

Machine Learning Review

M. Soleymani

Sharif University of Technology

Fall 2017

Some slides have been adapted from Fei Fei Li lectures, cs231n, Stanford 2017

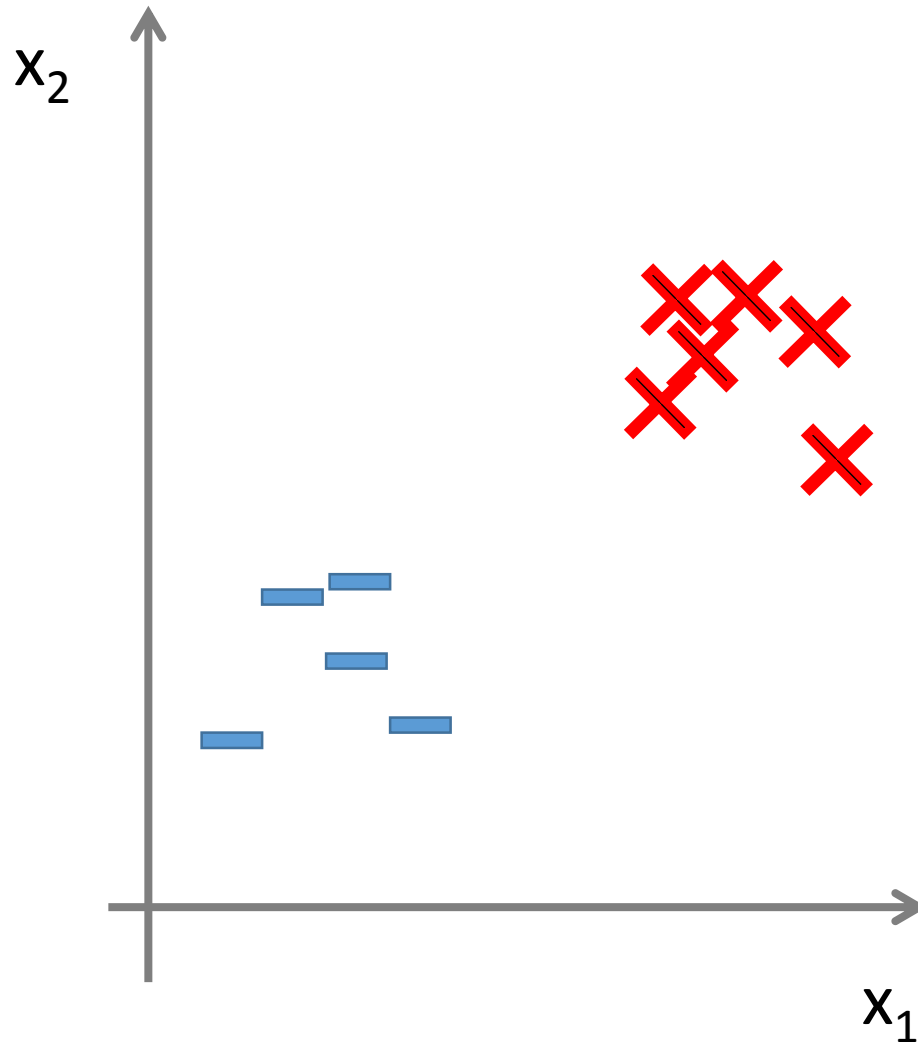
Types of ML problems

- Supervised learning (regression, classification)
 - predicting a target variable for which we get to see examples.
- Unsupervised learning
 - revealing structure in the observed data
- Reinforcement learning
 - partial (indirect) feedback, no explicit guidance
 - Given rewards for a sequence of moves to learn a policy and utility functions

Components of (Supervised) Learning

- Unknown target function: $f: \mathcal{X} \rightarrow \mathcal{Y}$
 - Input space: \mathcal{X}
 - Output space: \mathcal{Y}
- Training data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- We use **training set** to find the function that can also predict output on the **test set**

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	×

Supervised Learning: Regression vs. Classification

- Supervised Learning
 - **Regression**: predict a continuous target variable
 - E.g., $y \in [0,1]$
 - **Classification**: predict a discrete target variable
 - E.g., $y \in \{1,2, \dots, C\}$

Regression Example

- Housing price prediction

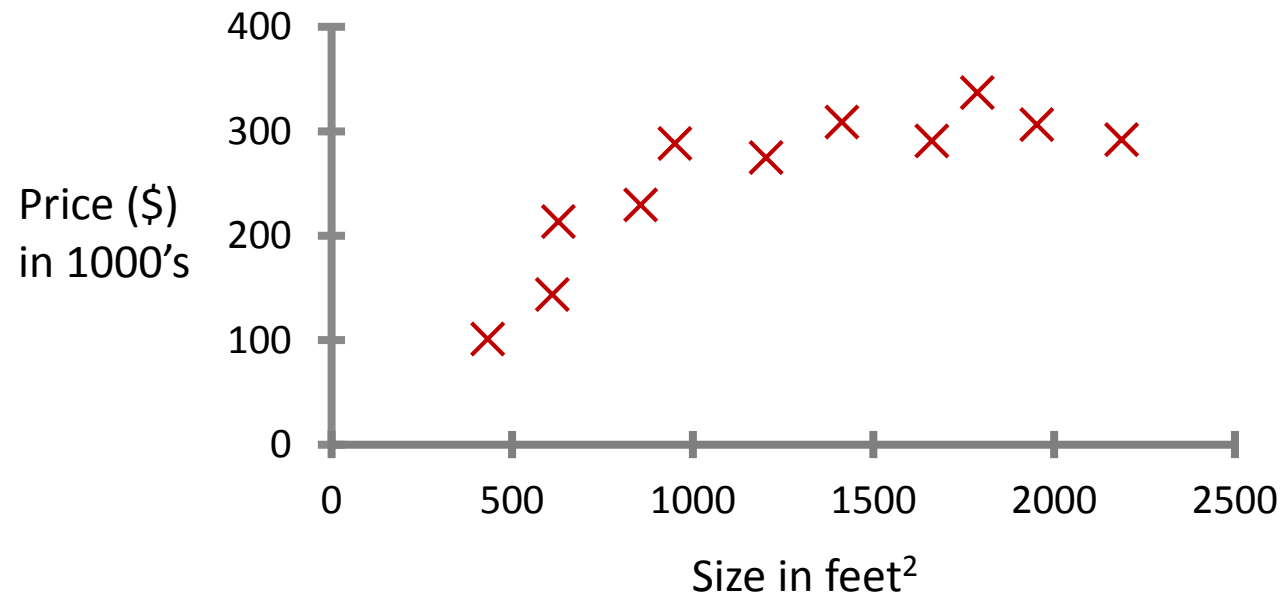
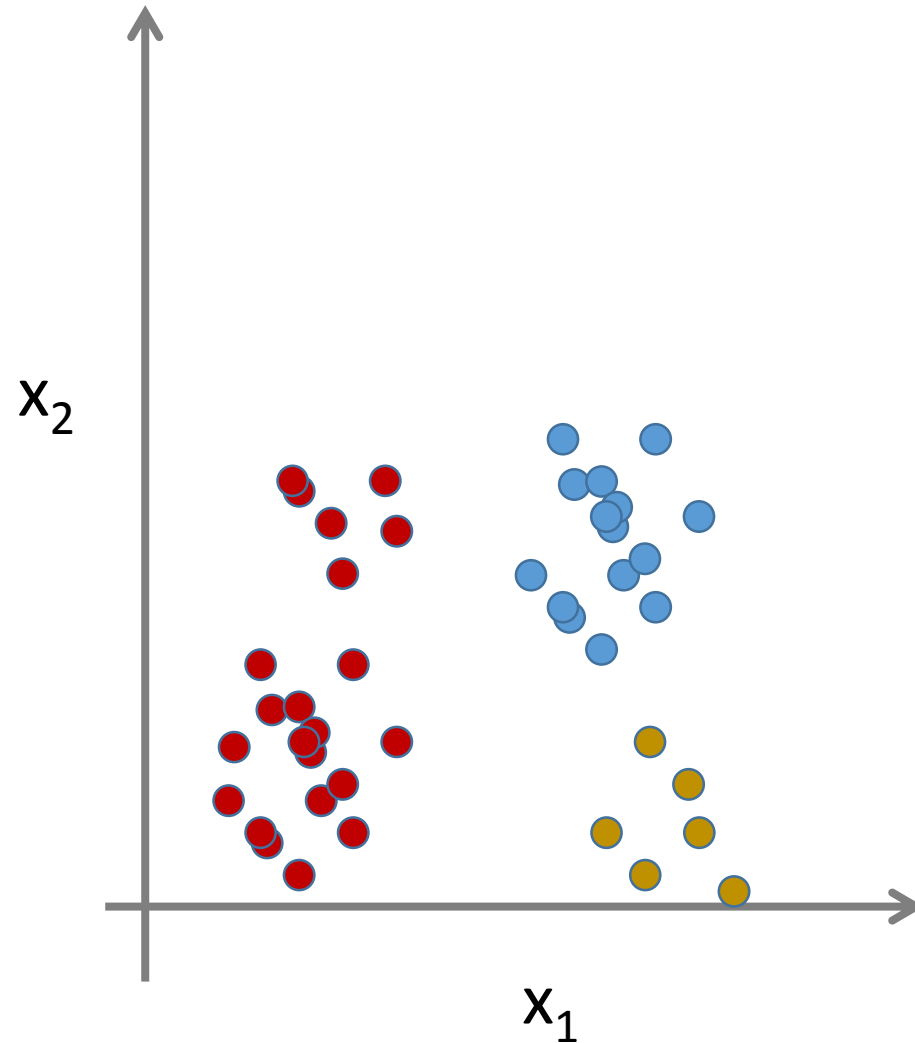


Figure adopted from slides of Andrew Ng,
Machine Learning course, Stanford.

