

Gradient Descent

W

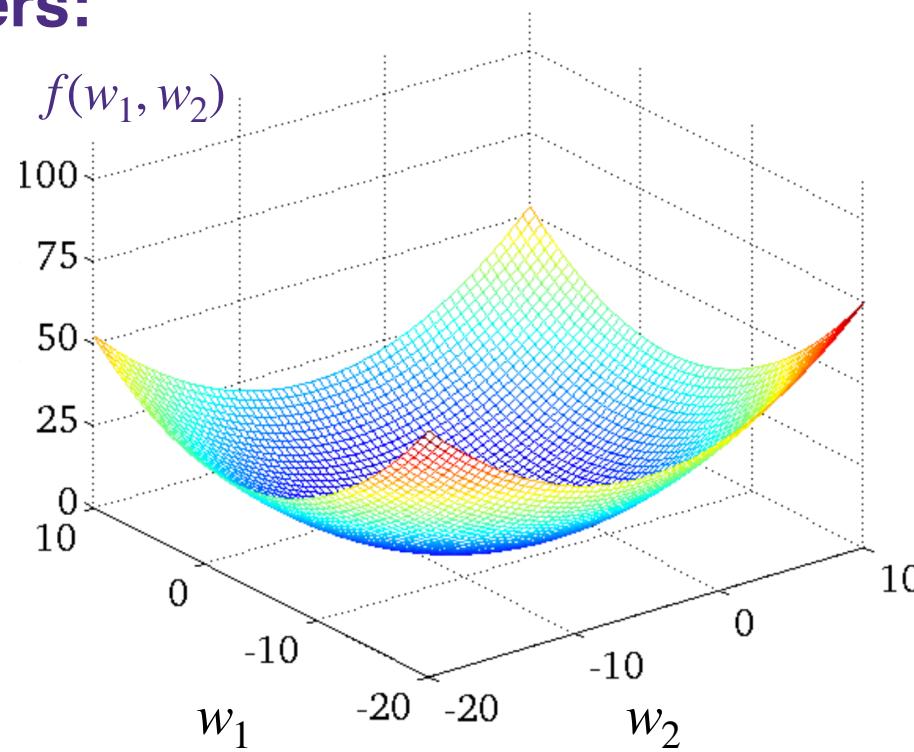
Running example: linear regression

- Given data:

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- Learning a model's parameters:

$$\hat{w}_{\text{LS}} = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|\mathbf{y} - \mathbf{X}w\|_2^2}_{f(w)}$$



Gradient descent

Example of a general non-convex $f(w)$

Initialize: $w_0 = 0$

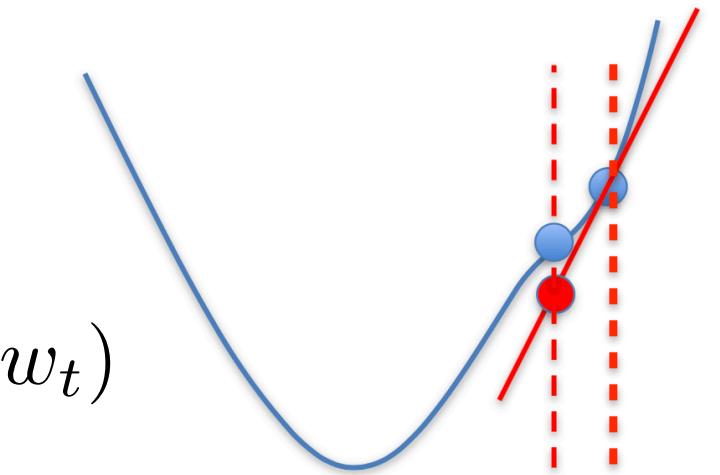
for $t = 1, 2, \dots$

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$



$$\eta \times \nabla f(w_t)$$

Learning rate or step-size,
a hyper-parameter to be chosen by the analyst



Gradient descent for linear regression

Initialize: $w_0 = 0$

for $t = 1, 2, \dots$

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

For linear regression, we have

$$\hat{w}_{\text{LS}} = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|y - Xw\|_2^2}_{f(w)}$$

Gradient descent for linear regression

Initialize: $w_0 = 0$

for $t = 1, 2, \dots$

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

For linear regression, we have

$$\hat{w}_{\text{LS}} = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|y - Xw\|_2^2}_{f(w)}$$

$$\nabla f(w_t) = -2X^T(y - Xw_t)$$

$$w_{t+1} = w_t + \eta 2X^T(y - Xw_t) = (I - 2\eta X^T X)w_t + 2\eta X^T y$$

Let the least-squares solution be $w^* = (X^T X)^{-1} X^T y$

$$\begin{aligned} w_{t+1} - w^* &= (I - 2\eta X^T X)w_t + 2\eta X^T y - w^* \\ &= (I - 2\eta X^T X)(w_t - w^*) + 2\eta X^T y - 2\eta X^T X w^* \\ &= (I - 2\eta X^T X)(w_t - w^*) \end{aligned}$$

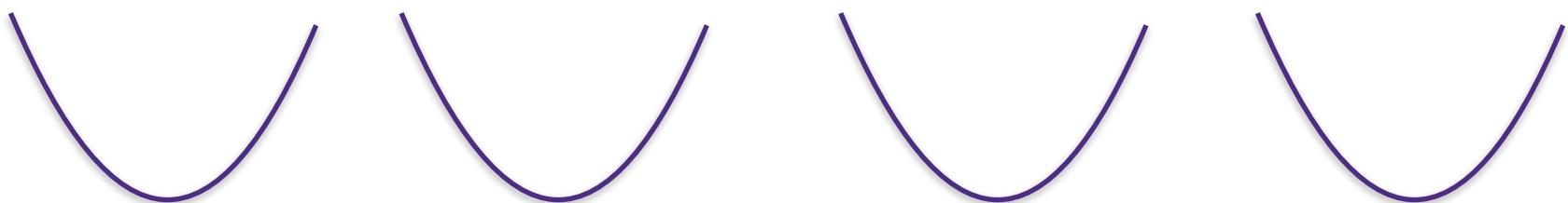
Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})(w_t - w^*)$$

Gradient descent for linear regression

$$\begin{aligned} w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* &= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})(w_t - w^*) \\ &= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^2(w_{t-1} - w^*) \\ &\quad \vdots \\ &= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1}(w_0 - w^*) \end{aligned}$$

In one dimension, $2\mathbf{X}^T \mathbf{X} = a$ is a scalar



$$0 < \eta < 1/a$$

$$\eta = 1/a$$

$$1/a < \eta < 2/a$$

$$\eta > 2/a$$

Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1}(w_0 - w^*)$$

In multi dimensions, **eigenvalues** of $\mathbf{X}^T \mathbf{X}$ are important
(you will see why I say $\mathbf{X}^T \mathbf{X}$ instead of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ in couple of slides)

Let the eigenvalue decomposition of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ be $Q^{-1}DQ$

Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1} (w_0 - w^*)$$

In multi dimensions, **eigenvalues** of $\mathbf{X}^T \mathbf{X}$ are important

(you will see why I say $\mathbf{X}^T \mathbf{X}$ instead of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ in couple of slides)

Let the eigenvalue decomposition of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ be $Q^{-1} D Q$

$$\begin{aligned} \text{Then, } w_{t+1} - w^* &= (Q^{-1} D Q)^{t+1} (w_0 - w^*) \\ &= \underbrace{Q^{-1} D Q Q^{-1} D Q \dots Q^{-1} D Q}_{t+1 \text{ times}} (w_0 - w^*) \end{aligned}$$

$$= Q^{-1} D^{t+1} Q (w_0 - w^*)$$

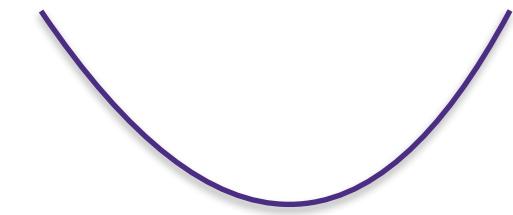
$$Q(w_{t+1} - w^*) = D^{t+1} Q (w_0 - w^*)$$

$$Q(w_{t+1} - w^*) \; = \; D^{t+1} Q\left(w_0 - w^*\right)$$

Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1}(w_0 - w^*)$$

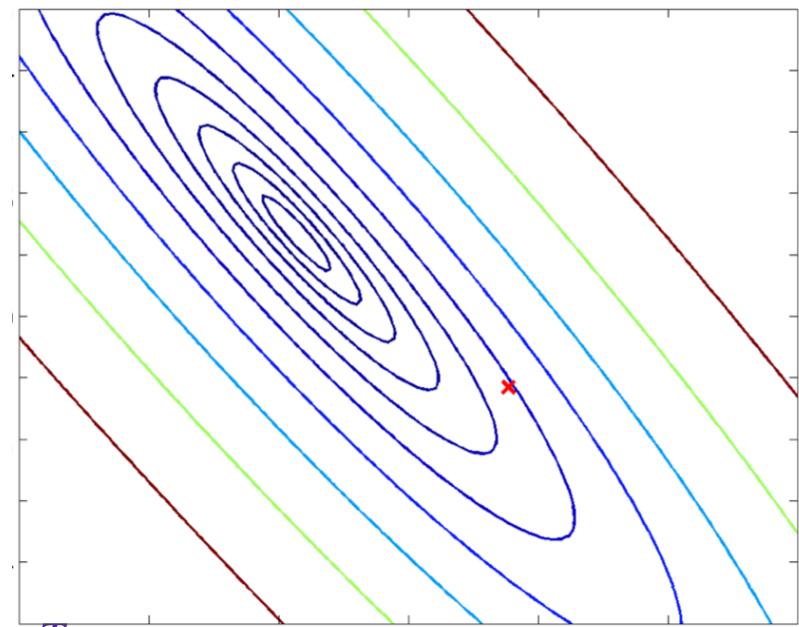
$$Q(w_{t+1} - w^*) = D^{t+1} Q(w_0 - w^*)$$



In direction q_1
 $0 < \lambda_1(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})$



In direction q_2
 $-1 < \lambda_2(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}) < 0$



Gradient descent for linear regression

Recall that in each eigen direction

$$q_i^T(w_{t+1} - w^*) = (\lambda_i(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}))^{t+1} q_i^T(w_0 - w^*)$$

We want the error to decay fast in all directions, whose bottleneck is the largest and the smallest eigen values:

$$\lambda_{\min}(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}) \text{ and } \lambda_{\max}(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})$$



