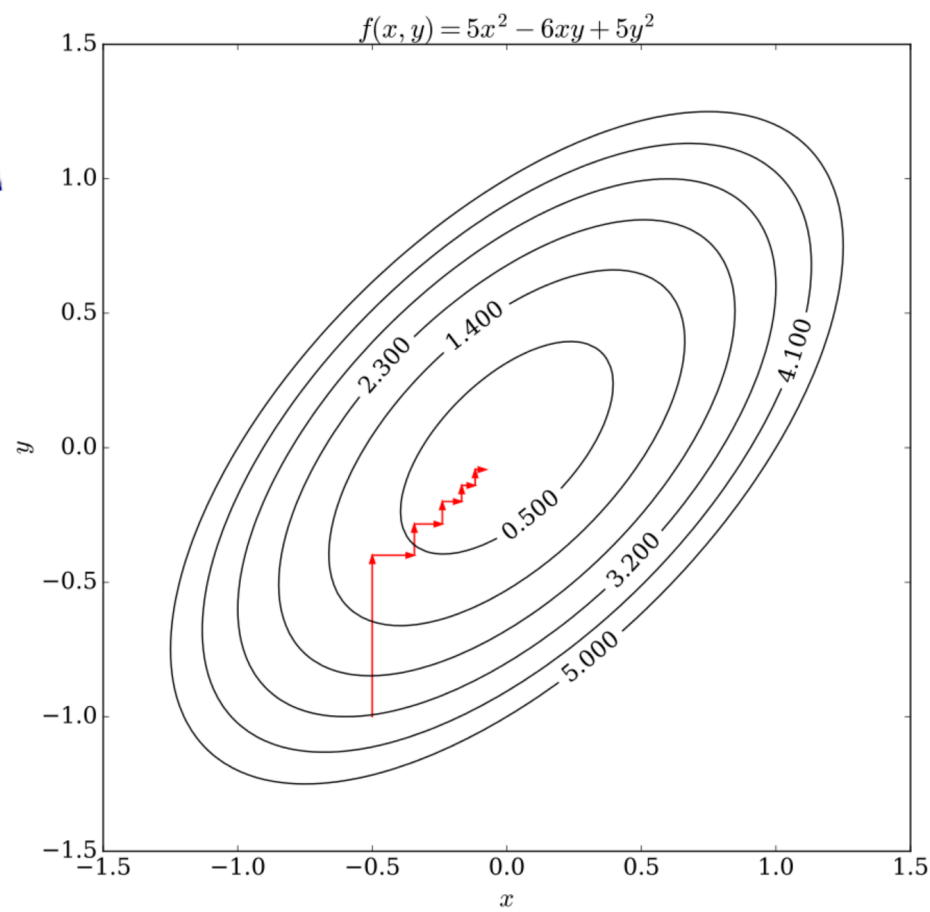
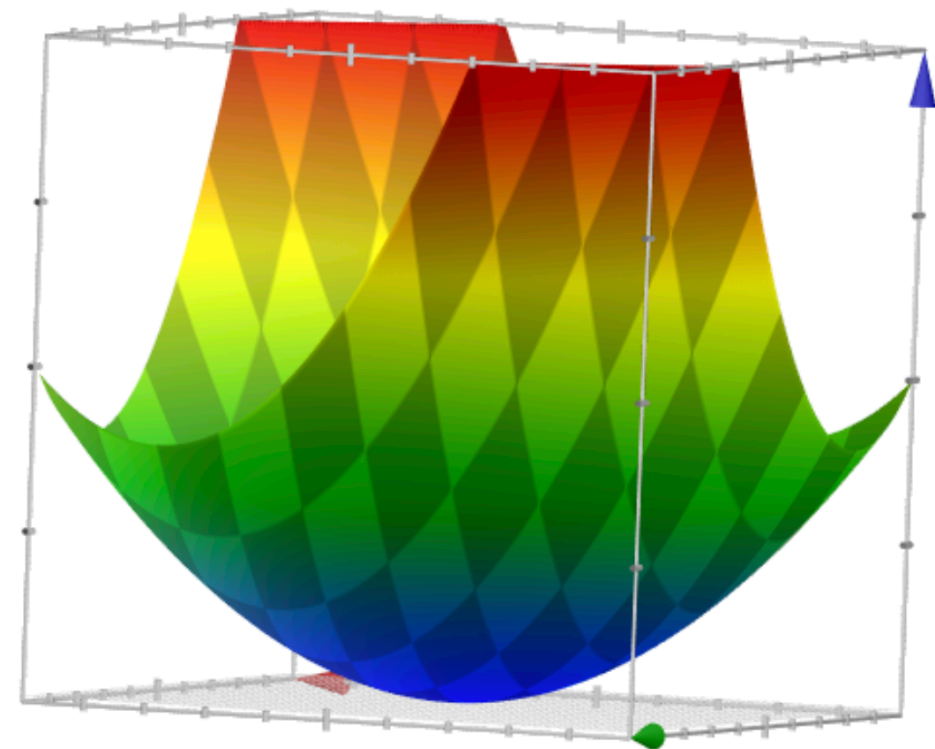


Coordinate Descent



Optimization: how do we solve Lasso?

