

Linear Regression

CE-717: Machine Learning
Sharif University of Technology

M. Soleymani

Fall 2016

Topics

- ▶ **Linear regression**
 - ▶ Error (cost) function
 - ▶ Optimization
 - ▶ Generalization

Regression problem

- ▶ The goal is to make (real valued) predictions given features
- ▶ Example: predicting house price from 3 attributes

Size (m^2)	Age (year)	Region	Price (10^6T)
100	2	5	500
80	25	3	250
...

Learning problem

- ▶ Selecting a **hypothesis space**

- ▶ Hypothesis space: a set of mappings from feature vector to target

- ▶ **Learning (estimation)**: optimization of a cost function

- ▶ Based on the training set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ and a cost function we find (an estimate) $f \in F$ of the target function

- ▶ **Evaluation**: we measure how well \hat{f} generalizes to unseen examples

Learning problem

- ▶ Selecting a **hypothesis space**

- ▶ Hypothesis space: a set of mappings from feature vector to target

- ▶ **Learning (estimation)**: optimization of a cost function

- ▶ Based on the training set $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ and a cost function we find (an estimate) $f \in F$ of the target function

- ▶ **Evaluation**: we measure how well \hat{f} generalizes to unseen examples

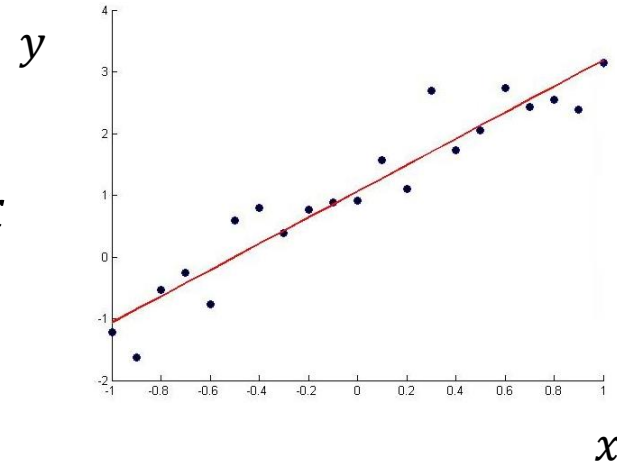
Hypothesis space

- ▶ Specify the class of functions (e.g., linear)
- ▶ We begin by the class of linear functions
 - ▶ easy to extend to generalized linear and so cover more complex regression functions

Linear regression: hypothesis space

▶ Univariate

$$f : \mathbb{R} \rightarrow \mathbb{R} \quad f(x; \mathbf{w}) = w_0 + w_1 x$$



▶ Multivariate

$$f : \mathbb{R}^d \rightarrow \mathbb{R} \quad f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$ are parameters we need to set.

Learning problem

- ▶ Selecting a **hypothesis space**

- ▶ Hypothesis space: a set of mappings from feature vector to target

- ▶ **Learning (estimation)**: optimization of a cost function

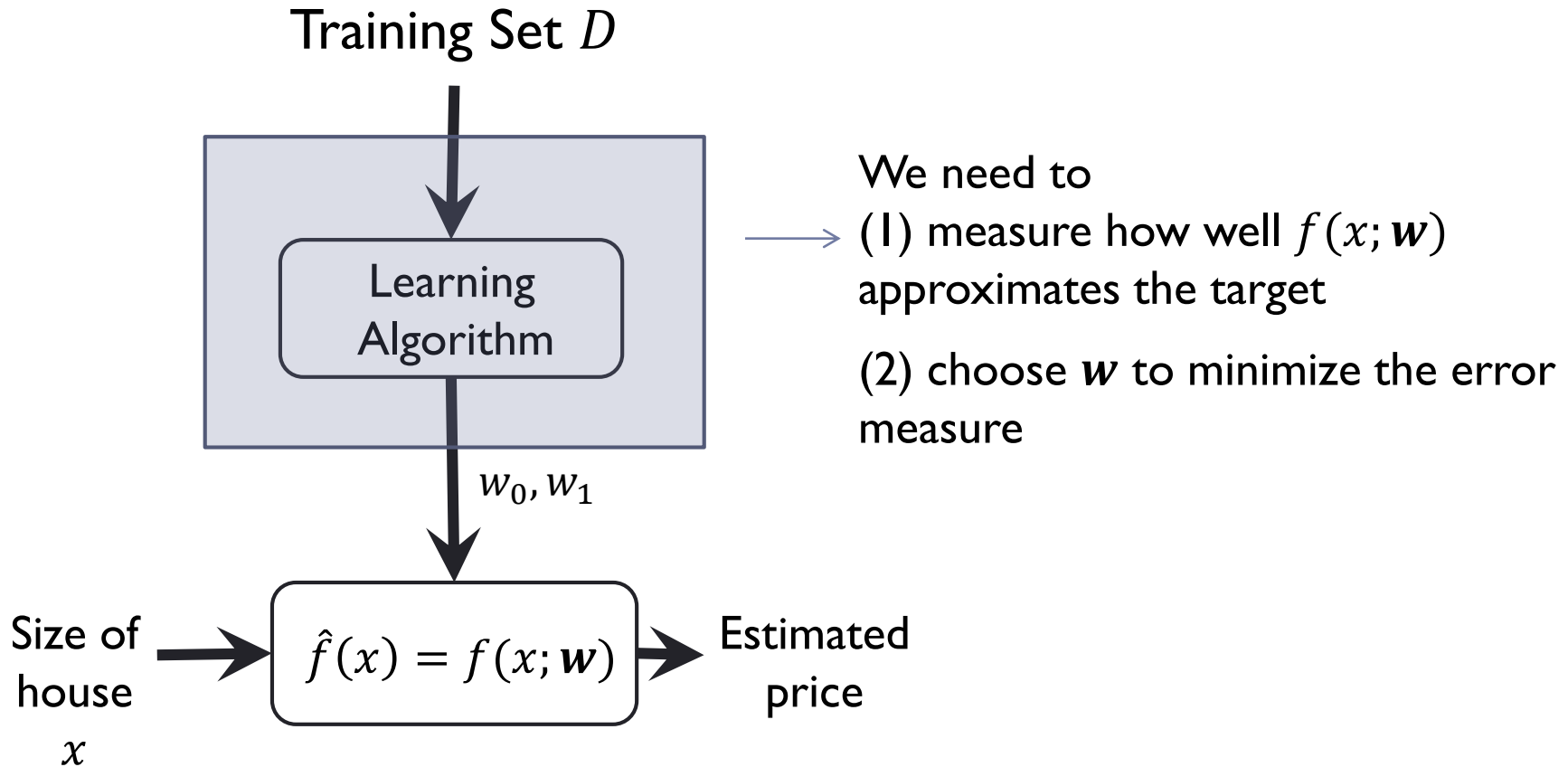
- ▶ Based on the training set $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$ and a cost function we find (an estimate) $f \in F$ of the target function

- ▶ **Evaluation**: we measure how well \hat{f} generalizes to unseen examples

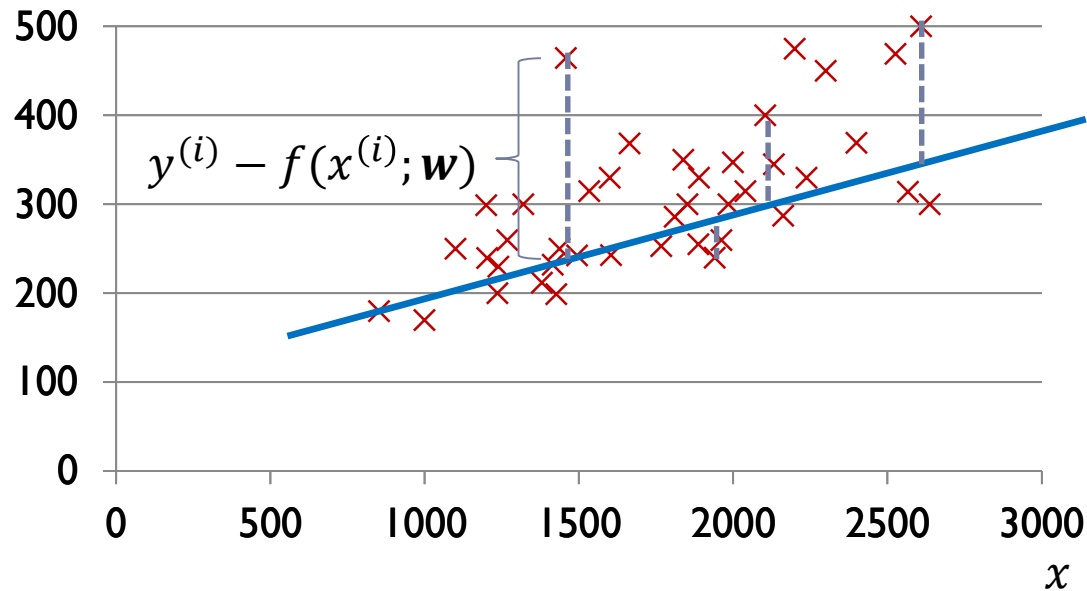
Learning algorithm

- ▶ Select how to measure the error (i.e. prediction loss)
- ▶ Find the minimum of the resulting error or cost function

Learning algorithm

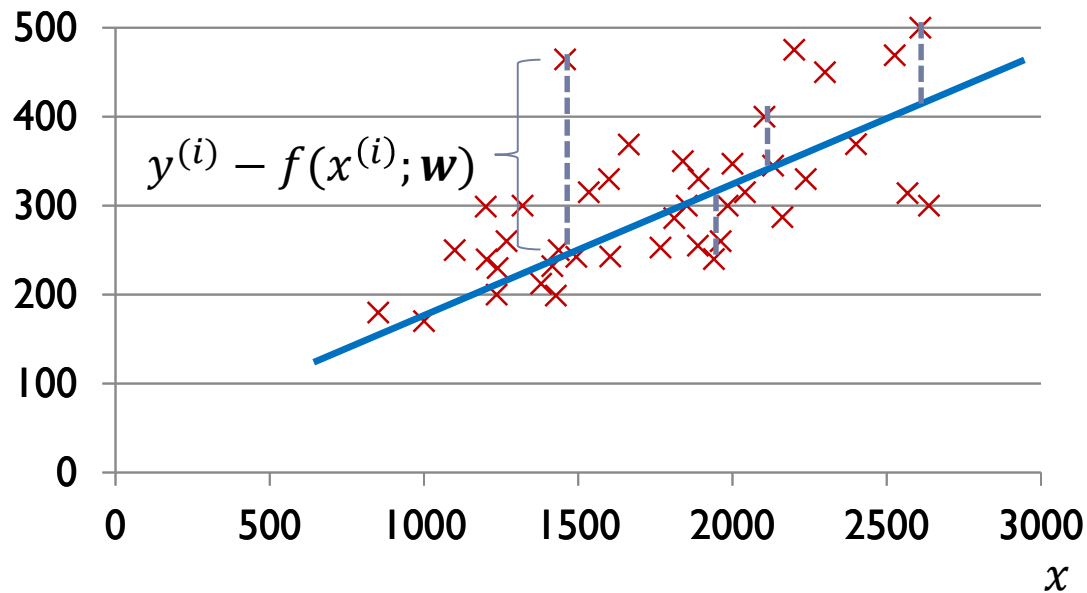


How to measure the error



$$\text{Squared error: } \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$

Linear regression: univariate example



Cost function:

$$\begin{aligned} J(\mathbf{w}) &= \sum_{i=1}^n (y^{(i)} - f(x; \mathbf{w}))^2 \\ &= \sum_{i=1}^n (y^{(i)} - w_0 - w_1 x^{(i)})^2 \end{aligned}$$

Regression: squared loss

- ▶ In the SSE cost function, we used squared error as the prediction loss:

$$\text{Loss}(y, \hat{y}) = (y - \hat{y})^2 \quad \hat{y} = f(\mathbf{x}; \mathbf{w})$$

- ▶ Cost function (based on the training set):

$$\begin{aligned} J(\mathbf{w}) &= \sum_{i=1}^n \text{Loss} \left(y^{(i)}, f(\mathbf{x}^{(i)}; \mathbf{w}) \right) \\ &= \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2 \end{aligned}$$