We want to choose the learning rate η such that

$$-1 \ll \lambda_d(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}) \leq \dots \leq \lambda_1(\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}) \ll 1$$

Gradient descent for linear regression

Recall that in each eigen direction

$$q_i^T(w_{t+1} - w^*) = \lambda_i(\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X})^{t+1} q_i^T(w_0 - w^*)$$

I will not prove the facts that $\|q_i\|_2 = 1$ and $q_i^T q_j = 0$

Claim: $\lambda_i(\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X}) = 1 - 2\eta\lambda_i(\mathbf{X}^T\mathbf{X})$

Claim: $\lambda_i(\mathbf{X}^T\mathbf{X}) \geq 0$



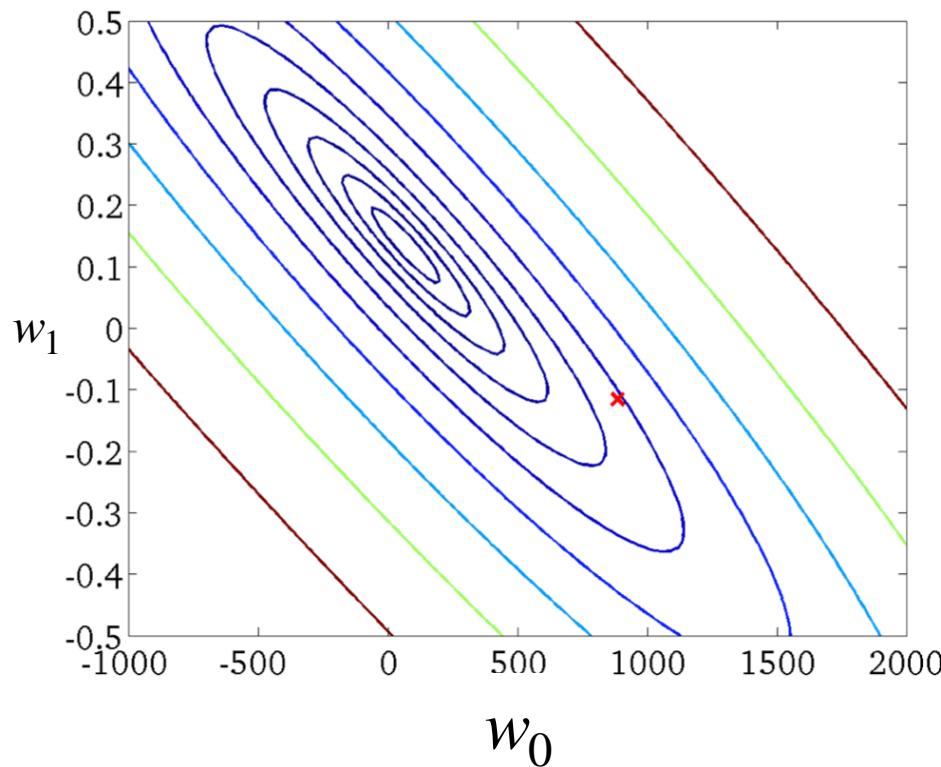
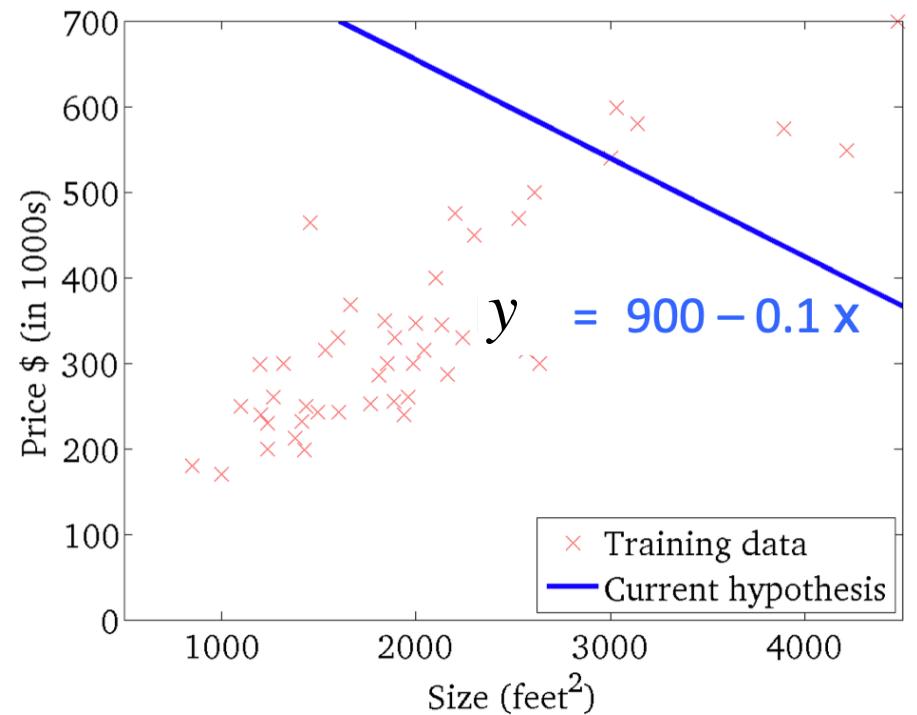
$$\eta = 0$$

$$1 - 2\eta\lambda_{\max}(\mathbf{X}^T\mathbf{X}) = -1$$

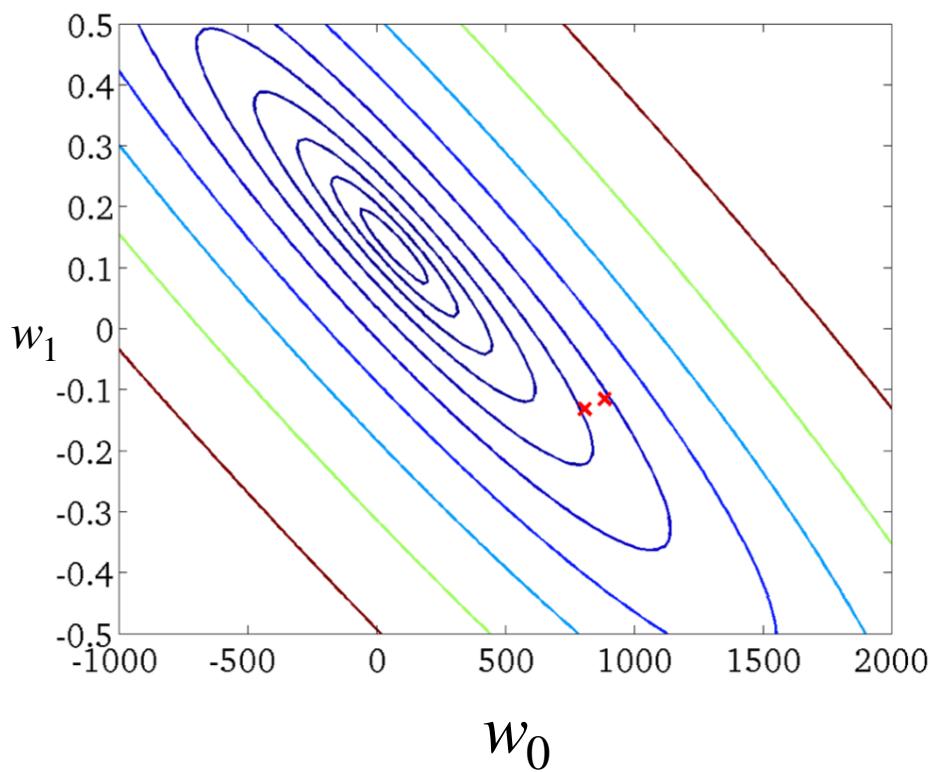
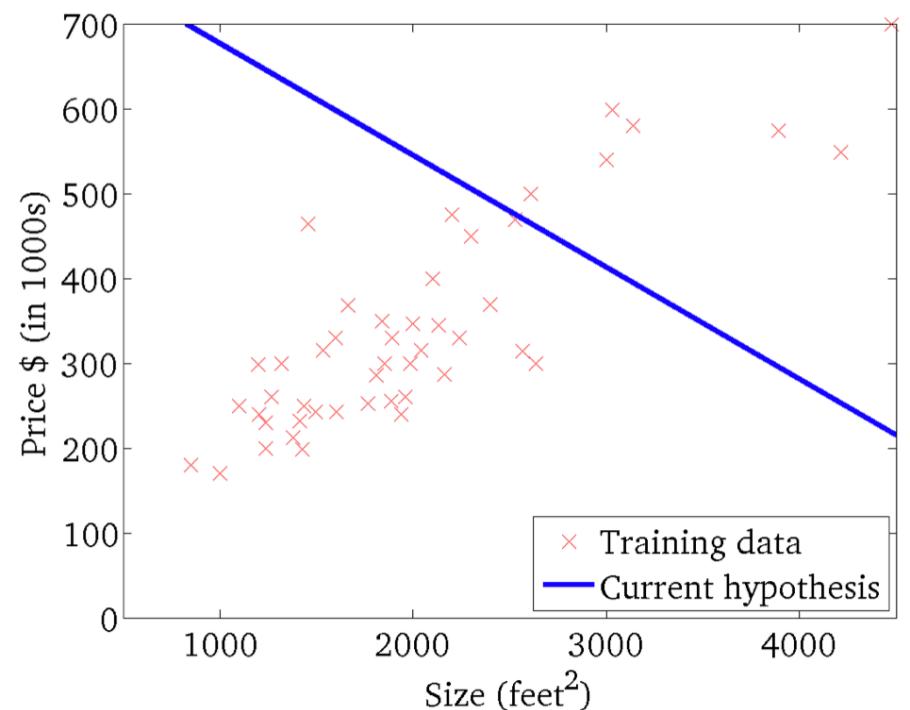