- among many methods to find the solution, we will learn **coordinate descent method**
- as an illustrating example, we show coordinate descent updates on finding the minimum of $f(x, y) = 5x^2 - 6xy + 5y^2$



How do we solve Lasso: $\min_w \mathcal{L}(w) + \lambda \|w\|_1$?

- Coordinate descent

- input: training data S_{train} , max # of iterations T
- initialize: $w^{(0)} = \mathbf{0} \in \mathbb{R}^d$
- for $t = 1, \dots, T$
 - for $j = 1, \dots, d$
 - fix $w_1^{(t)}, \dots, w_{j-1}^{(t)}$ and $w_{j+1}^{(t-1)}, \dots, w_d^{(t-1)}$, and

$$w_j^{(t)} \leftarrow \arg \min_{w_j \in \mathbb{R}} \mathcal{L} \left(\begin{bmatrix} w_1^{(t)} \\ \vdots \\ w_{j-1}^{(t)} \\ w_j \\ w_{j+1}^{(t-1)} \\ \vdots \\ w_d^{(t-1)} \end{bmatrix} \right) + \lambda \left\| \begin{bmatrix} w_1^{(t)} \\ \vdots \\ w_{j-1}^{(t)} \\ w_j \\ w_{j+1}^{(t-1)} \\ \vdots \\ w_d^{(t-1)} \end{bmatrix} \right\|_1$$

this is a one-dimensional optimization, which is much easier to solve

Coordinate descent for (un-regularized) linear regression

- let us understand what coordinate descent does on a simpler problem of linear least squares, which minimizes

$$\text{minimize}_w \mathcal{L}(w) = \|\mathbf{X}w - \mathbf{y}\|_2^2$$

- note that we know that the optimal solution is

$$\hat{w}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

so we do not need to run any optimization algorithm

- we are solving this problem with **coordinate descent** as a starting example
- the main challenge we address is, how do we update $w_j^{(t)}$?
- let us derive an **analytical rule** for updating $w_j^{(t)}$

Coordinate descent for (un-regularized) linear regression

Coordinate descent for (un-regularized) linear regression

- we will study the case $j = 1$, for now (other cases are almost identical)

- when updating $w_1^{(t)}$, recall that

$$w_1^{(t)} \leftarrow \arg \min \| \mathbf{X}w - \mathbf{y} \|_2^2$$

$$\text{where } w = [w_1, w_2^{(t-1)}, \dots, w_d^{(t-1)}]^T$$

- first step is to write the objective function in terms of the variable we are optimizing over, that is w_1 :

$$\mathcal{L}(w) = \left\| \mathbf{X}[:,1]w_1 + \mathbf{X}[:,2:d]w_{2:d} - \mathbf{y} \right\|_2^2$$

