



Linear regression, Gradient Descent

Iran University of Science and Technology

By: M. S. Tahaei, PhD.

Winter 2025

Outline

- Introduction to Learning
- Linear Regression
- Gradient Descent
- Generalized Linear Regression

A Definition of ML

- ▶ Tom Mitchell (1998): Well-posed learning problem
 - ▶ “A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E”.
- ▶ Using the observed data to make better decisions
 - ▶ Generalizing from the observed data

ML Definition: Example

- ▶ Consider an email program that learns how to filter spam according to emails you do or do not mark as spam.
 - ▶ T: Classifying emails as spam or not spam.
 - ▶ E: Watching you label emails as spam or not spam.
 - ▶ P: The number (or fraction) of emails correctly classified as spam/not spam.

The essence of machine learning

- A pattern exist
- We do not know it mathematically
- We have data on it

Example: Home Price

- ▶ Housing price prediction

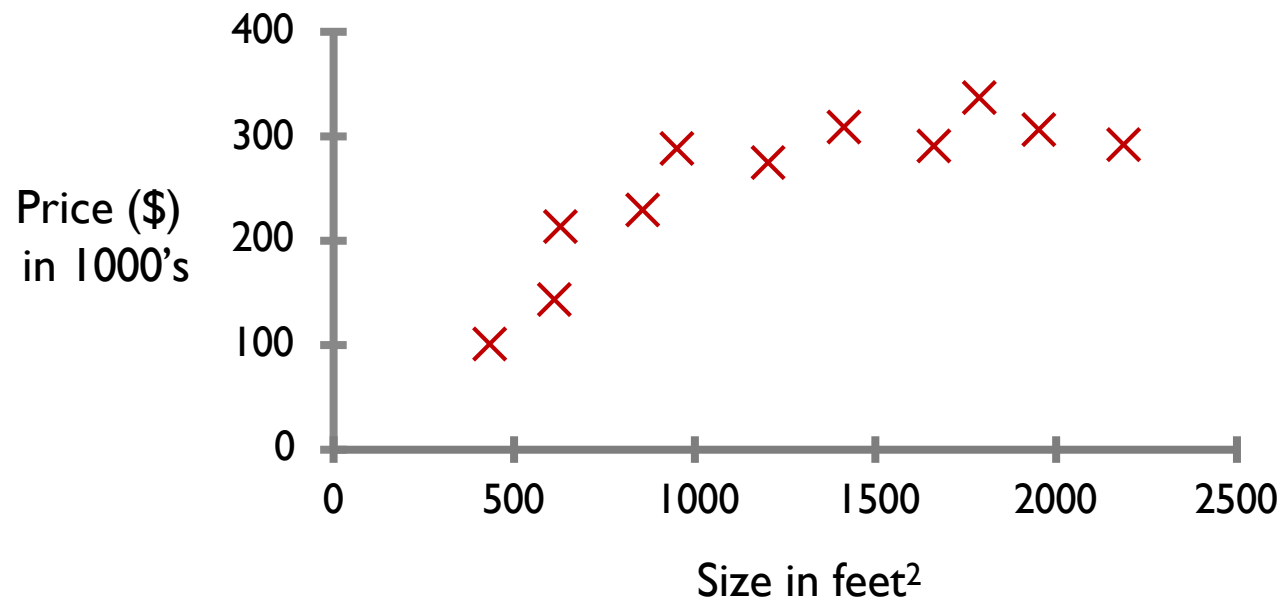


Figure adopted from slides of Andrew Ng,
Machine Learning course, Stanford.

Example: Bank loan

- ▶ Applicant form as the input:

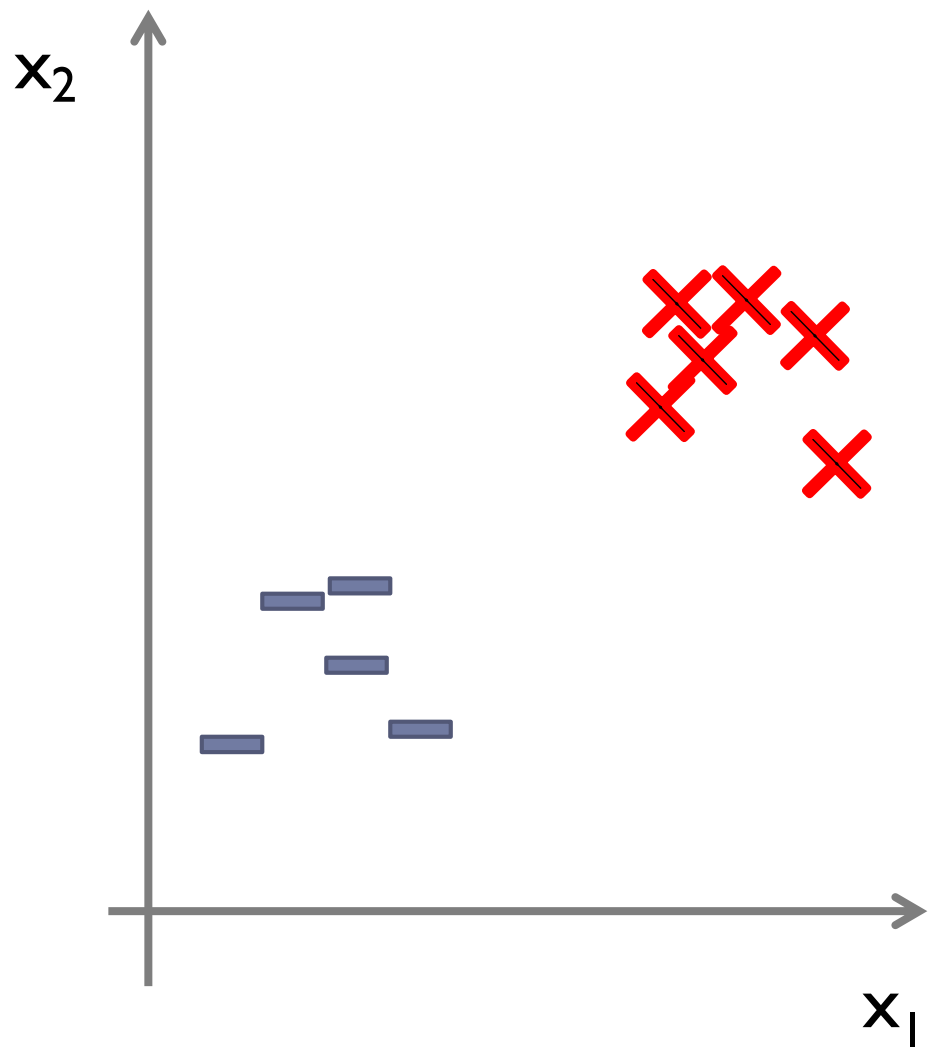
age	23 years
gender	male
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

- ▶ Output: approving or denying the request

Components of (Supervised) Learning

- Unknown target function: $f: \mathcal{X} \rightarrow \mathcal{Y}$
 - Input space: \mathcal{X}
 - Output space: \mathcal{Y}
- Training data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- Pick a formula $g: \mathcal{X} \rightarrow \mathcal{Y}$ that approximates the target function f
 - selected from a set of hypotheses \mathcal{H}

Training data: Example



Training data

x_1	x_2	y	
0.9	2.3	1	—
3.5	2.6	1	—
2.6	3.3	1	—
2.7	4.1	1	—
1.8	3.9	1	—
6.5	6.8	-1	×
7.2	7.5	-1	×
7.9	8.3	-1	×
6.9	8.3	-1	×
8.8	7.9	-1	×
9.1	6.2	-1	×

Solution Components

- ▶ **Learning model** composed of:
 - ▶ Learning algorithm
 - ▶ Hypothesis set
- ▶ Perceptron example

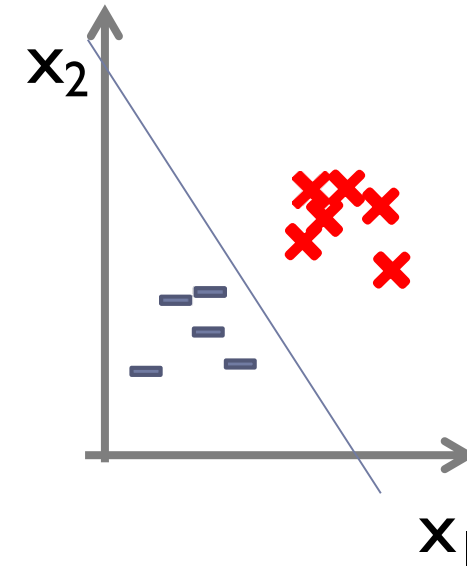
Perceptron classifier

- ▶ Input $\mathbf{x} = [x_1, \dots, x_d]$
- ▶ Classifier:
 - ▶ If $\sum_{i=1}^d w_i x_i > \text{threshold}$ then output 1
 - ▶ else output -1
- ▶ The linear formula $g \in \mathcal{H}$ can be written:

$$g(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^d \mathbf{w}_i x_i + \mathbf{w}_0 \right)$$

If we add a coordinate $x_0 = 1$ to the input:

$$g(\mathbf{x}) = \text{sign} \left(\sum_{i=0}^d \mathbf{w}_i x_i \right) \quad \xrightarrow{\text{Vector form}} \quad g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$



Perceptron learning algorithm: linearly separable data

- ▶ Give the training data $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})$
- ▶ **Misclassified** data $(\mathbf{x}^{(n)}, y^{(n)})$:
 $\text{sign}(\mathbf{w}^T \mathbf{x}^{(n)}) \neq y^{(n)}$

Repeat

Pick a **misclassified** data $(\mathbf{x}^{(n)}, y^{(n)})$ from training data and update \mathbf{w} :

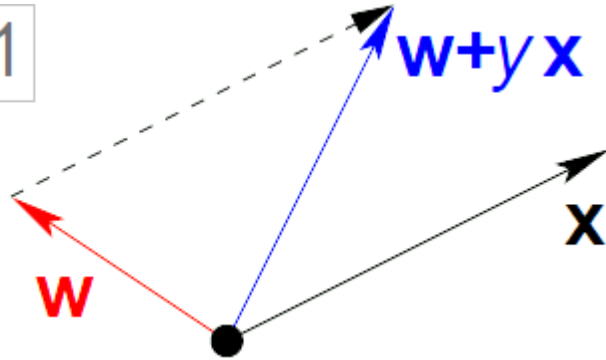
$$\mathbf{w} = \mathbf{w} + y^{(n)} \mathbf{x}^{(n)}$$

Until all training data points are correctly classified by g

Perceptron learning algorithm: Example of weight update

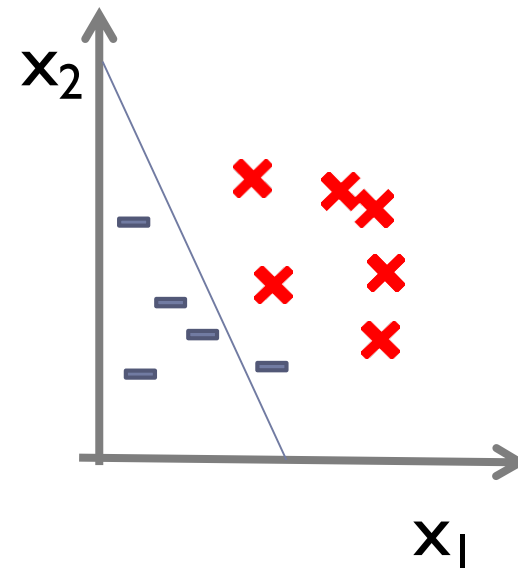
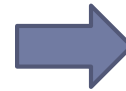
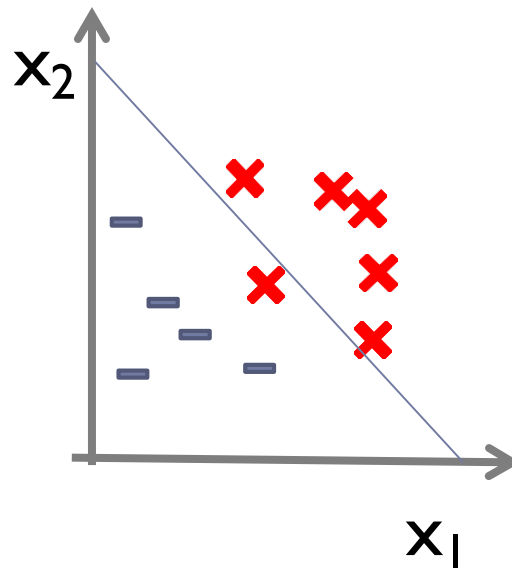
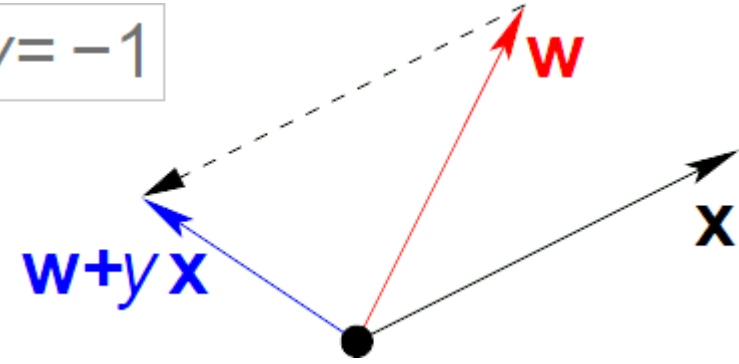
Correct Label

$$y = +1$$



Correct Label

$$y = -1$$



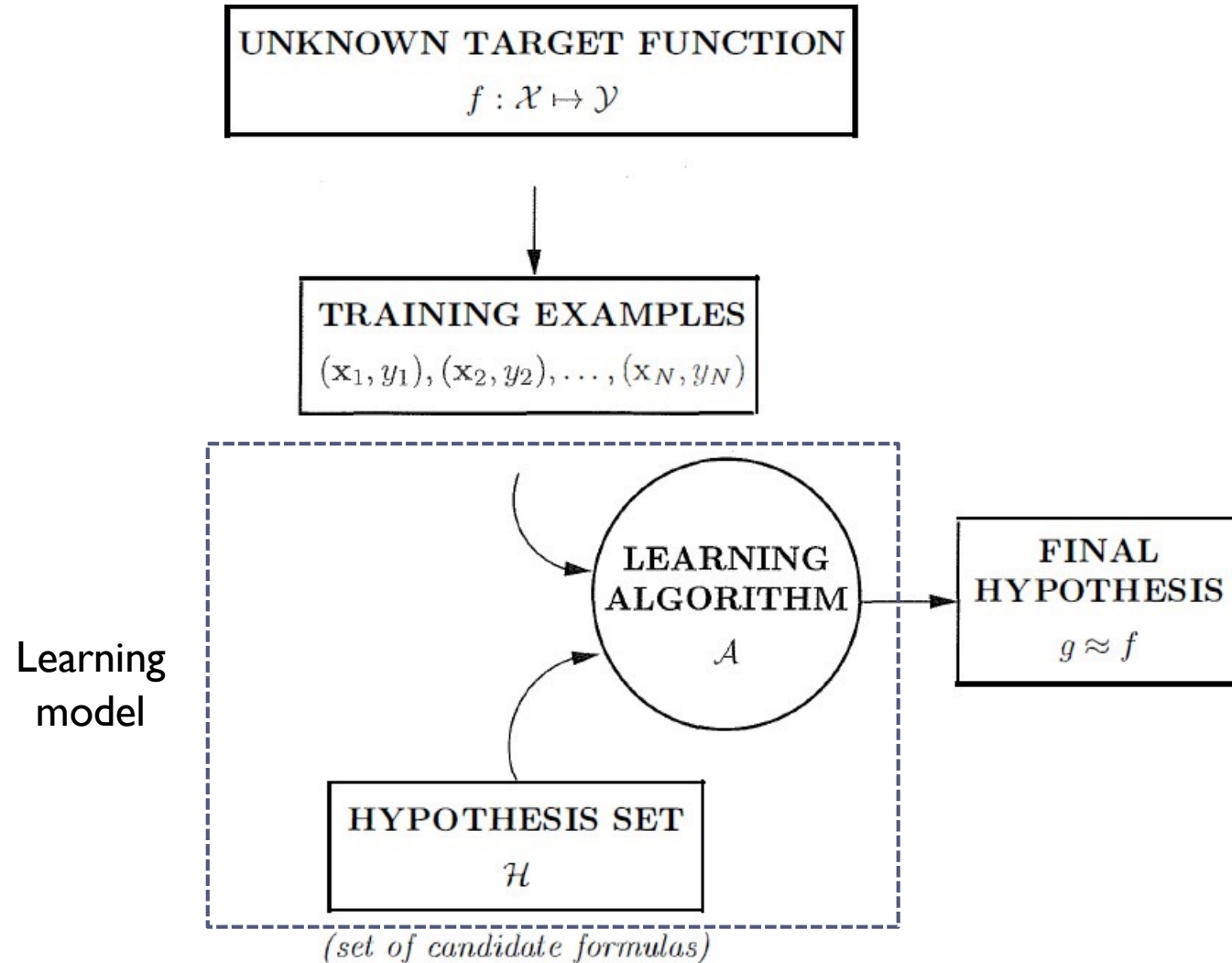
Experience (E) in ML

- ▶ Basic premise of learning:
 - ▶ “Using a set of observations to uncover an underlying process”
- ▶ We have different types of (getting) observations in different types or paradigms of ML methods

Main Steps of Learning Tasks

- ▶ Selection of hypothesis set (or model specification)
 - ▶ Which class of models (mappings) should we use for our data?
- ▶ Learning: find mapping f (from hypothesis set) based on the training data
 - ▶ Which notion of error should we use? (loss functions)
 - ▶ Optimization of loss function to find mapping f
- ▶ Evaluation: how well f generalizes to yet unseen examples
 - ▶ How do we ensure that the error on future data is minimized? (generalization)

Components of (Supervised) Learning



Linear regression, Cost Function and Gradient Descent

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
...