Supervised Learning vs. Unsupervised Learning

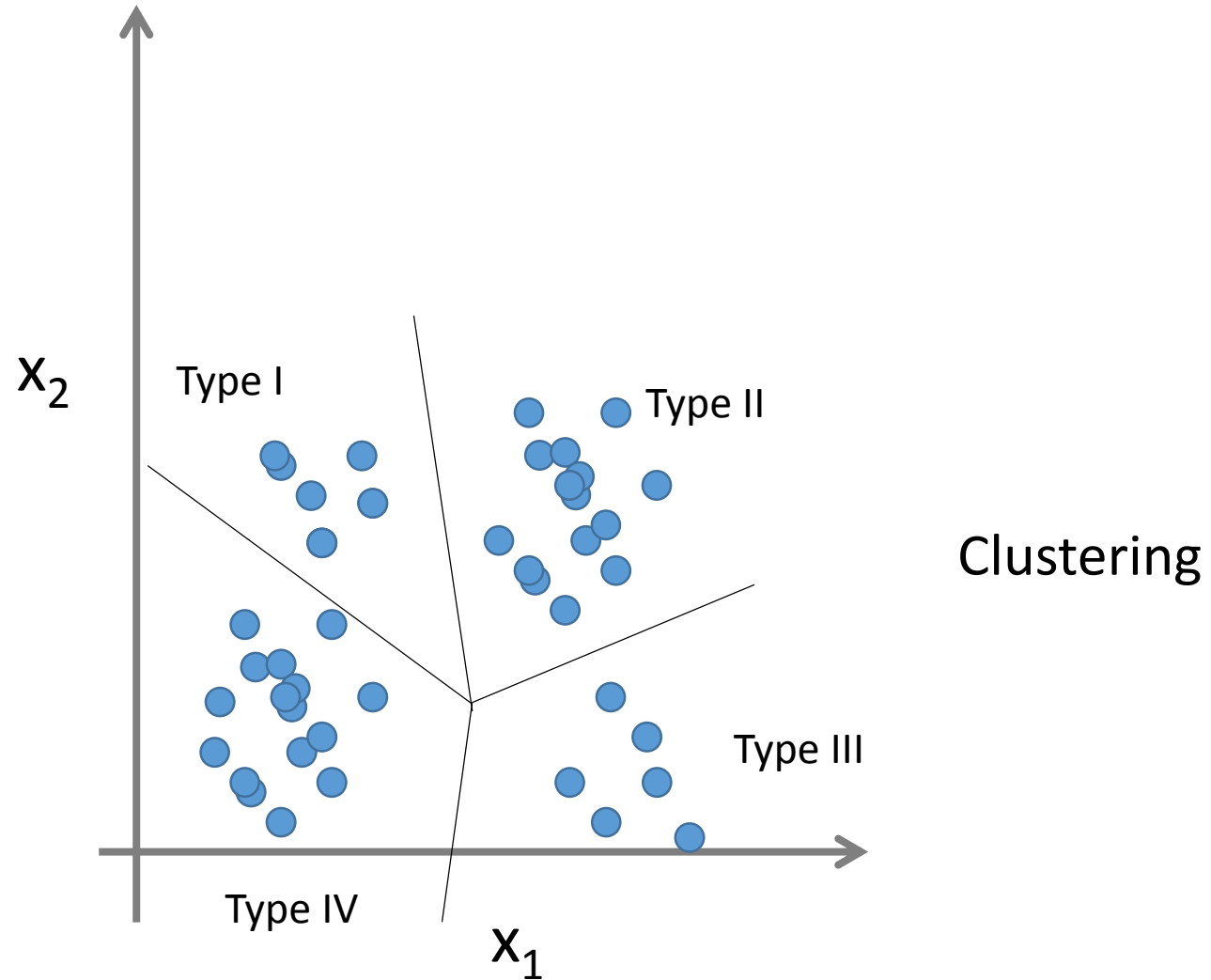
- Supervised learning
 - Given: Training set
 - labeled set of N input-output pairs $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - Goal: learning a mapping from \mathbf{x} to y
- Unsupervised learning
 - Given: Training set
 - $\{\mathbf{x}^{(i)}\}_{i=1}^N$
 - Goal: find groups or structures in the data
 - Discover the intrinsic structure in the data

Supervised Learning: Samples



Classification

Unsupervised Learning: Samples



Reinforcement Learning

- Provides only an indication as to whether an action is correct or not

Data in supervised learning:

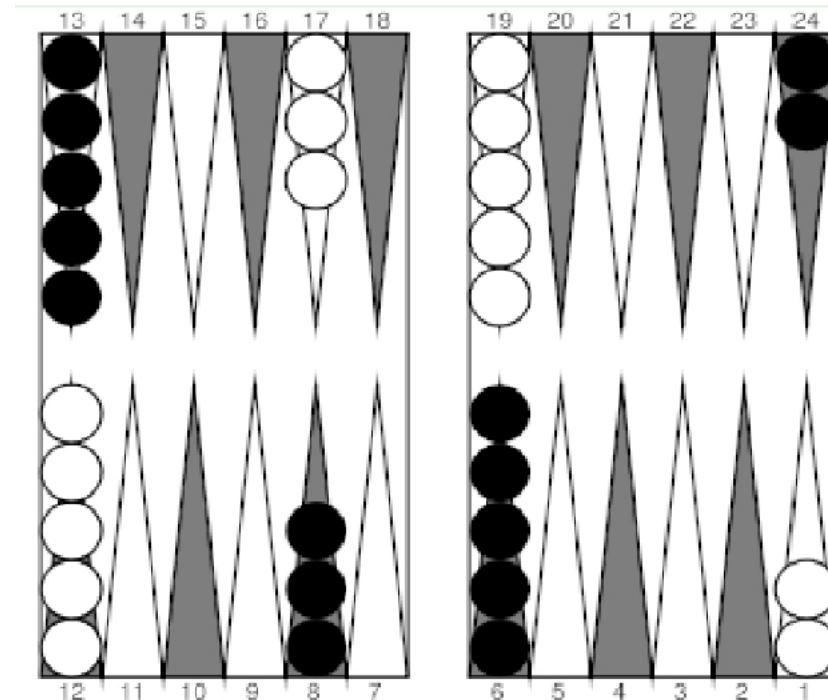
(input, correct output)

Data in Reinforcement Learning:

(input, some output, a grade of reward for this output)

Reinforcement Learning

- Typically, we need to get a sequence of decisions
 - it is usually assumed that reward signals refer to the entire sequence



Components of (Supervised) Learning

- Unknown target function: $f: \mathcal{X} \rightarrow \mathcal{Y}$
 - Input space: \mathcal{X}
 - Output space: \mathcal{Y}
- Training data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- We use **training set** to find the function that can also predict output on the **test set**

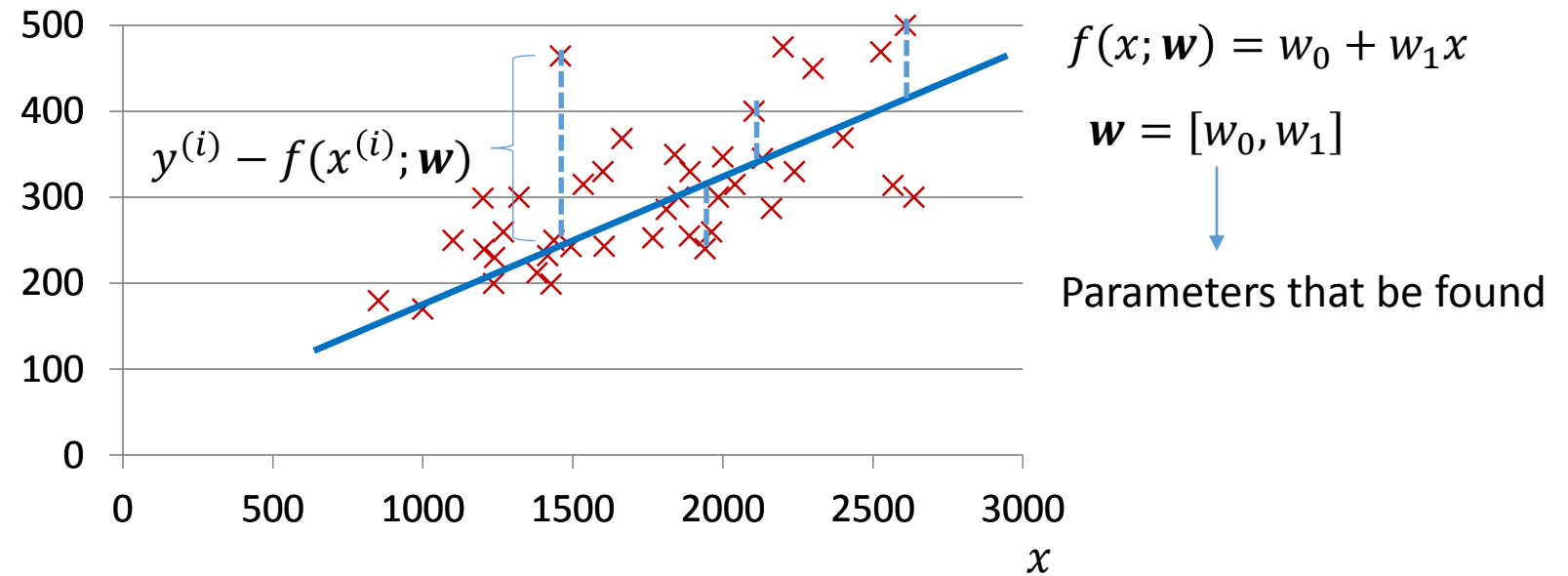
Generalization

- We don't intend to memorize data but need to figure out the pattern.
- A core objective of learning is to generalize from the experience.
 - Generalization: ability of a learning algorithm to perform accurately on new, unseen examples after having experienced.

(Typical) Steps of solving supervised learning problem

- Select the hypothesis space
 - A class of parametric models that map each input vector, \mathbf{x} , into a predicted output y .
- Define a **loss function** that quantifies how much undesirable is each parameter vector across the training data.
- Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)
- Evaluate the obtained model

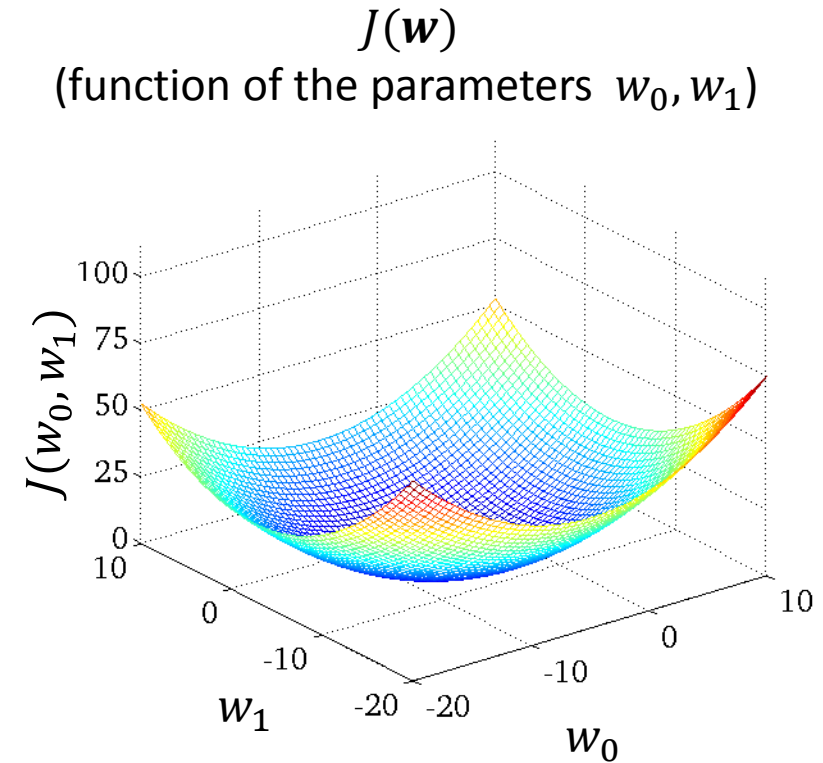
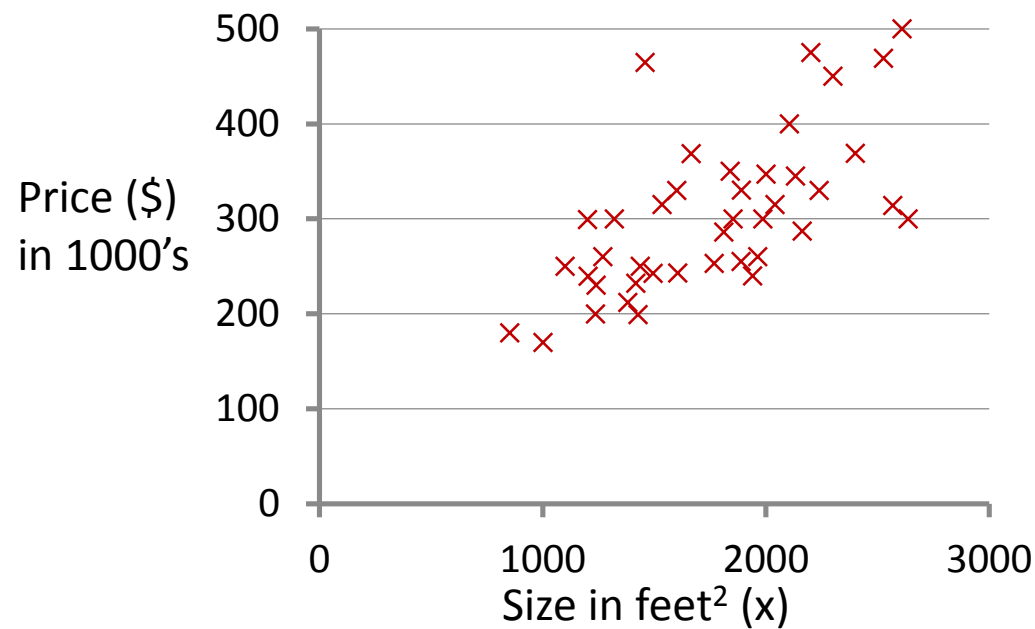
Linear regression: square error loss function