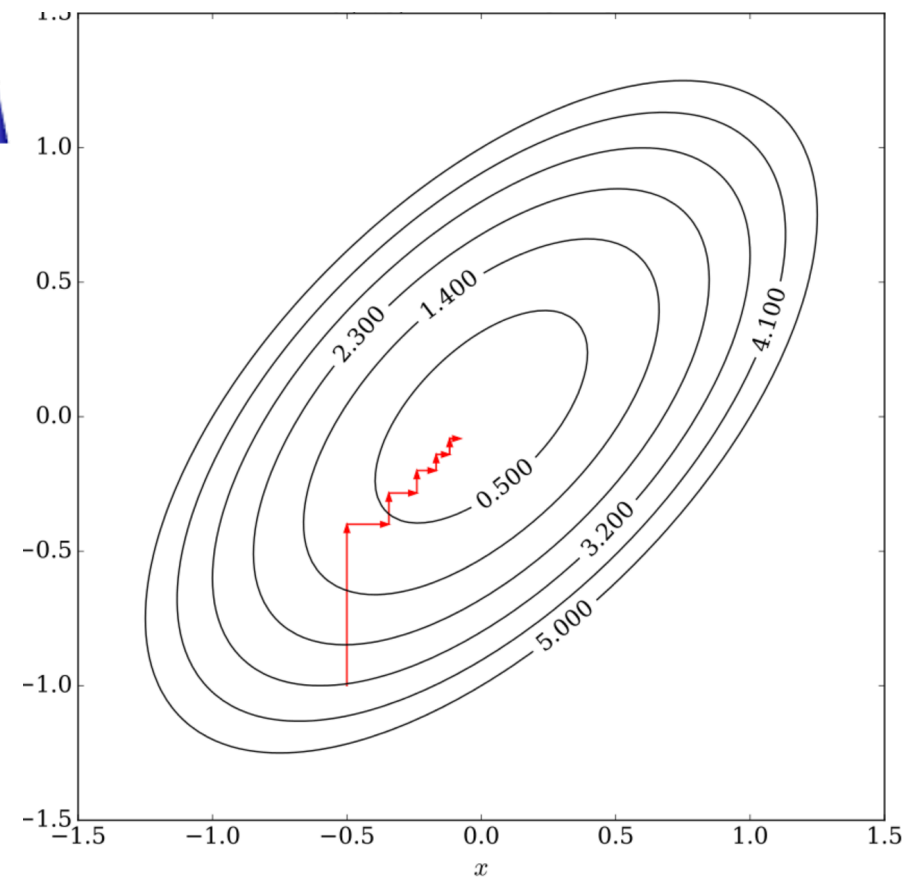
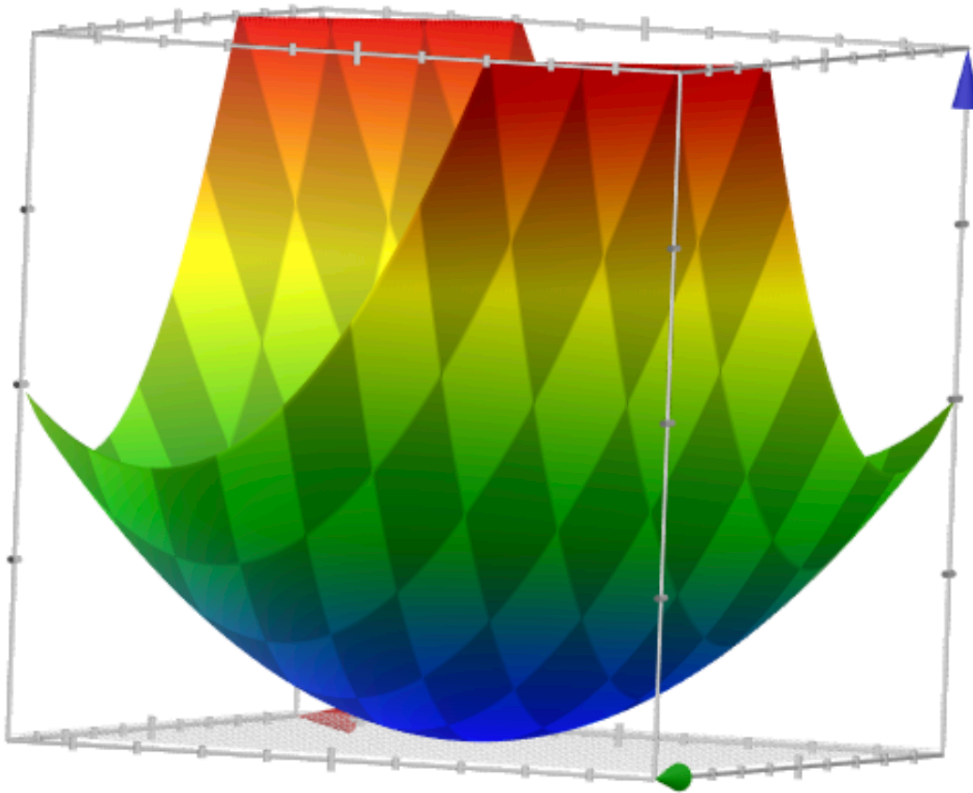
$$\text{where } w_{2:d} = [w_2^{(t-1)}, \dots, w_d^{(t-1)}]^T$$

$$\begin{bmatrix} \mathbf{X}[:,1] & | & \mathbf{X}[:,2:d] \end{bmatrix} \begin{bmatrix} w_1 \\ w_{2:d} \end{bmatrix} - \mathbf{y} = \begin{bmatrix} \mathbf{X}[:,1] \end{bmatrix} w_1 + \left(\begin{bmatrix} \mathbf{X}[:,2:d] \end{bmatrix} w_{2:d} - \mathbf{y} \right)$$

- we know from linear least squares that the minimizer is

$$w_1^{(t)} \leftarrow (\mathbf{X}[:,1]^T \mathbf{X}[:,1])^{-1} \mathbf{X}[:,1]^T (\mathbf{y} - \mathbf{X}[:,2:d]w_{2:d})$$

- Coordinate descent applied to a quadratic loss



Coordinate descent for Lasso

- let us apply coordinate descent on Lasso, which minimizes
$$\text{minimize}_w \mathcal{L}(w) + \lambda \|w\|_1 = \|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda \|w\|_1$$
- the goal is to derive an **analytical rule** for updating $w_j^{(t)}$'s
- let us first write the update rule explicitly for $w_1^{(t)}$
 - first step is to write the loss in terms of w_1

$$\left\| \mathbf{X}[:,1]w_1 - (\mathbf{y} - \mathbf{X}[:,2:d]w_{2:d}) \right\|_2^2 + \lambda \left(|w_1| + \underbrace{\|w_{2:d}\|_1}_{\text{constant}} \right)$$