$$\eta = \frac{1}{\lambda_{\max}(\mathbf{X}^T\mathbf{X})}$$

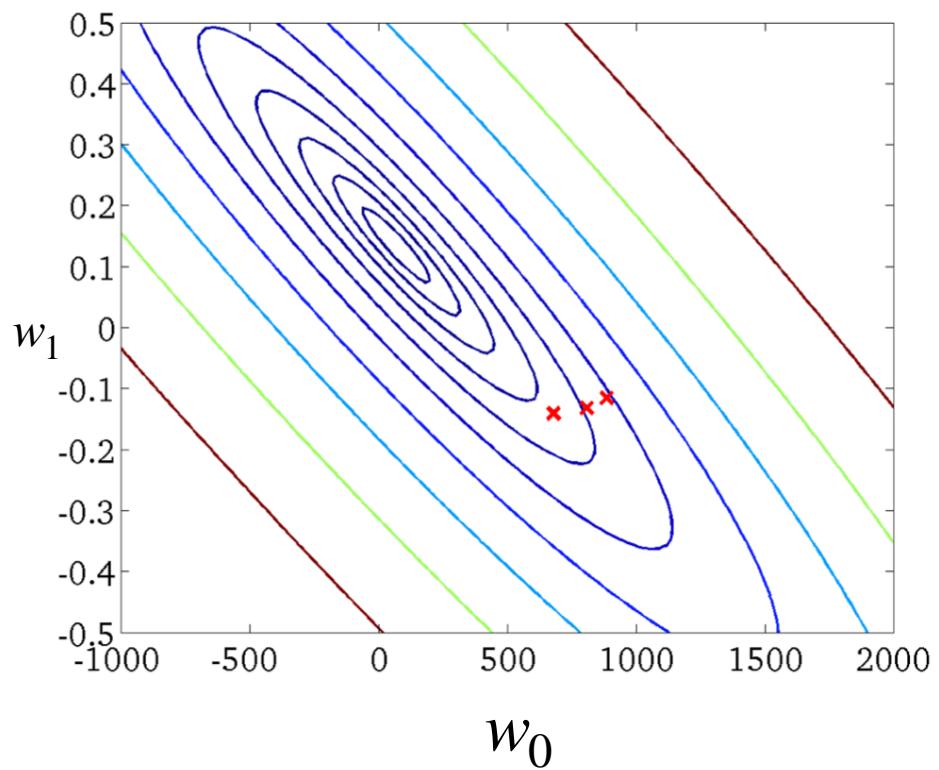
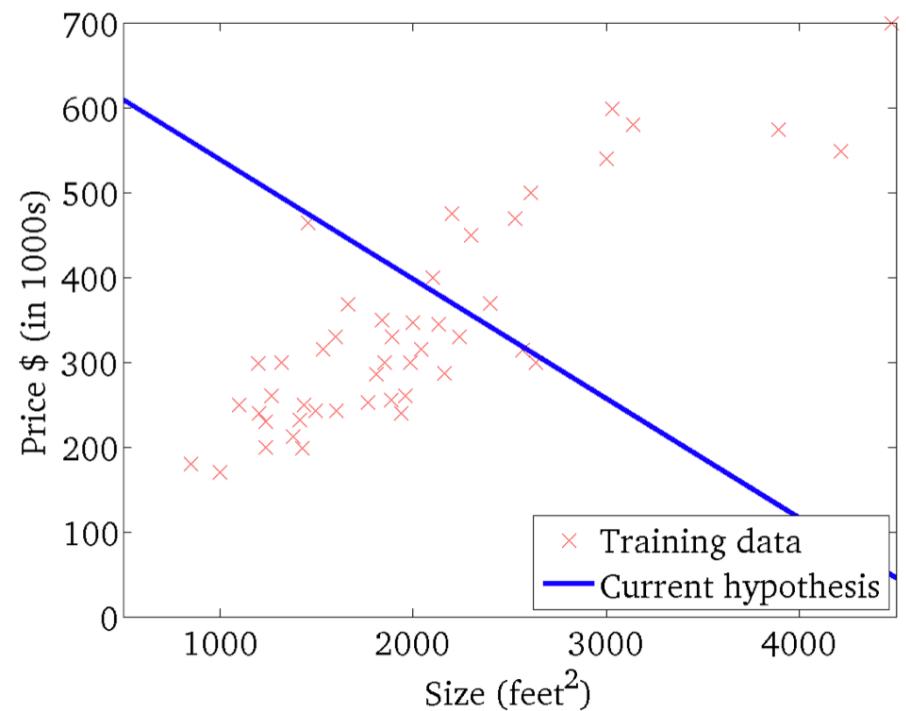

$$y = w_0 + w_1 x$$



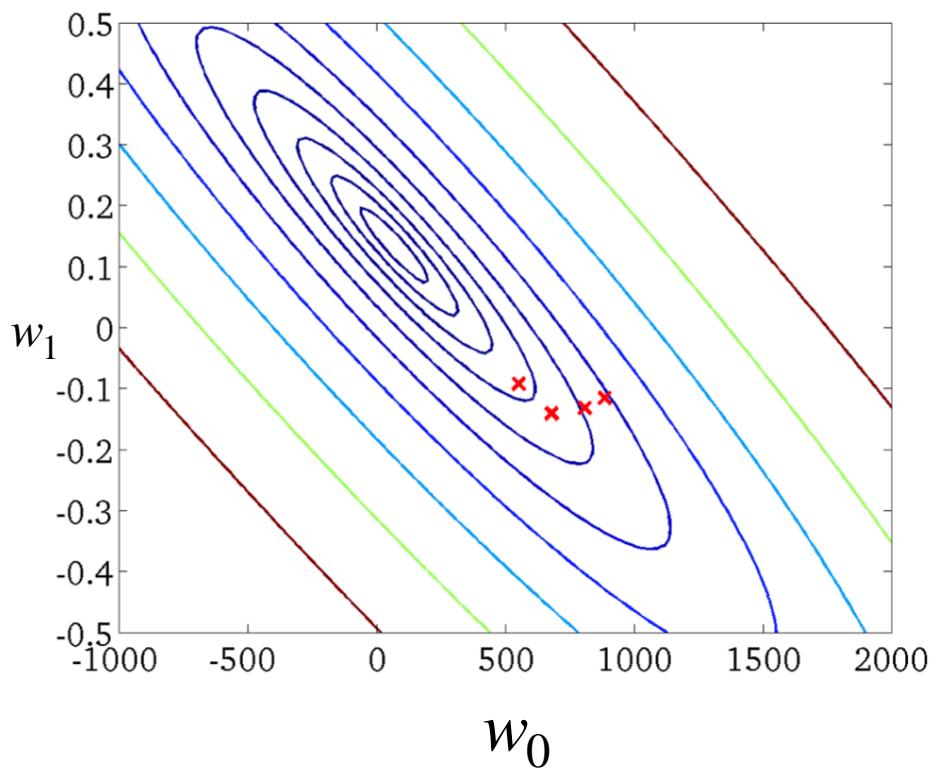
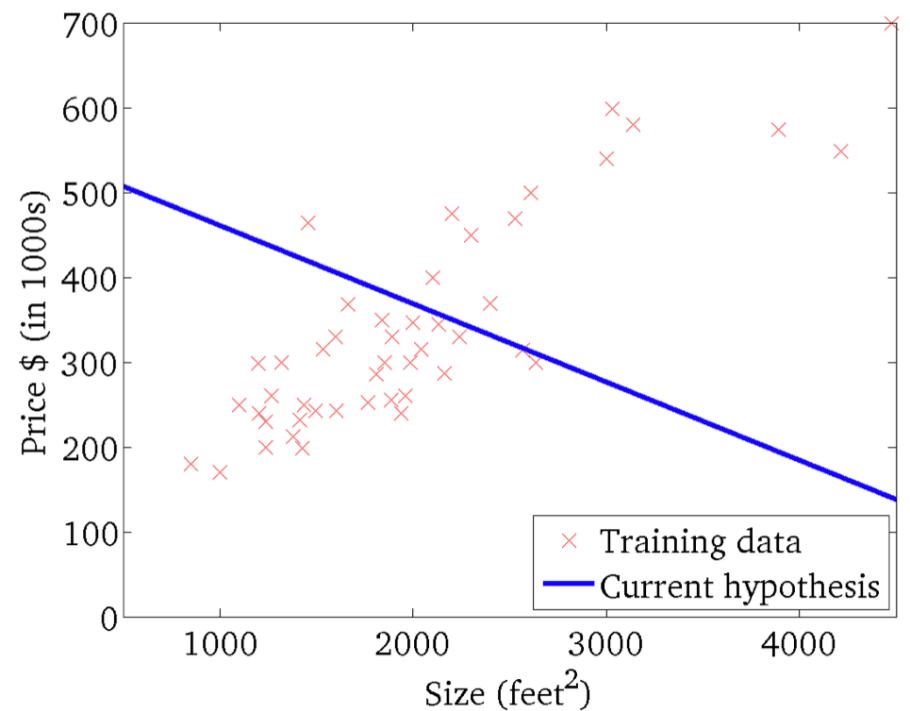

$$y = w_0 + w_1 x$$



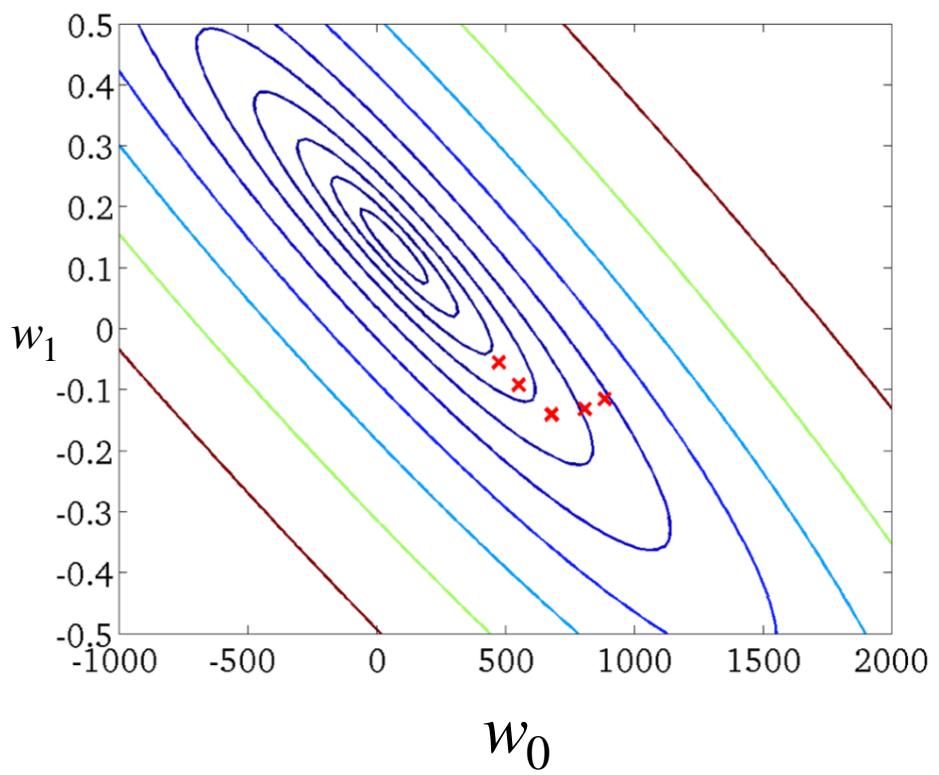
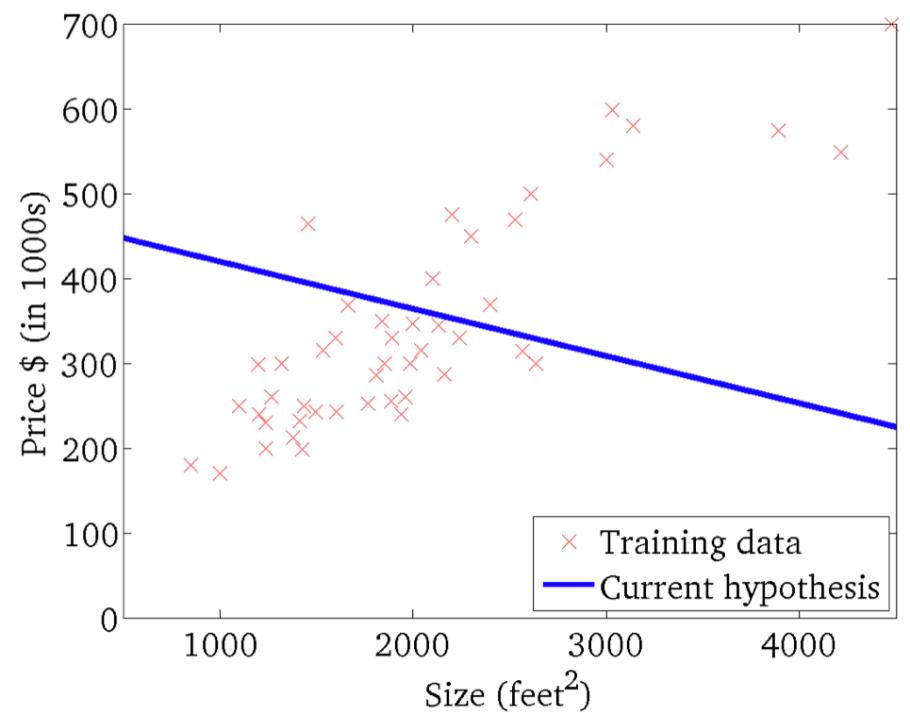