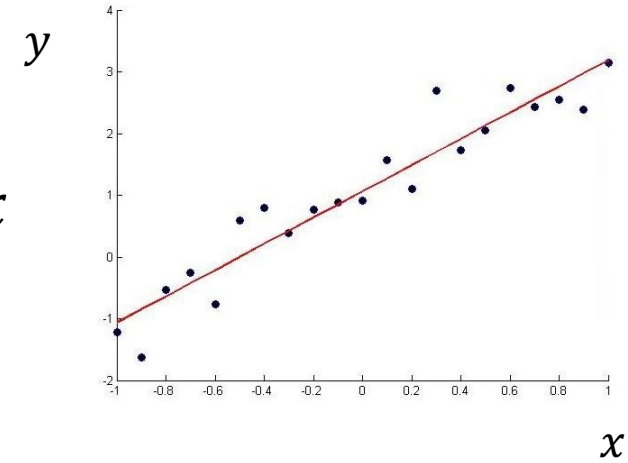
Learning problem

- ▶ 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 f generalizes to unseen examples

Linear regression: hypothesis space

- ▶ Univariate

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



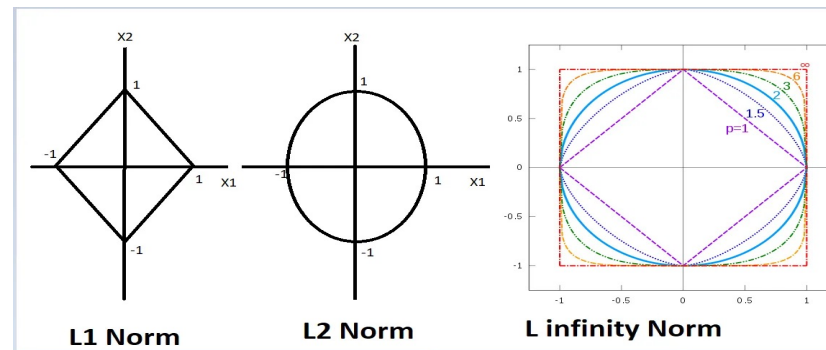
- ▶ Multivariate

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

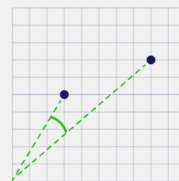
$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$ are parameters we need to set.

Learning algorithm and distance metrics

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

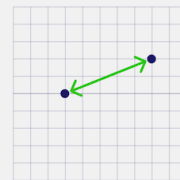


Distance Metrics in Vector Search



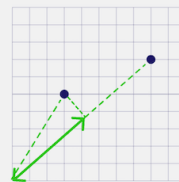
Cosine Distance

$$1 - \frac{A \cdot B}{||A|| ||B||}$$



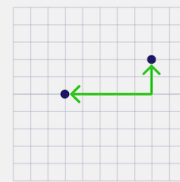
Squared Euclidean
(L2 Squared)

$$\sum_{i=1}^n (x_i - y_i)^2$$



Dot Product

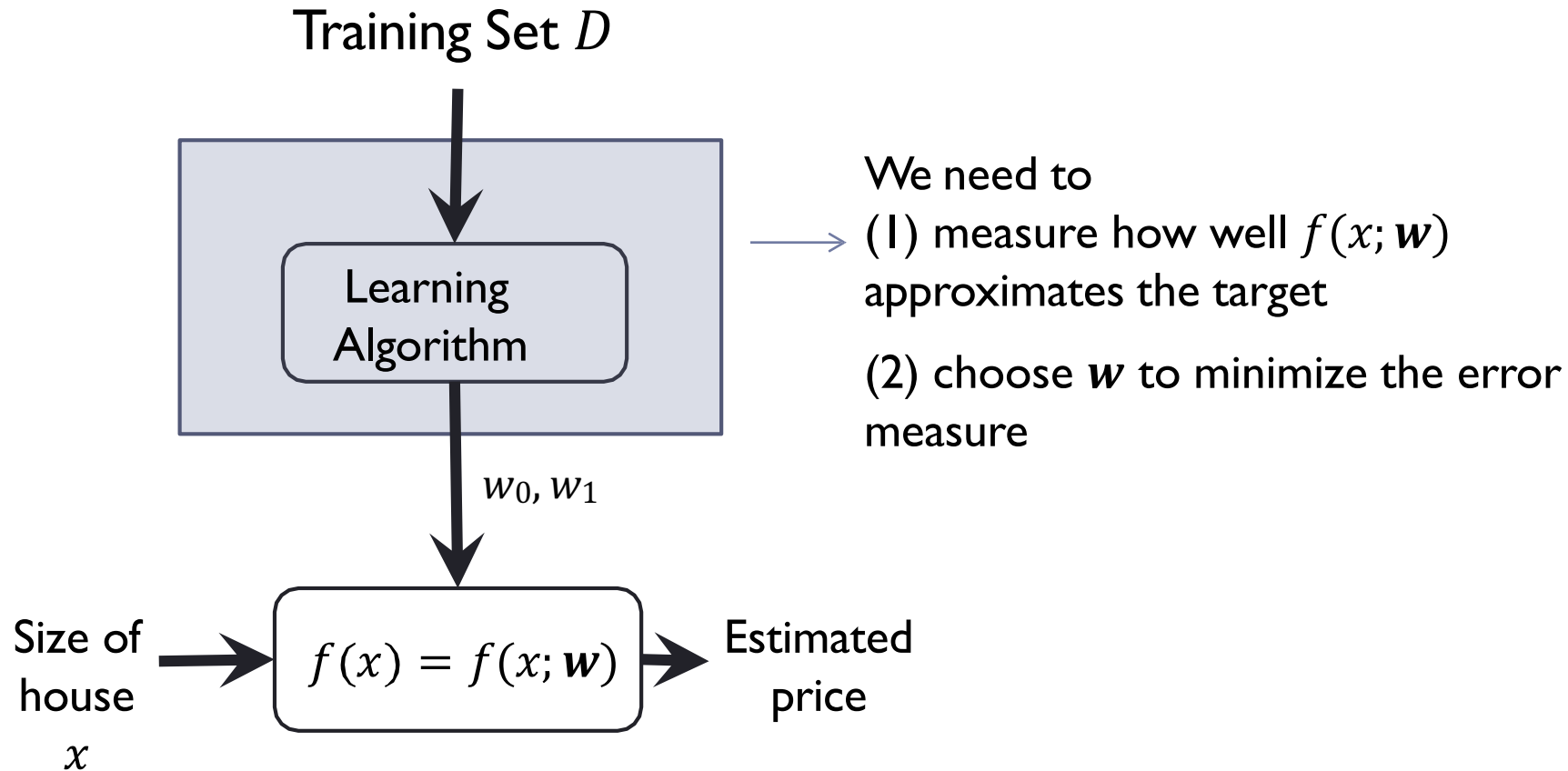
$$A \cdot B = \sum_{i=1}^n A_i B_i$$



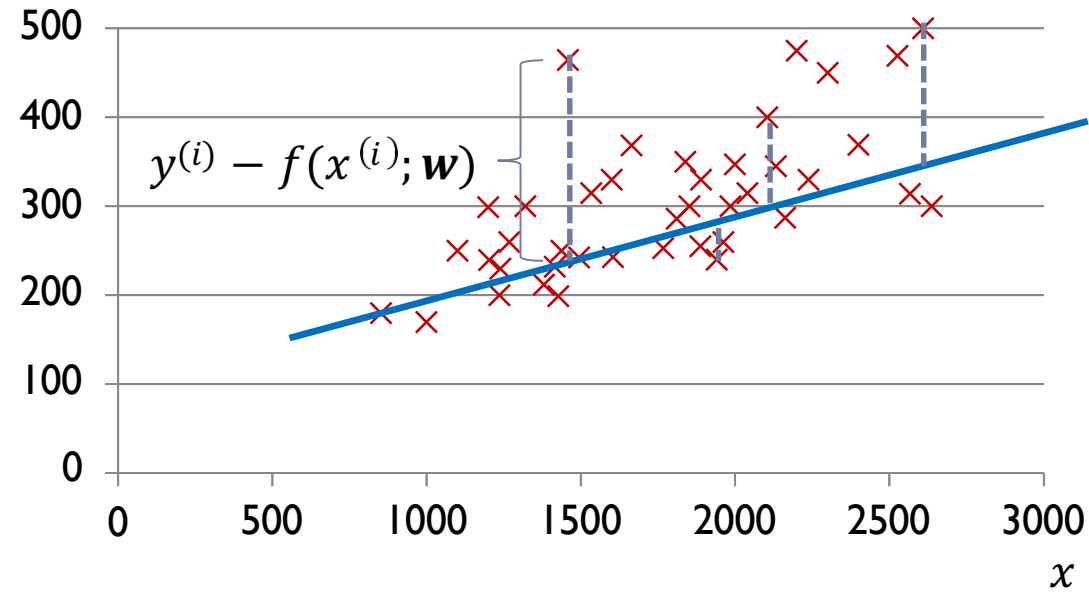
Manhattan (L1)

$$\sum_{i=1}^n |x_i - y_i|$$

Learning algorithm

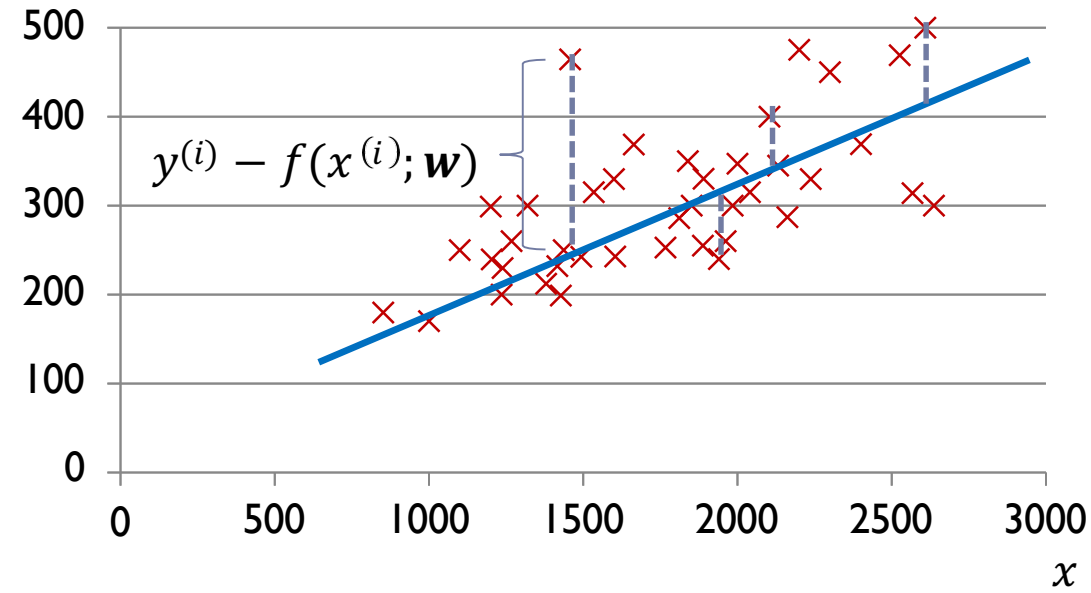


How to measure the error



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:

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

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

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

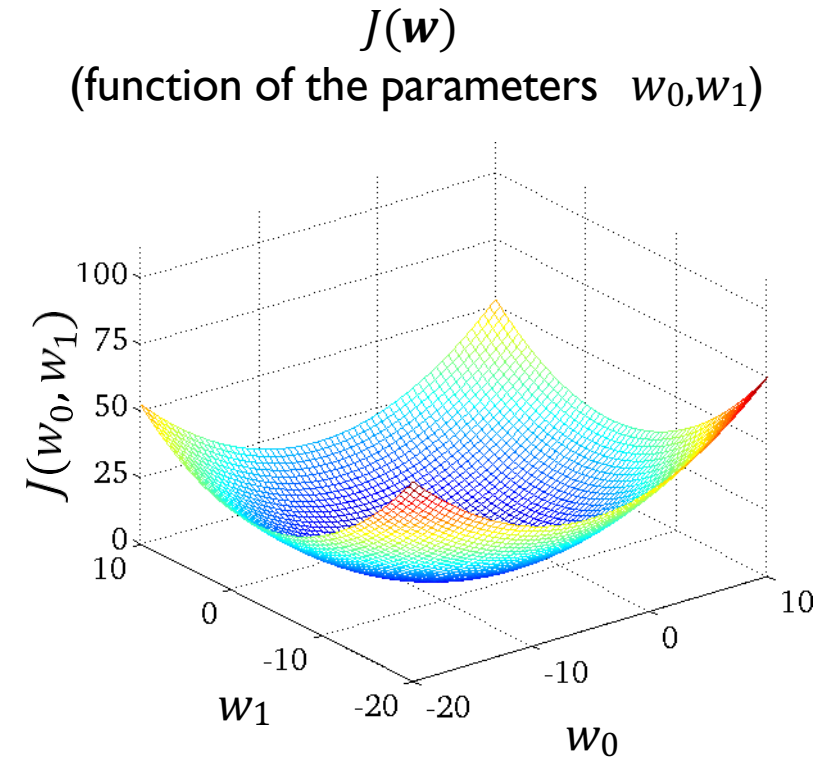
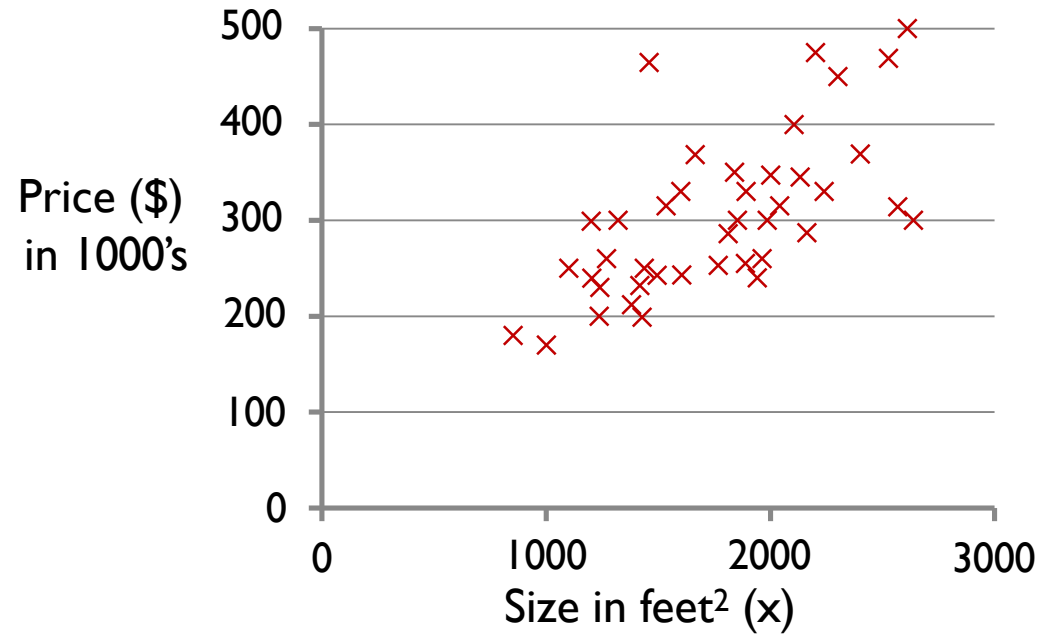
- ▶ Minimizing sum (or mean) of squared errors is a common approach in curve fitting, neural network, etc.