- hence, the coordinate descent update boils down to

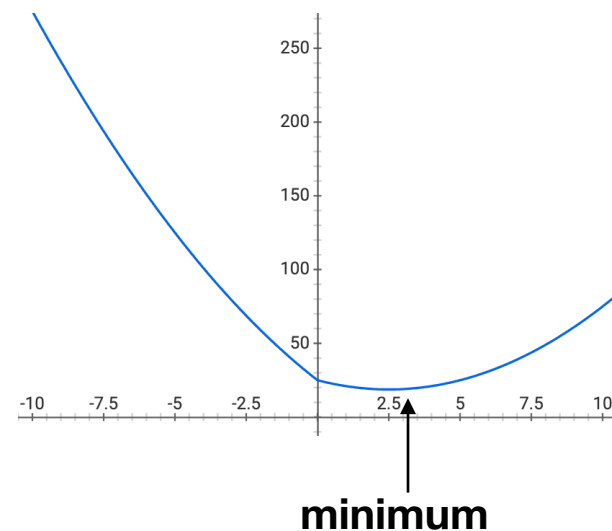
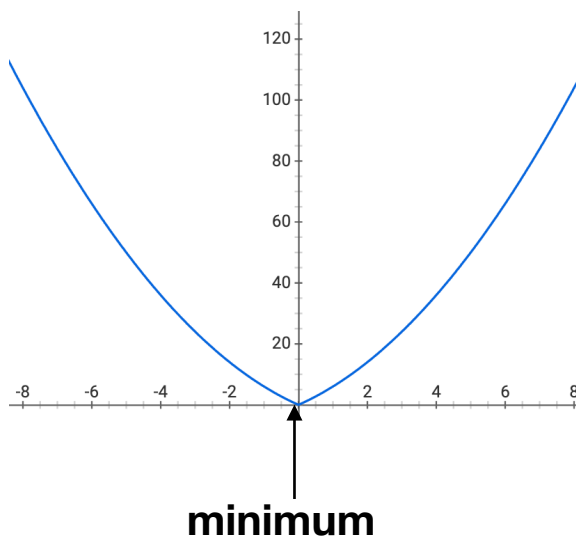
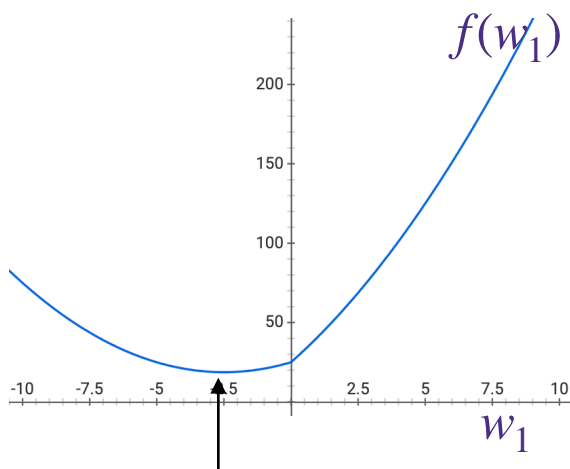
$$w_1^{(t)} \leftarrow \arg \min_{w_1} \underbrace{\left\| \mathbf{X}[:,1]w_1 - (\mathbf{y} - \mathbf{X}[:,2:d]w_{2:d}) \right\|_2^2 + \lambda |w_1|}_{f(w_1)}$$

Convexity

- to find the minimizer of $f(w_1)$, let's study some properties
- for simplicity, we represent the objective function as

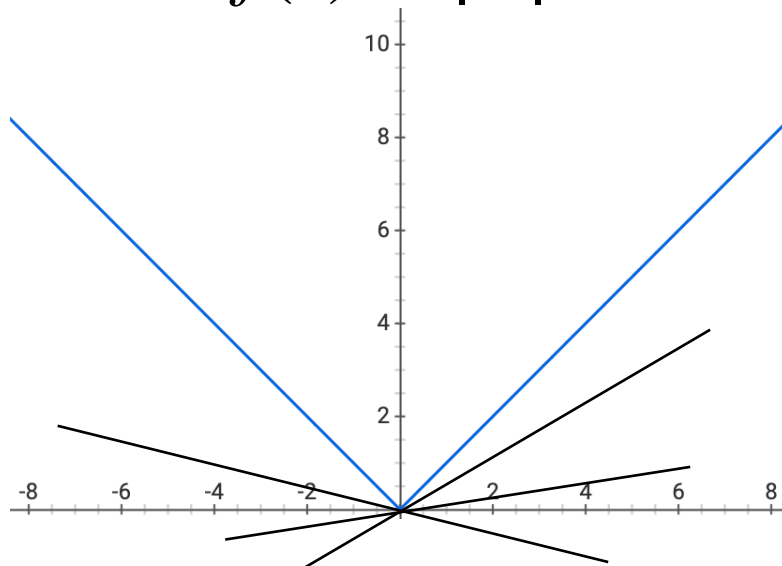
$$f(w_1) = (aw_1 - b)^2 + \lambda |w_1|$$

- this function is
 - **convex**, and
 - **non-differentiable**
- depending on the values of a and b, the function looks like one of the three below