$$y = w_0 + w_1 x$$



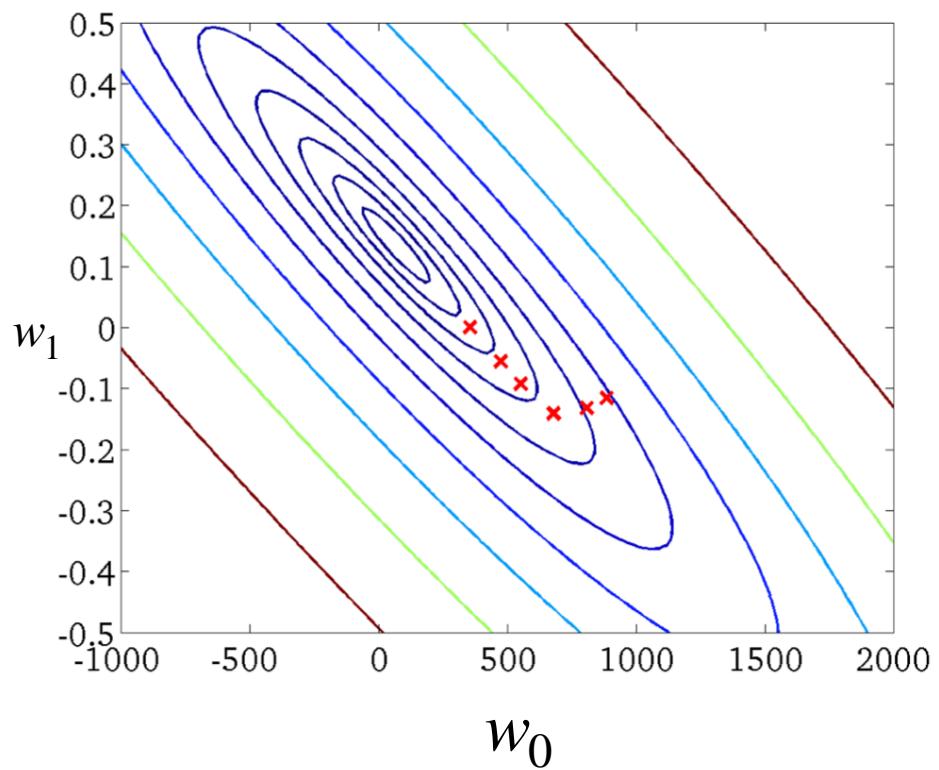
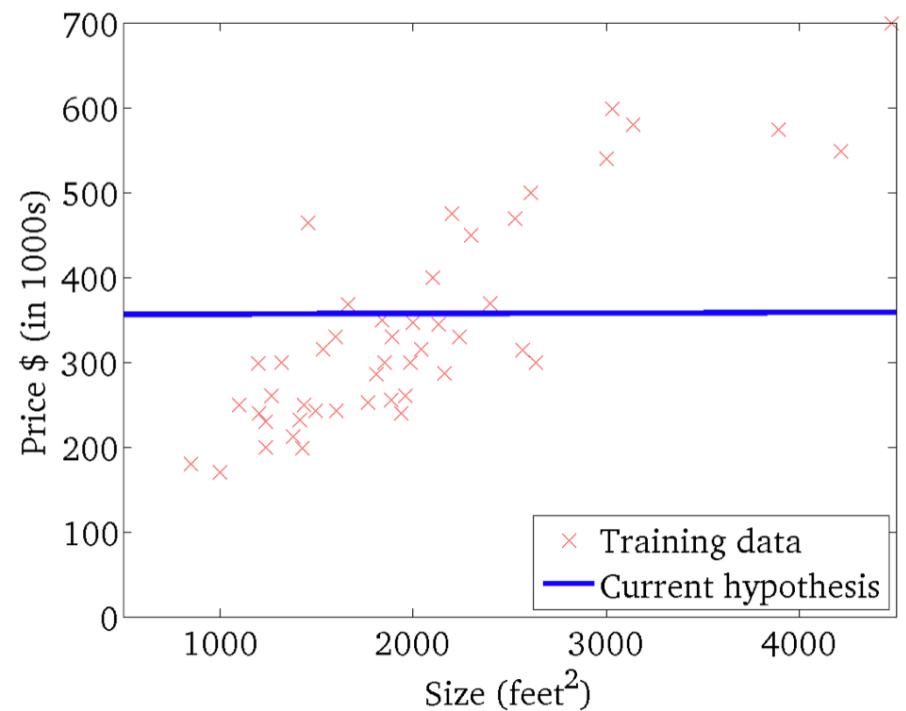

$$y = w_0 + w_1 x$$



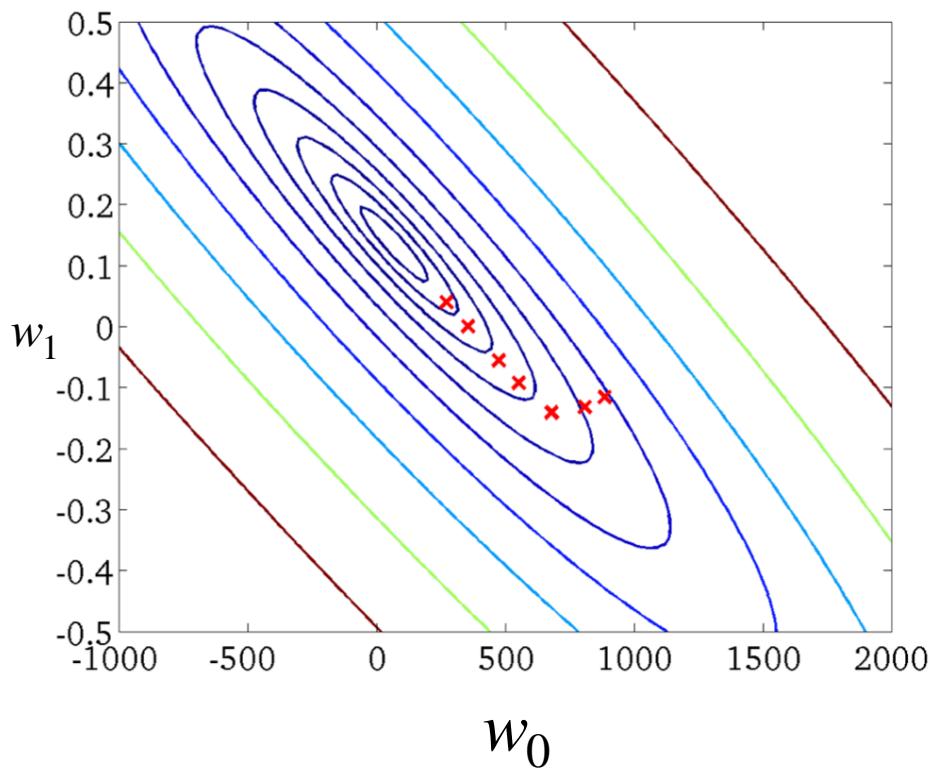
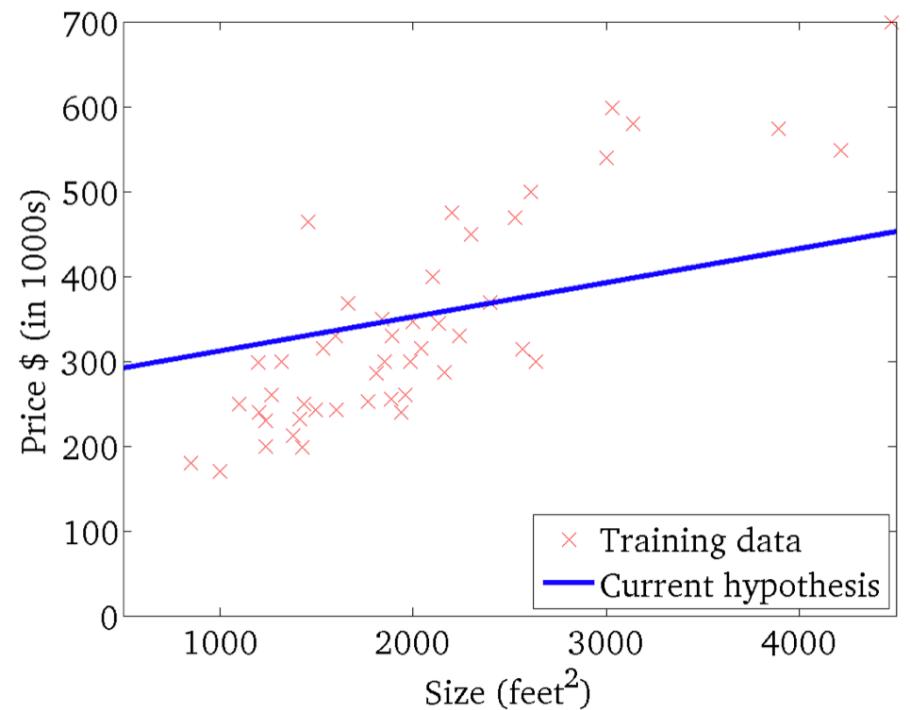