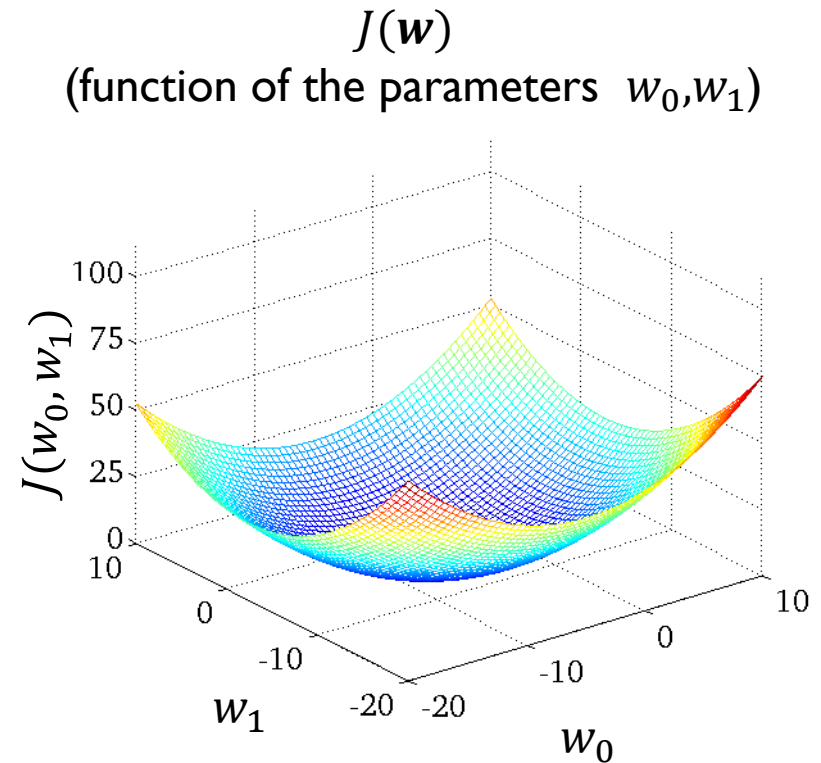
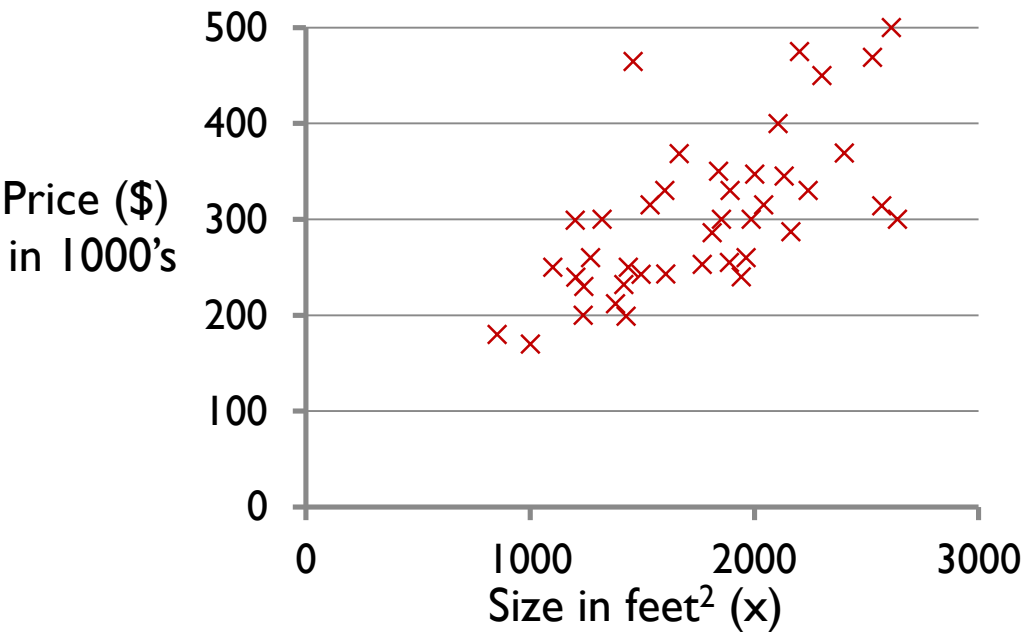
- ▶ Minimizing sum (or mean) of squared errors is a common approach in curve fitting, neural network, etc.

Sum of Squares Error (SSE) cost function

$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2$$

- ▶ $J(\mathbf{w})$: sum of the squares of the prediction errors on the training set
- ▶ We want to find the best regression function $f(\mathbf{x}^{(i)}; \mathbf{w})$
 - ▶ equivalently, the best \mathbf{w}
- ▶ Minimize $J(\mathbf{w})$
 - ▶ Find optimal $\hat{f}(\mathbf{x}) = f(\mathbf{x}; \hat{\mathbf{w}})$ where $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$

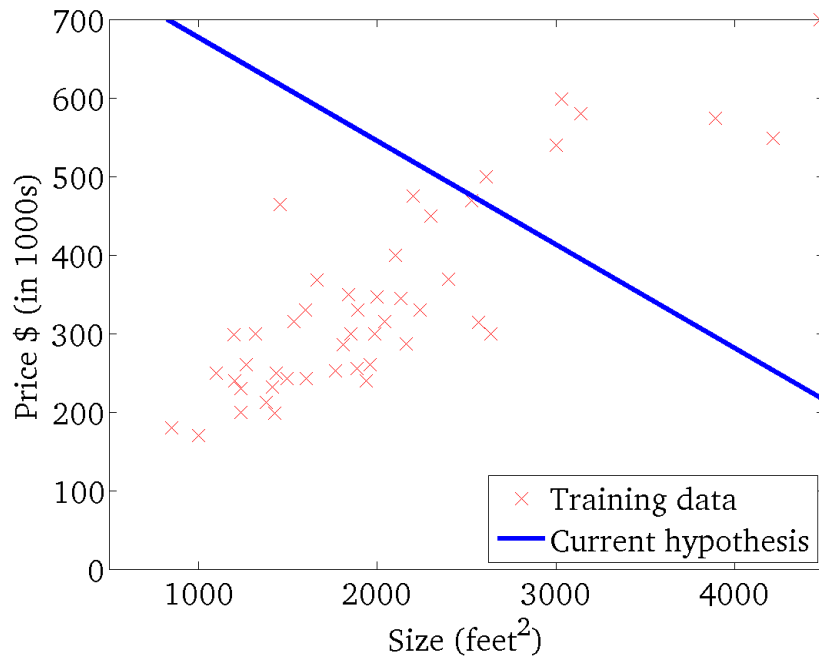
Cost function: univariate example



Cost function: univariate example

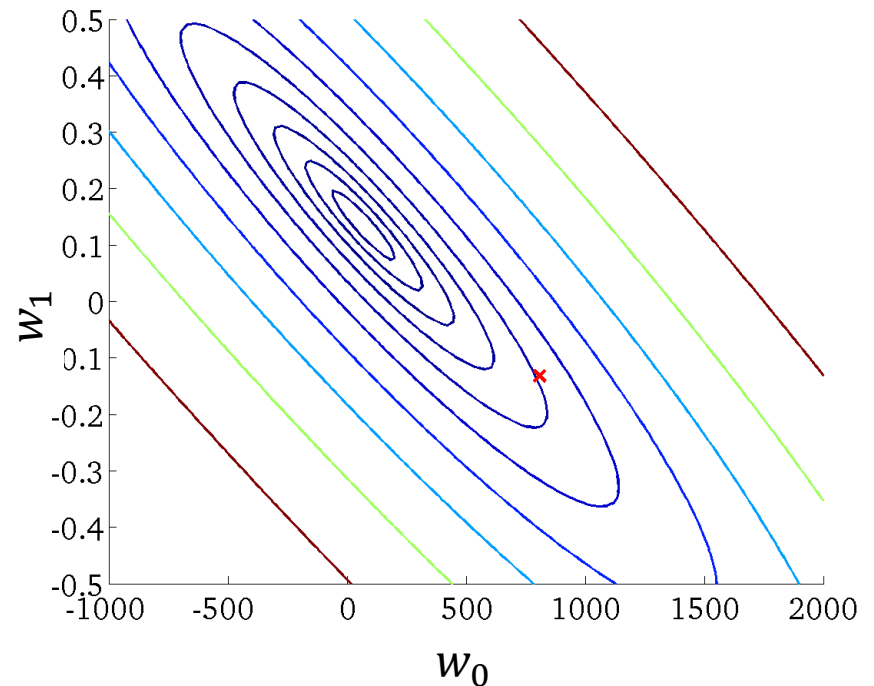
$$f(x; w_0, w_1) = w_0 + w_1 x$$

(for fixed w_0, w_1 , this is a function of x)



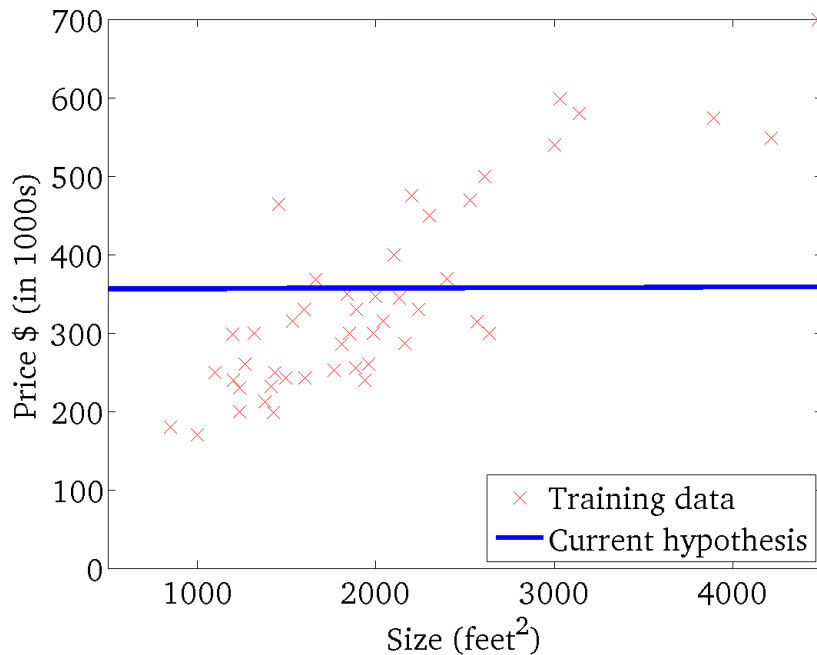
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



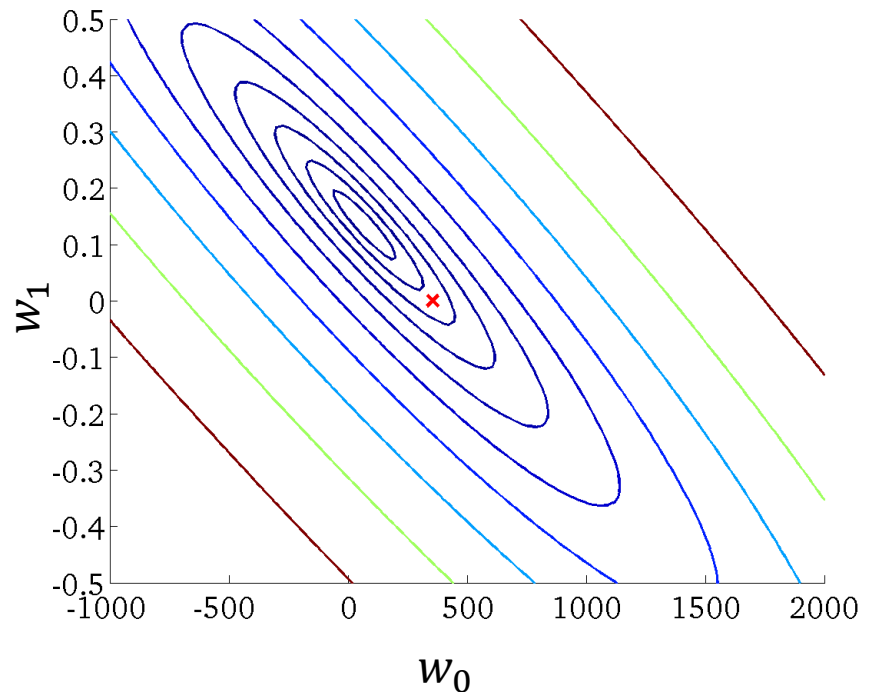
Cost function: univariate example

$$f(x; w_0, w_1) = w_0 + w_1 x$$



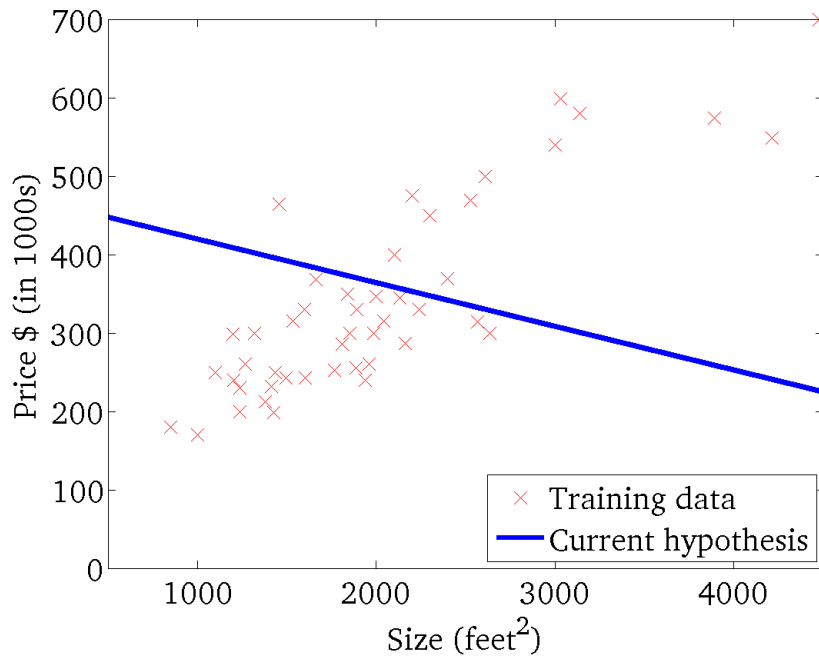
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