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 $f(\mathbf{x}) = f(\mathbf{x}; \mathbf{w}^*)$ where $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$

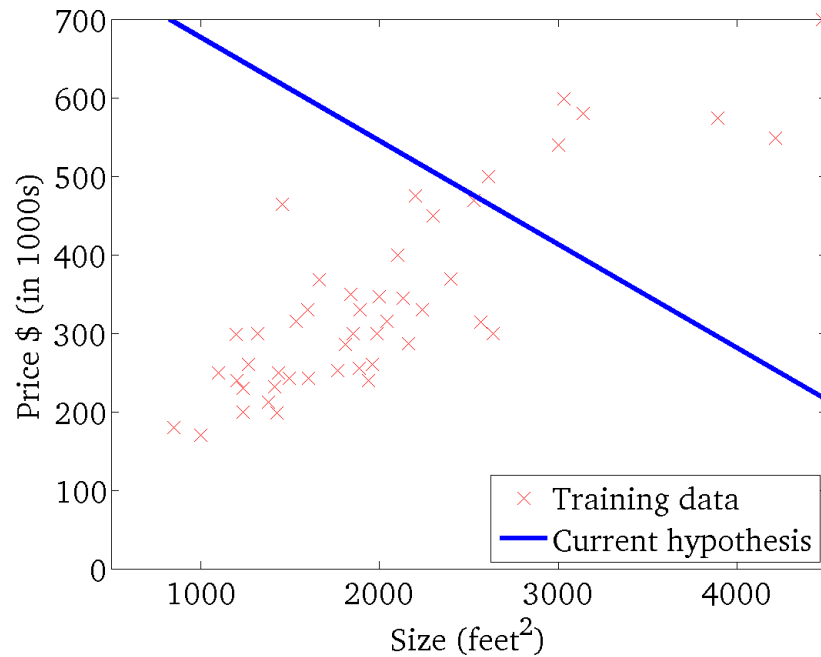
Cost function: univariate example



Cost function: univariate example

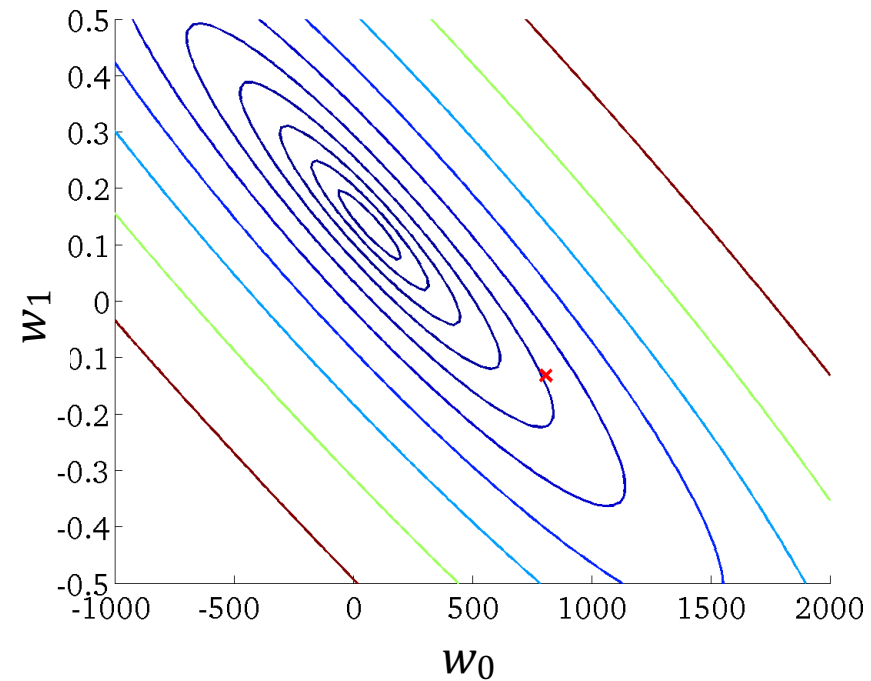
$$f(x; w_0, w_1) = w_0 + w_1 x$$

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



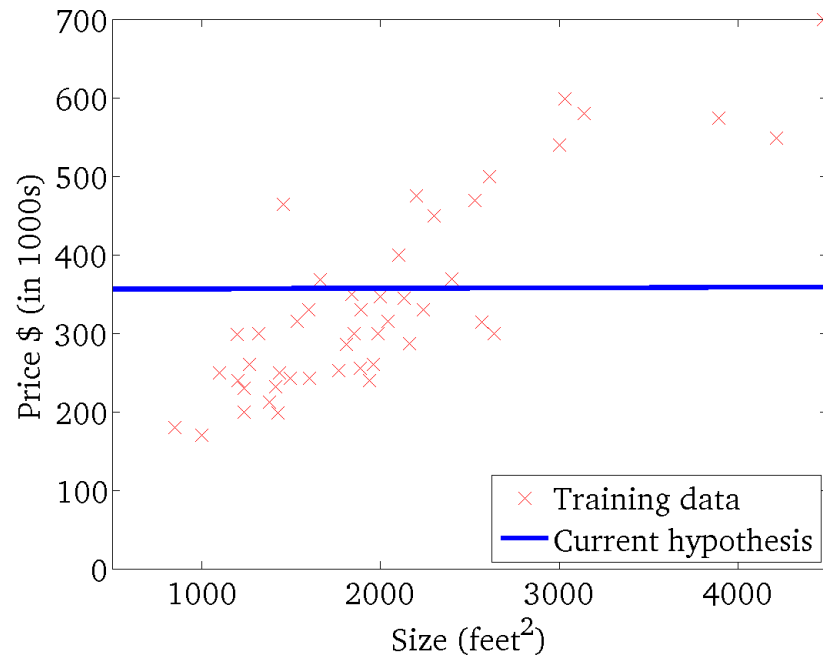
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



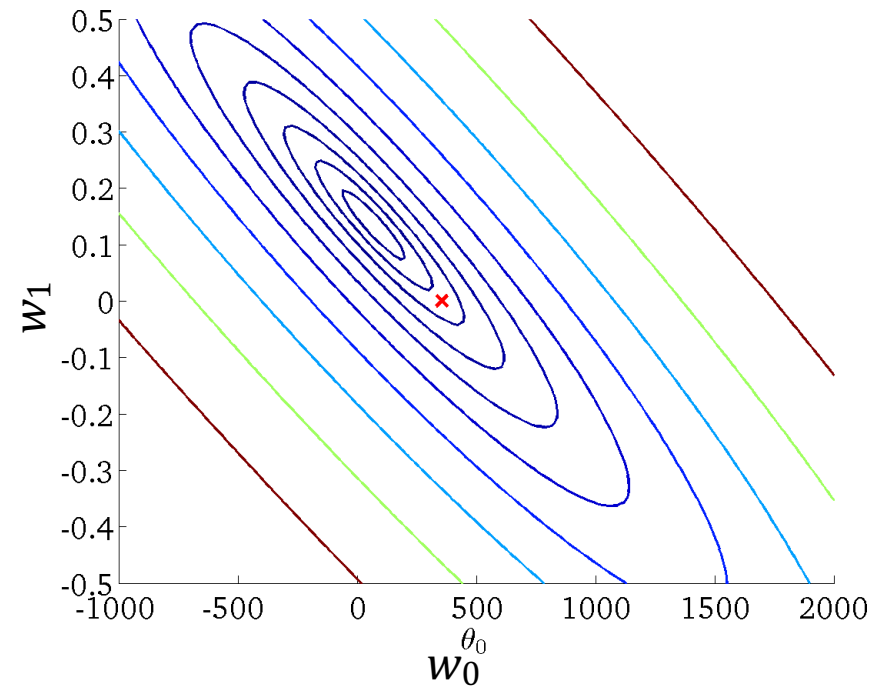
Cost function: univariate example

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



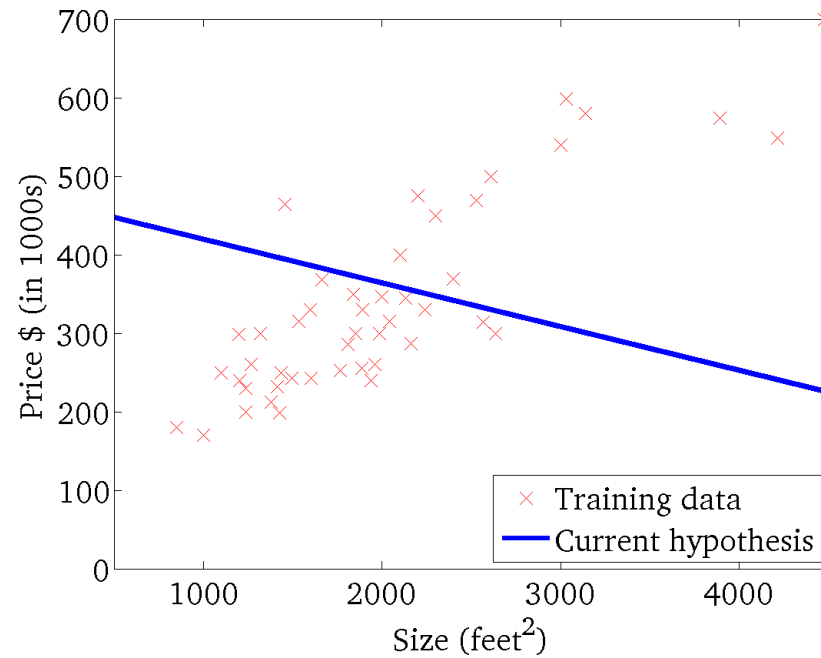
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



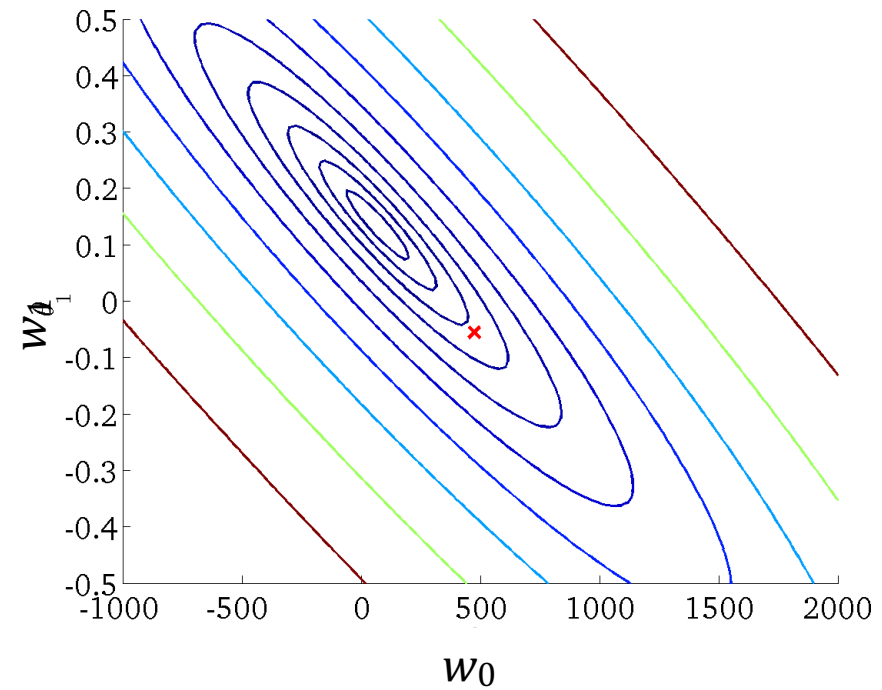
Cost function: univariate example

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



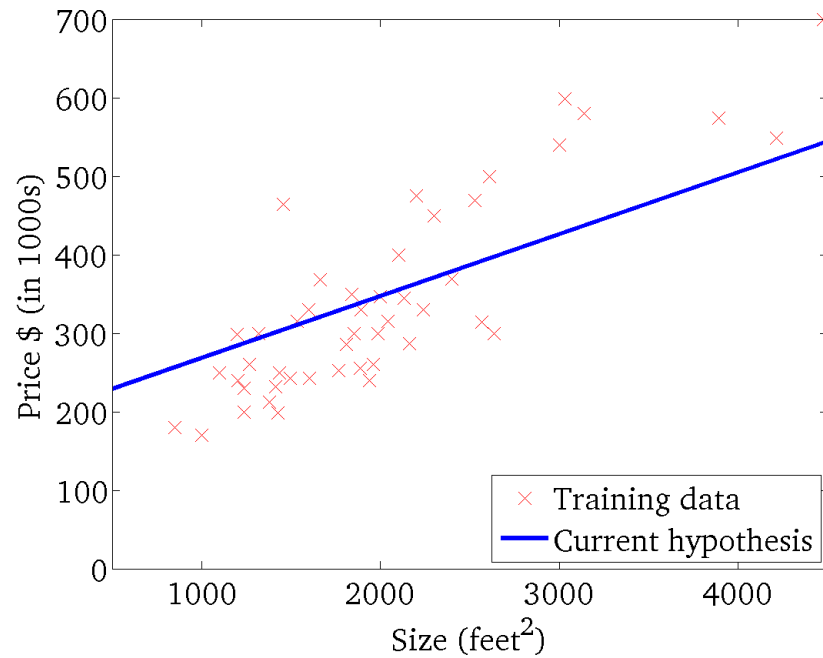
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



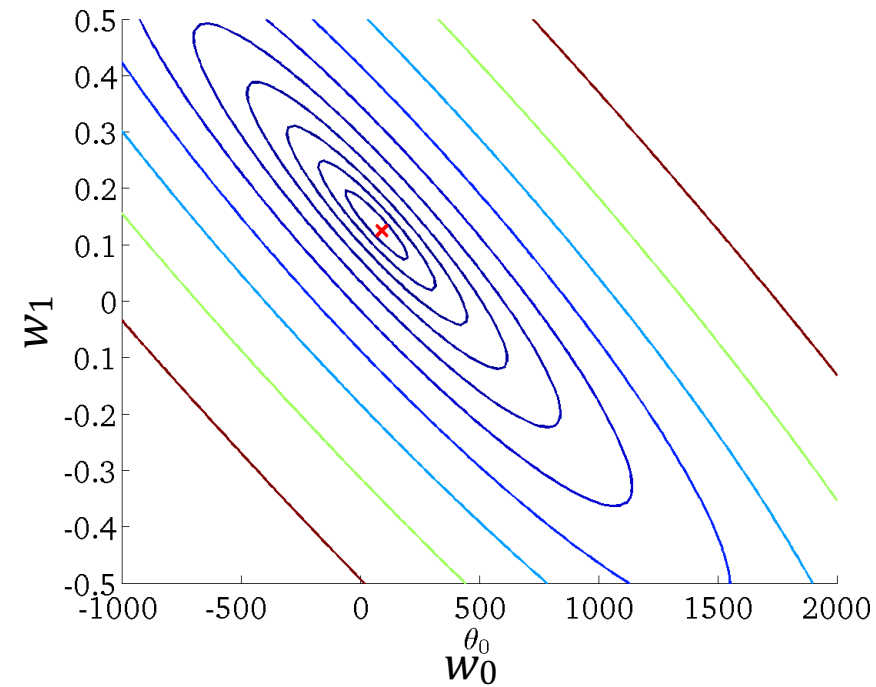
Cost function: univariate example

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



$$J(w_0, w_1)$$

(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

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

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

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

w_1 can be obtained by setting the partial derivative of the SSE to 0 and solving for w_1 , ultimately resulting in:

Exercise!



$$w_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

Cost function: multivariate for various features

- ▶ We have to minimize the empirical squared loss:

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

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

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

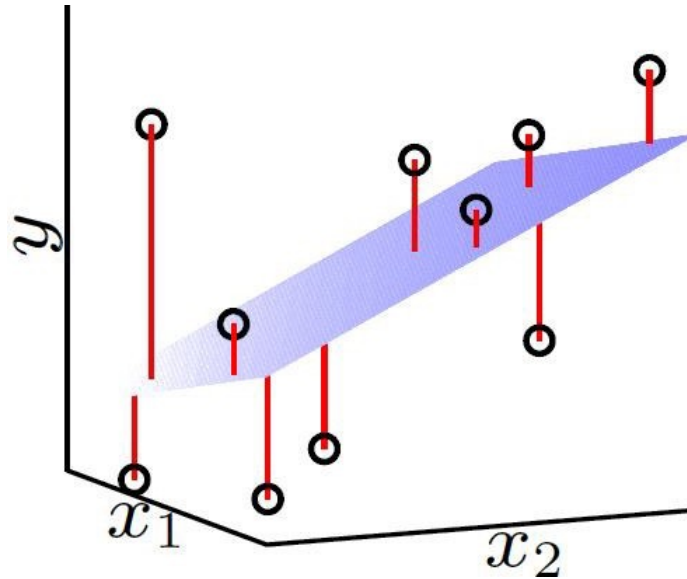
$$\mathbf{w} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^{d+1}} J(\mathbf{w})$$

$$\mathbf{w} \in \mathbb{R}^{d+1}$$

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41.0	880.0	129.0	322.0	126.0	8.3252	452600.0	NEAR BAY
1	-122.22	37.86	21.0	7099.0	1106.0	2401.0	1138.0	8.3014	358500.0	NEAR BAY
2	-122.24	37.85	52.0	1467.0	190.0	496.0	177.0	7.2574	352100.0	NEAR BAY
3	-122.25	37.85	52.0	1274.0	235.0	558.0	219.0	5.6431	341300.0	NEAR BAY
4	-122.25	37.85	52.0	1627.0	280.0	565.0	259.0	3.8462	342200.0	NEAR BAY
...
20635	-121.09	39.48	25.0	1665.0	374.0	845.0	330.0	1.5603	78100.0	INLAND
20636	-121.21	39.49	18.0	697.0	150.0	356.0	114.0	2.5568	77100.0	INLAND
20637	-121.22	39.43	17.0	2254.0	485.0	1007.0	433.0	1.7000	92300.0	INLAND
20638	-121.32	39.43	18.0	1860.0	409.0	741.0	349.0	1.8672	84700.0	INLAND
20639	-121.24	39.37	16.0	2785.0	616.0	1387.0	530.0	2.3886	89400.0	INLAND

20640 rows × 10 columns

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

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

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

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

Minimizing cost function

Optimal linear weight vector (for SSE cost function):

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

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

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

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

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

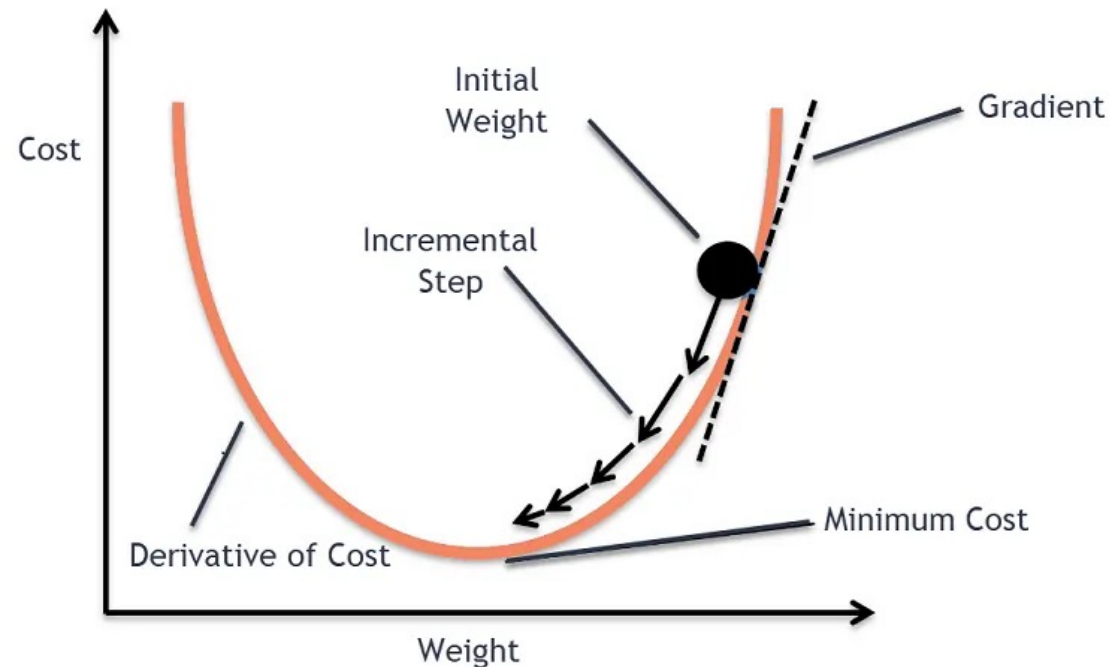
$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

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

Another approach for optimizing the sum squared error

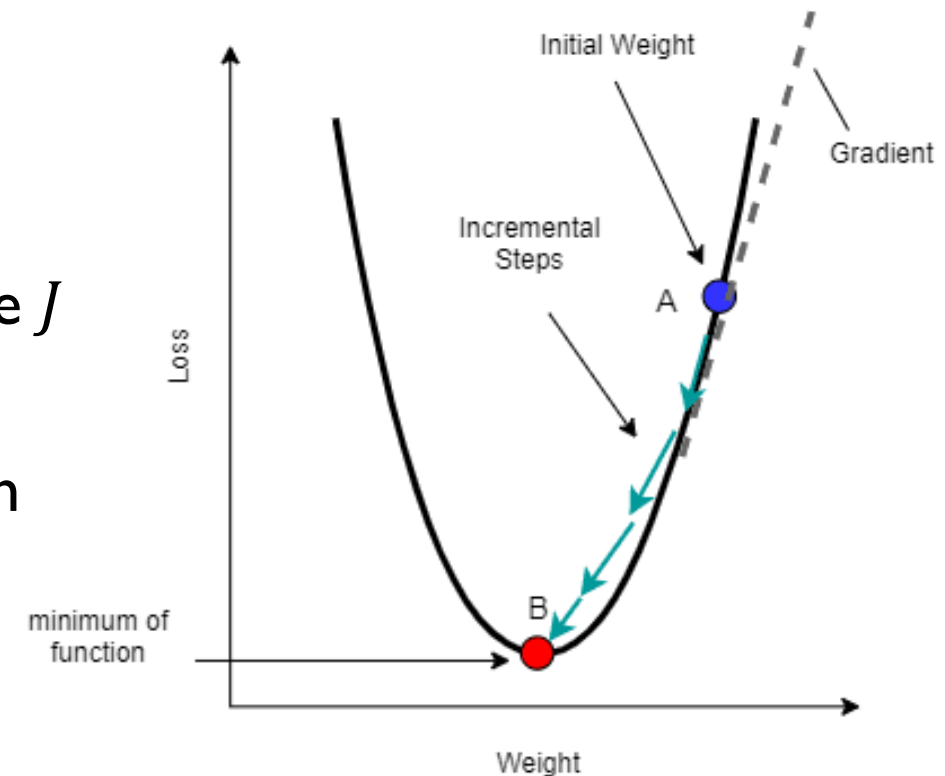
- Iterative approach for solving the following optimization problem:

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



Gradient descent

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



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

Gradient descent

- ▶ Minimize $J(\mathbf{w})$

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

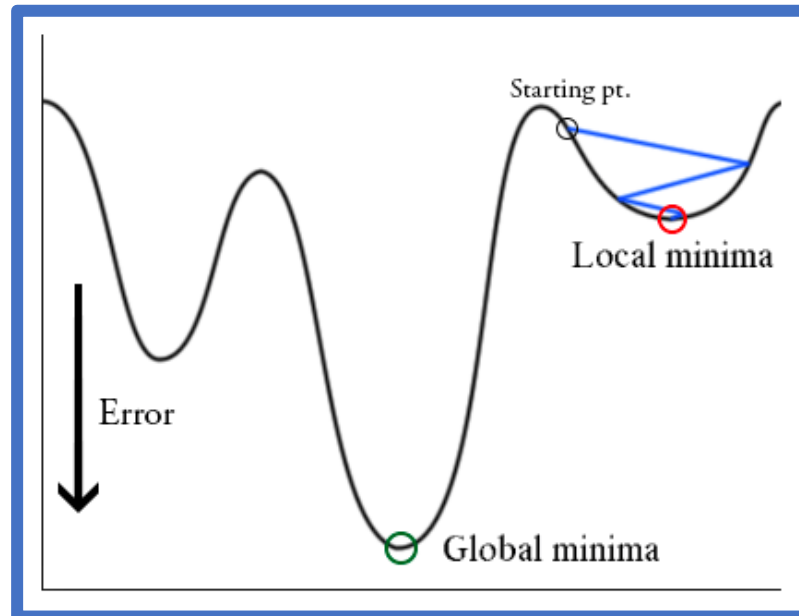
**Step size
(Learning rate
parameter)**

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

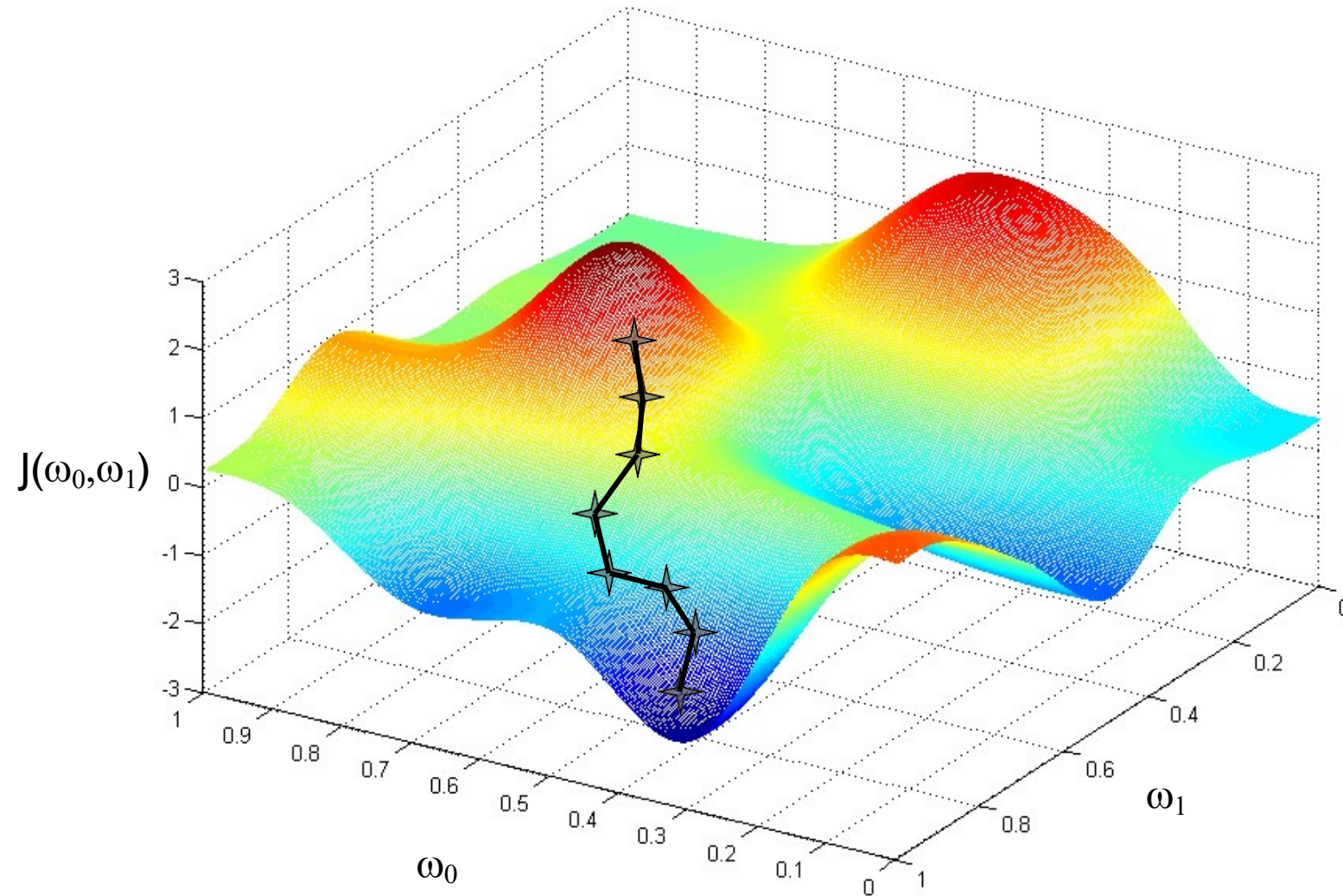
- ▶ If η is small enough, then $J(\mathbf{w}^{t+1}) \leq J(\mathbf{w}^t)$.
- ▶ η can be allowed to change at every iteration as η_t .

Gradient descent

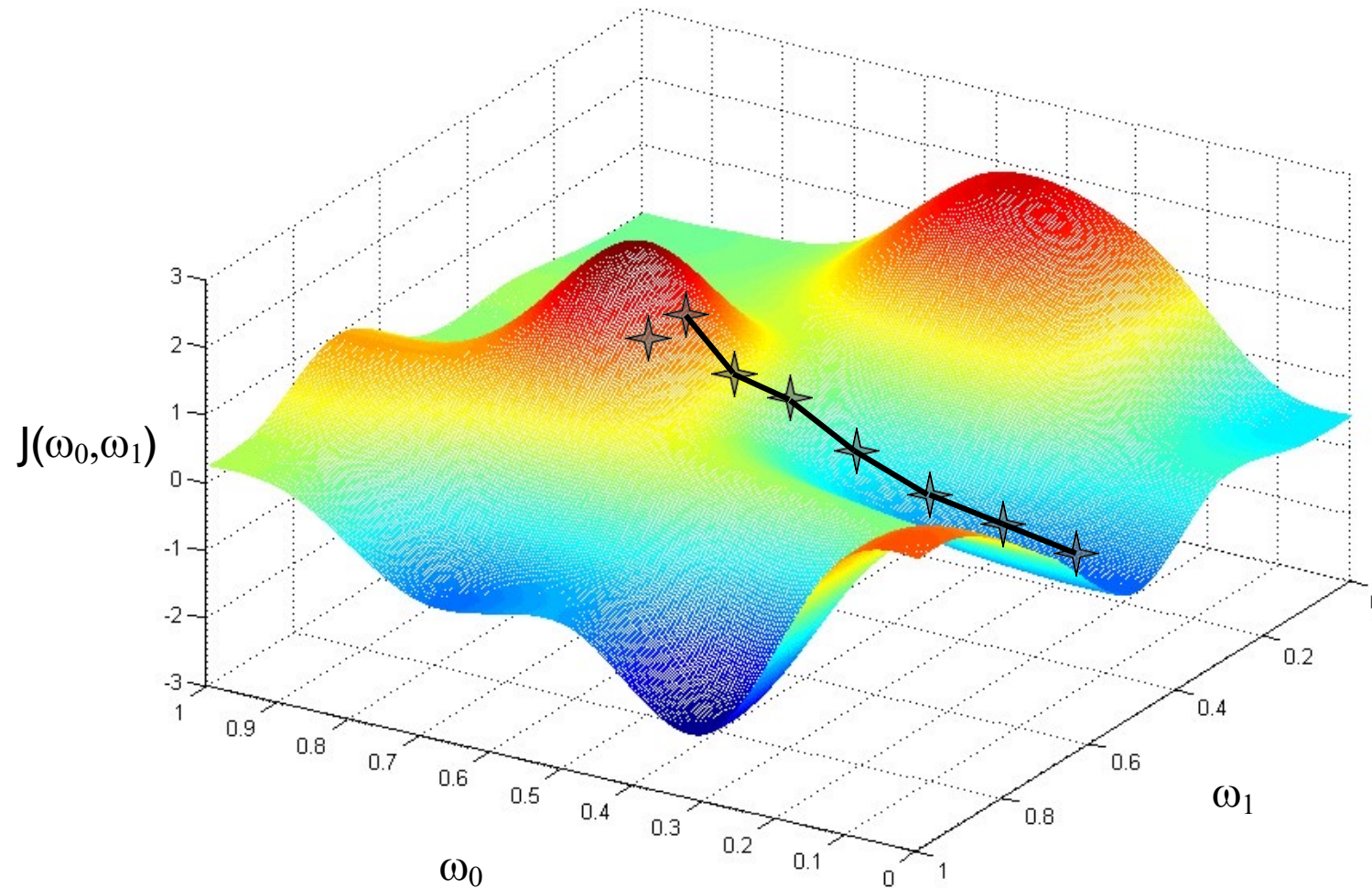
- ▶ Local minima problem
- ▶ However, when J is convex, all local minima are also global minima \Rightarrow gradient descent can converge to the global solution.