$$y = w_0 + w_1 x$$



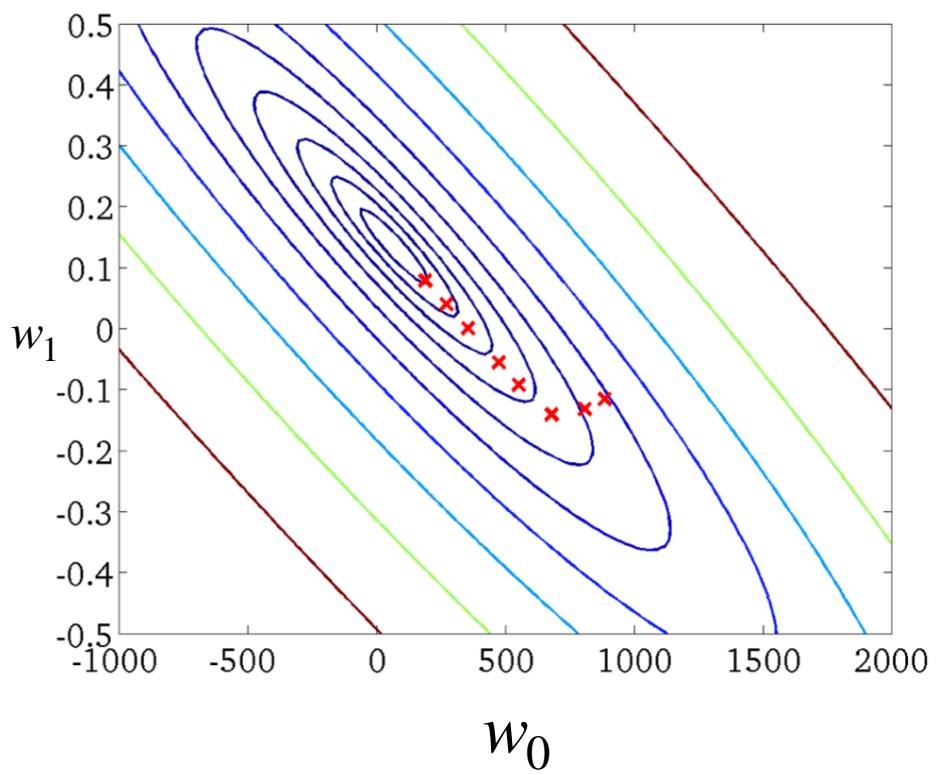
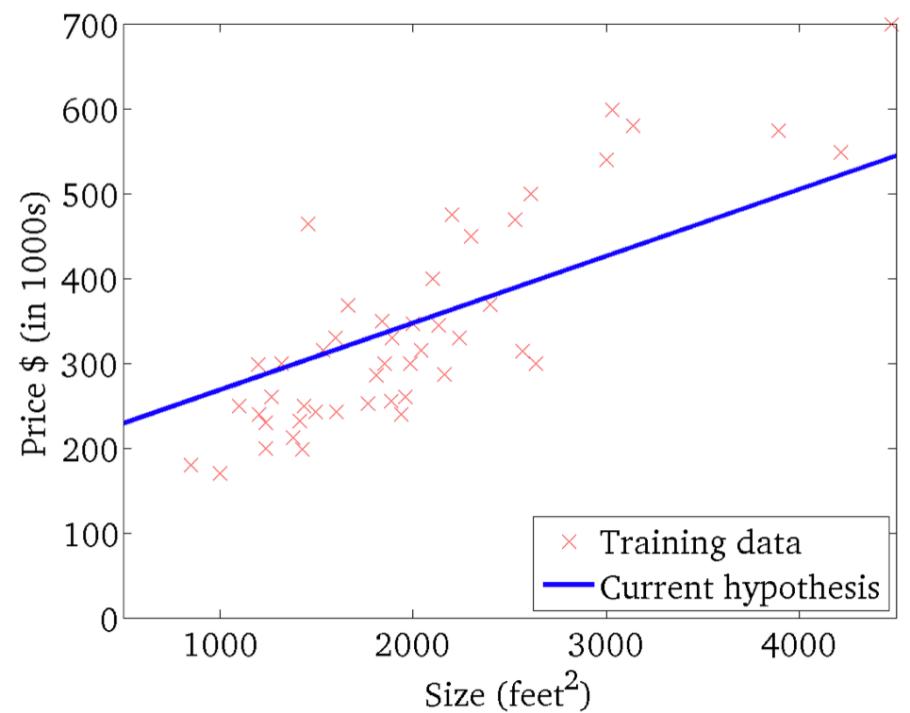

$$y = w_0 + w_1 x$$



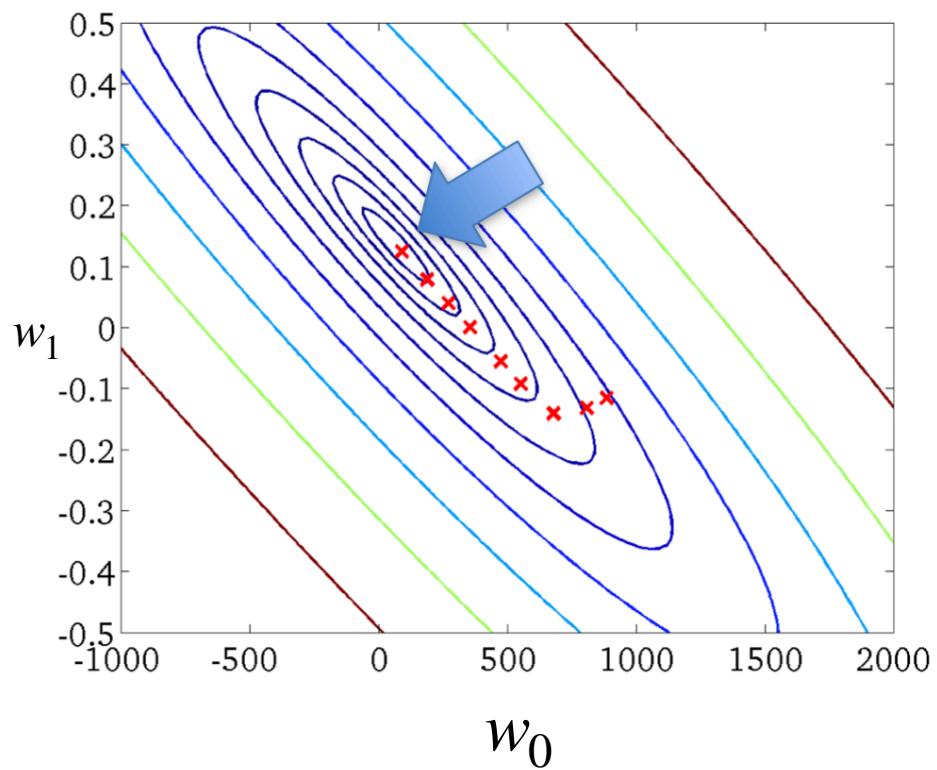
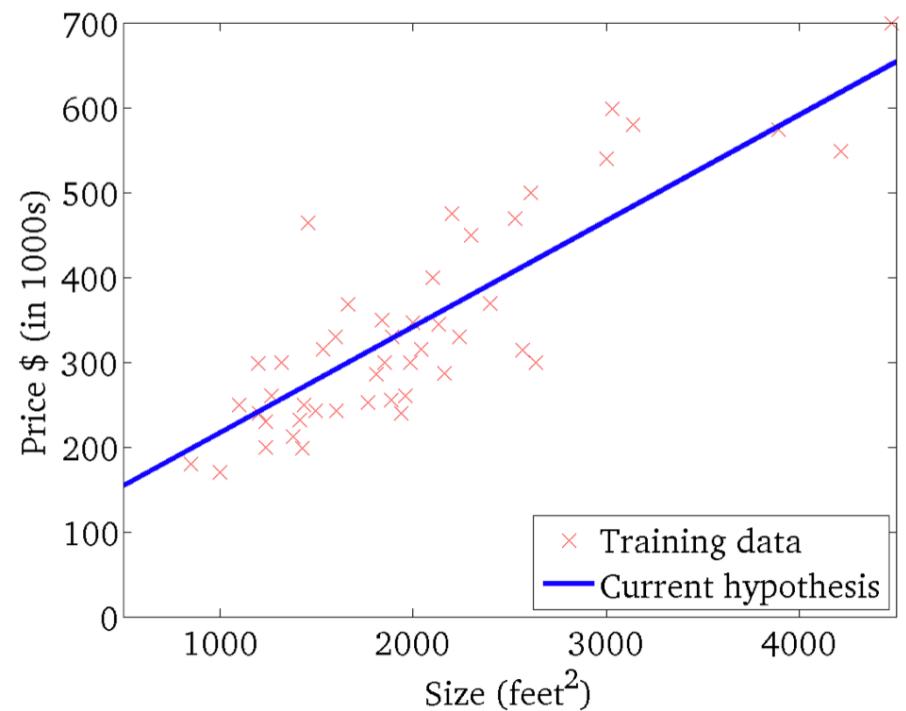

$$y = w_0 + w_1 x$$




$$y = w_0 + w_1 x$$




$$y = w_0 + w_1 x$$



Gradient descent for logistic regression

Loss function: Conditional Likelihood

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$$

$$\widehat{w}_{MLE} = \arg \max_w \prod_{i=1}^n P(y_i|x_i, w) \quad P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

$$= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))$$

$$\nabla f(w) = \sum_{i=1}^n \frac{1}{1 + \exp(-y_i x_i^T w)} \exp(-y_i x_i^T w) (-y_i x_i)$$

Gradient descent

- For smooth functions, with a (small enough) fixed step size η

- if $f(w)$ is convex,

$$f(w_t) - f(w^*) \leq \frac{\|w_0 - w^*\|_2^2}{2\eta t}$$

- if $f(w)$ is μ -strongly convex,

$$f(w_t) - f(w^*) \leq (1 - \eta\mu)^t (f(w_0) - f(w^*))$$

- What can we do for non-smooth function $f(w)$?