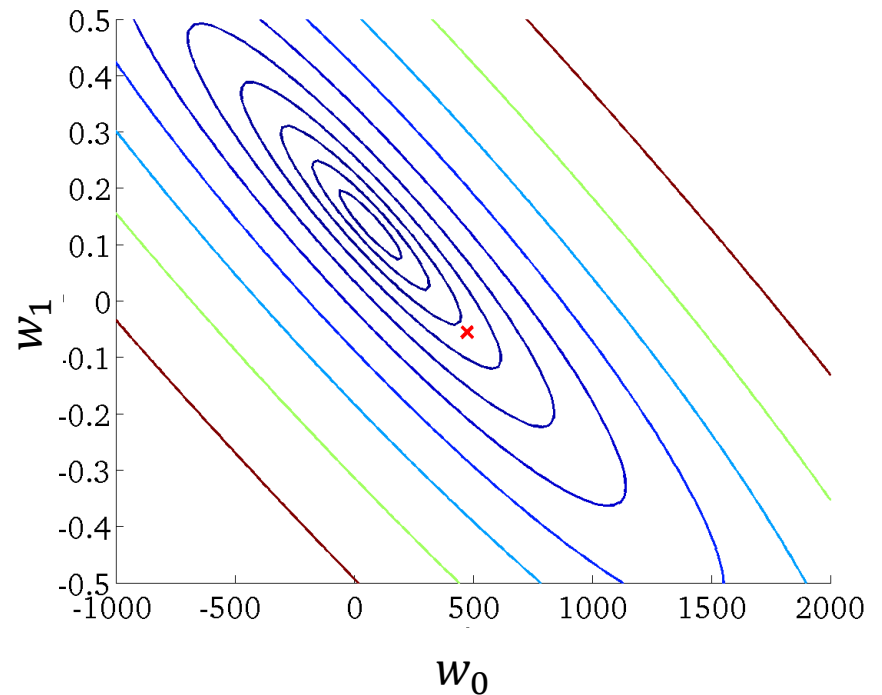
Cost function: univariate example

$$f(x; w_0, w_1) = w_0 + w_1 x$$



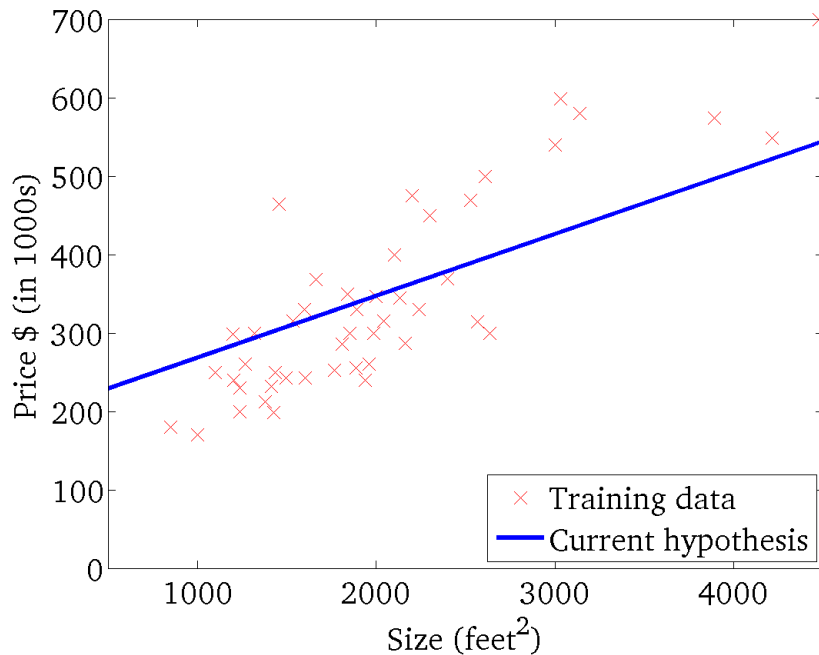
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



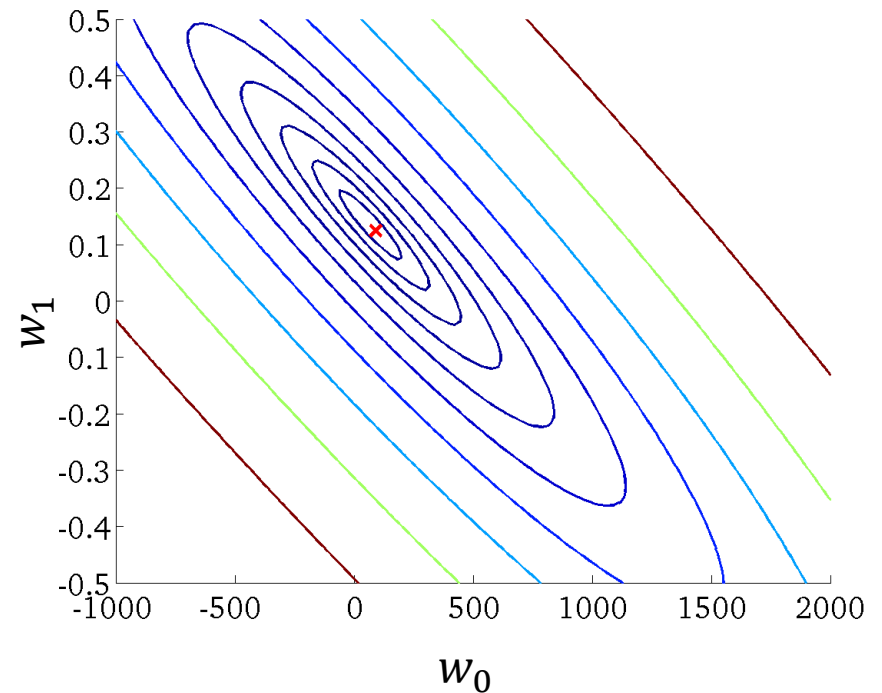
Cost function: univariate example

$$f(x; w_0, w_1) = w_0 + w_1 x$$



$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



Cost function optimization: univariate

$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - w_0 - w_1 x^{(i)})^2$$

- ▶ Necessary conditions for the “optimal” parameter values:

$$\frac{\partial J(\mathbf{w})}{\partial w_0} = 0$$

$$\frac{\partial J(\mathbf{w})}{\partial w_1} = 0$$

Optimality conditions: univariate

$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - w_0 - w_1 x^{(i)})^2$$

$$\frac{\partial J(\mathbf{w})}{\partial w_1} = \sum_{i=1}^n 2(y^{(i)} - w_0 - w_1 x^{(i)})(-x^{(i)}) = 0$$

$$\frac{\partial J(\mathbf{w})}{\partial w_0} = \sum_{i=1}^n 2(y^{(i)} - w_0 - w_1 x^{(i)})(-1) = 0$$

- ▶ A systems of 2 linear equations

Cost function: multivariate

- ▶ We have to minimize the empirical squared loss:

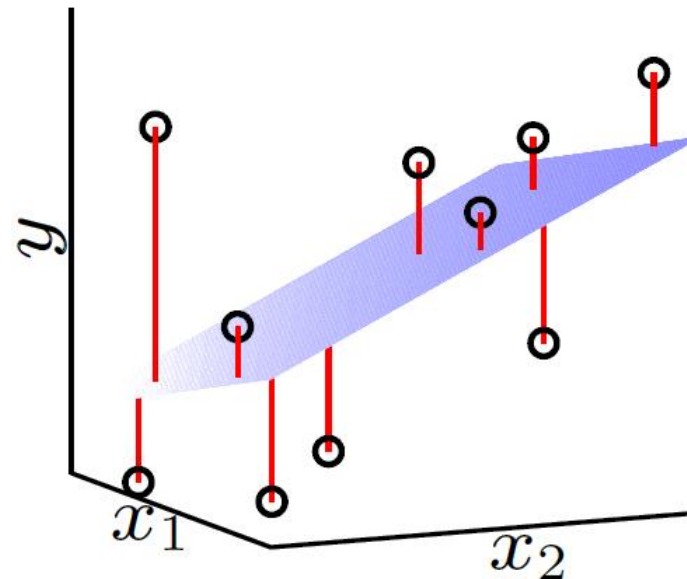
$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2$$

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

$$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} J(\mathbf{w})$$

Cost function and optimal linear model



- ▶ Necessary conditions for the “optimal” parameter values:

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \mathbf{0}$$

- ▶ A system of $d + 1$ linear equations

Cost function: matrix notation

$$\begin{aligned} J(\mathbf{w}) &= \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2 = \\ &= \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)} \right)^2 \end{aligned}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

Minimizing cost function

Optimal linear weight vector (for SSE cost function):

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

$$\nabla_{\mathbf{w}}J(\mathbf{w}) = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\nabla_{\mathbf{w}}J(\mathbf{w}) = \mathbf{0} \Rightarrow \mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

Minimizing cost function

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$$

$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

\mathbf{X}^\dagger is pseudo inverse of \mathbf{X}

Another approach for optimizing the sum squared error

- ▶ Iterative approach for solving the following optimization problem:

$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2$$

Review:

Iterative optimization of cost function

- ▶ Cost function: $J(\mathbf{w})$
- ▶ Optimization problem: $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$

- ▶ Steps:
 - ▶ Start from \mathbf{w}^0
 - ▶ Repeat
 - ▶ Update \mathbf{w}^t to \mathbf{w}^{t+1} in order to reduce J
 - ▶ $t \leftarrow t + 1$
 - ▶ until we hopefully end up at a minimum

Review:

Gradient descent

- ▶ First-order optimization algorithm to find $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$
 - ▶ Also known as "steepest descent"
- ▶ In each step, takes steps proportional to the negative of the gradient vector of the function at the current point \mathbf{w}^t :

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \gamma_t \nabla J(\mathbf{w}^t)$$

- ▶ $J(\mathbf{w})$ decreases fastest if one goes from \mathbf{w}^t in the direction of $-\nabla J(\mathbf{w}^t)$
- ▶ Assumption: $J(\mathbf{w})$ is defined and differentiable in a neighborhood of a point \mathbf{w}^t

Gradient ascent takes steps proportional to (the positive of) the gradient to find a local maximum of the function

Review:

Gradient descent

► Minimize $J(\mathbf{w})$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^t)$$

Step size
(Learning rate parameter)

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \left[\frac{\partial J(\mathbf{w})}{\partial w_1}, \frac{\partial J(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_d} \right]$$