Problem of gradient descent with non-convex cost functions



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

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

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

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

Gradient descent for SSE cost function

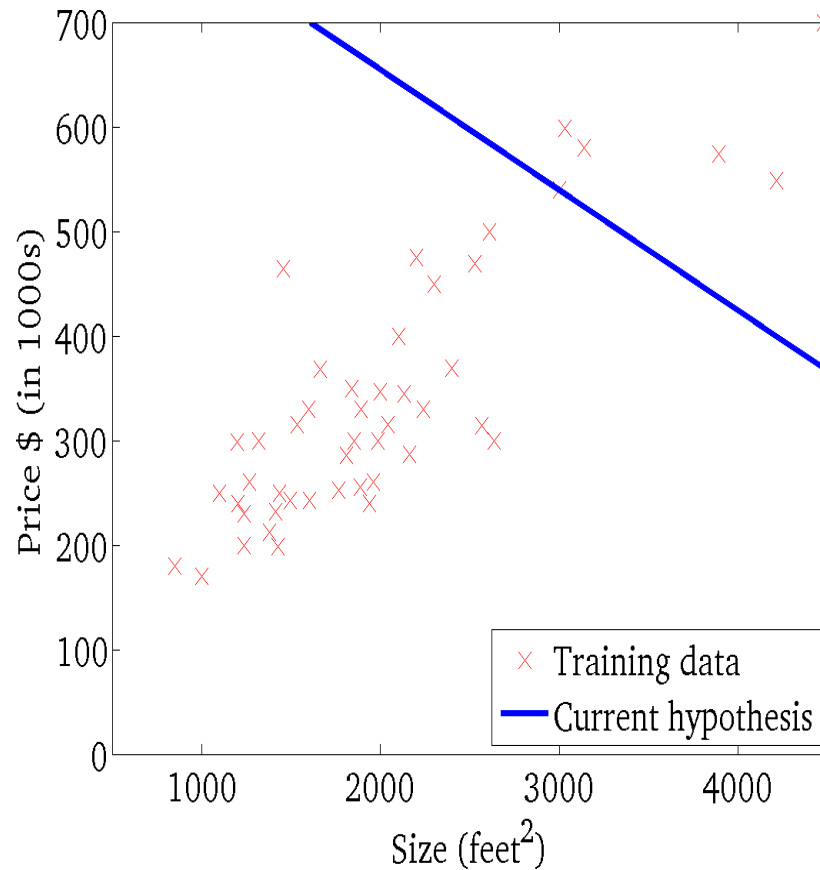
- ▶ Weight update rule: $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$

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

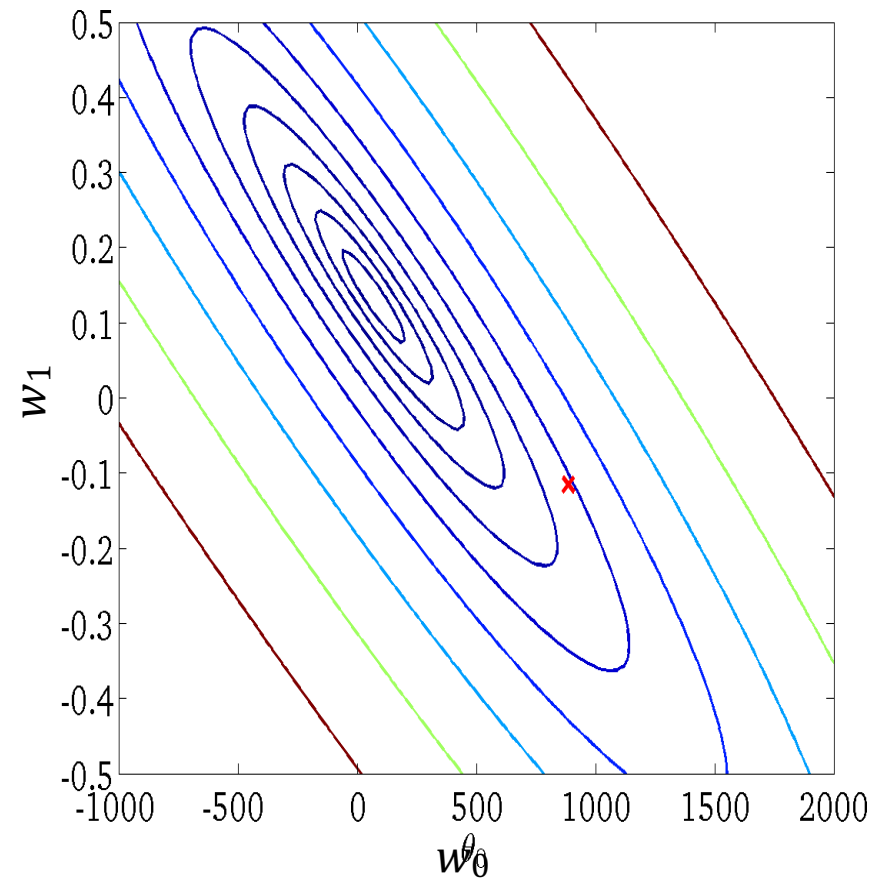
Batch mode: each step
considers all training data

- ▶ 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.

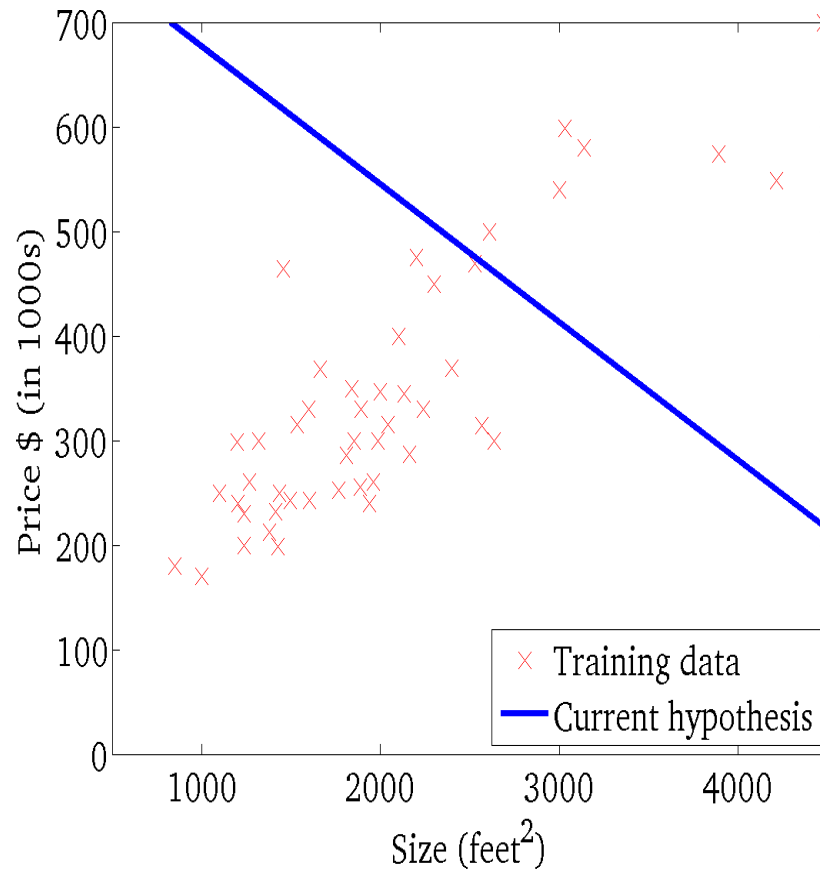
$$f(x; w_0, w_1) = w_0 + w_1 x$$



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

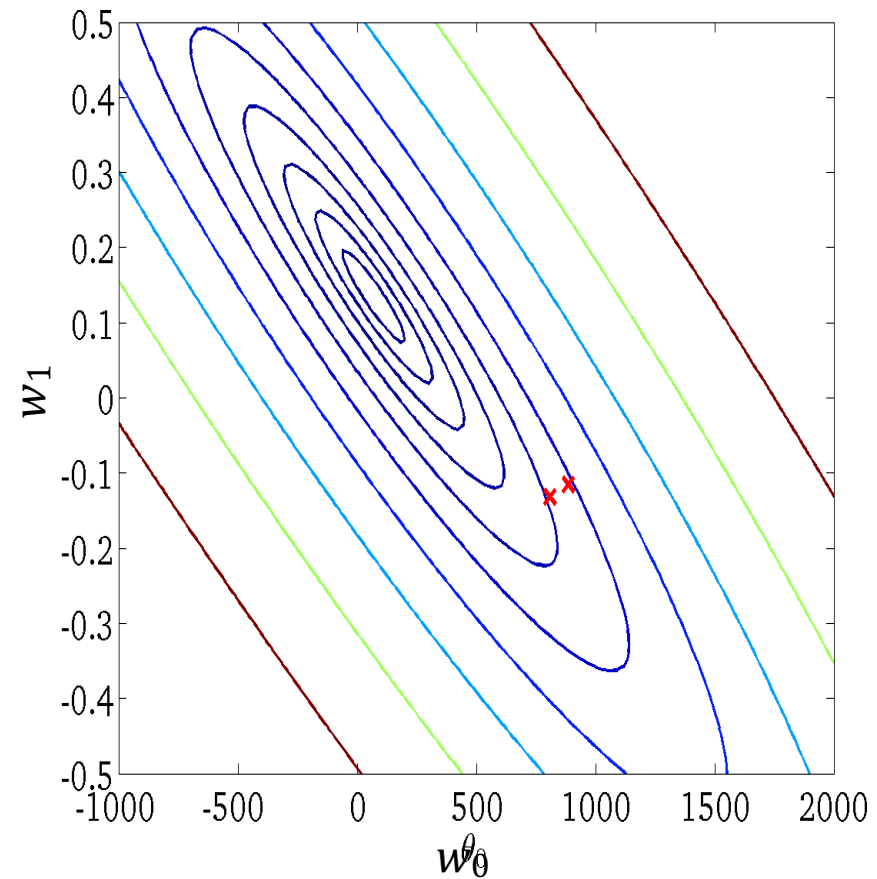


$$f(x; w_0, w_1) = w_0 + w_1 x$$

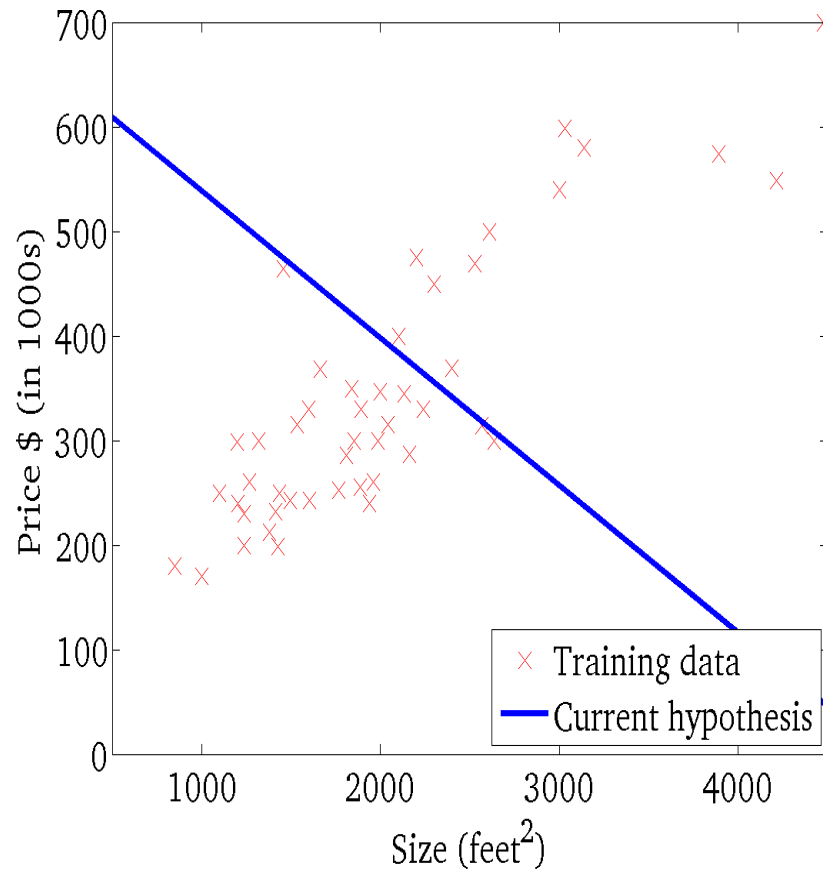


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

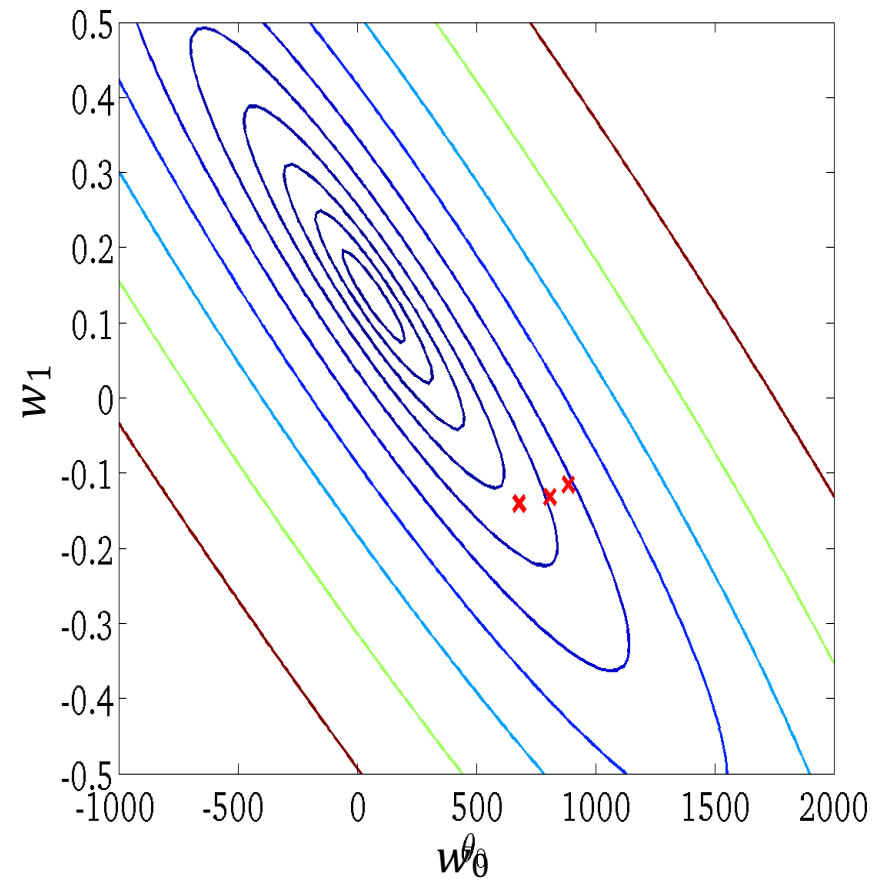


$$f(x; w_0, w_1) = w_0 + w_1 x$$



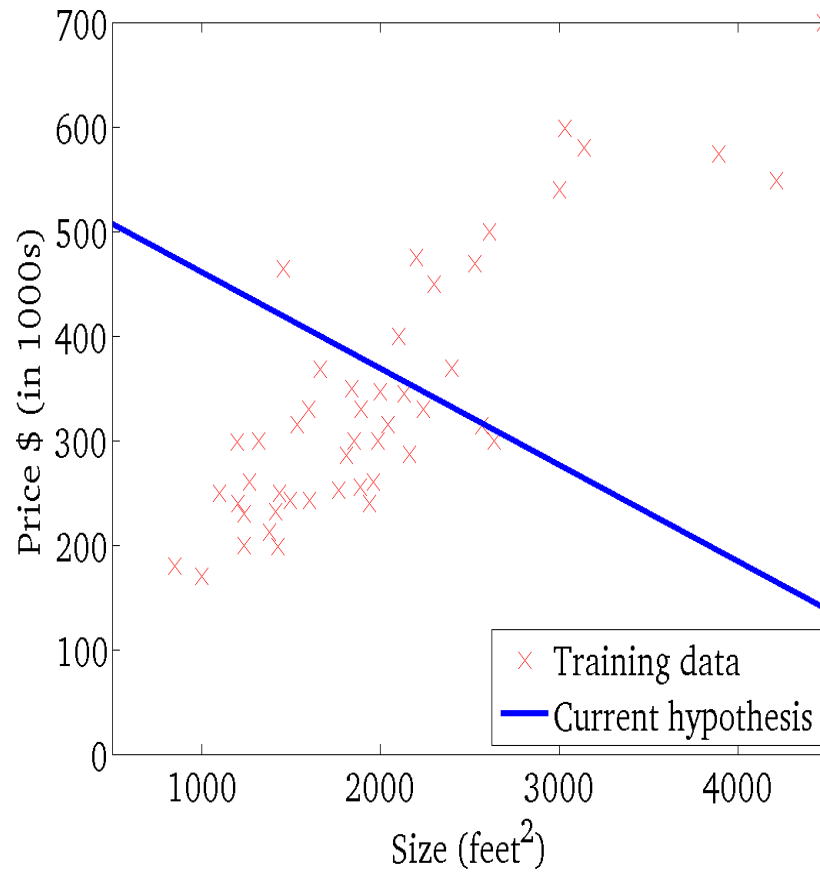
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