Cost function:

$$J(\mathbf{w}) = \sum_{i=1}^n (y^{(i)} - (w_0 + w_1x^{(i)}))^2$$

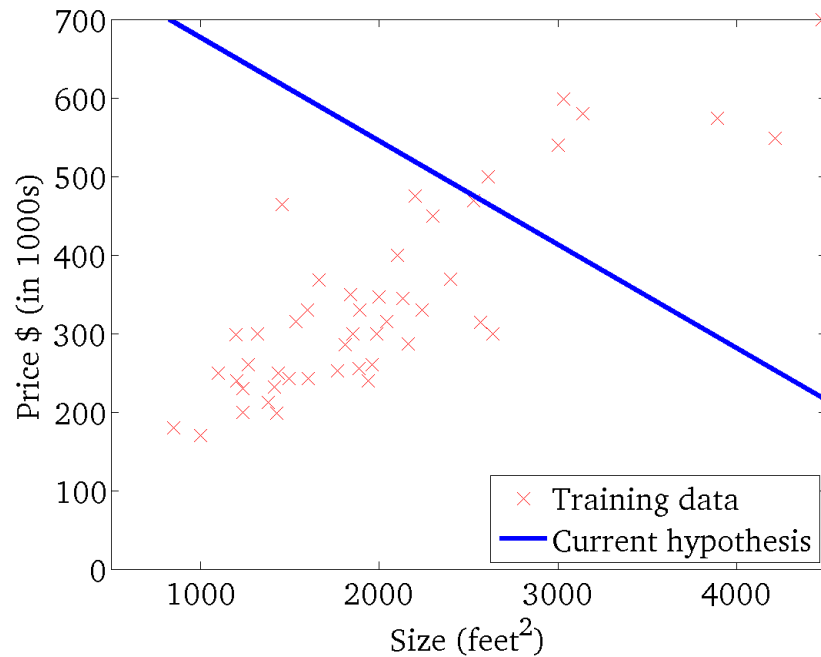
Cost function: example



Cost function: example

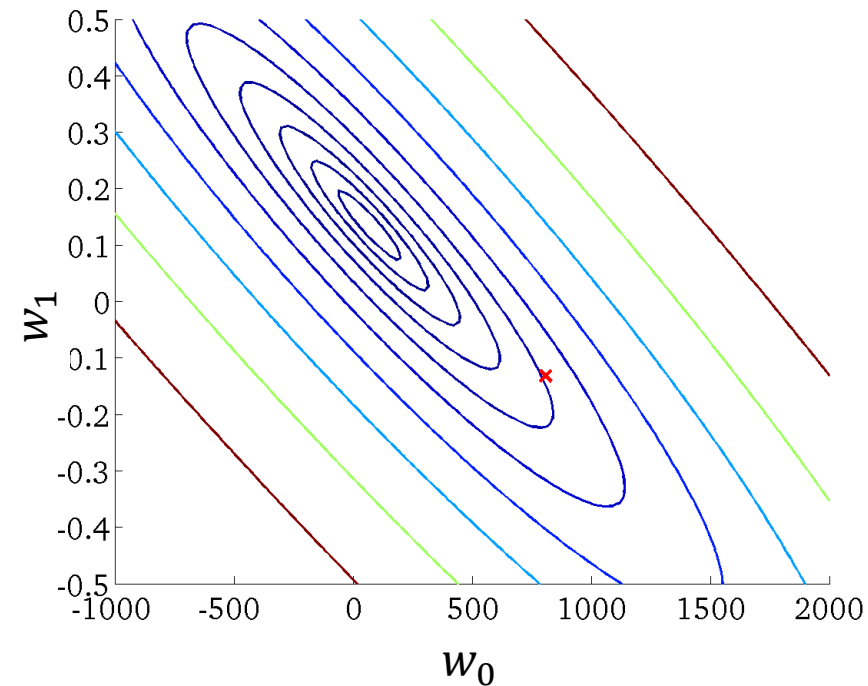
$$f(x; w_0, w_1) = w_0 + w_1 x$$

(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)

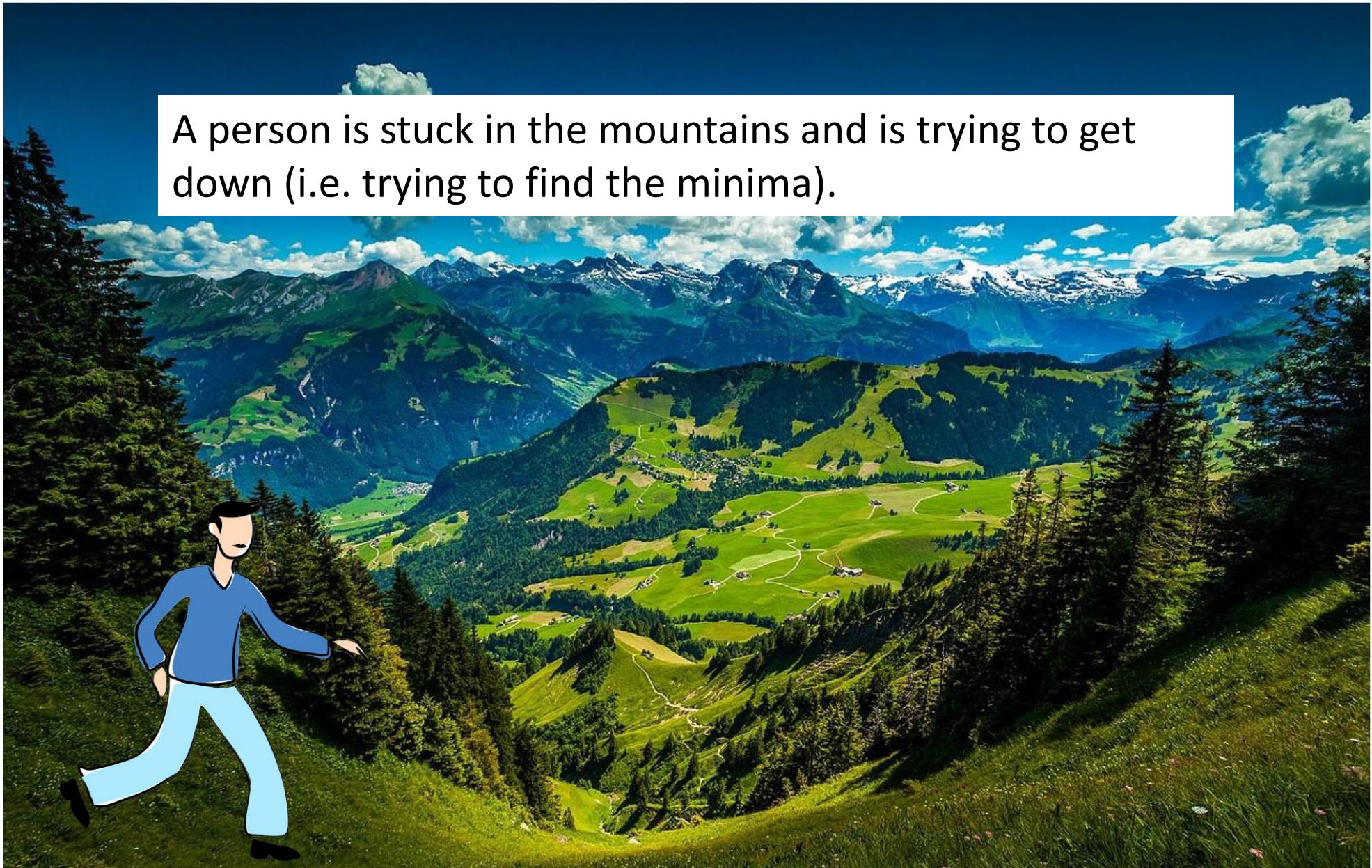


Review: Iterative optimization of cost function

- 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

How to optimize parameters?

A person is stuck in the mountains and is trying to get down (i.e. trying to find the minima).



Follow up the slope

The steepness of the hill represents the slope of the surface at that point.



How to compute the slope?

- In 1-dimension, the derivative of a function:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- the slope of the error surface can be calculated by taking the derivative of the error function at that point
- In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension
- The direction of steepest descent is the **negative gradient**

Gradient descent (or 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


Learning rate: The amount of time he travels before taking another measurement is the learning rate of the algorithm.

Gradient descent

- Minimize $J(\mathbf{w})$

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

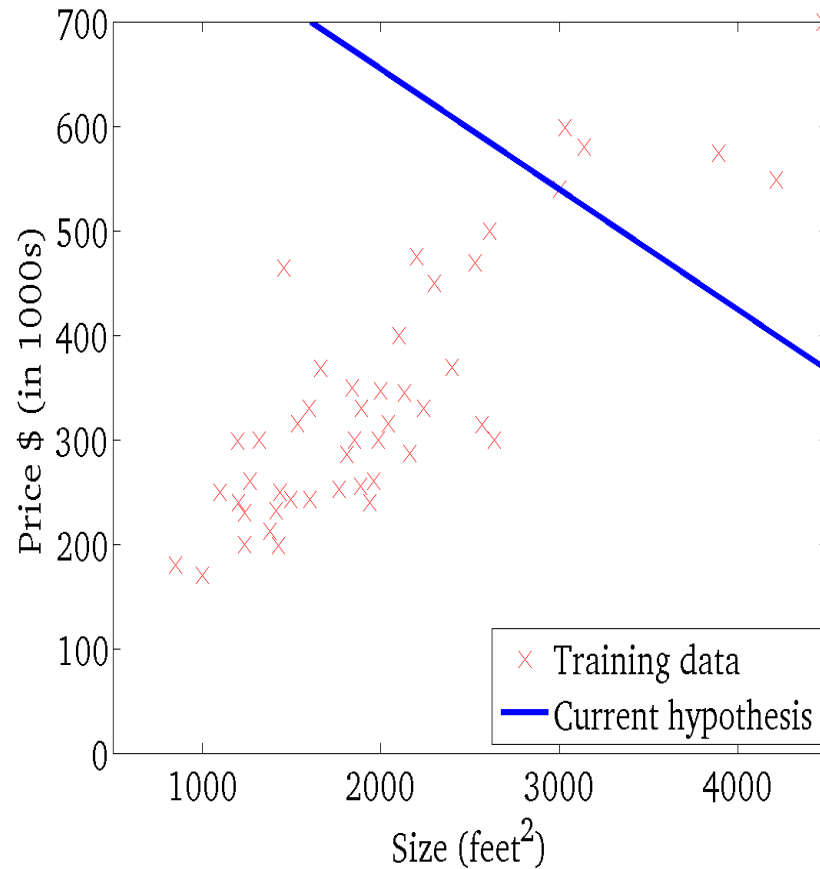
Step size
(Learning rate parameter)