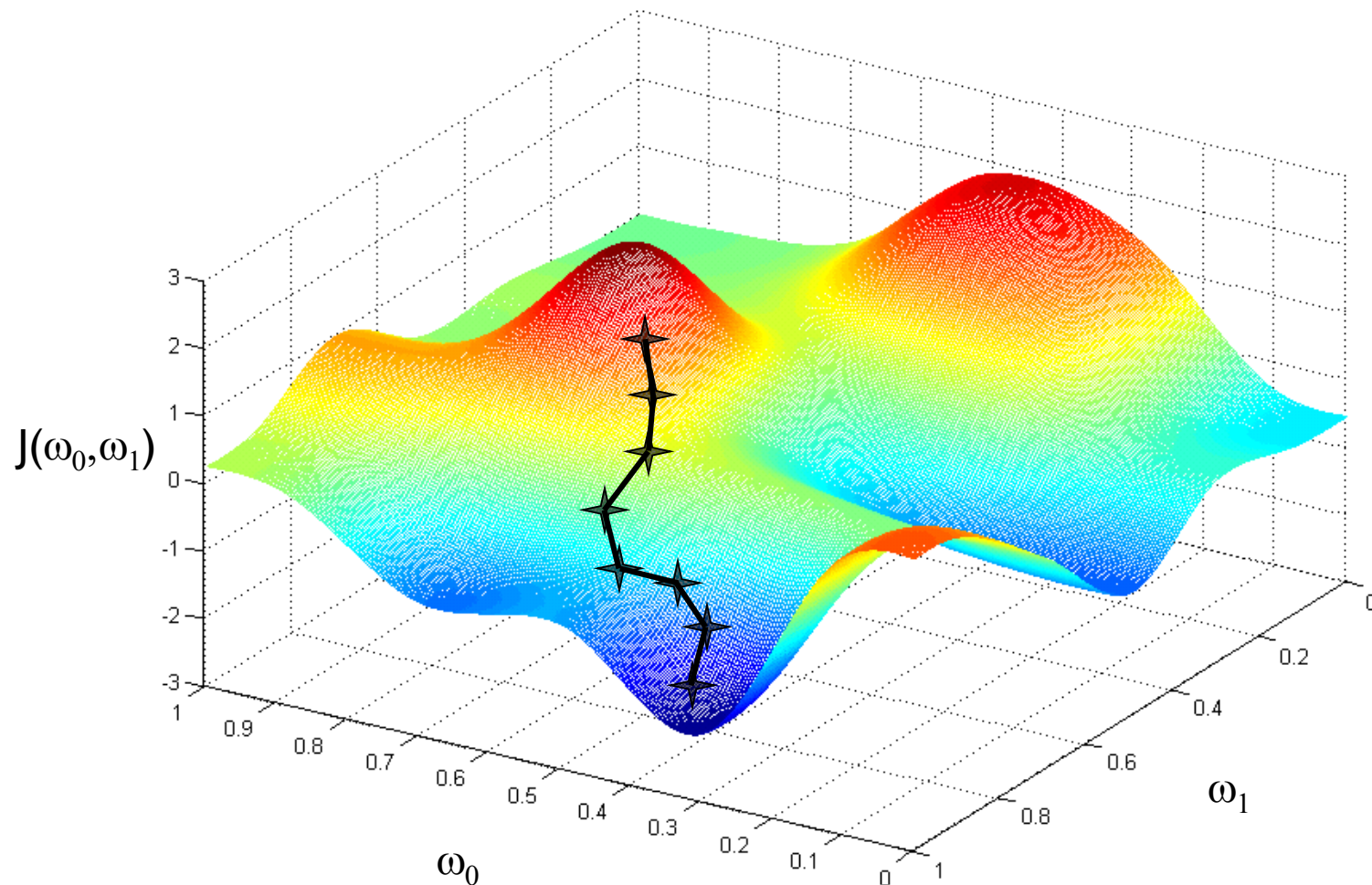
- If η is small enough, then $J(\mathbf{w}^{t+1}) \leq J(\mathbf{w}^t)$.
- η can be allowed to change at every iteration as η_t .

Review:

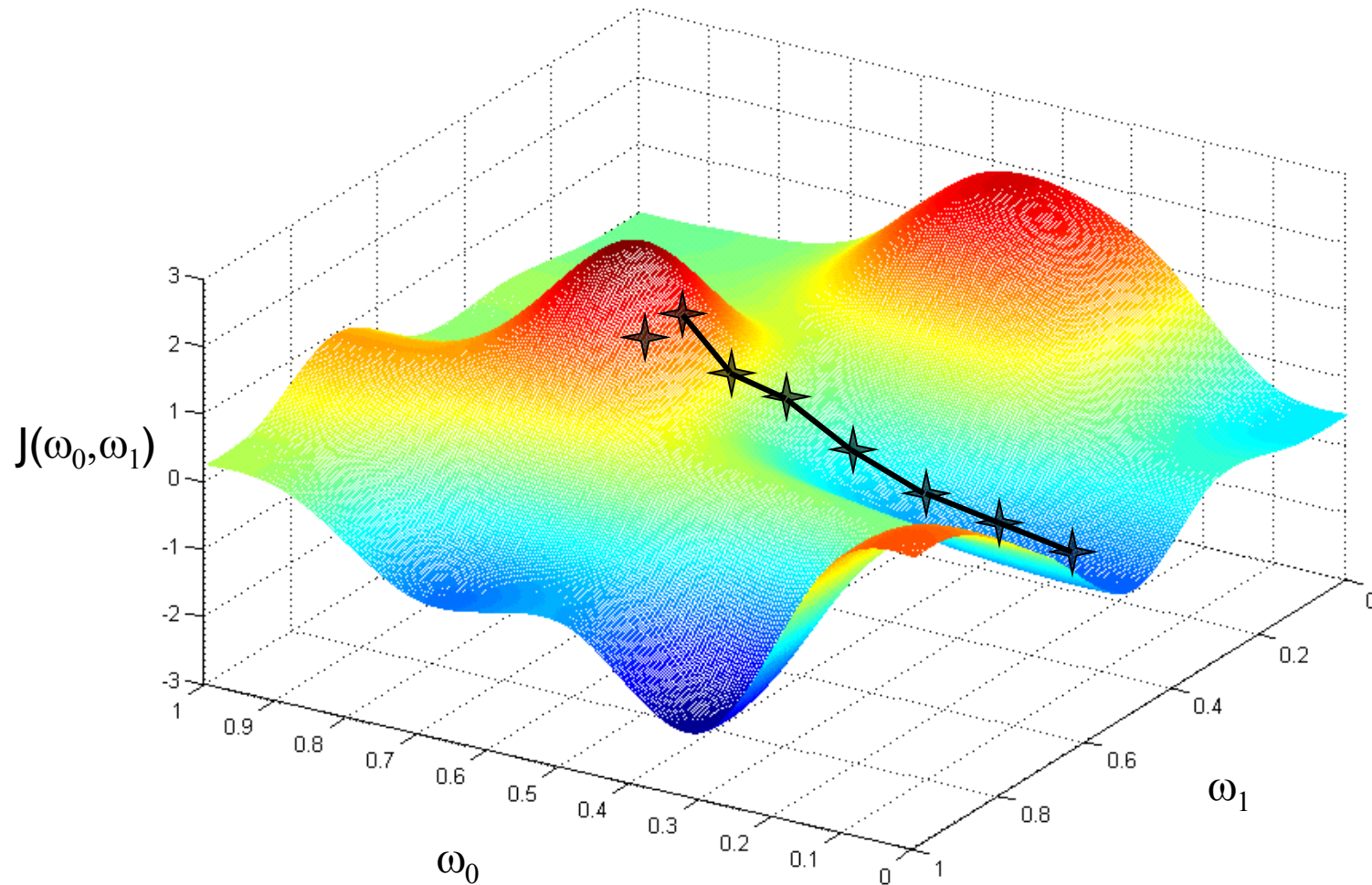
Gradient descent disadvantages

- ▶ Local minima problem
- ▶ However, when J is convex, all local minima are also global minima \Rightarrow gradient descent can converge to the global solution.

Review: Problem of gradient descent with non-convex cost functions



Review: Problem of gradient descent with non-convex cost functions



Gradient descent for SSE cost function

- ▶ Minimize $J(\mathbf{w})$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^t)$$

- ▶ $J(\mathbf{w})$: Sum of squares error

$$J(\mathbf{w}) = \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$

- ▶ Weight update rule for $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^{tT} \mathbf{x}^{(i)} \right) \mathbf{x}^{(i)}$$

Gradient descent for SSE cost function

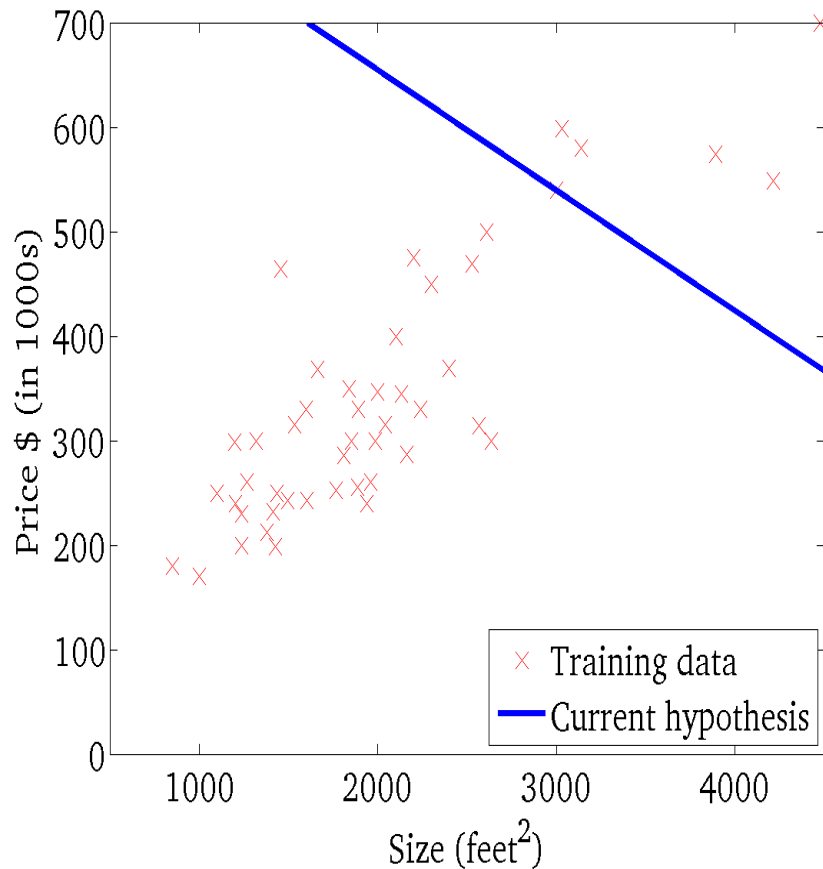
- ▶ Weight update rule: $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta \sum_{i=1}^n (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}) \mathbf{x}^{(i)}$$

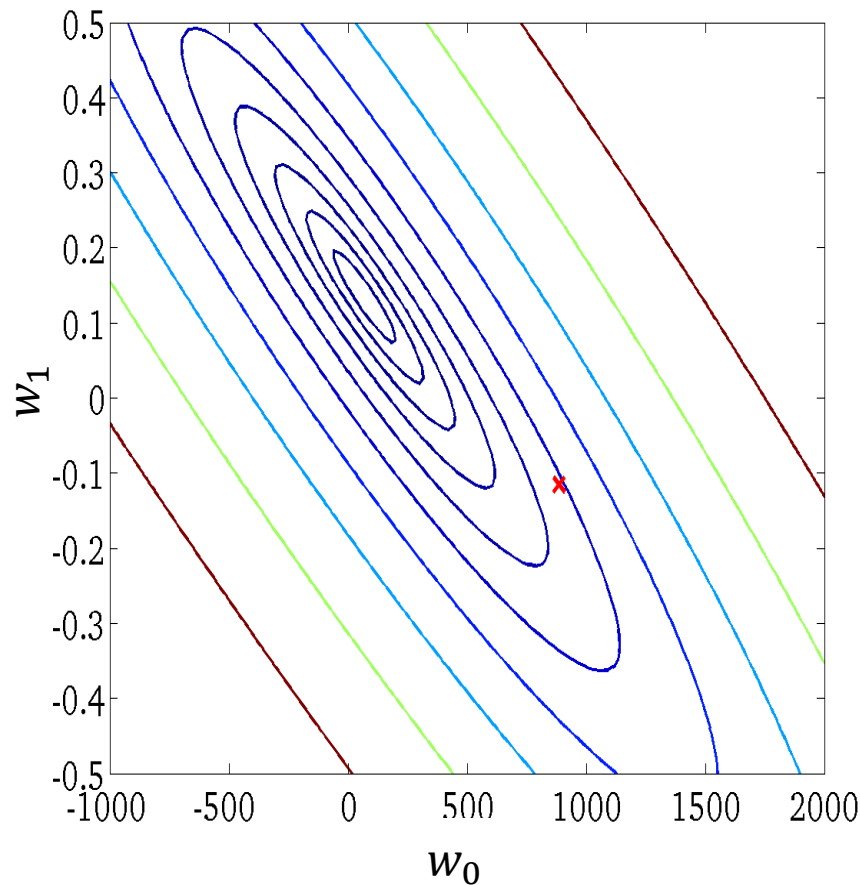
Batch mode: each step
considers all training data

- ▶ η : too small \rightarrow gradient descent can be slow.
- ▶ η : too large \rightarrow gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

