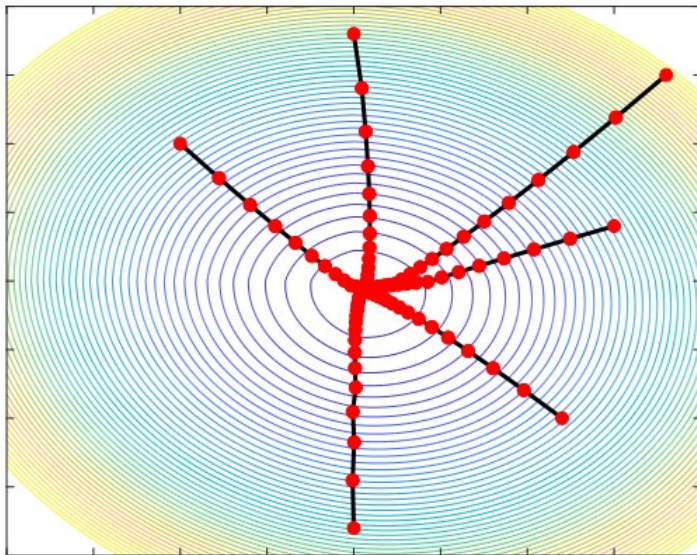
$$\begin{aligned} &\text{if } \alpha > \beta \\ &\text{if } \alpha \in [-\beta, \beta] \\ &\text{if } \alpha < -\beta \end{aligned}$$

# Boston Housing Data: Solution Path



# Subgradient Method vs CD

- $y = 1 \cdot x - 1 + z$
- $x, z$  are independent  $\mathcal{N}(0, 1)$
- $\lambda = 100$
- $\eta = 0.1$



# Final Thoughts on Lasso

- CD is typically much faster to converge than the subgradient method
- CD also has no tuning parameters
- Least angle regression (LARS): Solve Lasso for all  $\lambda$  in one algorithm
- Selected features are not stable: slight perturbation of training data can lead to significant changes to sparsity pattern

# List of Methods

- Gradient descent / SGD
  - Subgradient method / SSM
  - Coordinate descent
  - Newton's method
  - Quasi-Newton methods
  - Majorize-minimize
  - ADMM
  - . . .
- 
- Convergence guarantees exist for all of these methods, but are beyond our time constraints