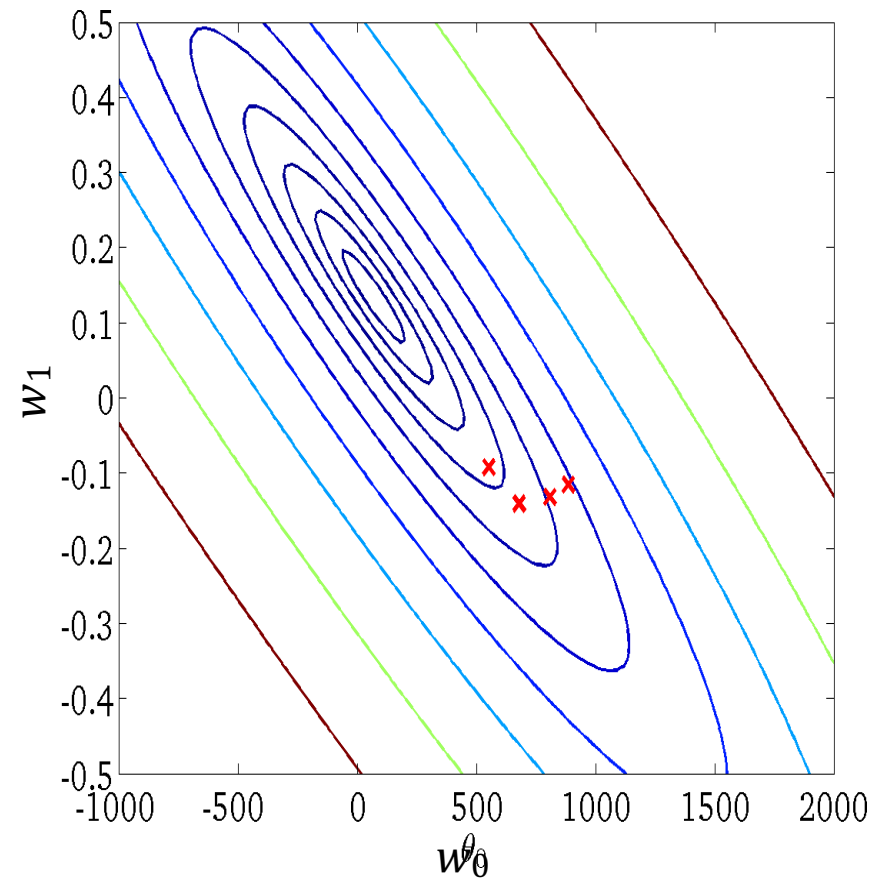
This example has been adopted from: Prof. Ng's slides (ML Online Course, Stanford)

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



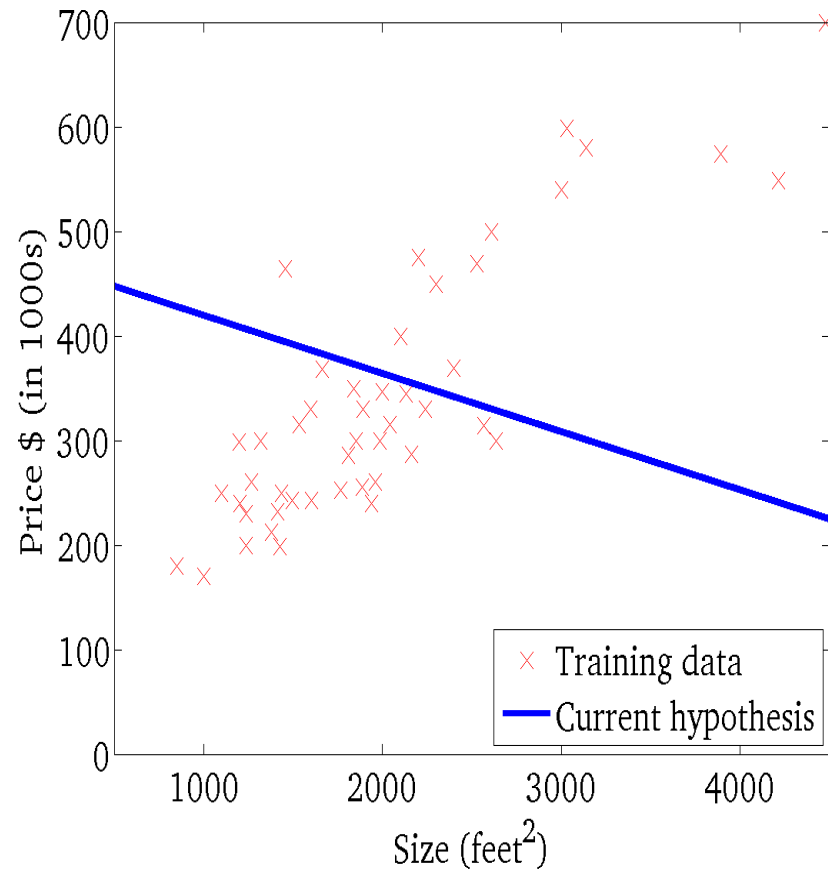
$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



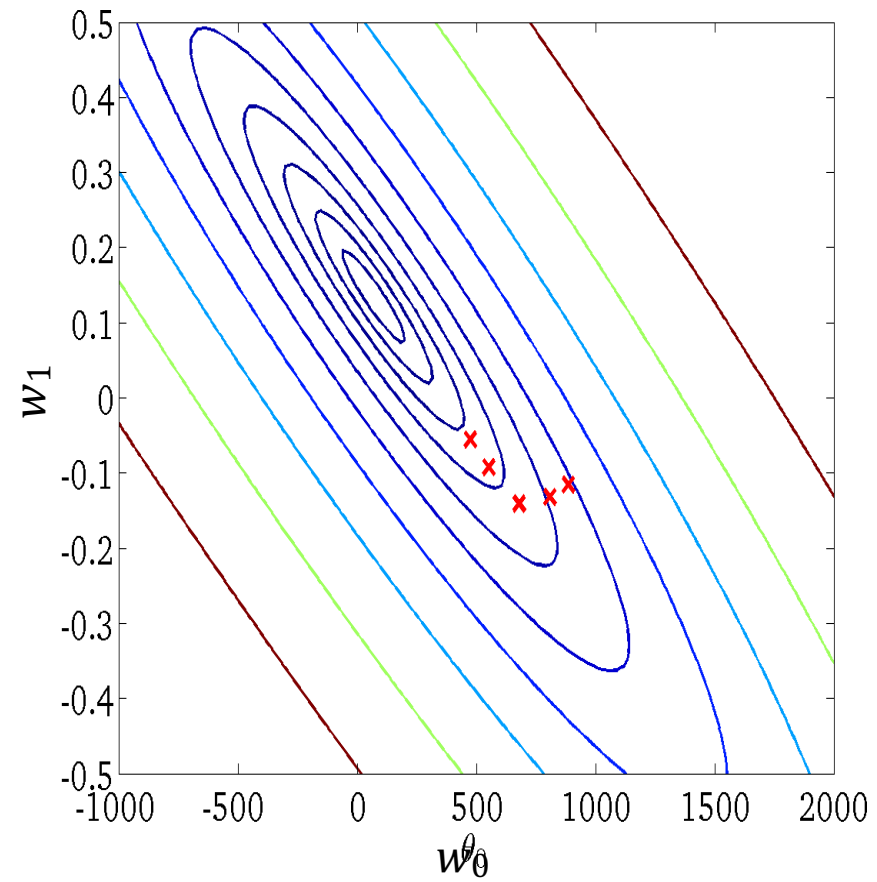
This example has been adopted from: Prof. Ng's slides (ML Online Course, Stanford)