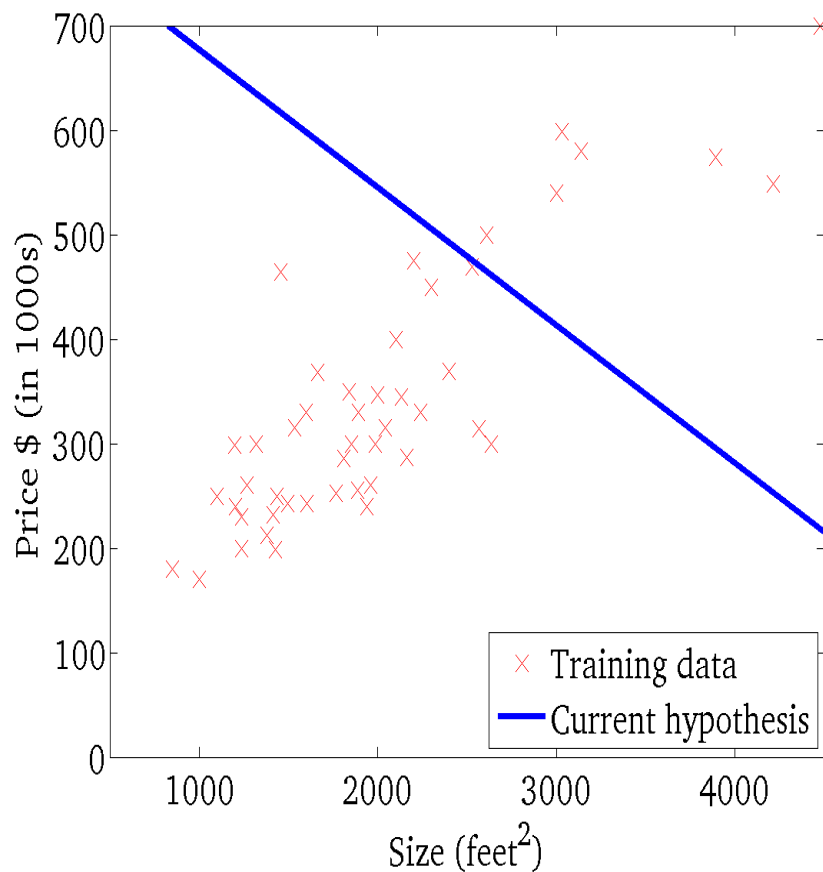
$$f(x; w_0, w_1) = w_0 + w_1 x$$



$J(w_0, w_1)$
(function of the parameters w_0, w_1)