Convexity

$$f(x) = |x|$$



- for a **non-differentiable** function, gradient is not defined at some points, for example at $x = 0$ for $f(x) = |x|$
- at such points, **sub-gradient** plays the role of gradient
 - sub-gradient at a differentiable point is the same as the gradient
 - sub-gradient at a non-differentiable point is a set of vector satisfying

$$\partial f(x) = \{ g \in \mathbb{R}^d \mid f(y) \geq f(x) + g^T(y - x), \text{ for all } y \in \mathbb{R}^d \}$$

- for example, $\partial |x| = \begin{cases} +1 & \text{for } x > 0 \\ [-1, 1] & \text{for } x = 0 \\ -1 & \text{for } x < 0 \end{cases}$

Computing the sub-gradient

$$w_1^{(t)} = \arg \min_{w_1} \underbrace{\left\| \mathbf{X}[:,1]w_1 - (\mathbf{y} - \mathbf{X}[:,2:d]w_{-1}) \right\|_2^2}_{f(w_1)} + \lambda |w_1|$$

Computing the sub-gradient

$$w_1^{(t)} = \arg \min_{w_1} \underbrace{\left\| \mathbf{X}[:, 1]w_1 - (\mathbf{y} - \mathbf{X}[:, 2 : d]w_{-1}) \right\|_2^2 + \lambda |w_1|}_{f(w_1)}$$

- this is $f(w_1) = (aw_1 - b)^2 + \lambda |w_1| + \text{constants}$, with

- $a = \sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}$, and

- $b = \frac{\mathbf{X}[:, 1]^T (\mathbf{y} - \mathbf{X}[:, 2 : d]w_{-1})}{\sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}}$

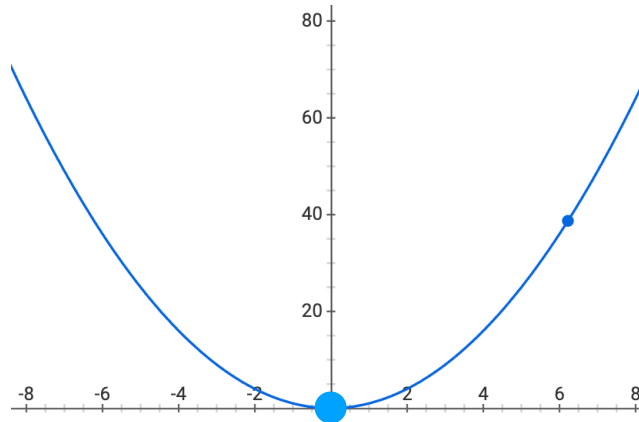
- $f(w_1)$ is non-differentiable, and its sub-gradient is

$$\partial f(w_1) = (2a(aw_1 - b) + \lambda \partial |w_1|$$

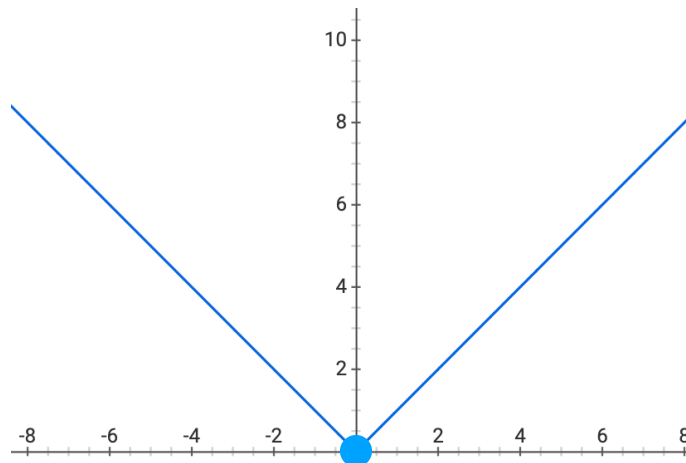
$$= \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

Convexity

- for convex differentiable functions, the minimum is achieved at points where gradient is zero



- for convex non-differentiable functions, the minimum is achieved at points where sub-gradient includes zero



Computing the sub-gradient

- the minimizer $w_1^{(t)}$ is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

Computing the sub-gradient

- the minimizer $w_1^{(t)}$ is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

Computing the sub-gradient

- the minimizer $w_1^{(t)}$ is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

Computing the sub-gradient

- considering all three cases, we get the following update rule by setting the sub-gradient to zero

$$w_1^{(t)} \leftarrow \begin{cases} \frac{b}{a} - \frac{\lambda}{2a^2} & \text{for } 2ab > \lambda \\ 0 & \text{for } -\lambda \leq 2ab \leq \lambda \\ \frac{b}{a} + \frac{\lambda}{2a^2} & \text{for } \lambda < -2ab \end{cases}$$

How do we find the minimizer?