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

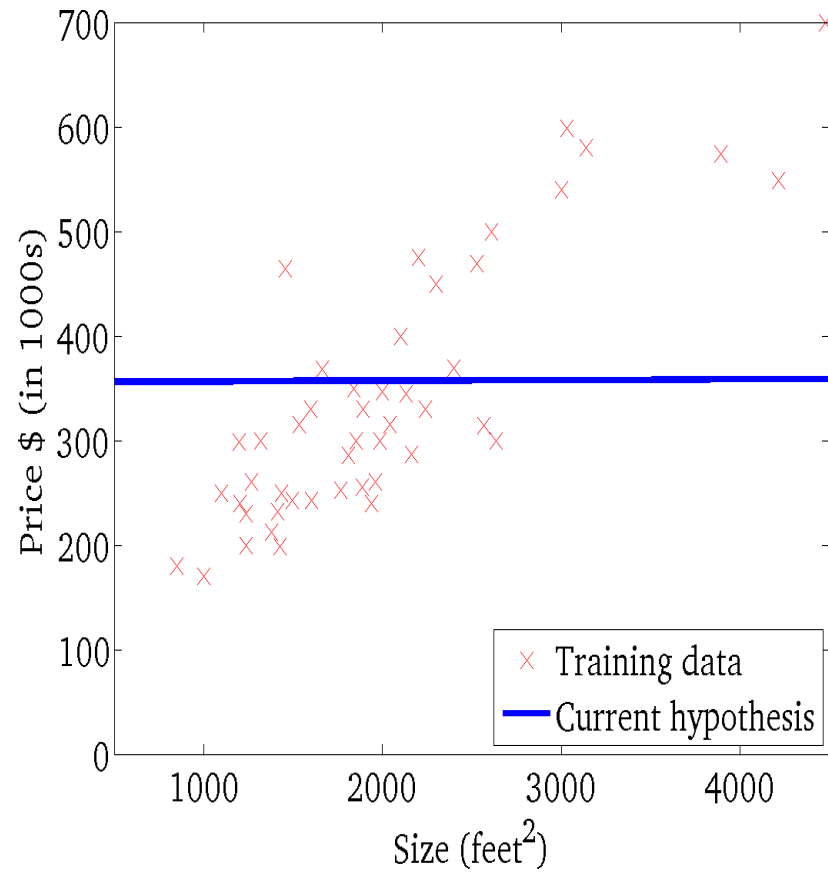


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

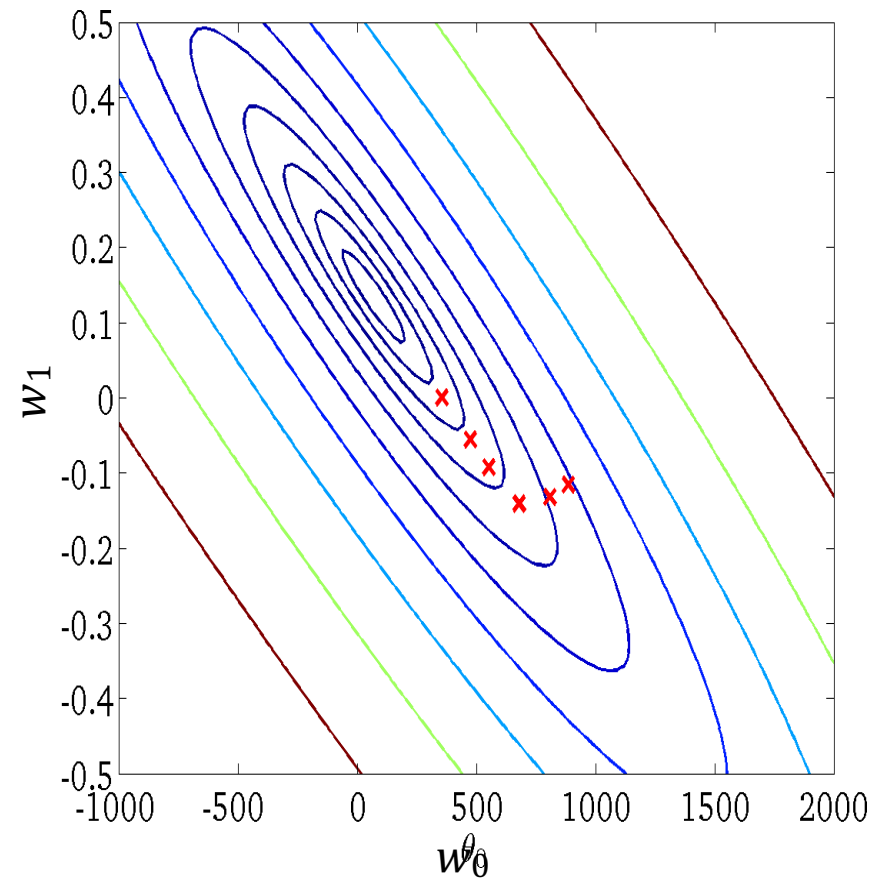


$$f(x; w_0, w_1) = w_0 + w_1 x$$

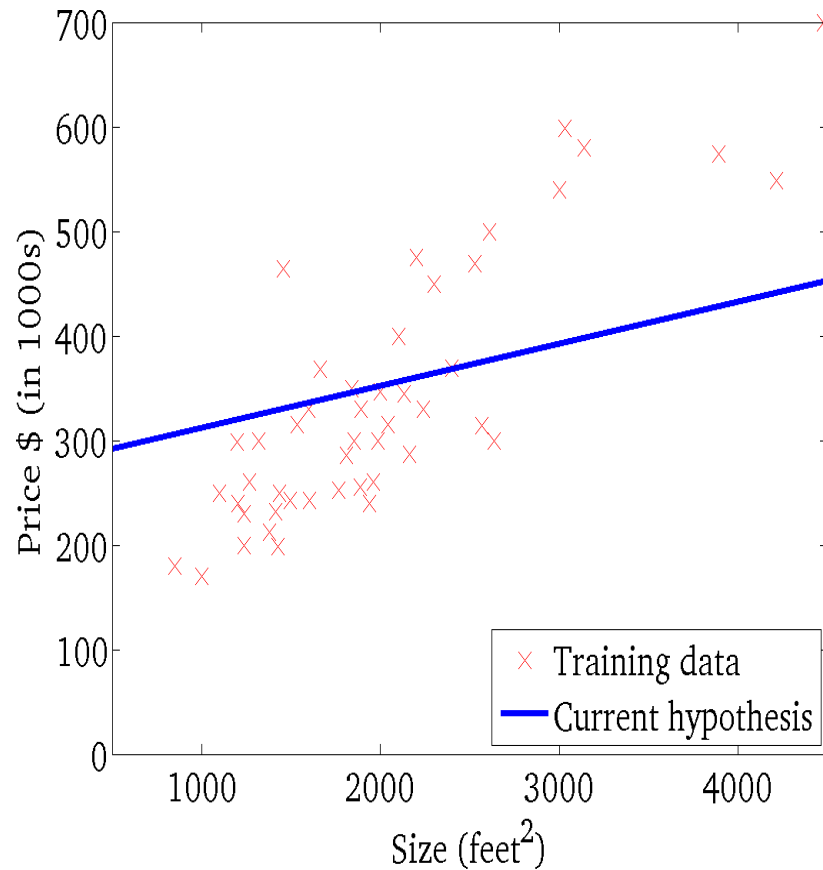


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

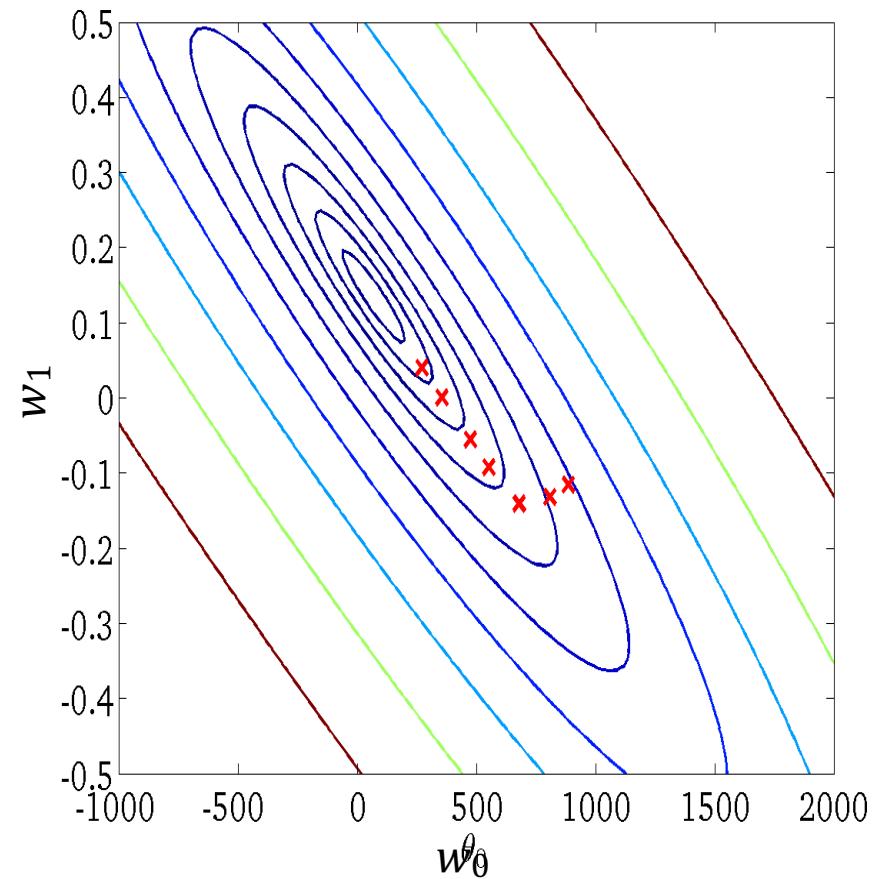


$$f(x; w_0, w_1) = w_0 + w_1 x$$

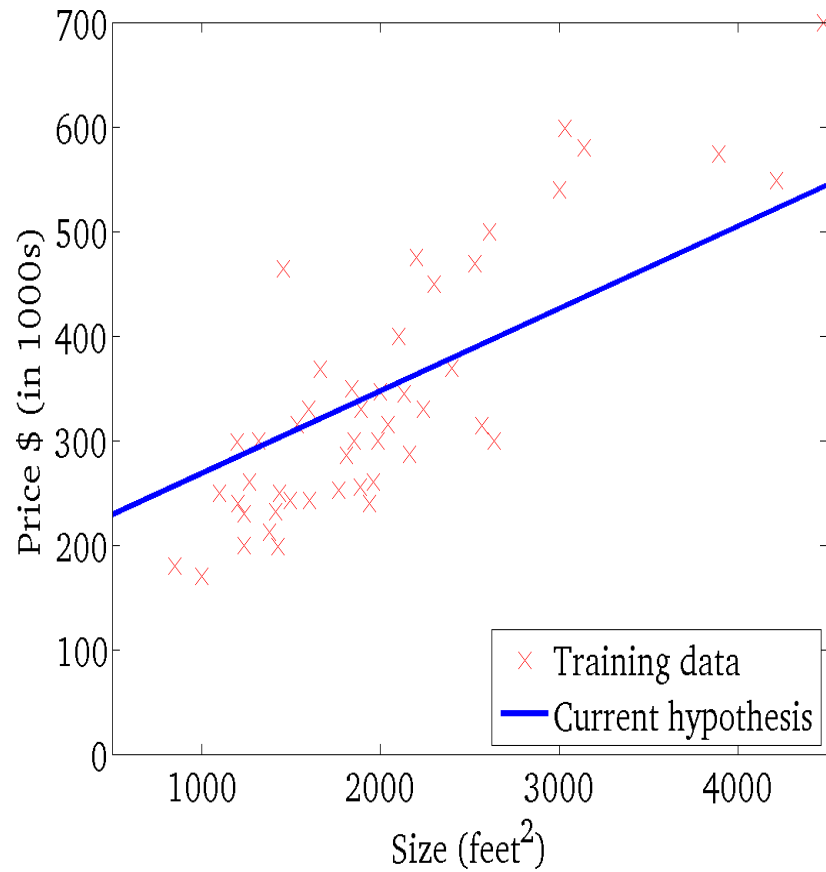


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

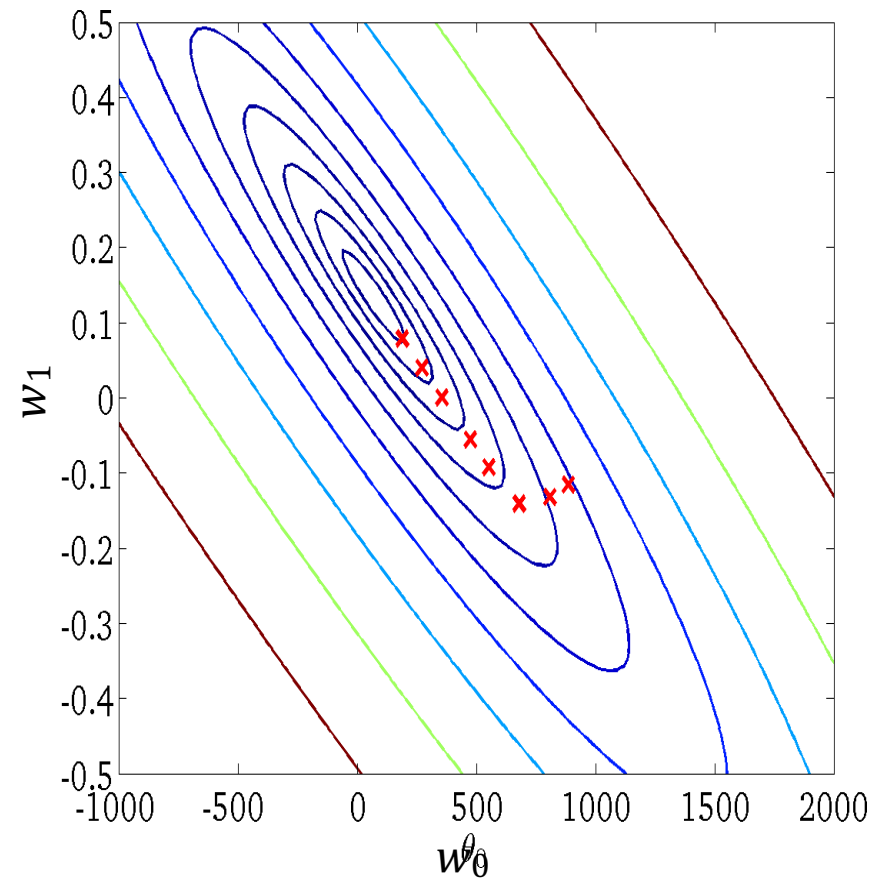


$$f(x; w_0, w_1) = w_0 + w_1 x$$

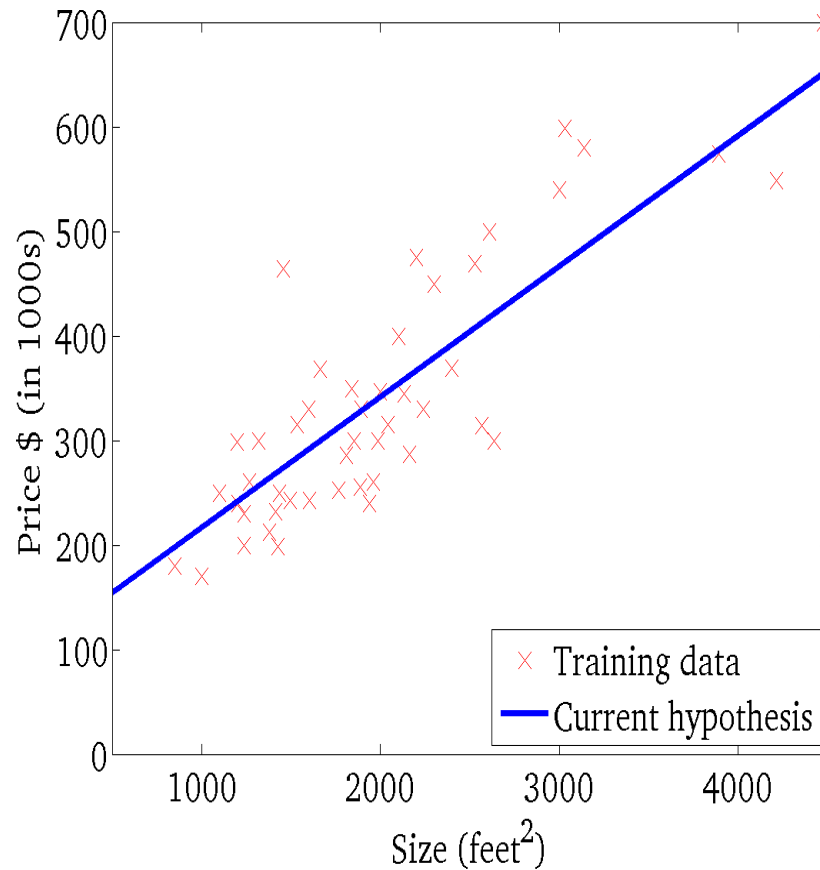


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

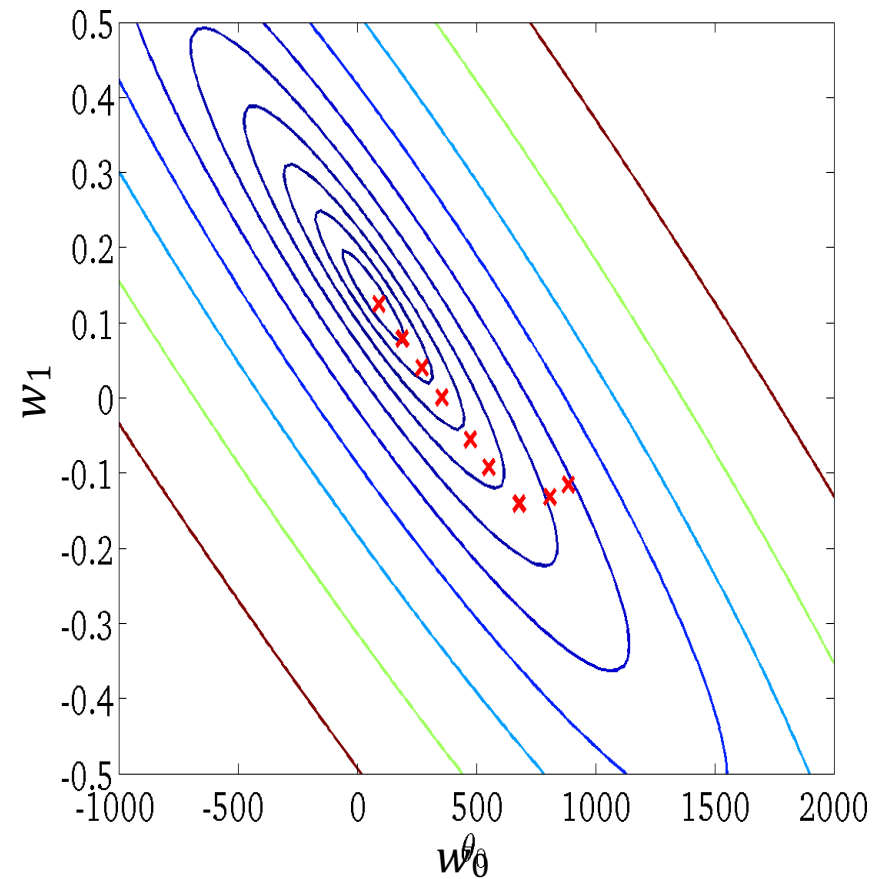


$$f(x; w_0, w_1) = w_0 + w_1 x$$



$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



Stochastic gradient descent

- ▶ Example: Linear regression with SSE cost function

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

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

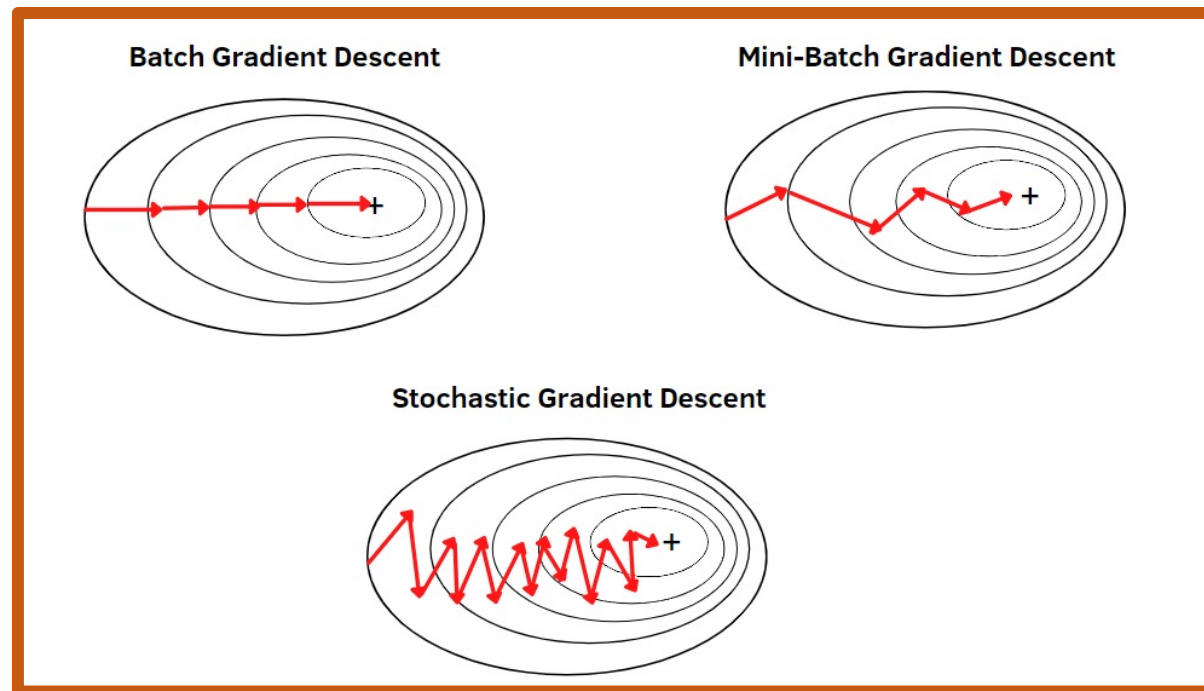
$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}) \mathbf{x}^{(i)}$$