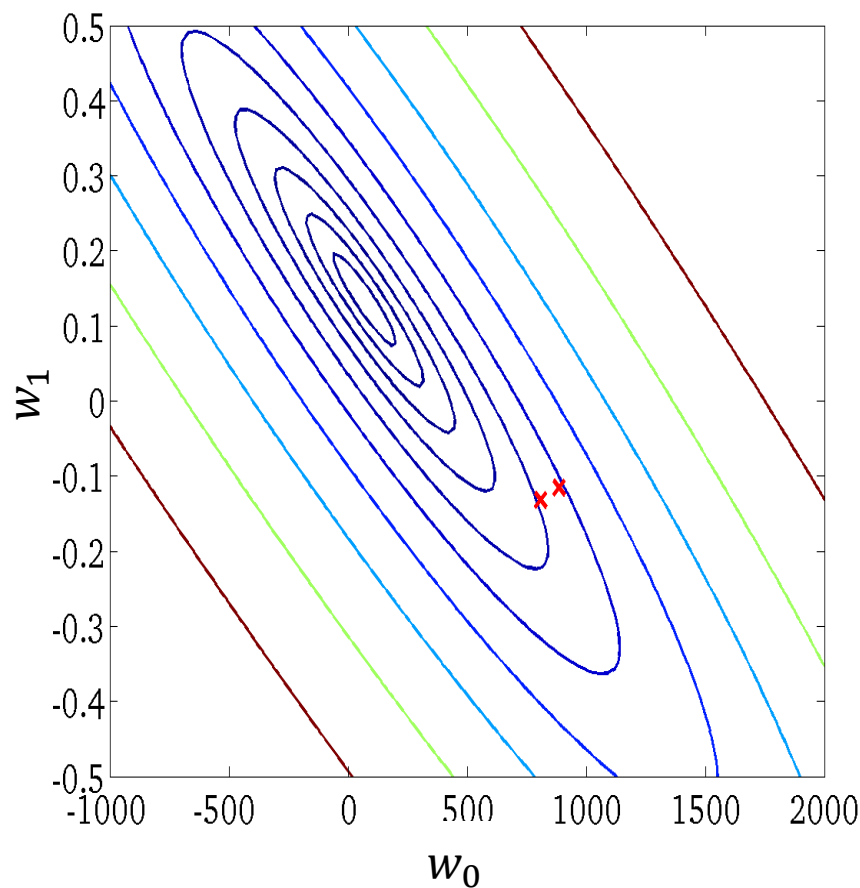


$$f(x; w_0, w_1) = w_0 + w_1 x$$

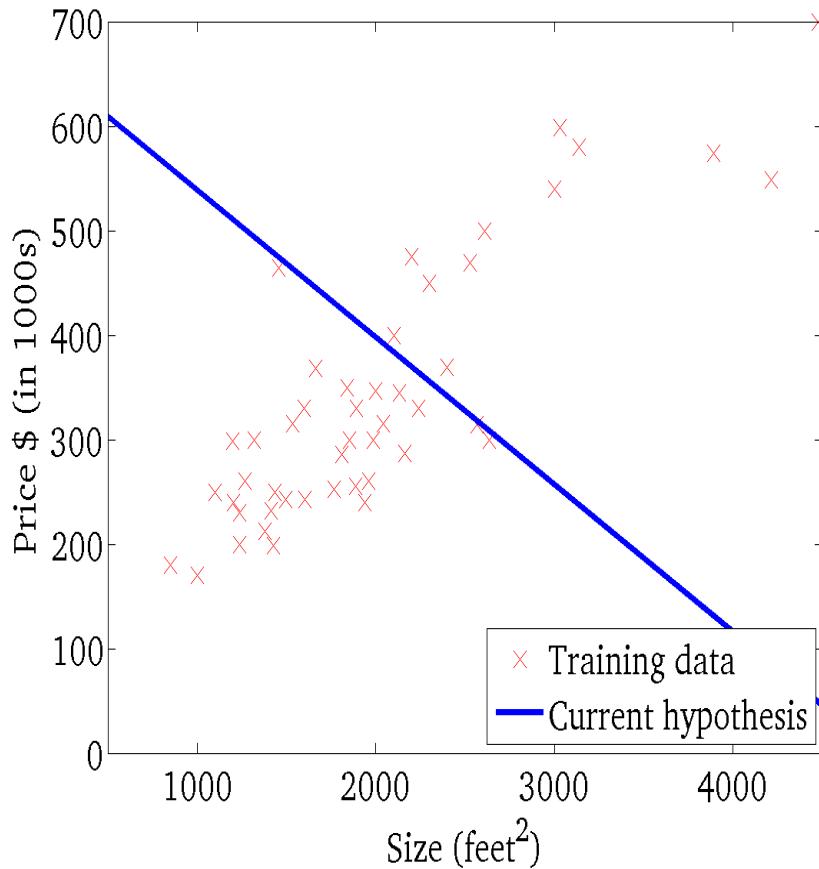


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

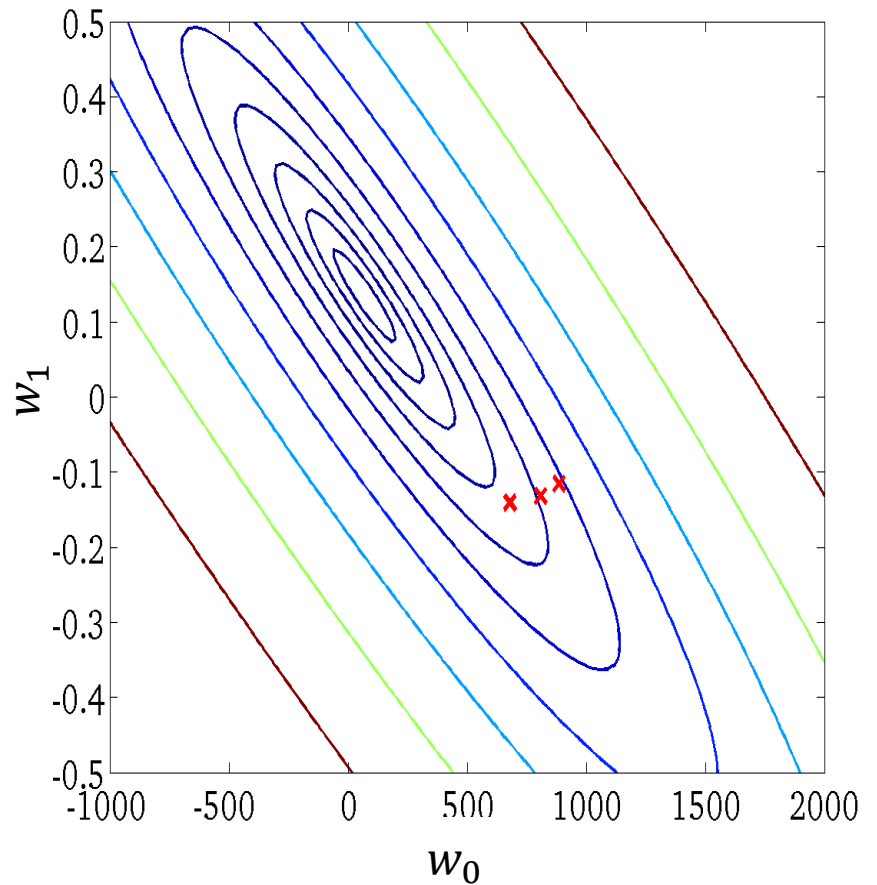


$$f(x; w_0, w_1) = w_0 + w_1 x$$

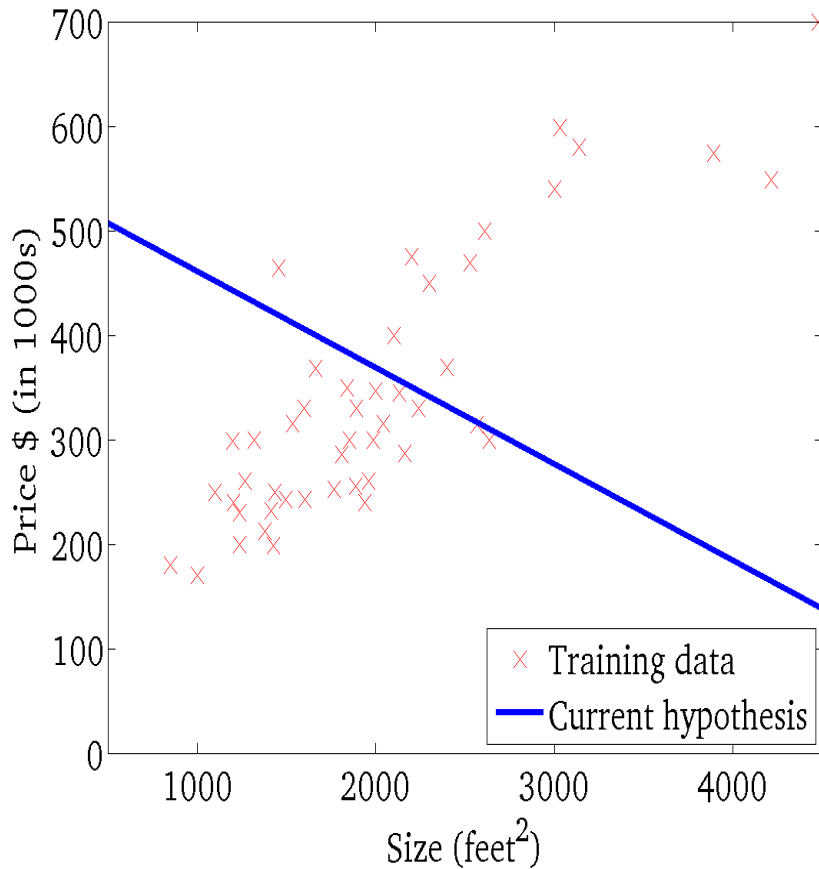


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

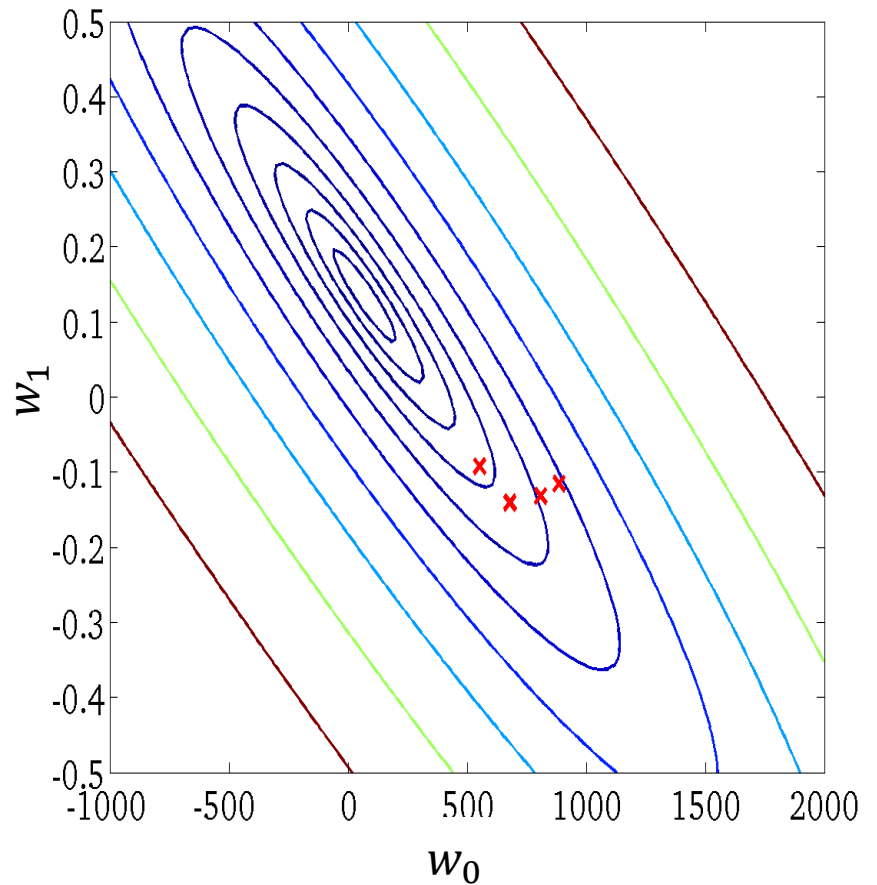


$$f(x; w_0, w_1) = w_0 + w_1 x$$

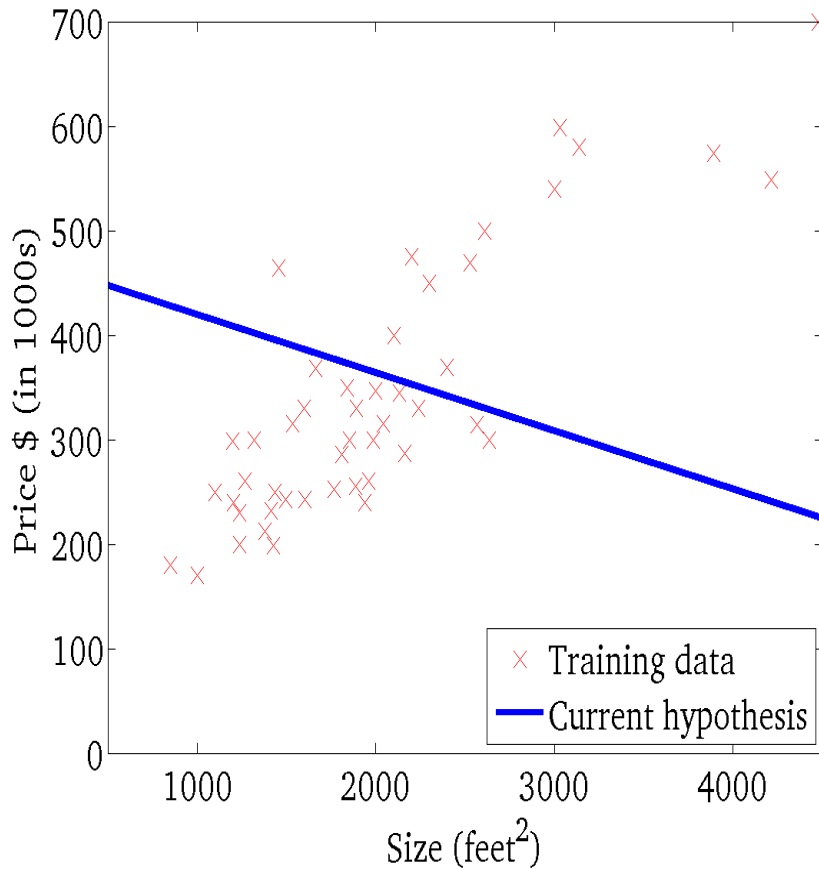


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

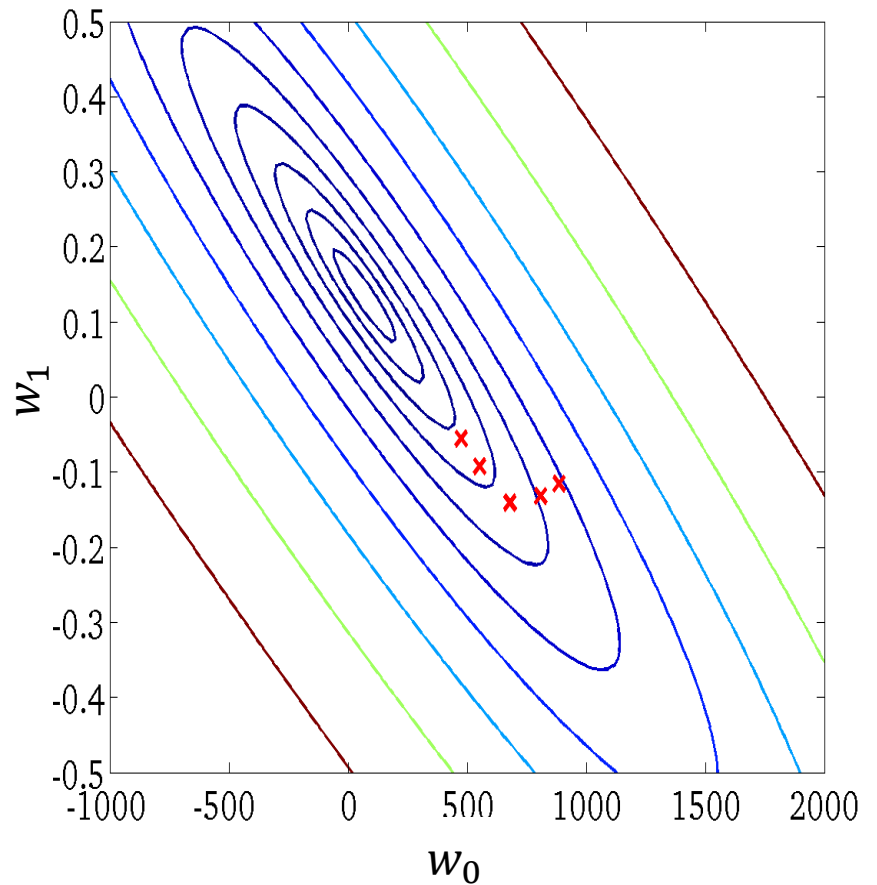


$$f(x; w_0, w_1) = w_0 + w_1 x$$

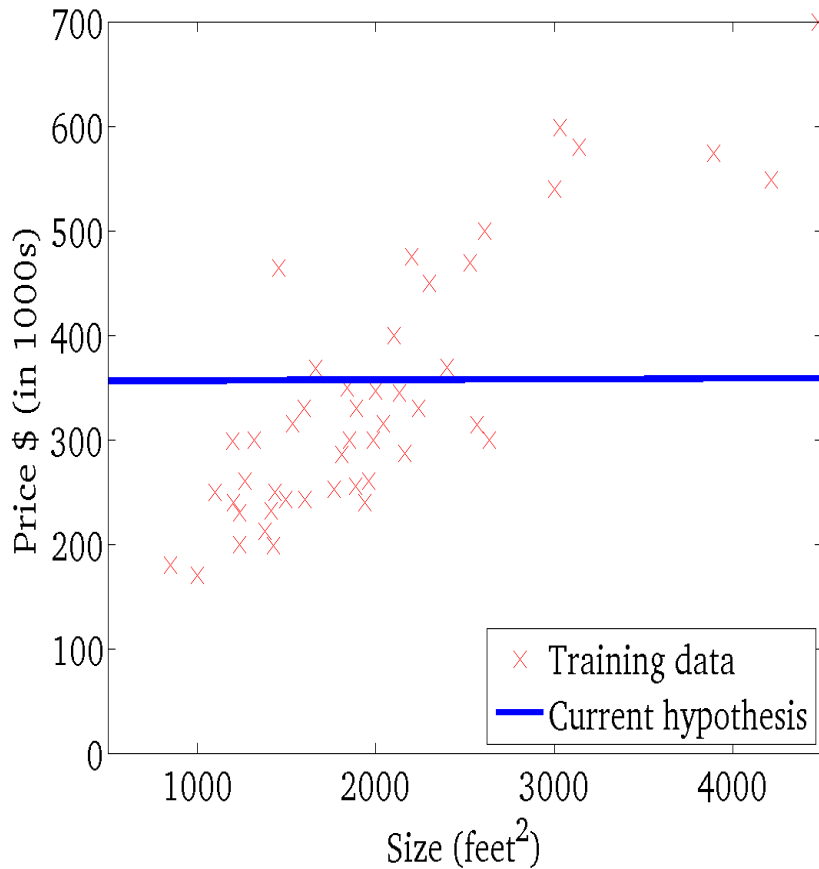


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

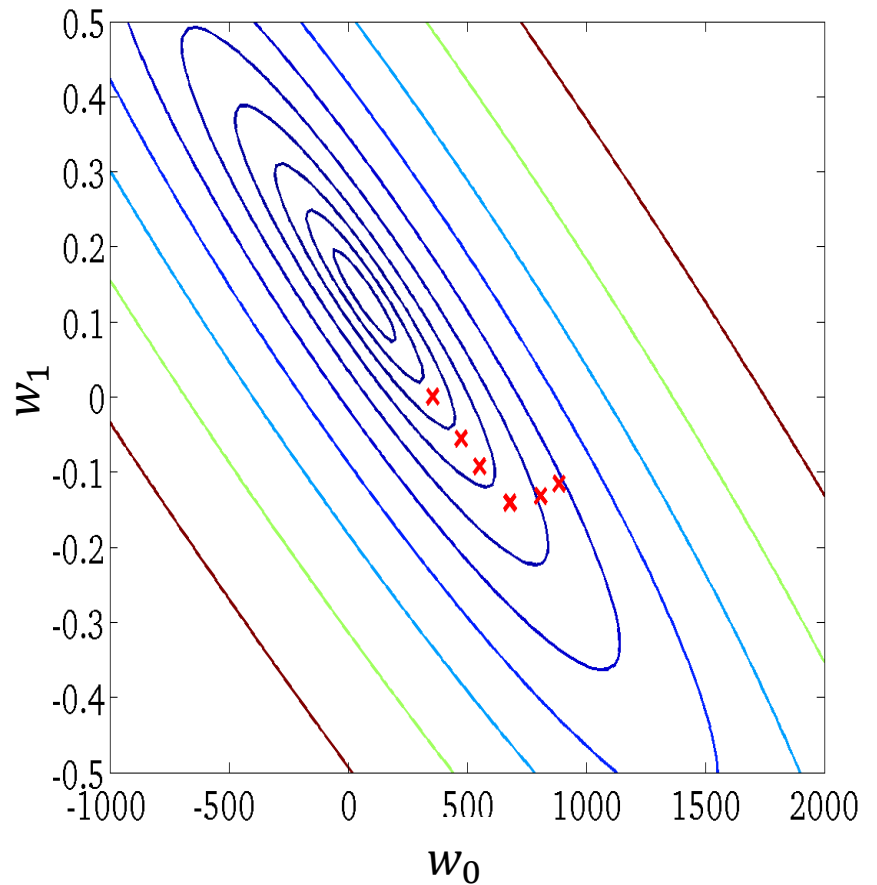


$$f(x; w_0, w_1) = w_0 + w_1 x$$

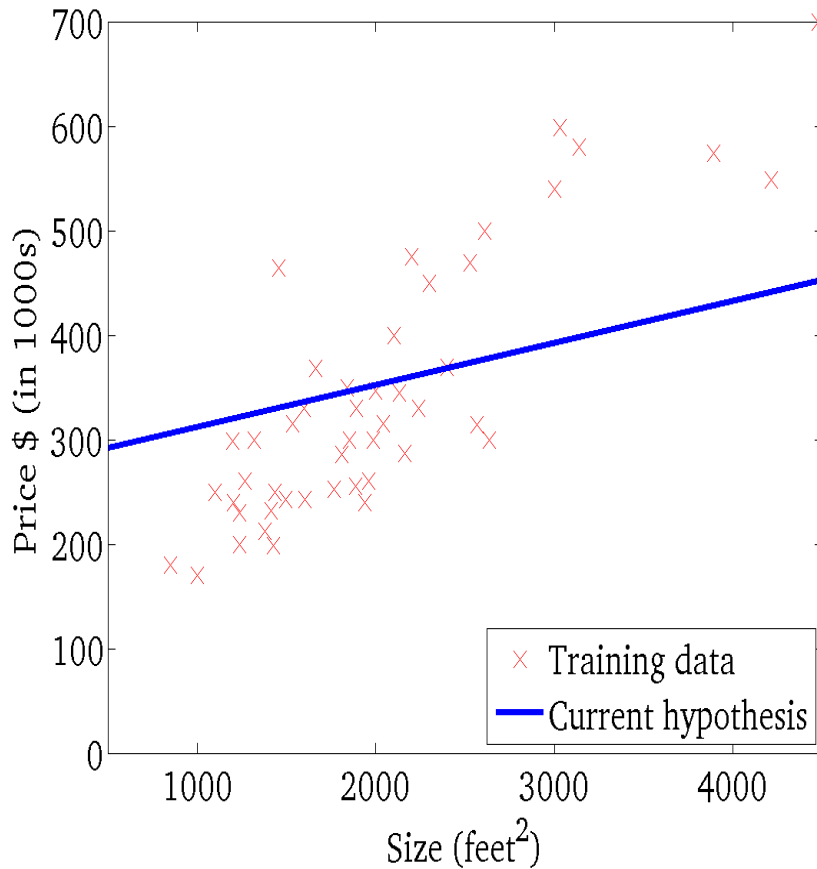


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

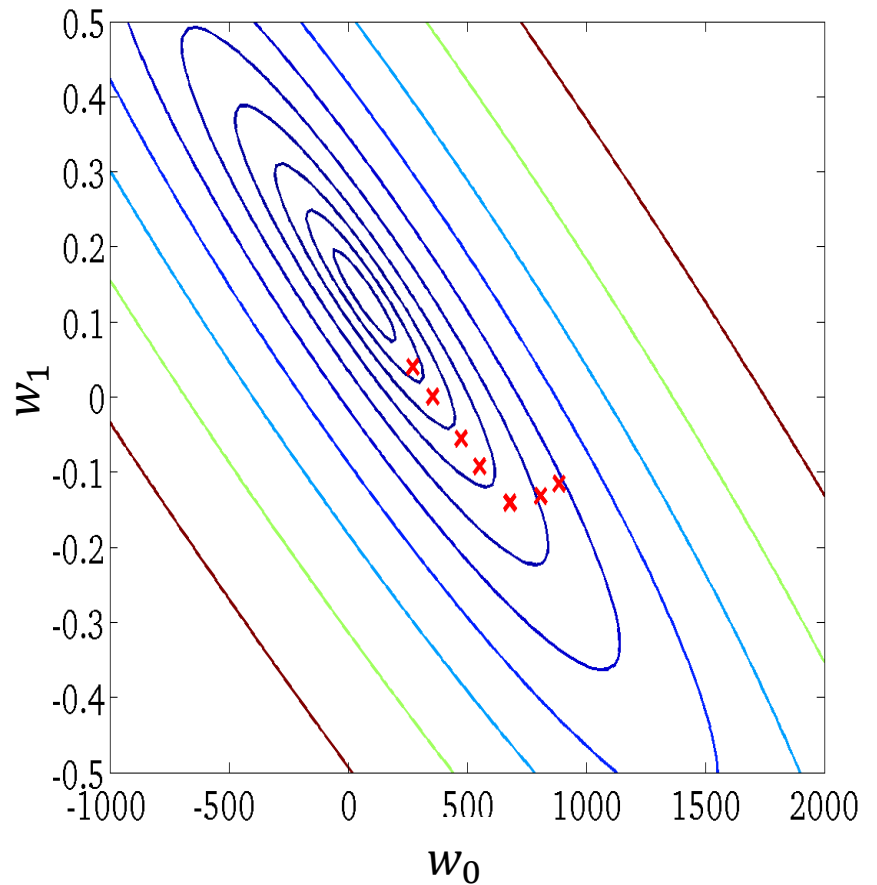