- the minimizer $w_1^{(t)}$ is when zero is included in the sub-gradient

$$\partial f(w_1) = \begin{cases} 2a(aw_1 - b) + \lambda & \text{for } w_1 > 0 \\ [-2ab - \lambda, -2ab + \lambda] & \text{for } w_1 = 0 \\ 2a(aw_1 - b) - \lambda & \text{for } w_1 < 0 \end{cases}$$

- case 1:

- $2a(aw_1 - b) + \lambda = 0$ for some $w_1 > 0$

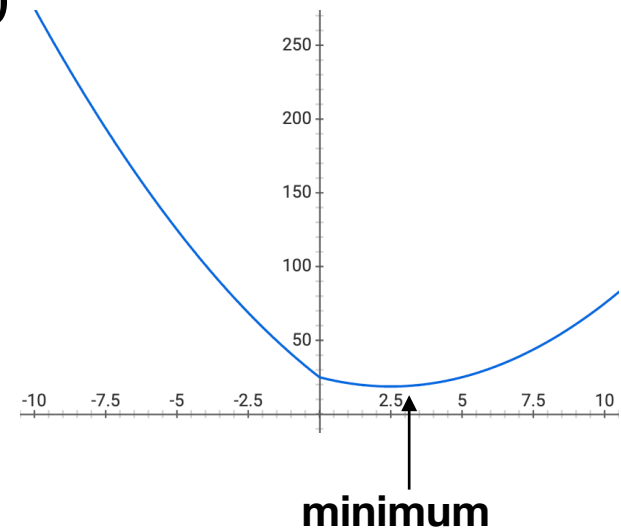
- this happens when

$$w_1 = \frac{-\lambda + 2ab}{2a^2} > 0$$

- hence,

$$w_1^{(t)} \leftarrow \frac{b}{a} - \frac{\lambda}{2a^2},$$

if $\lambda < 2ab$



- case 2:

- $2a(aw_1 - b) - \lambda = 0$ for some $w_1 < 0$

- this happens when

$$w_1 = \frac{\lambda + 2ab}{2a^2} < 0$$

- hence,

$$w_1^{(t)} \leftarrow \frac{b}{a} + \frac{\lambda}{2a^2},$$

if $\lambda < -2ab$

- case 3:

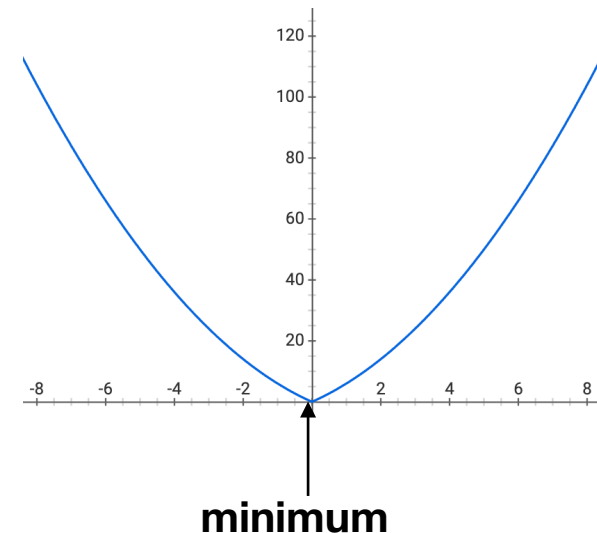
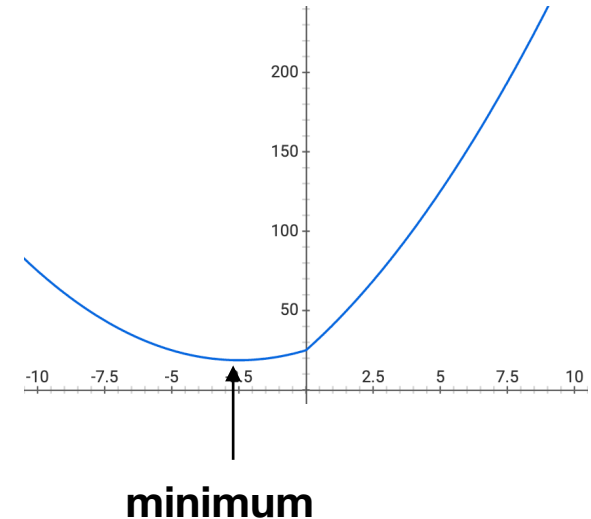
- $0 \in [-2ab - \lambda, -2ab + \lambda]$

- and $w_1 = 0$

- hence,

$$w_1^{(t)} \leftarrow 0,$$

if $-\lambda \leq 2ab \leq \lambda$

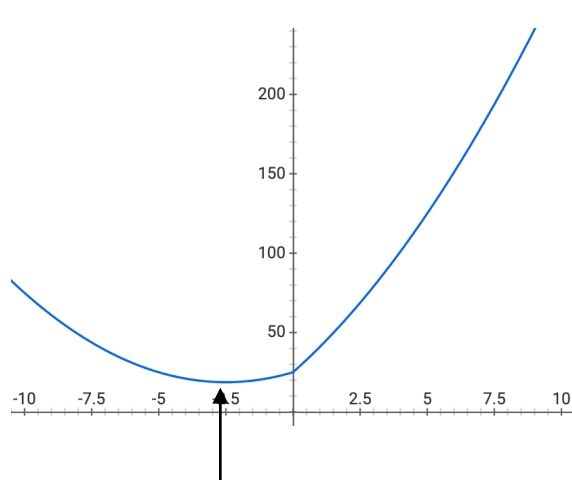


Coordinate descent on Lasso

- considering all three cases, we get the following update rule by setting the sub-gradient to zero