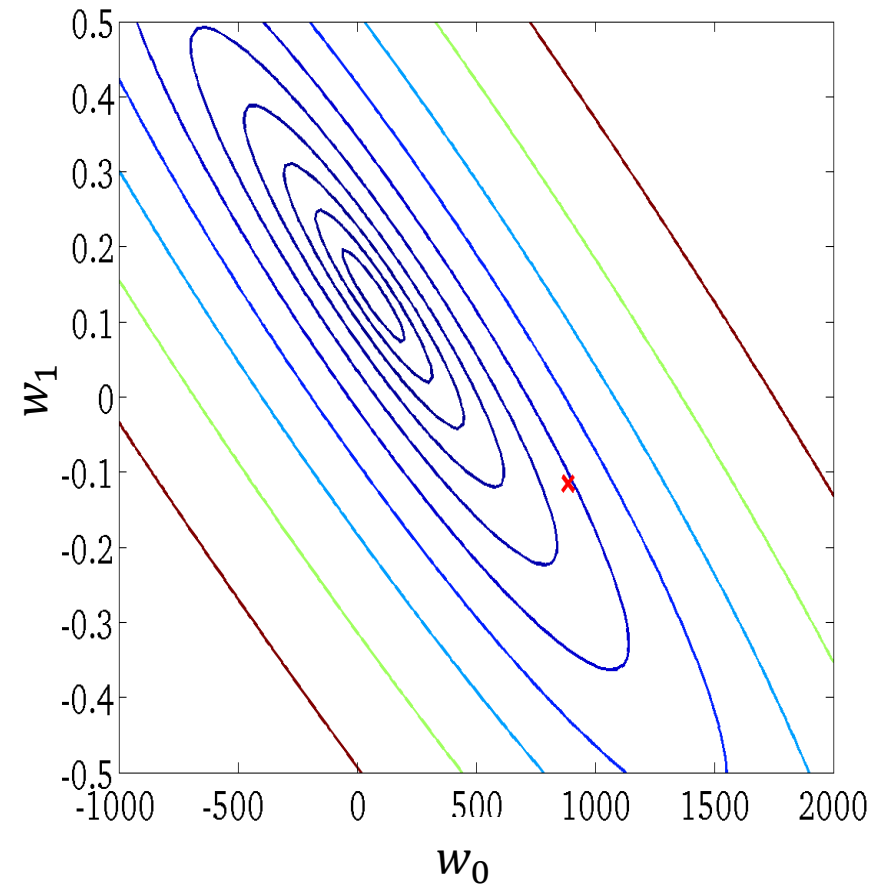
$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \left[\frac{\partial J(\mathbf{w})}{\partial w_0}, \frac{\partial J(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_d} \right]$$

- 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 .

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

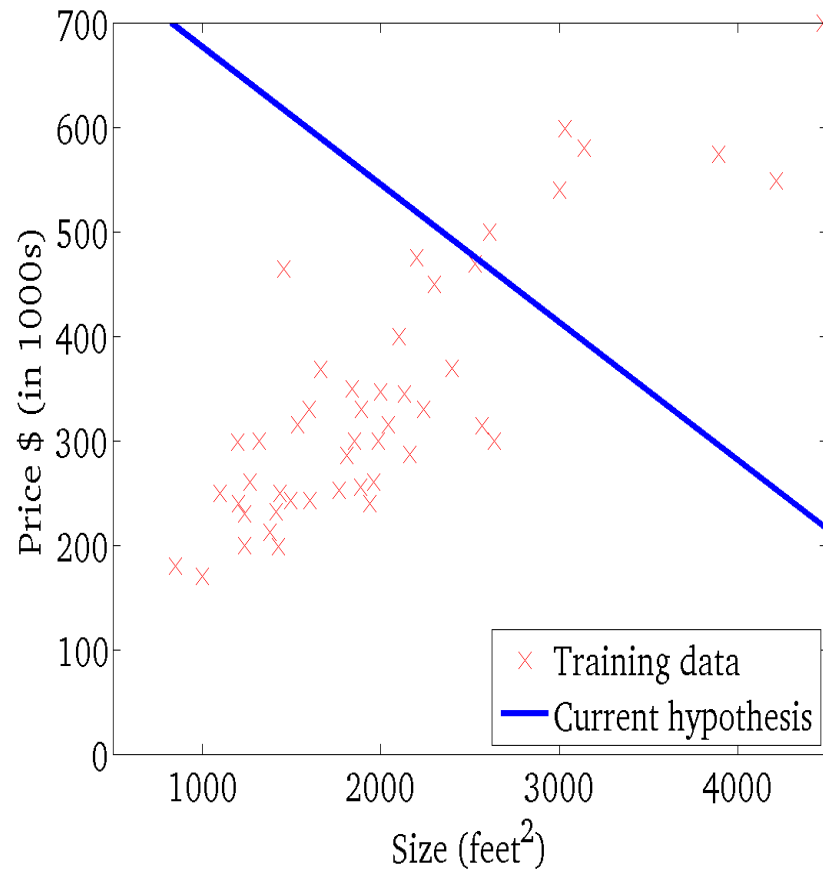


$J(w_0, w_1)$
(function of the parameters w_0, w_1)



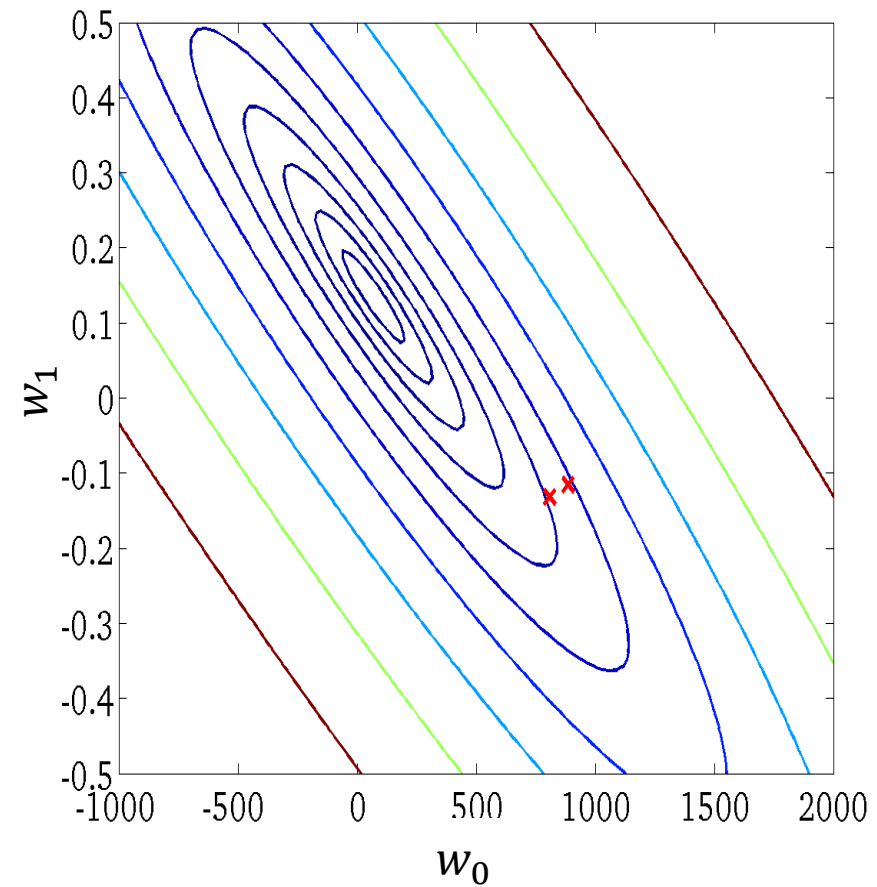
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \sum_{i=1}^N (\mathbf{w}^{tT} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$$