$$f(x; w_0, w_1) = w_0 + w_1 x$$

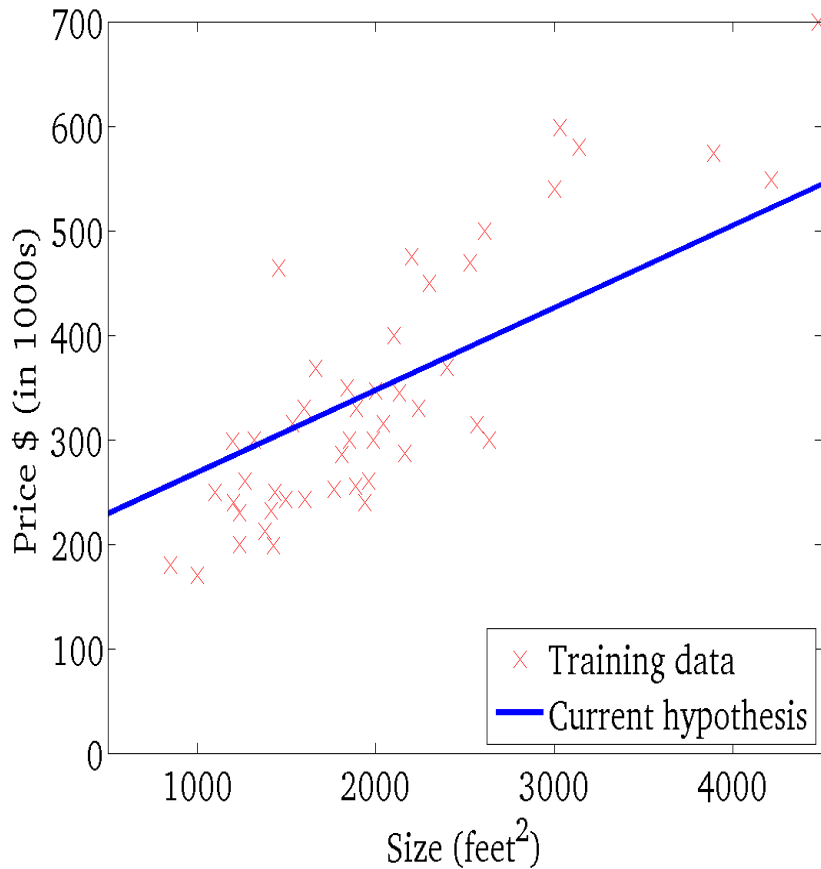


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

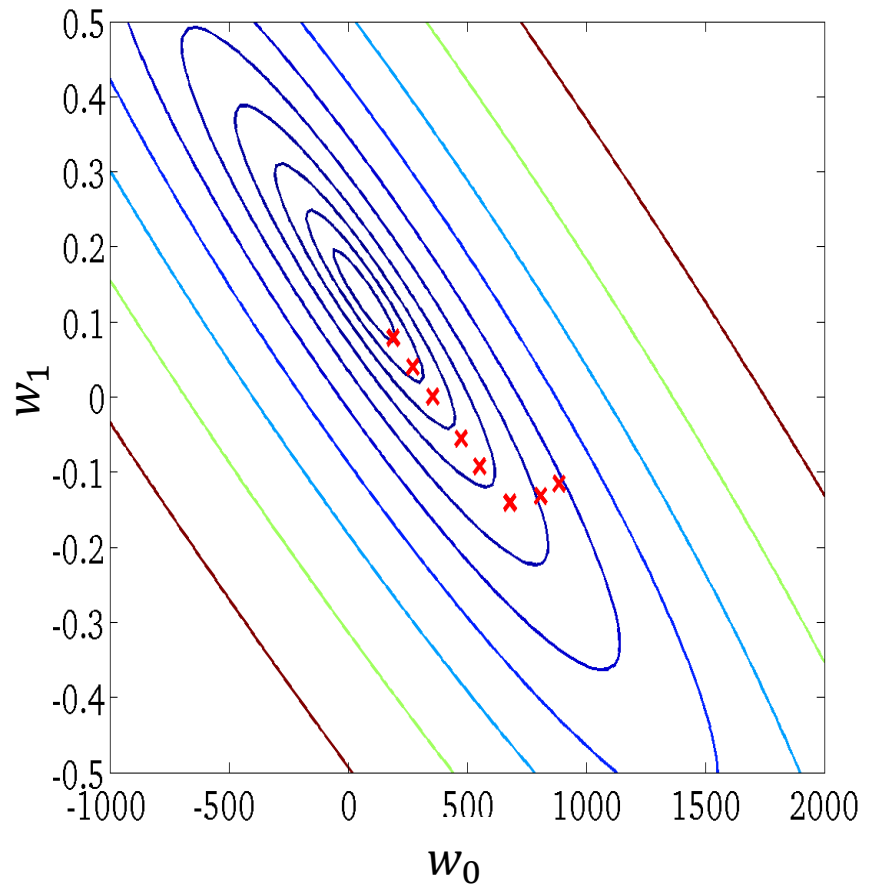


$$f(x; w_0, w_1) = w_0 + w_1 x$$

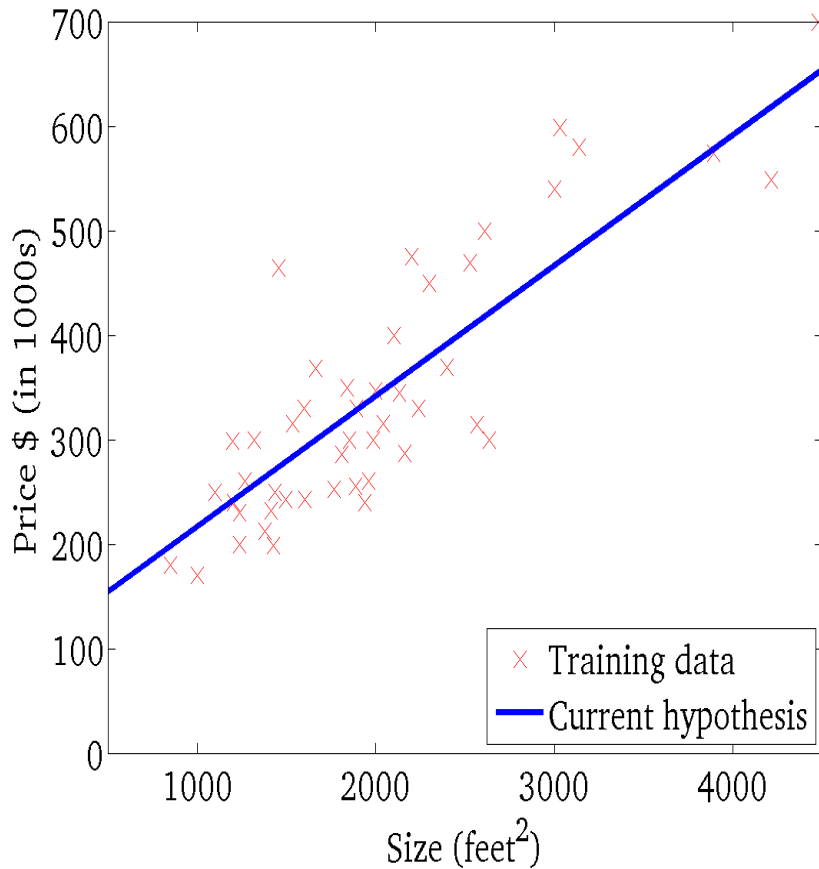


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

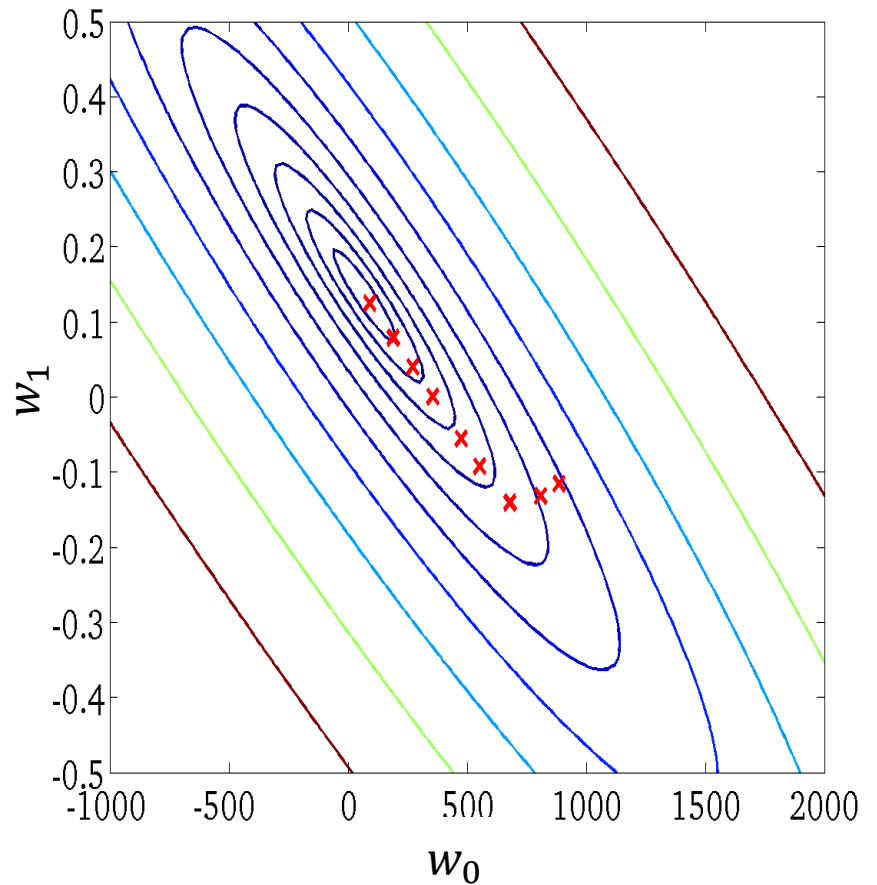


$$f(x; w_0, w_1) = w_0 + w_1 x$$



$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



Stochastic gradient descent

- ▶ Batch techniques process the entire training set in one go
 - ▶ thus they can be computationally costly for large data sets.
- ▶ Stochastic gradient descent: when the cost function can comprise a sum over data points:

$$J(\mathbf{w}) = \sum_{i=1}^n J^{(i)}(\mathbf{w})$$

- ▶ Update after presentation of $(\mathbf{x}^{(i)}, y^{(i)})$:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J^{(i)}(\mathbf{w})$$