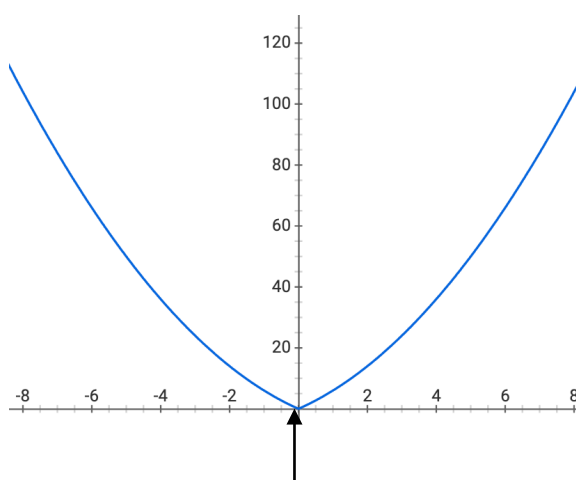
$$w_1^{(t)} \leftarrow \begin{cases} \frac{b}{a} - \frac{\lambda}{2a^2} & \text{for } 2ab > \lambda \\ 0 & \text{for } -\lambda \leq 2ab \leq \lambda \\ \frac{b}{a} + \frac{\lambda}{2a^2} & \text{for } \lambda < -2ab \end{cases}$$

• where $a = \sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}$, and $b = \frac{\mathbf{X}[:, 1]^T (\mathbf{y} - \mathbf{X}[:, 2:d] w_{-1})}{\sqrt{\mathbf{X}[:, 1]^T \mathbf{X}[:, 1]}}$

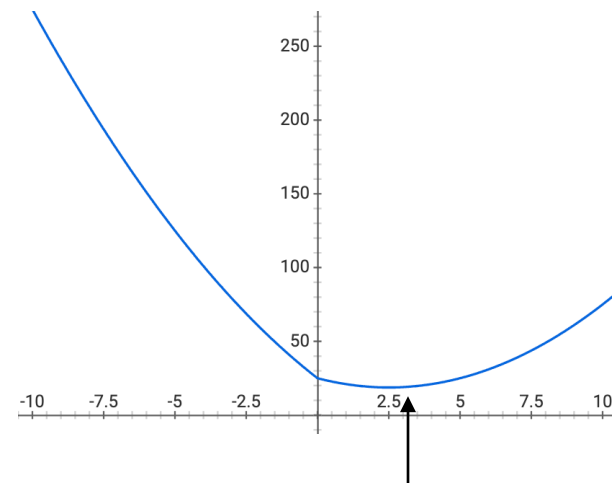


20

minimum



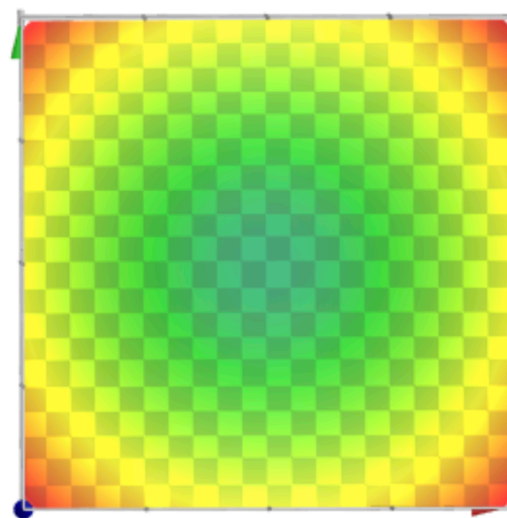
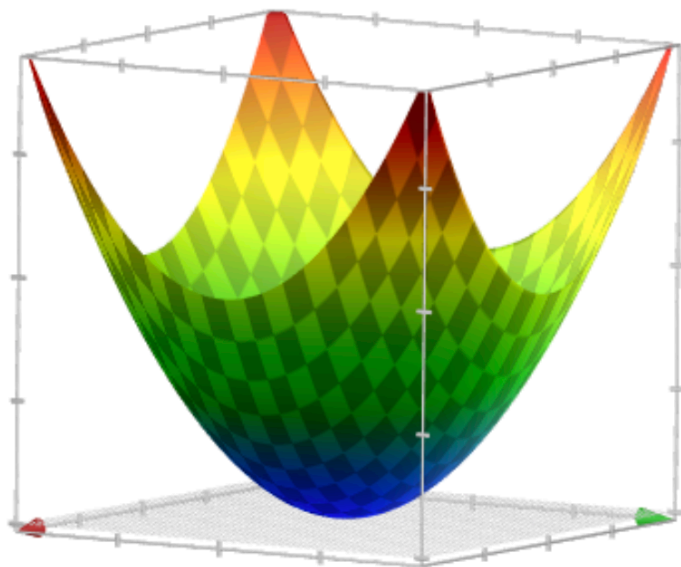
minimum



minimum

When does coordinate descent work?