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

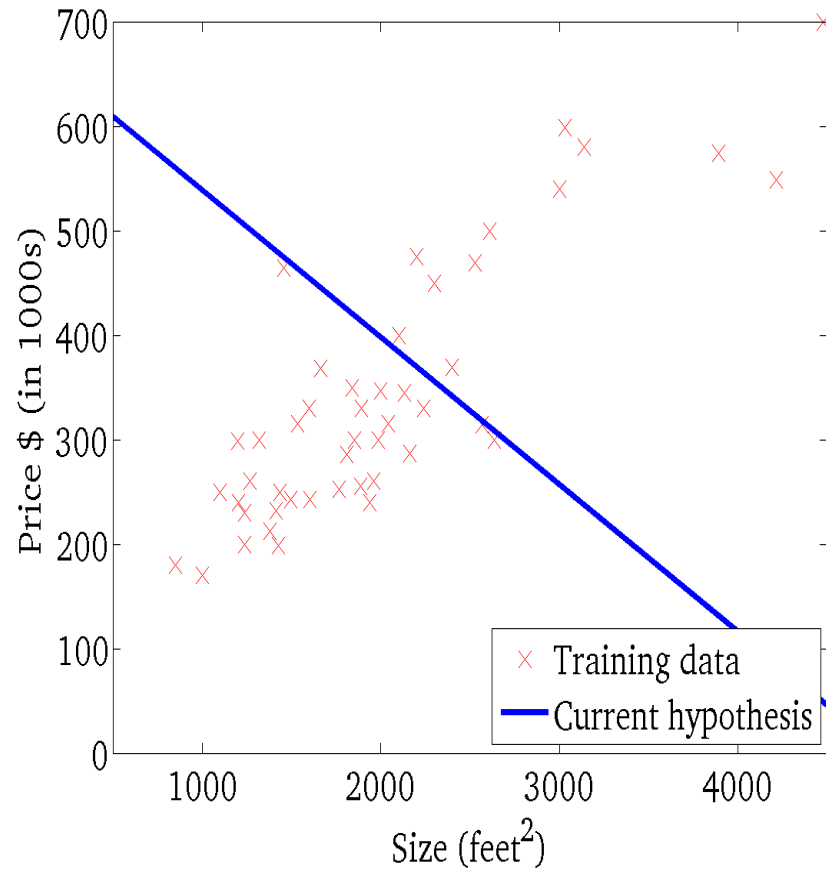


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

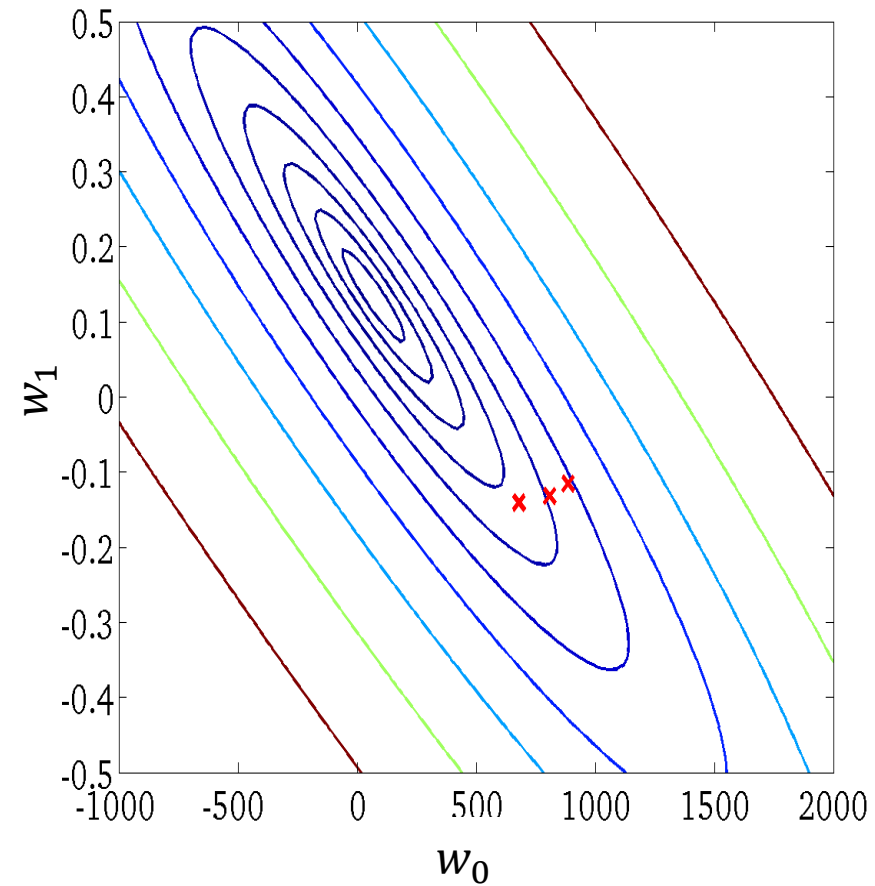


$$f(x; w_0, w_1) = w_0 + w_1 x$$

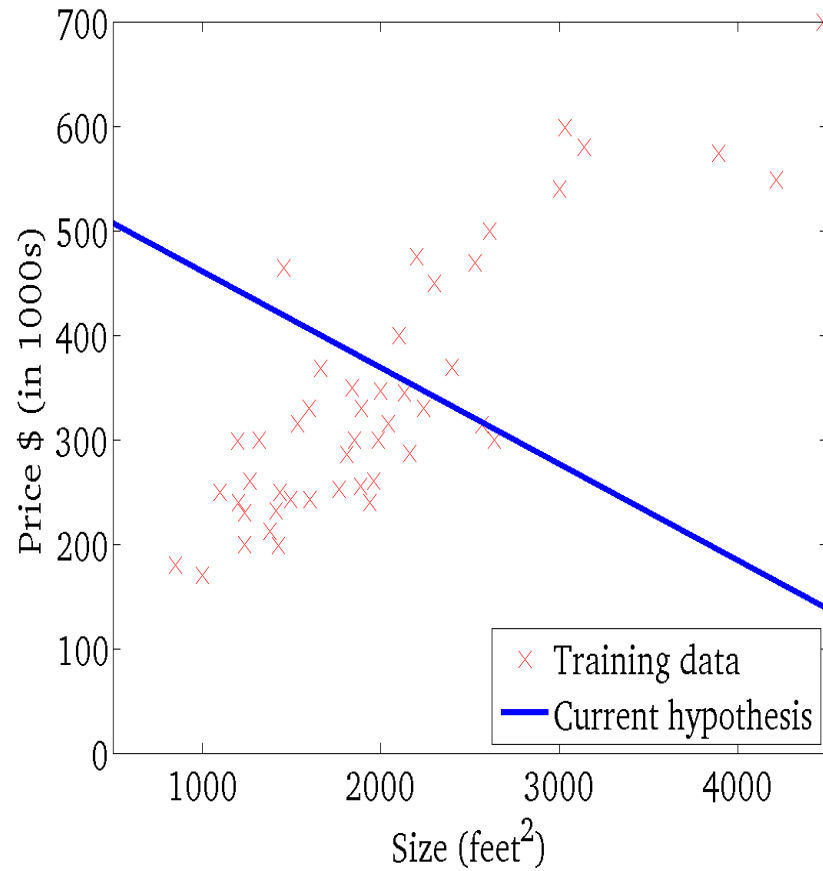


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

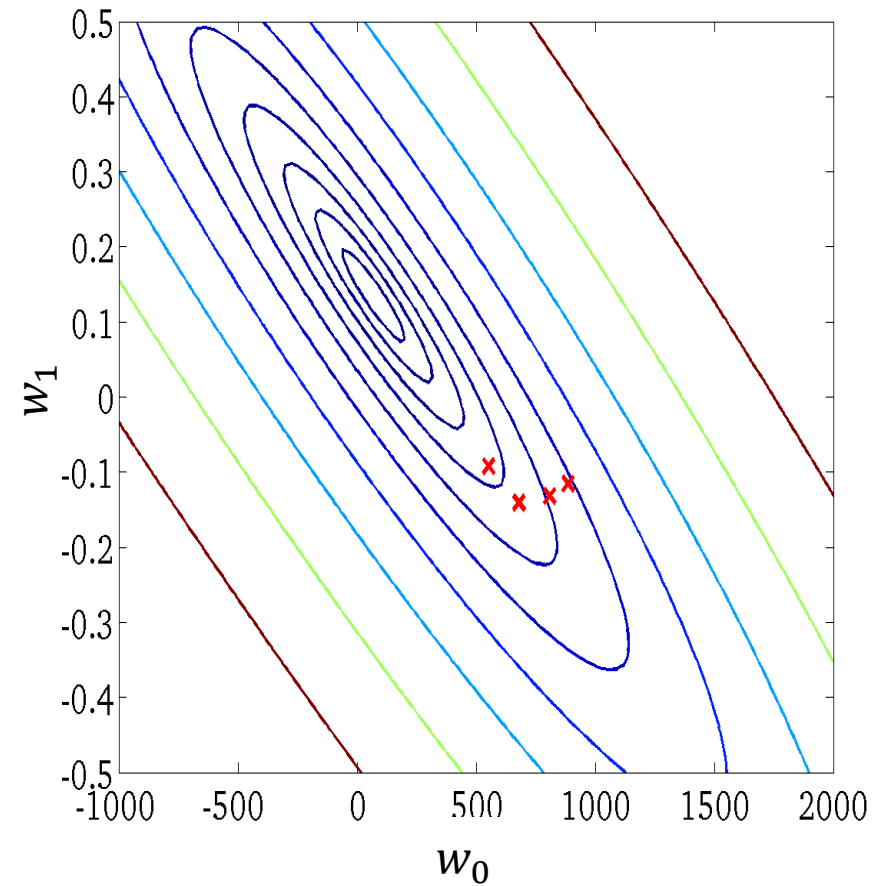


$$f(x; w_0, w_1) = w_0 + w_1 x$$

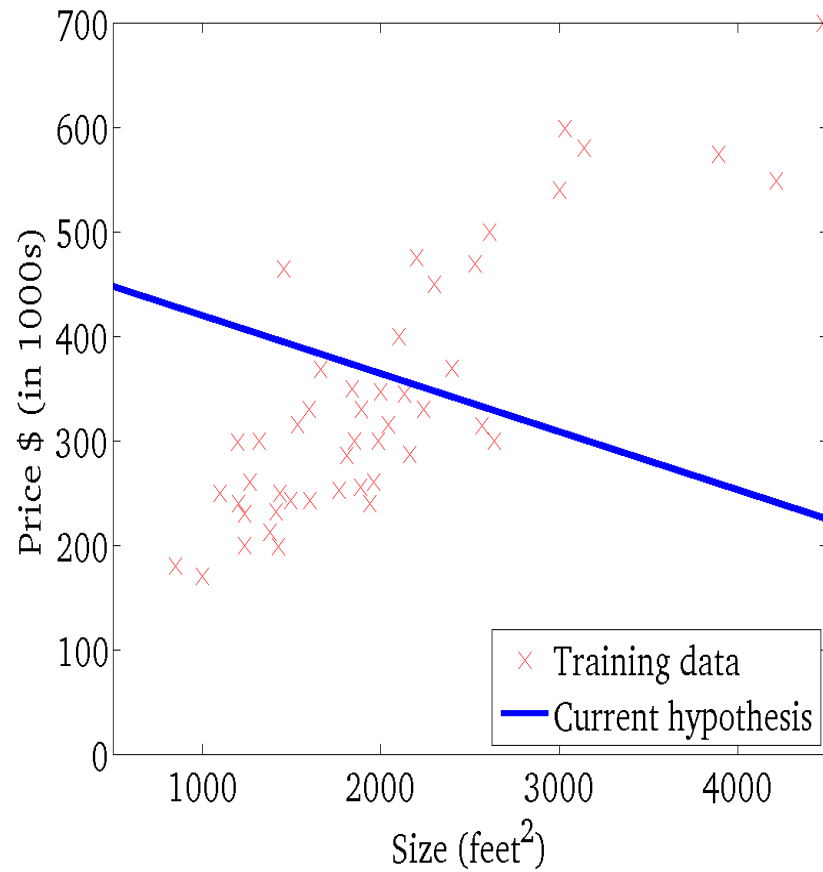


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

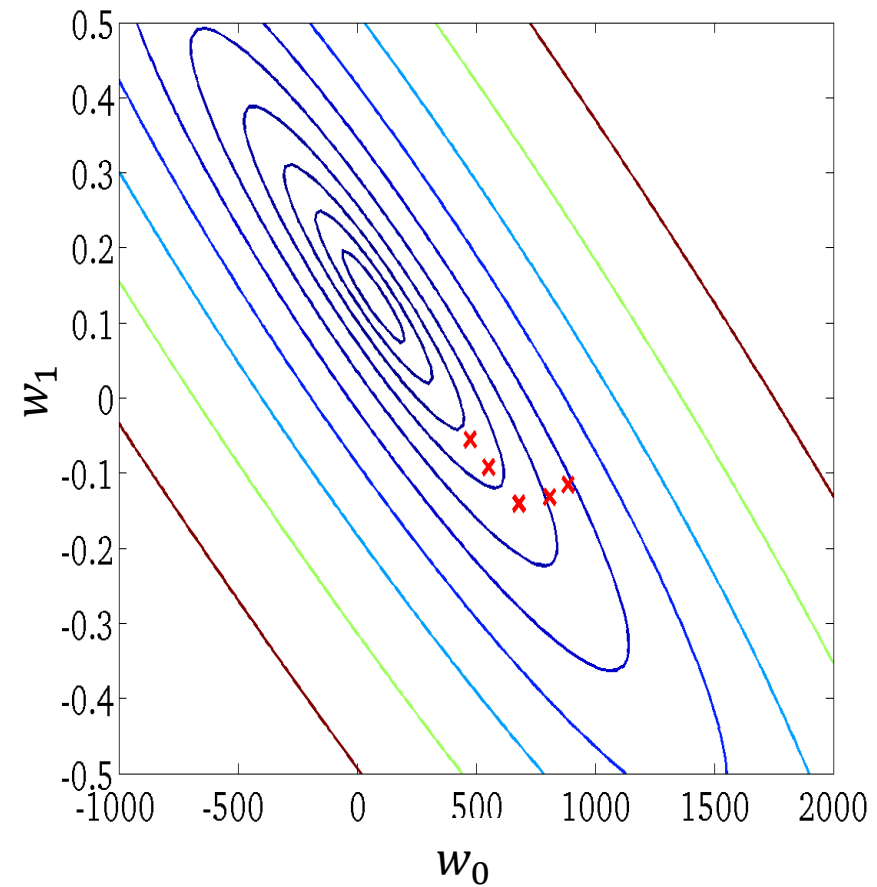


$$f(x; w_0, w_1) = w_0 + w_1 x$$

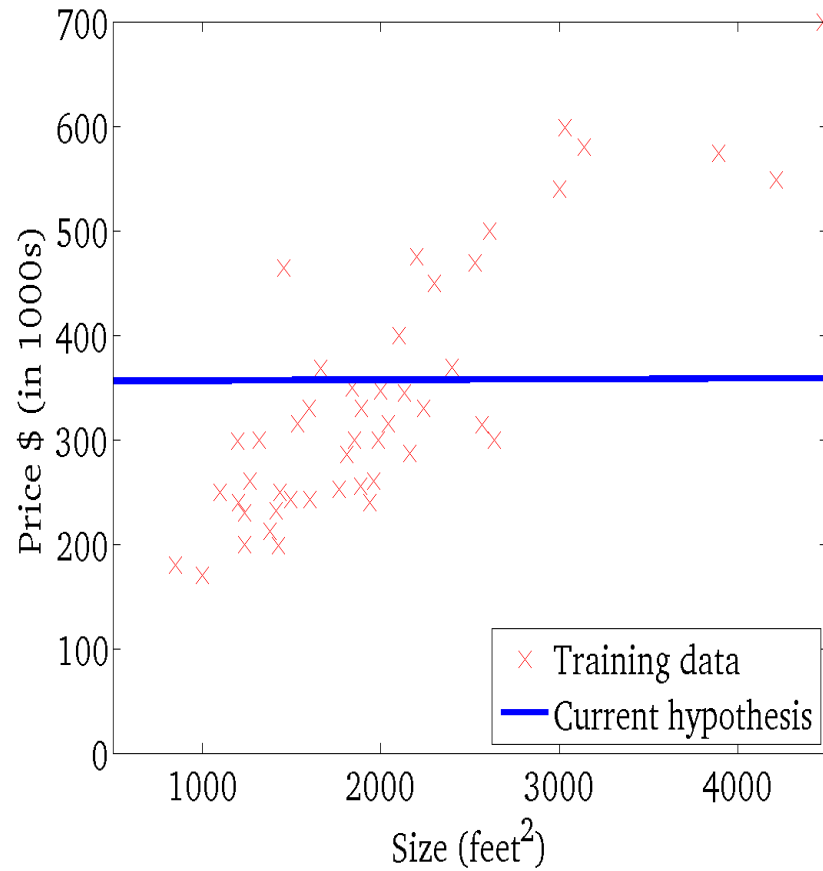


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

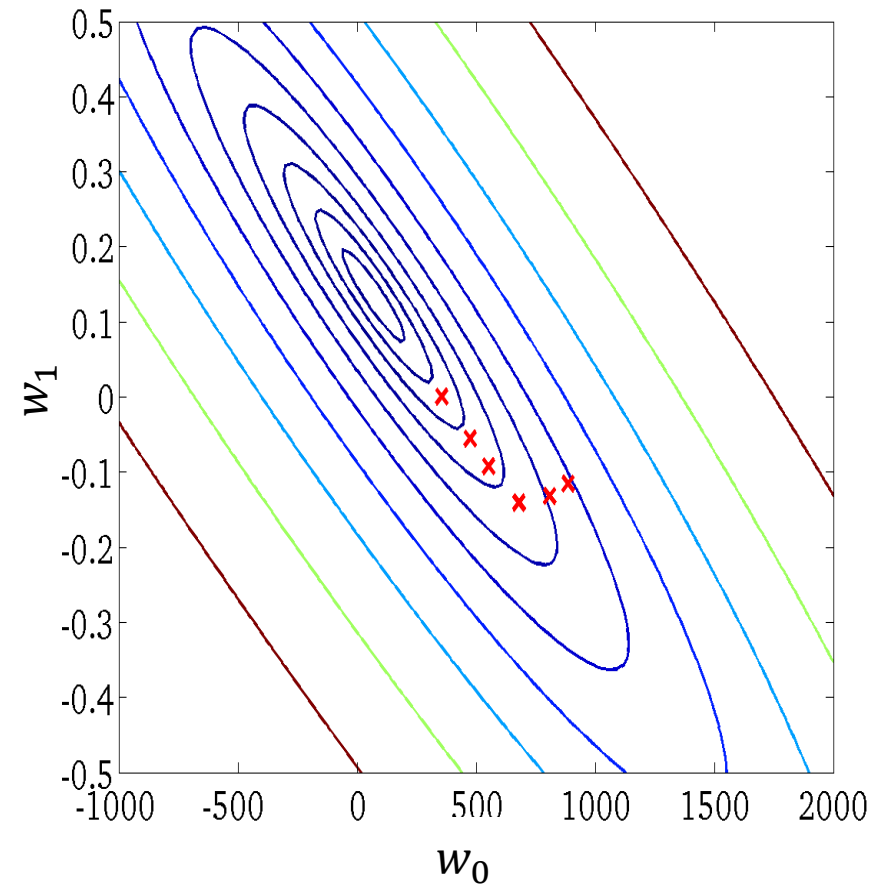


$$f(x; w_0, w_1) = w_0 + w_1 x$$

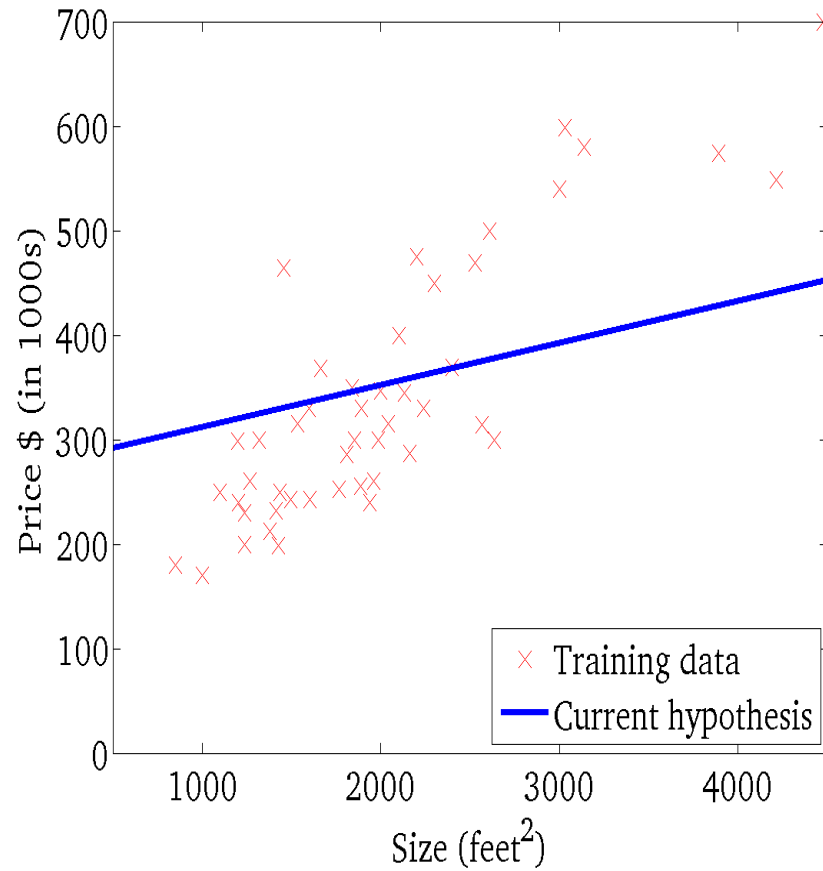


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