- for example, LASSO

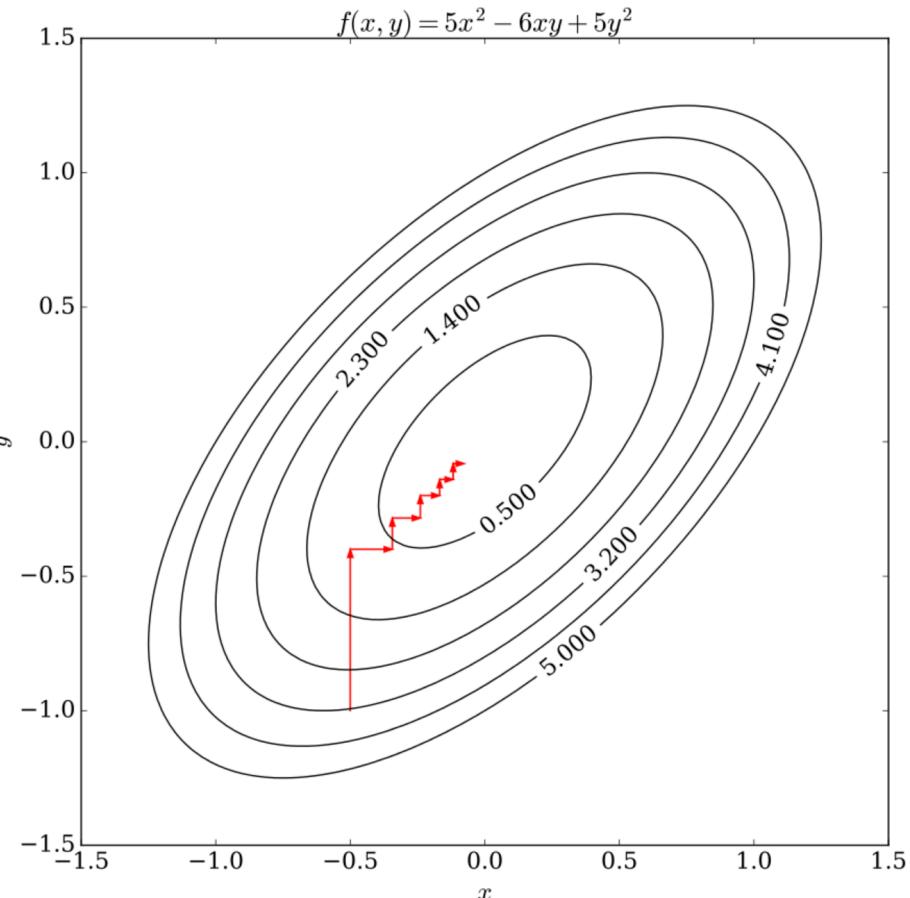
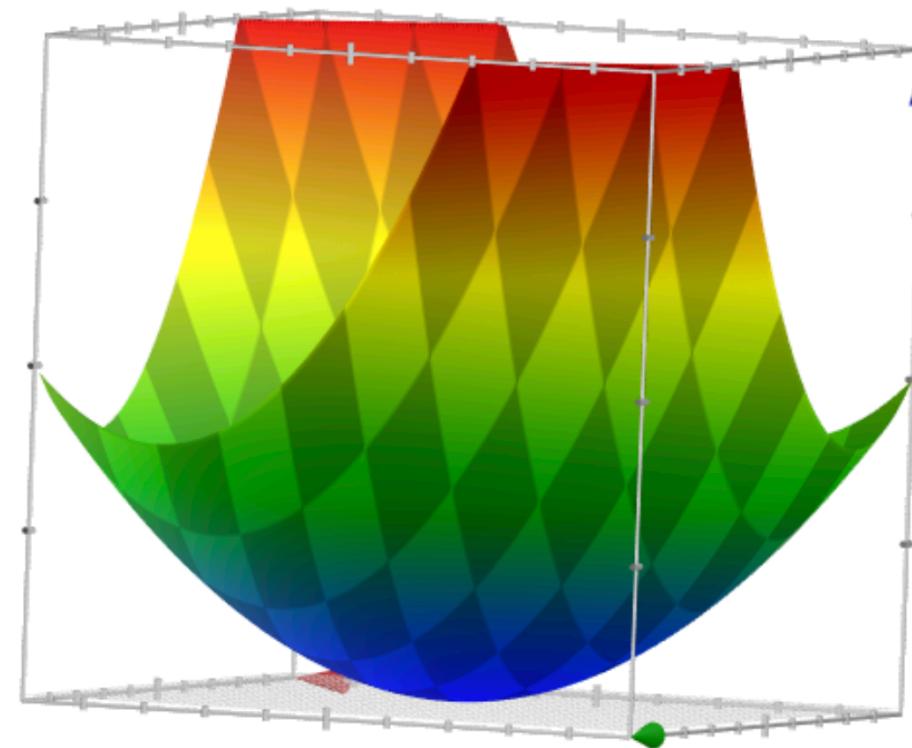
$$\hat{w}_{\text{Lasso}} = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|y - Xw\|_2^2 + \lambda \|w\|_1}_{f(w)}$$

Coordinate Descent

W

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 $f(x, y) = 5x^2 - 6xy + 5y^2$



Optimization: how do we solve Lasso?

- 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 for illustration purpose

- the main challenge we want to 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

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

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 :

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

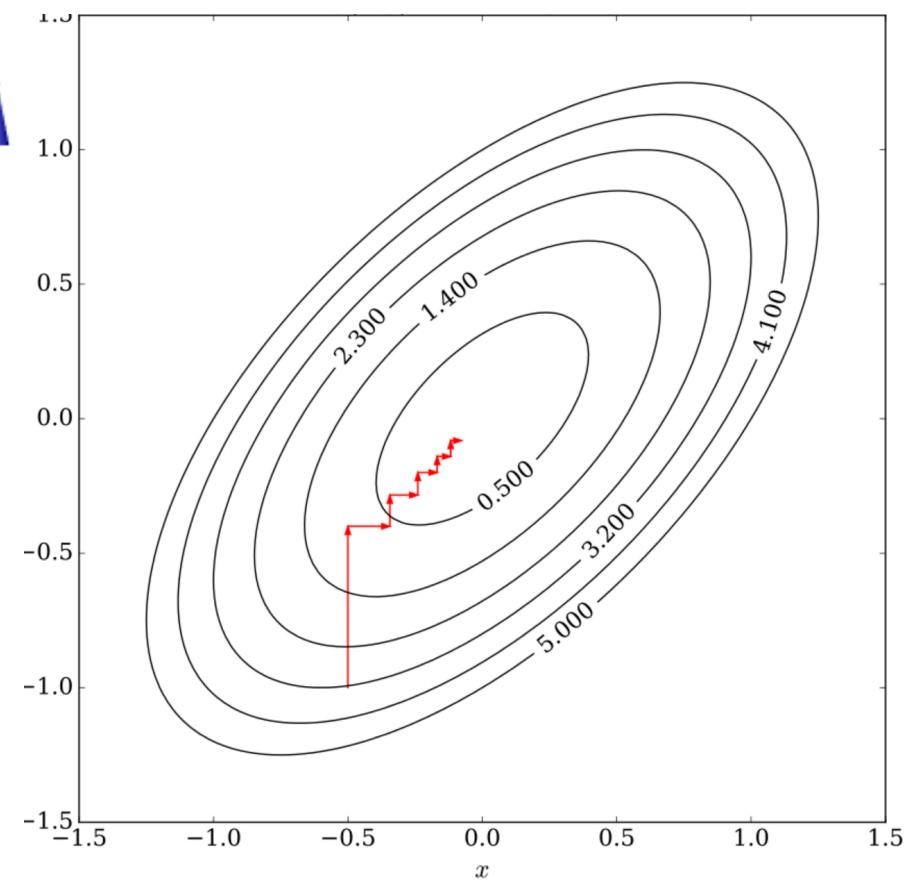
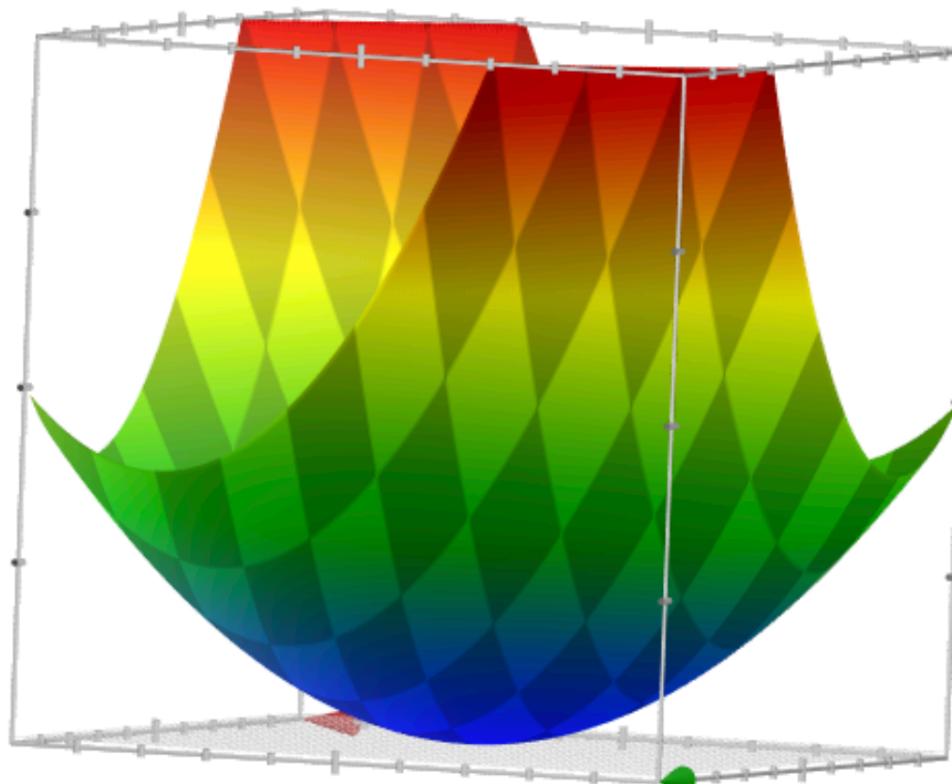
where $w_{-1} = [w_2^{(t-1)}, \dots, w_d^{(t-1)}]^T$