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

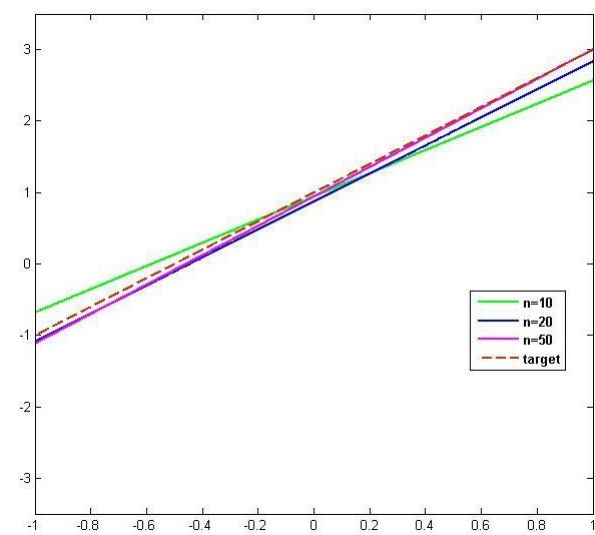
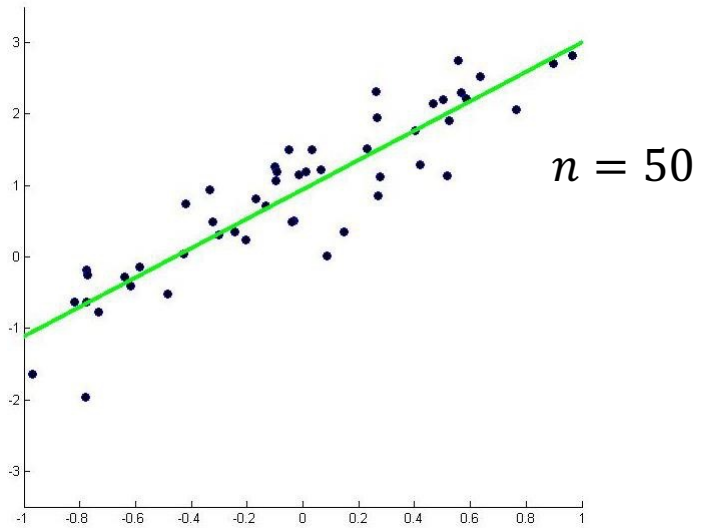
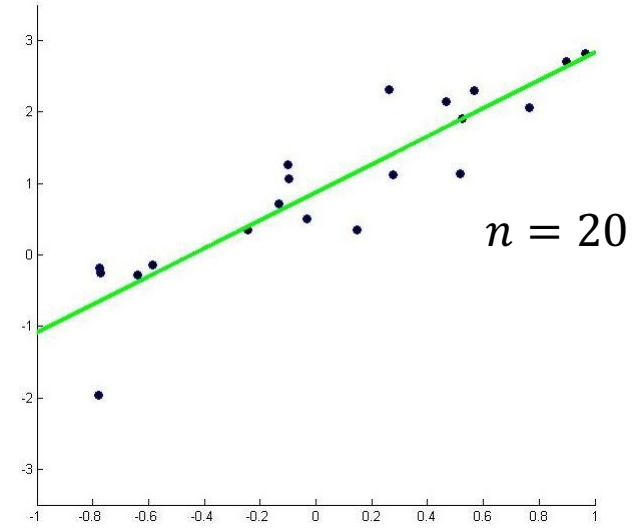
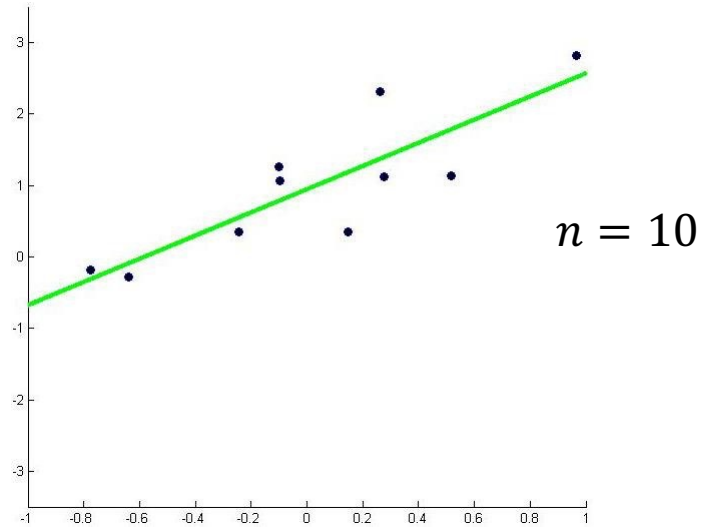
- ▶ Assumption: training and test examples are drawn independently at random from the same but unknown distribution.
- ▶ Each training/test example (\mathbf{x}, y) is a sample from joint probability distribution $P(\mathbf{x}, y)$, i.e., $(\mathbf{x}, y) \sim P$

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

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

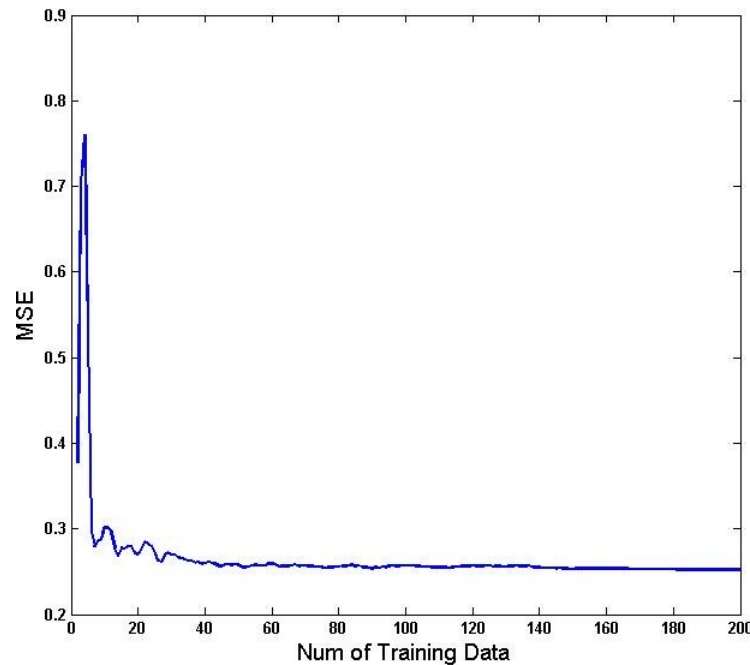
- ▶ We minimize empirical loss (on the training data) and expect to also find an acceptable expected loss
 - ▶ Empirical loss as a proxy for the performance over the whole distribution.

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:

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

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

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

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

Recall: Linear regression (squared loss)

- ▶ Linear regression functions

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

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

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

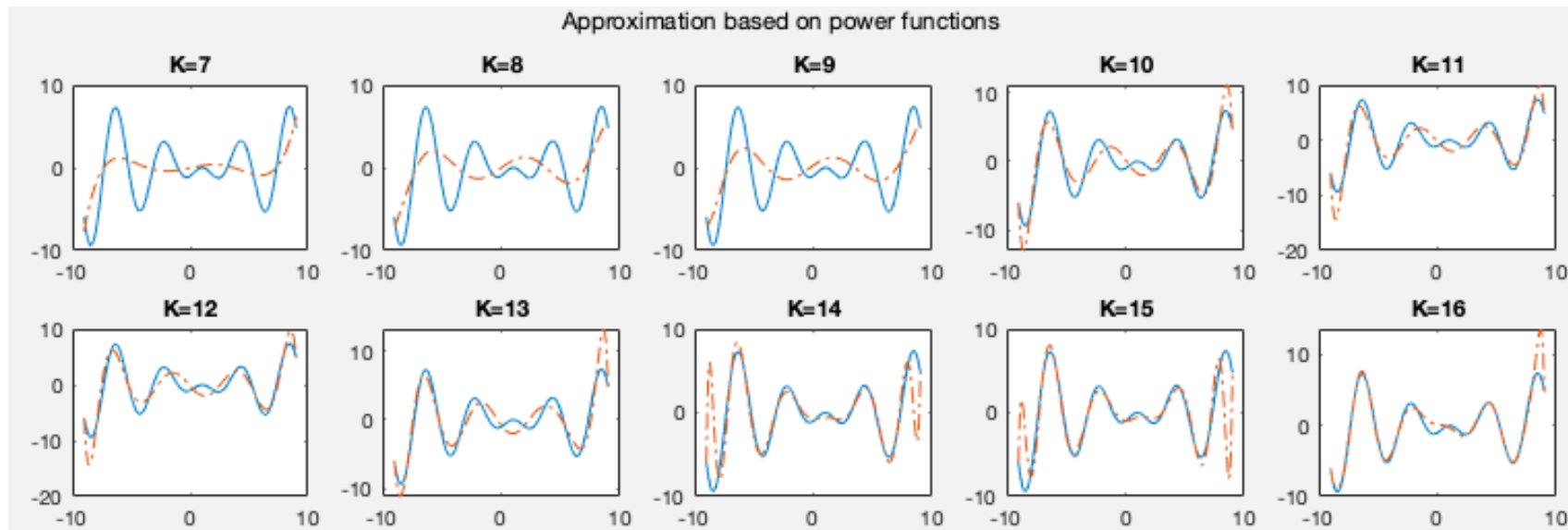
- ▶ Minimizing the squared loss for linear regression

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

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

Beyond linear regression

- ▶ How to extend the linear regression to non-linear functions?
 - ▶ Transform the data using basis functions
 - ▶ Learn a linear regression on the new feature vectors (obtained
 - by basis functions)



Beyond linear regression

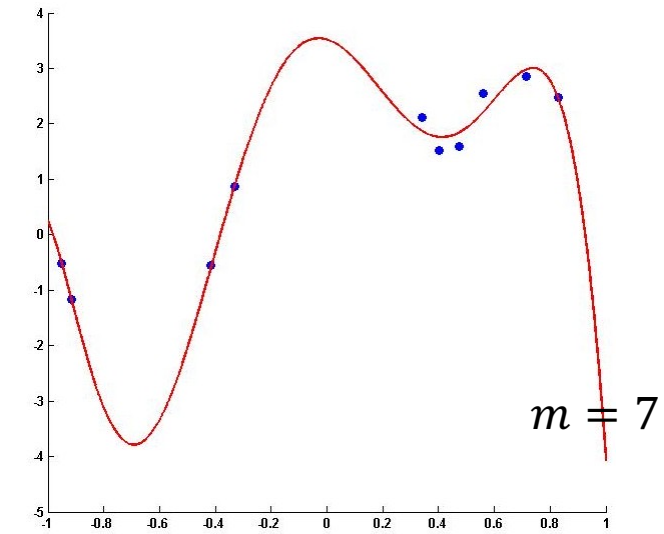
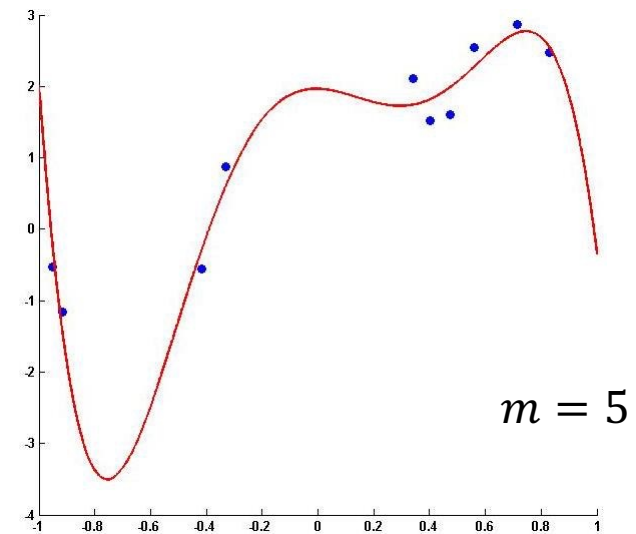
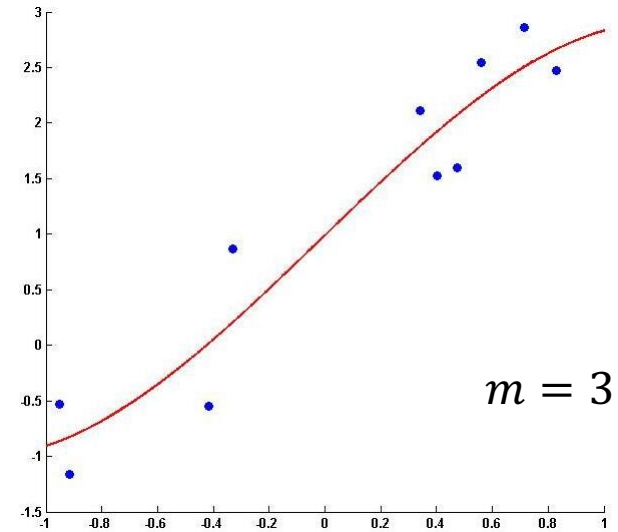
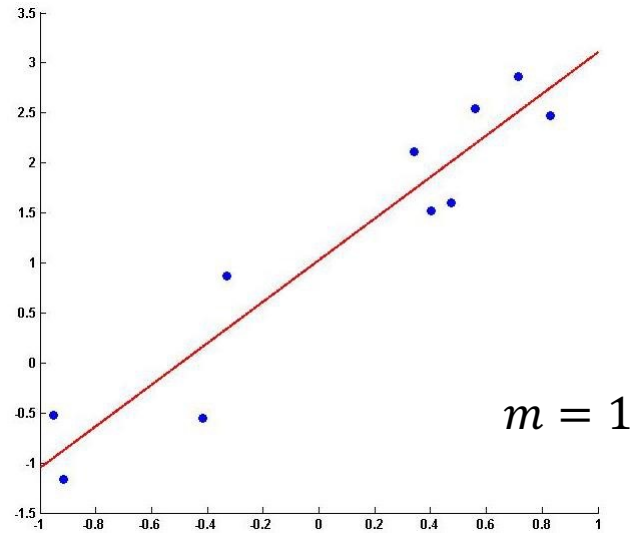
- ▶ m^{th} order polynomial regression (univariate $f : \mathbb{R} \rightarrow \mathbb{R}$)

$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{m-1}x^{m-1} + w_mx^m$$

- ▶ Solution: $\mathbf{w} = (\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T \mathbf{y}$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X}' = \begin{bmatrix} 1 & x^{(1)1} & x^{(1)2} & \dots & x^{(1)m} \\ 1 & x^{(2)1} & x^{(2)2} & \dots & x^{(2)m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x^{(n)1} & x^{(n)2} & \dots & x^{(n)m} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} \mathbf{w}_0 \\ \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_m \end{bmatrix}$$

Polynomial regression: example



Generalized linear

- ▶ Linear combination of fixed non-linear function of the input vector

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots w_m\phi_m(\mathbf{x})$$

$\{\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})\}$: set of basis functions (or features)

$$\phi_i(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$$

Basis functions: examples

- ▶ Linear

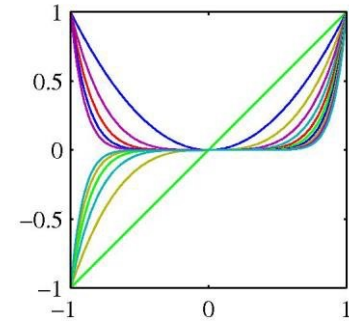
If $m = d$, $\phi_i(\mathbf{x}) = x_i$, $i = 1, \dots, d$, then

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + \dots + w_dx_d$$

- ▶ Polynomial (univariate)

If $\phi_i(x) = x^i$, $i = 1, \dots, m$, then

$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{m-1}x^{m-1} + w_mx^m$$



Generalized linear: optimization

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

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \boldsymbol{\Phi} = \begin{bmatrix} 1 & \phi_1(\mathbf{x}^{(1)}) & \cdots & \phi_m(\mathbf{x}^{(1)}) \\ 1 & \phi_1(\mathbf{x}^{(2)}) & \cdots & \phi_m(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}^{(n)}) & \cdots & \phi_m(\mathbf{x}^{(n)}) \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

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

Resource

- 1 C. M. Bishop, *Pattern Recognition and Machine Learning*.
- 2 Y. S. Abu-Mostafa, “Machine learning.” California Institute of Technology, 2012.
- 3 Machine Learning, Dr. Soleymani, Sharif University