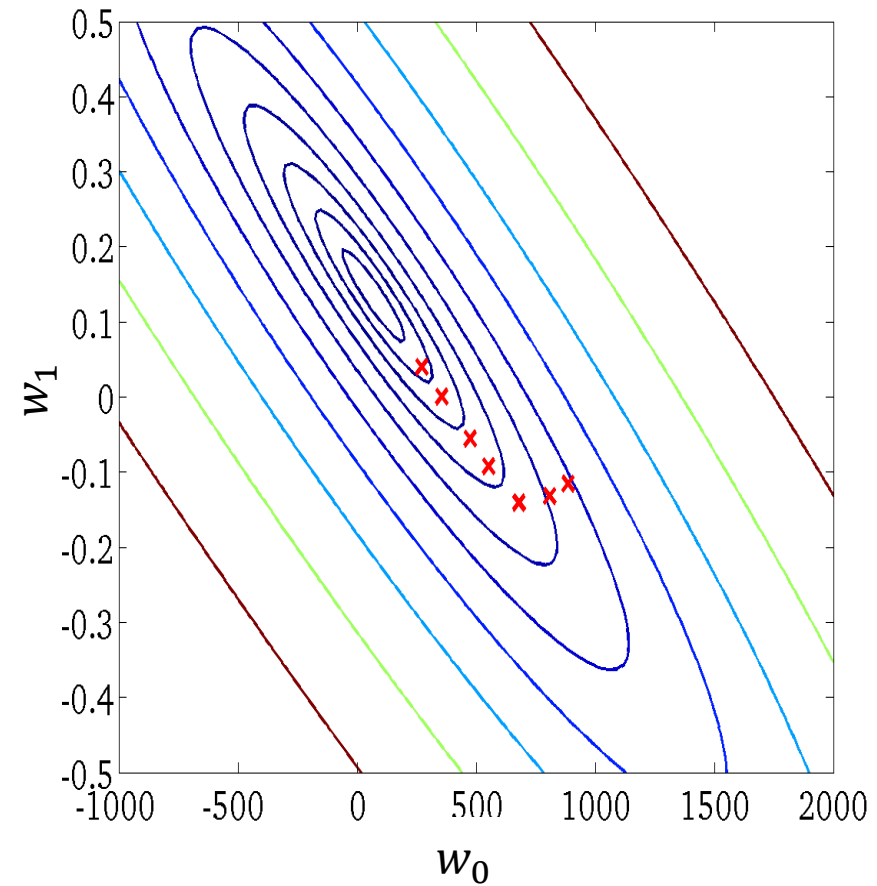


$$f(x; w_0, w_1) = w_0 + w_1 x$$

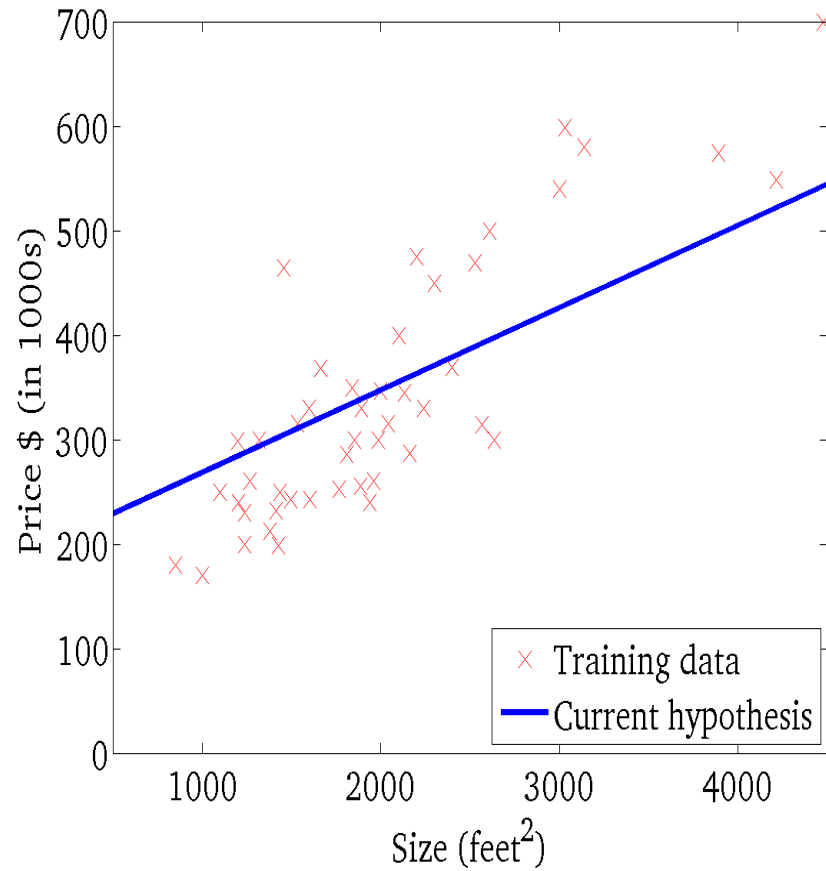


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

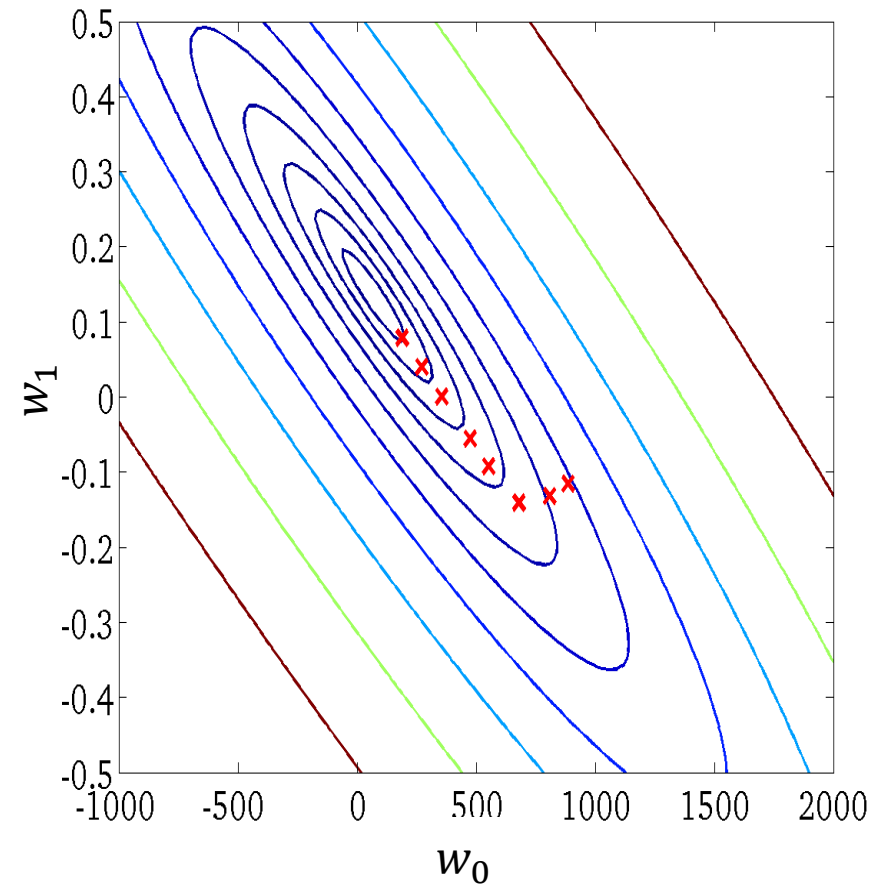


$$f(x; w_0, w_1) = w_0 + w_1 x$$

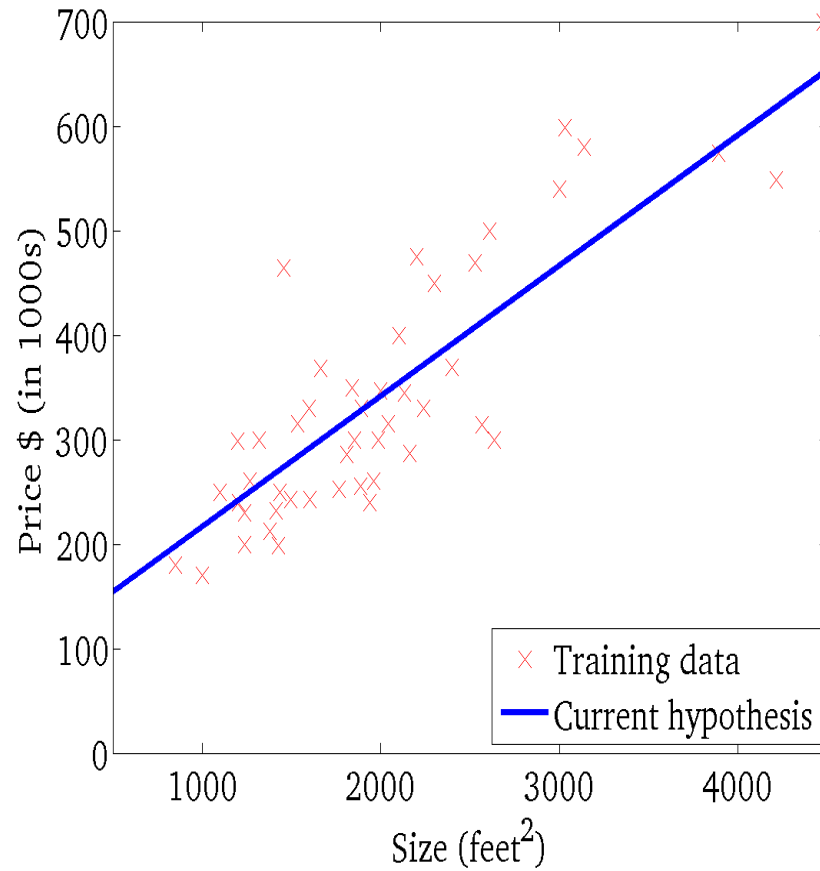


$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)

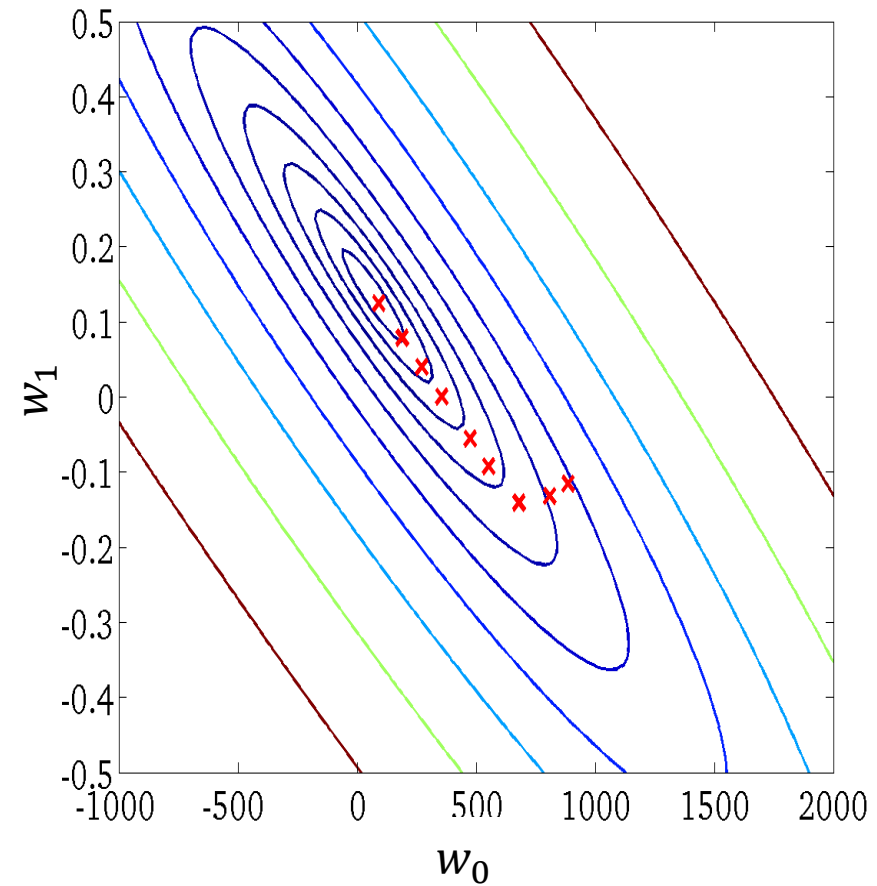


$$f(x; w_0, w_1) = w_0 + w_1 x$$



$$J(w_0, w_1)$$

(function of the parameters w_0, w_1)



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.

Stochastic gradient descent

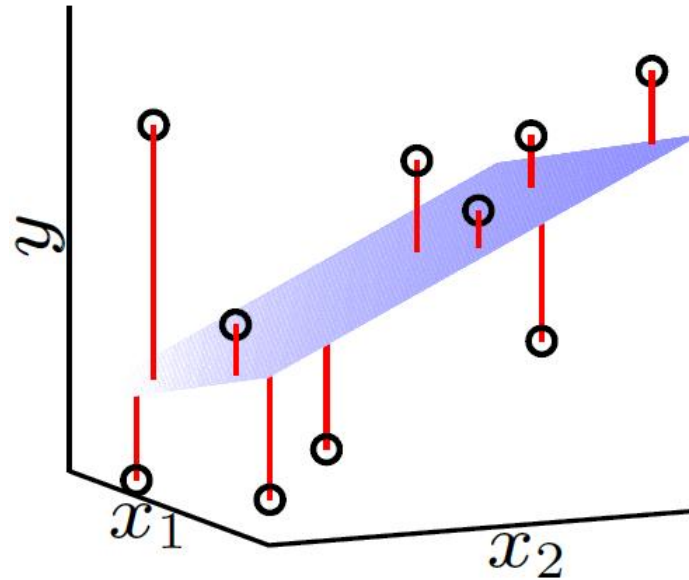
- Batch techniques process the entire training set in one go
 - thus they can be computationally costly for large data sets.
- Stochastic gradient descent: when the cost function can comprise a sum over data points:

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

- Update after presentation of a **mini-batch** S of data:

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

Linear model: multi-dimensional inputs



$$\begin{aligned} f(\mathbf{x}; \mathbf{w}) &= w_0 + w_1 x_1 + \cdots + w_d x_d \\ &= \mathbf{w}^T \mathbf{x} \end{aligned}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Generalized linear regression

- 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)

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

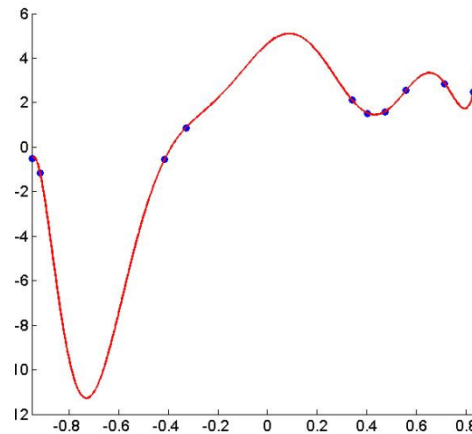
- Polynomial (univariate)

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

$$f(x; \mathbf{w}) = w_0 + w_1 x + \dots + w_{m-1} x^{m-1} + w_m x^m$$

Model complexity and overfitting

- With limited training data, models may achieve zero training error but a large test error.