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

- we know from linear least squares that the minimizer is

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

- 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

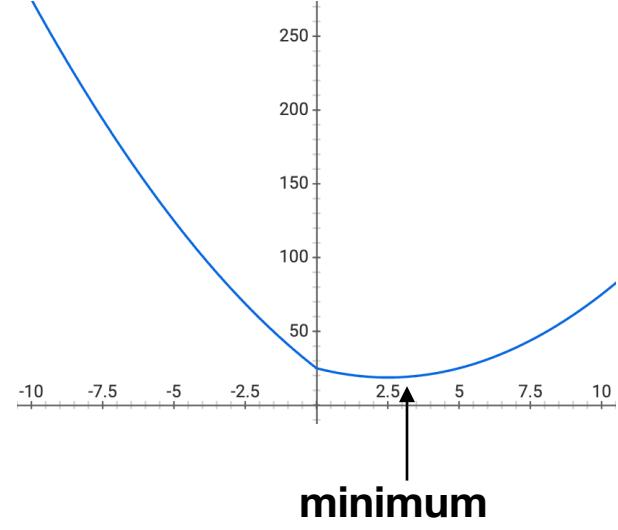
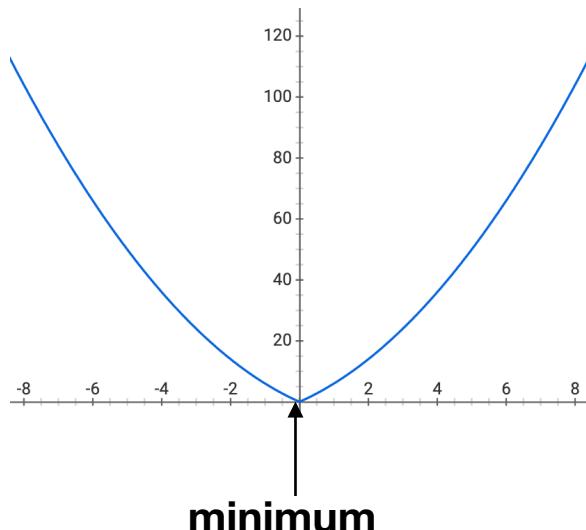
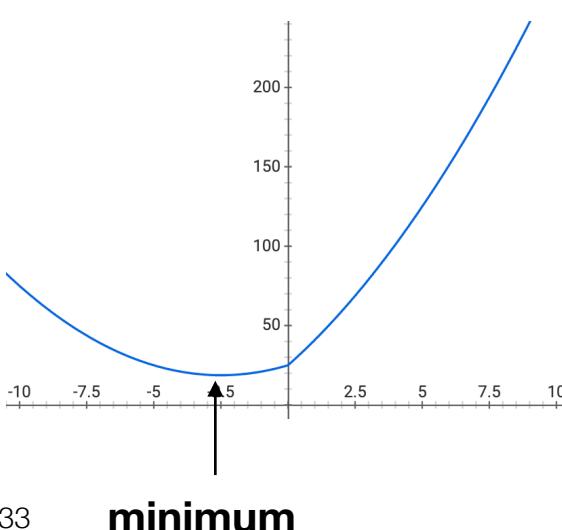
$$\left\| \mathbf{X}[:, 1]w_1 - (\mathbf{y} - \mathbf{X}[:, 2:d]w_{-1}) \right\|_2^2 + \lambda \left(|w_1| + \underbrace{\|w_{-1}\|_1}_{\text{constant}} \right)$$

- hence, the coordinate descent update boils down to

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

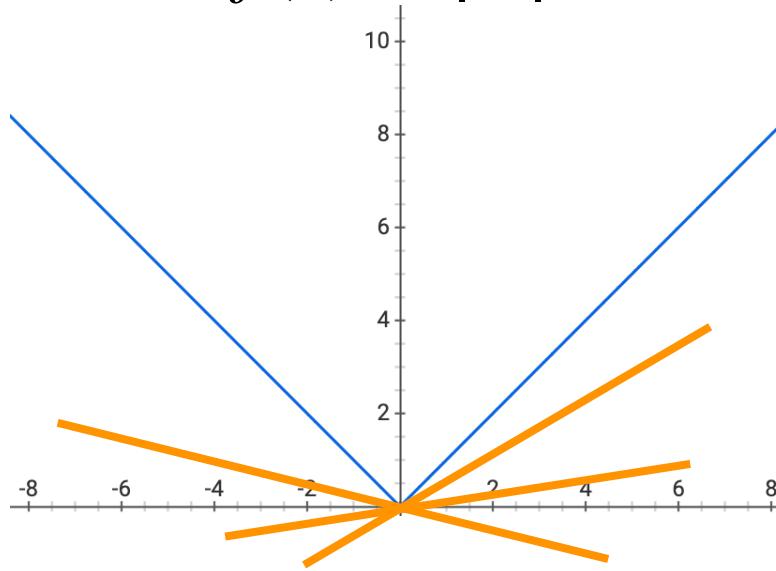
Convexity

- 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

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



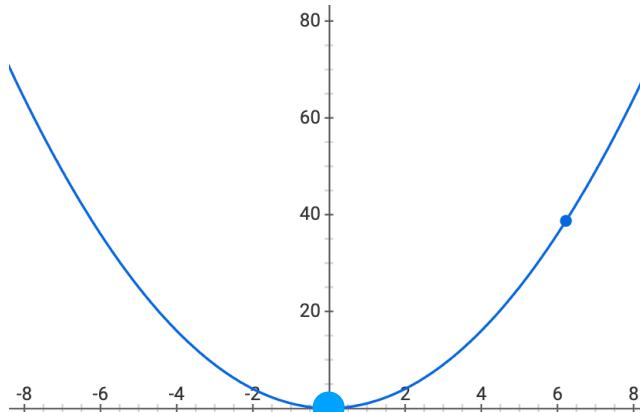
- for a **non-differentiable** function, gradient is not defined at some points, for example at $x = 0$ for $f(x) = |x|$
- at such points, **sub-gradient** plays the role of gradient
 - sub-gradient at a differentiable point is the same as the gradient
 - sub-gradient at a non-differentiable point is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

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

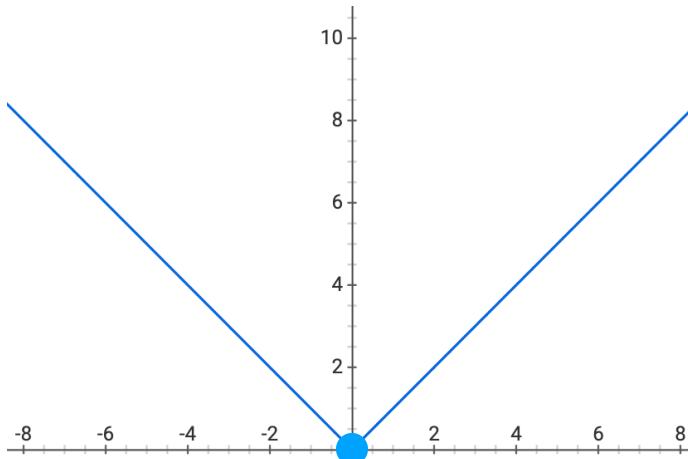
$$\bullet \text{ for example, } \partial |x| = \begin{cases} +1 & \text{for } x > 0 \\ [-1, 1] & \text{for } x = 0 \\ -1 & \text{for } x < 0 \end{cases}$$

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

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

Coordinate descent update on Lasso

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

- 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

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

Computing the minimizer

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

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

Computing the minimizer

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

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

Computing the minimizer

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

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

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 $w_1 > 0$

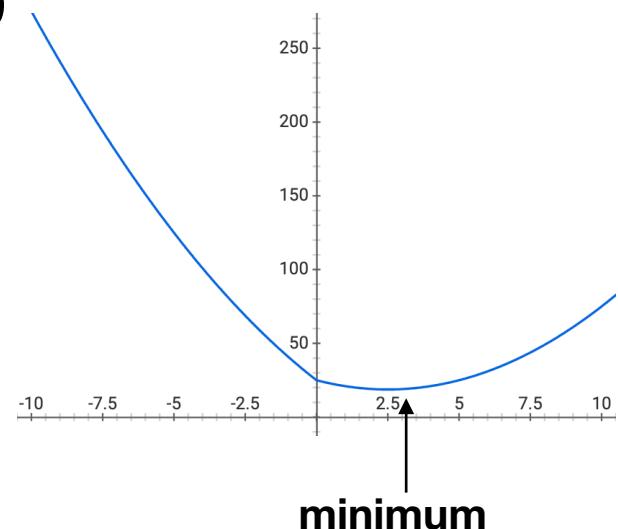
- this happens when

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

- 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:

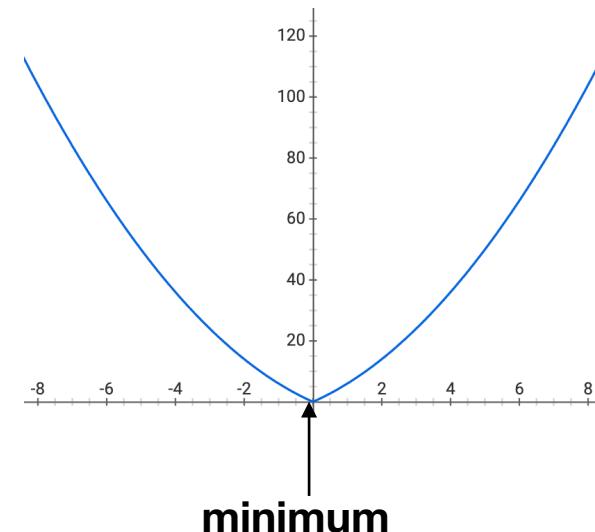
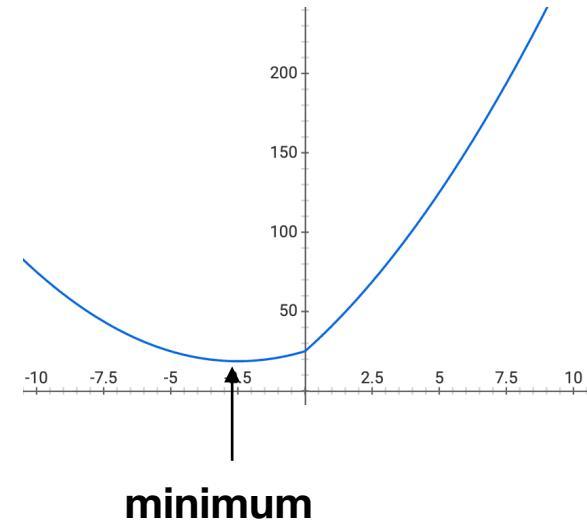
- $0 \in [-2ab - \lambda, -2ab + \lambda]$