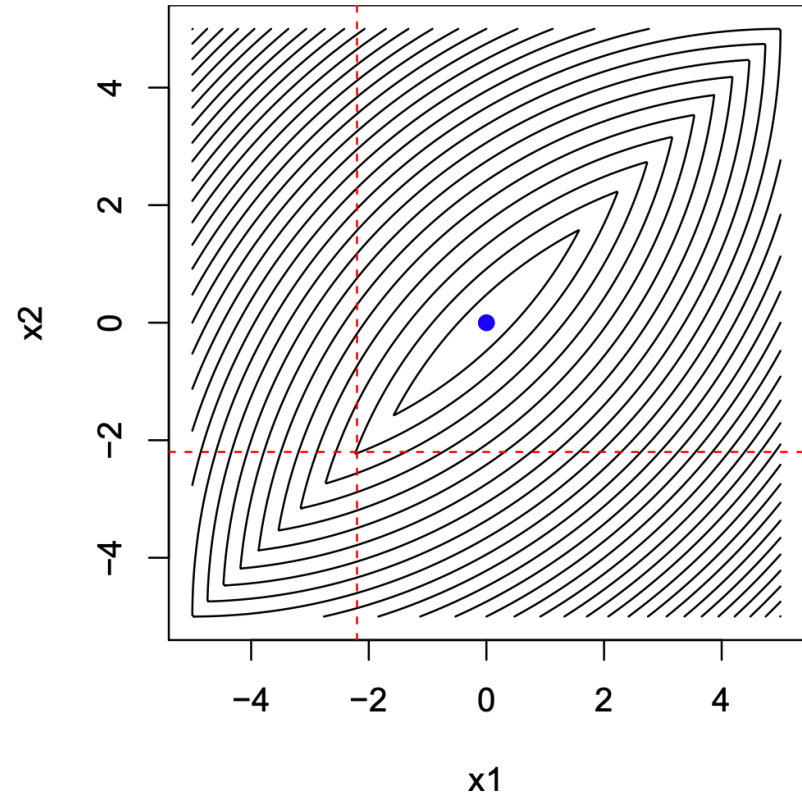
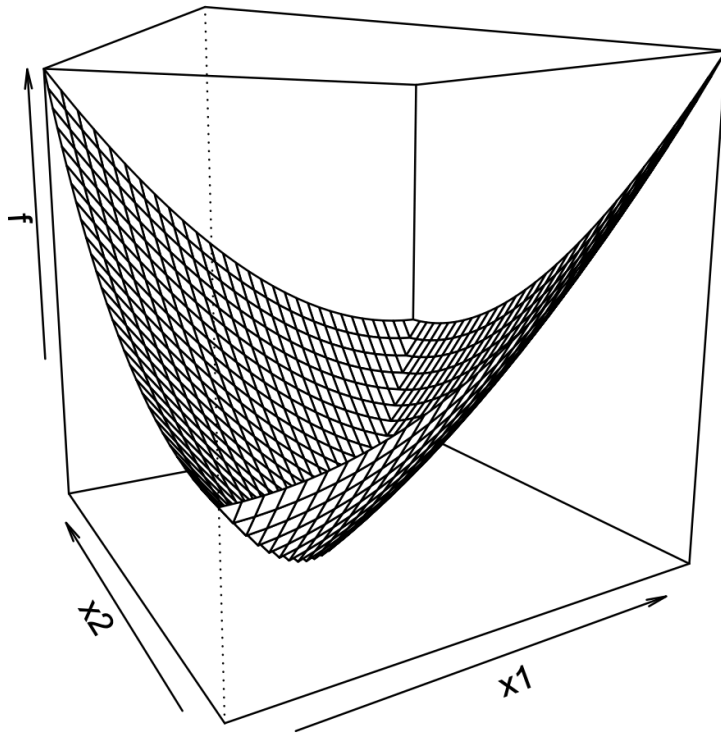
- Consider minimizing a **differentiable convex** function $f(x)$, then coordinate descent converges to the global minima



- when coordinate descent has stopped, that means $\frac{\partial f(x)}{\partial x_j} = 0$ for all $j \in \{1, \dots, d\}$
- this implies that the gradient $\nabla_x f(x) = 0$, which happens only at minimum

When does coordinate descent work?

- Consider minimizing a **non-differentiable convex** function $f(x)$, then coordinate descent can get stuck



When does coordinate descent work?

- then how can coordinate descent find optimal solution for Lasso?
- consider minimizing a **non-differentiable convex** function but has a structure of $f(x) = g(x) + \sum_{j=1}^d h_j(x_j)$, with differentiable convex function $g(x)$ and coordinate-wise non-differentiable convex functions $h_j(x_j)$'s, then coordinate descent converges to the global minima