Stochastic gradient descent

- ▶ Example: Linear regression with SSE cost function

$$J^{(i)}(\mathbf{w}) = (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J^{(i)}(\mathbf{w})$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}) \mathbf{x}^{(i)}$$

Least Mean Squares (LMS)

It is proper for sequential or online learning

Stochastic gradient descent: online learning

- ▶ Sequential learning is also appropriate for real-time applications
 - ▶ data observations are arriving in a continuous stream
 - ▶ and predictions must be made before seeing all of the data
- ▶ The value of η needs to be chosen with care to ensure that the algorithm converges

Evaluation and generalization

- ▶ Why minimizing the cost function (based on only training data) while we are interested in the performance on new examples?

$$\min_{\theta} \sum_{i=1}^n \text{Loss} \left(y^{(i)}, f(\mathbf{x}^{(i)}; \theta) \right) \longrightarrow \text{Empirical loss}$$

- ▶ **Evaluation:** After training, we need to measure how well the learned prediction function can predicts the target for unseen examples

Training and test performance

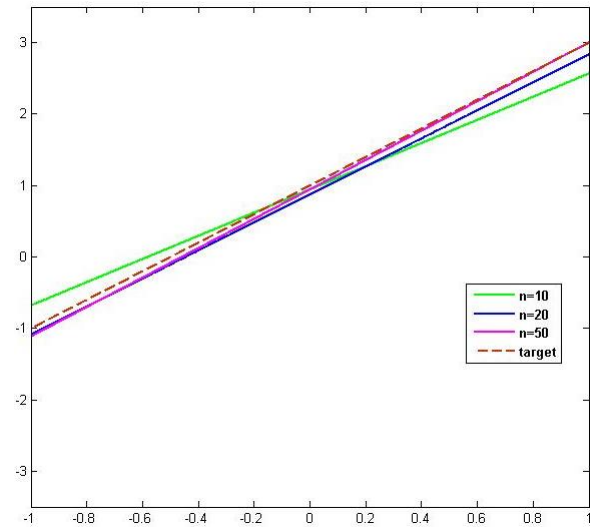
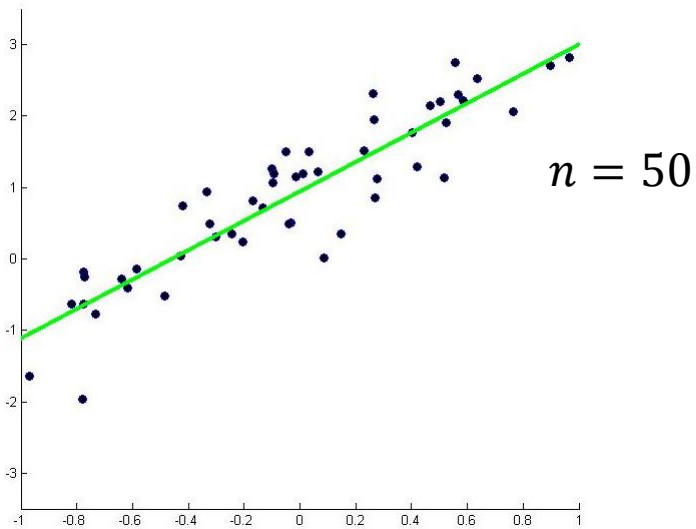
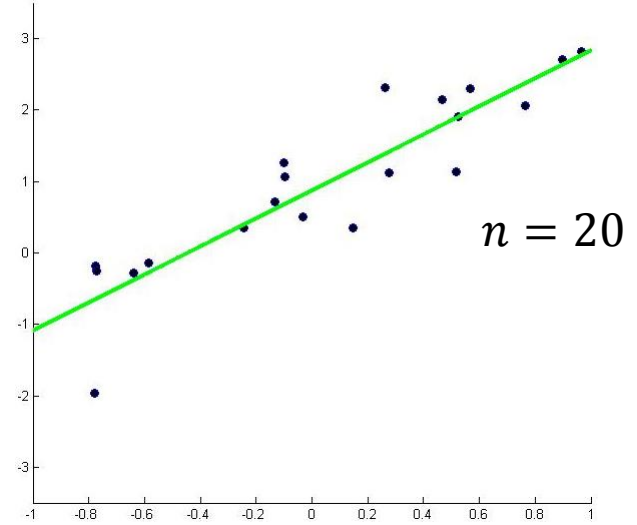
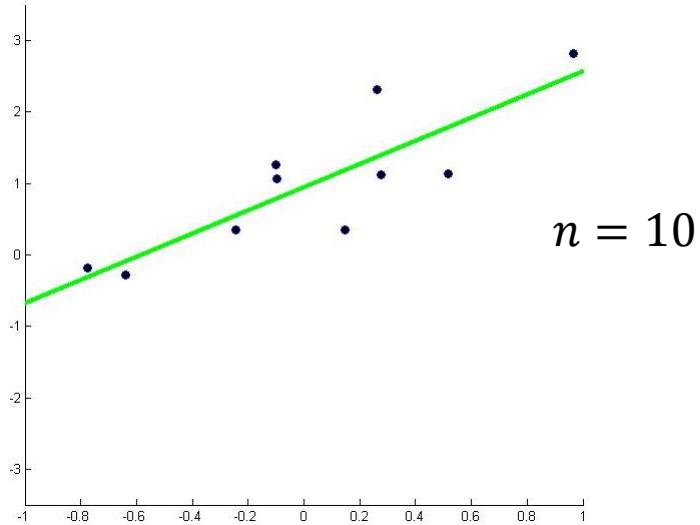
- ▶ **Assumption:** training and test examples are drawn independently at random from the same but unknown distribution.
 - ▶ Each training/test example (\mathbf{x}, y) is a sample from joint probability distribution $P(\mathbf{x}, y)$, i.e., $(\mathbf{x}, y) \sim P$

$$\text{Empirical (training) loss} = \frac{1}{n} \sum_{i=1}^n \text{Loss} \left(y^{(i)}, f(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \right)$$

$$\text{Expected (test) loss} = E_{\mathbf{x}, y} \{ \text{Loss}(y, f(\mathbf{x}; \boldsymbol{\theta})) \}$$

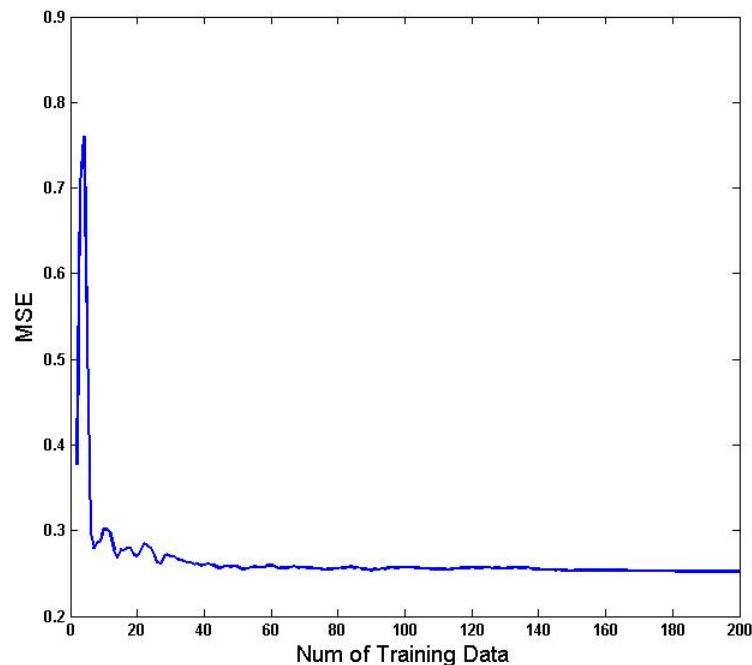
- ▶ We minimize empirical loss (on the training data) and expect to also find an acceptable expected loss
 - ▶ Empirical loss as a proxy for the performance over the whole distribution.

Linear regression: number of training data



Linear regression: generalization

- ▶ By increasing the number of training examples, will solution be better?
- ▶ Why the mean squared error does not decrease more after reaching a level?



Linear regression: types of errors

- ▶ Structural error: the error introduced by the limited function class (infinite training data):

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} E_{x,y} [(y - \mathbf{w}^T \mathbf{x})^2]$$

$$\text{Structural error: } E_{x,y} \left[(y - \mathbf{w}^{*T} \mathbf{x})^2 \right]$$

where $\mathbf{w}^* = (w_0^*, \dots, w_d^*)$ are the optimal linear regression parameters (infinite training data)

Linear regression: types of errors

- ▶ Approximation error measures how close we can get to the optimal linear predictions with limited training data:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} E_{\mathbf{x},y}[(y - \mathbf{w}^T \mathbf{x})^2]$$

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

$$\text{Approximation error: } E_{\mathbf{x}} \left[(\mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x})^2 \right]$$

Where $\hat{\mathbf{w}}$ are the parameter estimates based on a small training set (so themselves are random variables).

Linear regression: error decomposition

- ▶ The expected error can decompose into the sum of structural and approximation errors

$$\begin{aligned} E_{x,y}[(y - \hat{\mathbf{w}}^T \mathbf{x})^2] \\ = E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})^2] + E_x[(\mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x})^2] \end{aligned}$$

- ▶ Derivation

$$\begin{aligned} E_{x,y}[(y - \hat{\mathbf{w}}^T \mathbf{x})^2] &= E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x} + \mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x})^2] \\ &= E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})^2] + E_x[(\mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x})^2] \\ &\quad + 2E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})(\mathbf{w}^{*T} \mathbf{x} - \hat{\mathbf{w}}^T \mathbf{x})] \end{aligned}$$

Linear regression: error decomposition

- ▶ The expected error can decompose into the sum of structural and approximation errors

$$\begin{aligned} E_{x,y}[(y - \widehat{\mathbf{w}}^T \mathbf{x})^2] \\ = E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})^2] + E_x[(\mathbf{w}^{*T} \mathbf{x} - \widehat{\mathbf{w}}^T \mathbf{x})^2] \end{aligned}$$

- ▶ Derivation

$$\begin{aligned} E_{x,y}[(y - \widehat{\mathbf{w}}^T \mathbf{x})^2] &= E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x} + \mathbf{w}^{*T} \mathbf{x} - \widehat{\mathbf{w}}^T \mathbf{x})^2] \\ &= E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})^2] + E_x[(\mathbf{w}^{*T} \mathbf{x} - \widehat{\mathbf{w}}^T \mathbf{x})^2] \\ &+ 0 \end{aligned}$$

Note: Optimality condition for \mathbf{w}^* give us $E_{x,y}[(y - \mathbf{w}^{*T} \mathbf{x})\mathbf{x}] = 0$
since $\nabla_{\mathbf{w}} E_{x,y}[(y - \mathbf{w}^T \mathbf{x})^2]|_{\mathbf{w}^*} = 0$