- and $w_1 = 0$

- hence,

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

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

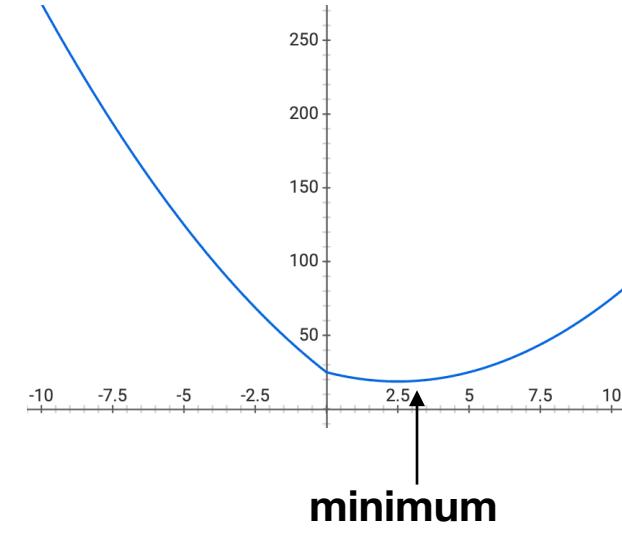
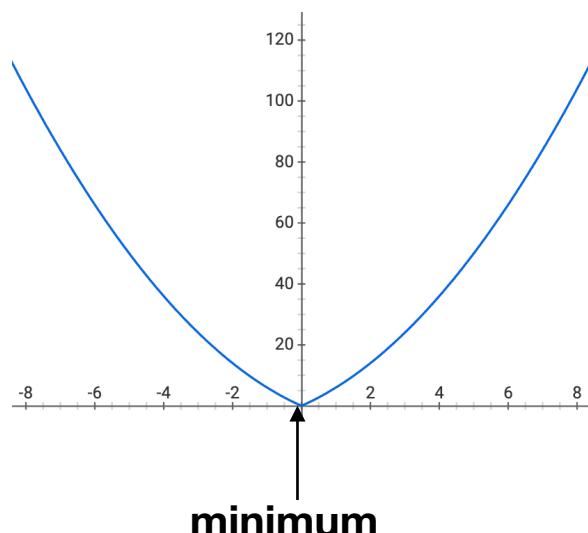
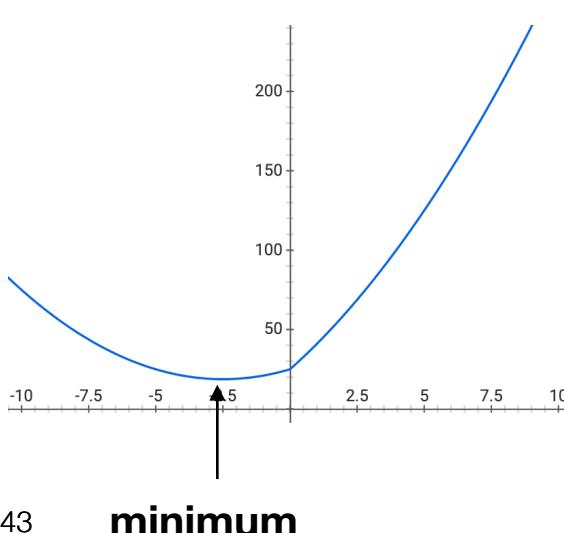


Coordinate descent on Lasso

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

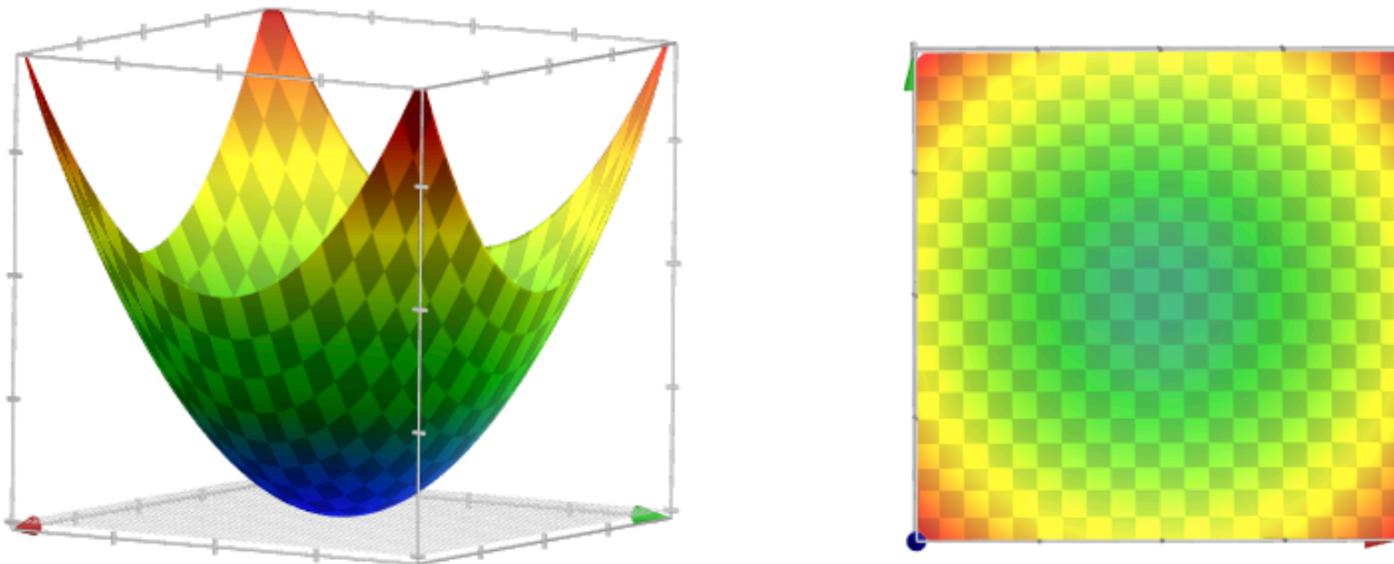
$$w_1^{(t)} \leftarrow \begin{cases} \frac{b}{a} - \frac{\lambda}{2a^2} & \text{for } 2ab > \lambda \\ 0 & \text{for } -\lambda \leq 2ab \leq \lambda \\ \frac{b}{a} + \frac{\lambda}{2a^2} & \text{for } \lambda < -2ab \end{cases}$$

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



When does coordinate descent work?

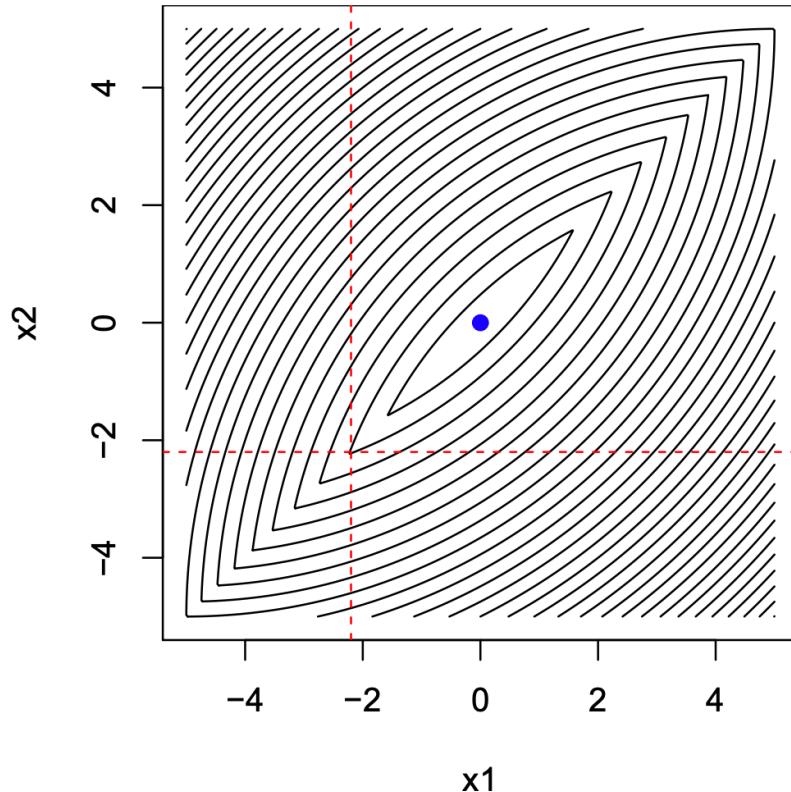
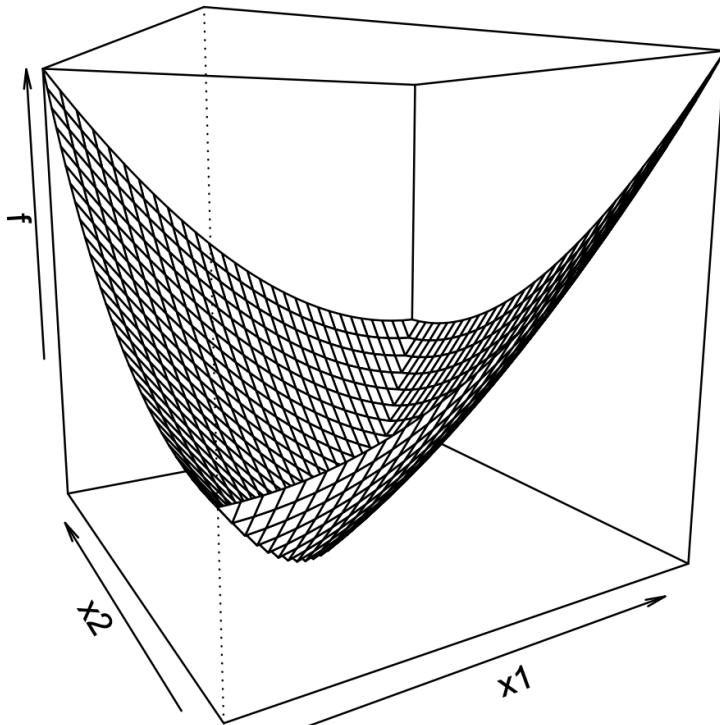
- Consider minimizing a **differentiable convex** function $f(x)$, then coordinate descent converges to the global minima



- when coordinate descent has stopped, that means
$$\frac{\partial f(x)}{\partial x_j} = 0 \text{ for all } j \in \{1, \dots, d\}$$
- 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