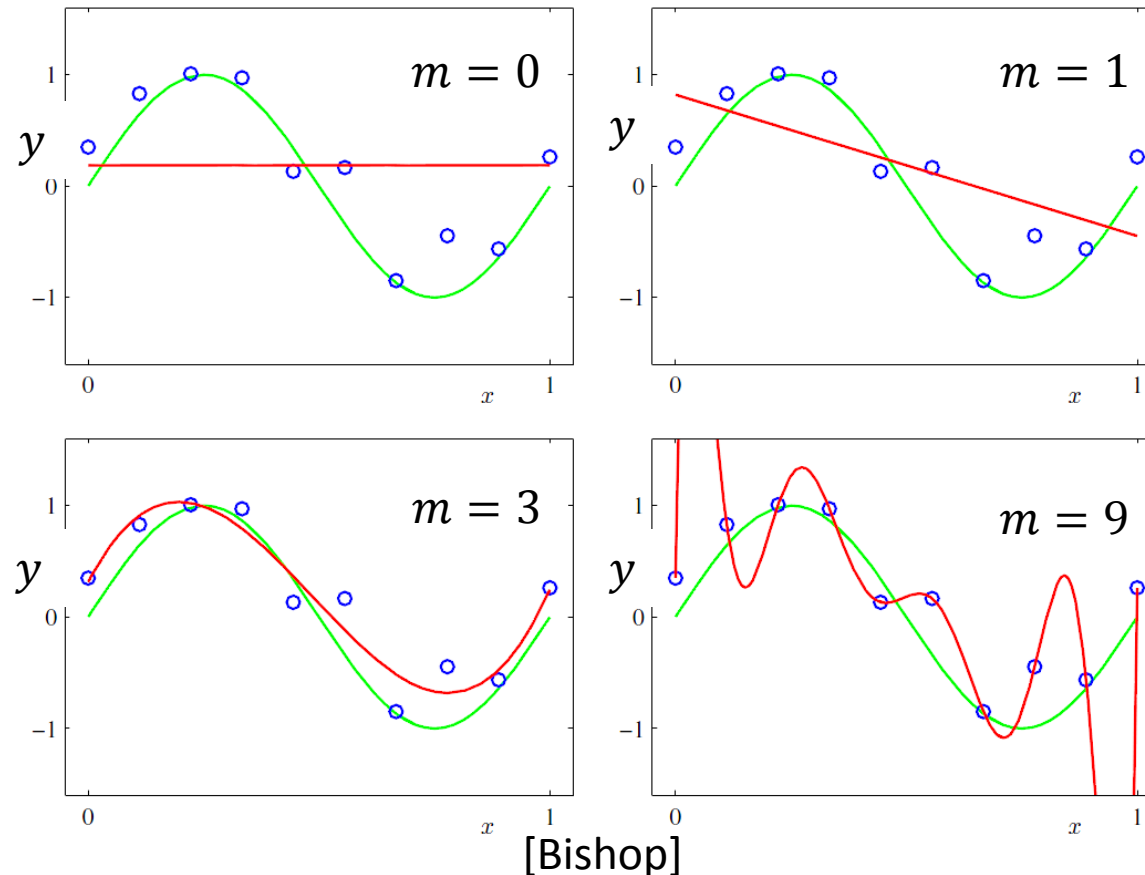
- Over-fitting: when the training loss no longer bears any relation to the test (generalization) loss.
 - Fails to generalize to unseen examples.

Over-fitting causes

- Model complexity
 - E.g., Model with a large number of parameters (degrees of freedom)
- Low number of training data
 - Small data size compared to the complexity of the model

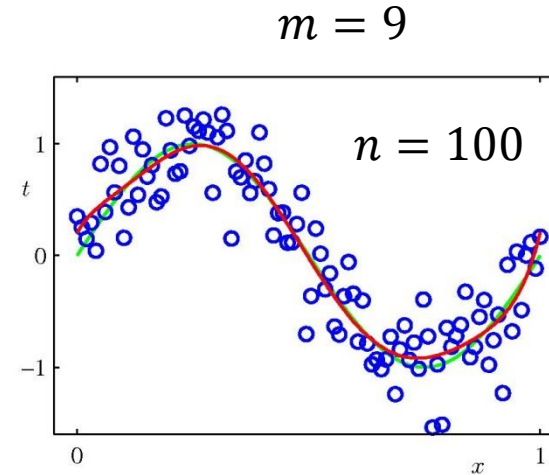
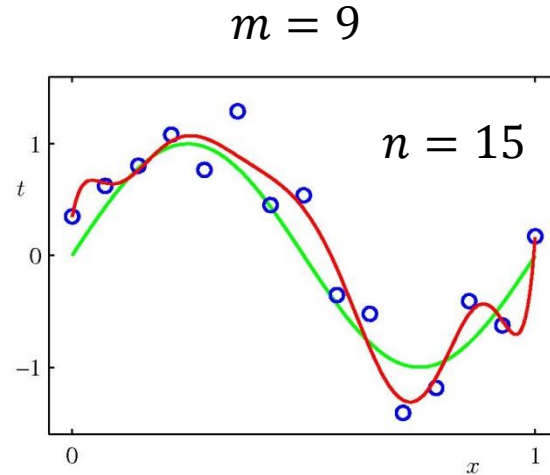
Model complexity

- Example:
 - Polynomials with larger m are becoming increasingly tuned to the random noise on the target values.



Number of training data & overfitting

- ▶ Over-fitting problem becomes less severe as the size of training data increases.



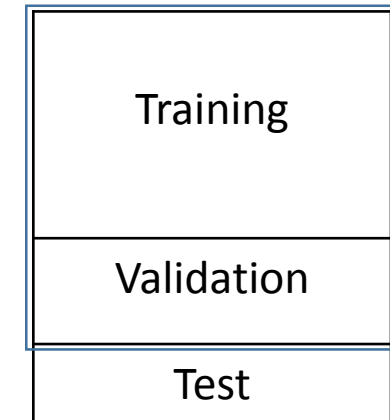
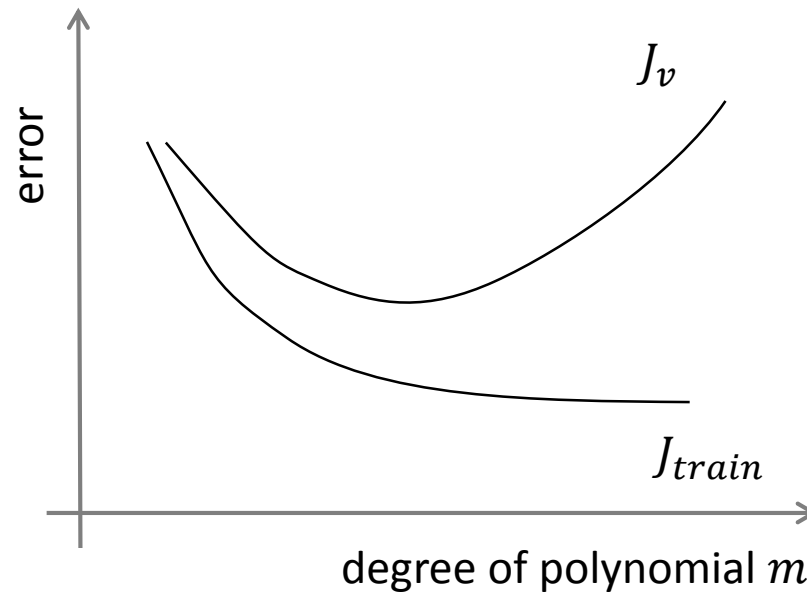
[Bishop]

Avoiding over-fitting

- Determine a suitable value for model complexity
 - **Simple hold-out method**
 - **Cross-validation**
- Regularization (Occam's Razor)
 - Explicit preference towards simpler models
 - Penalize for the model complexity in the objective function

Simple hold out: training, validation, and test sets

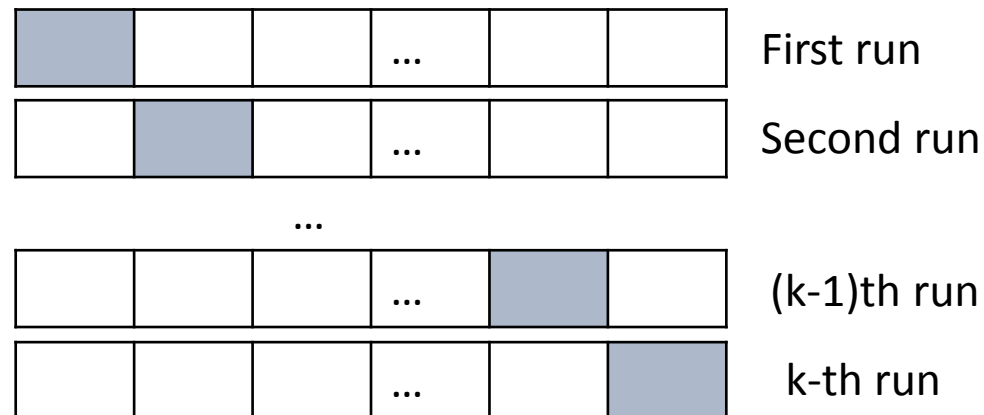
- Simple hold-out chooses the model (hyperparameters) that minimizes error on validation set.



- run on the test set once at the very end!

Cross-Validation (CV): Evaluation

- k -fold cross-validation steps:
 - Shuffle the dataset and randomly partition training data into k groups of approximately equal size
 - for $i = 1$ to k
 - Choose the i -th group as the held-out validation group
 - Train the model on all but the i -th group of data
 - Evaluate the model on the held-out group
 - Performance scores of the model from k runs are **averaged**.



Regularization

- Adding a penalty term in the cost function to discourage the coefficients from reaching large values.
- Ridge regression (weight decay):

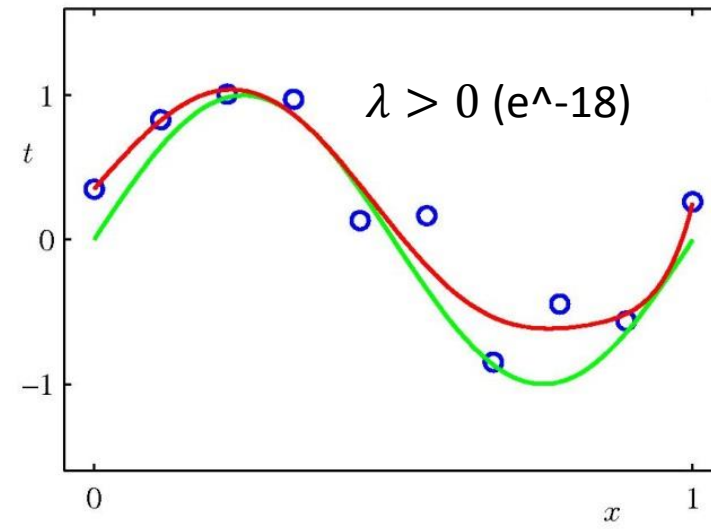
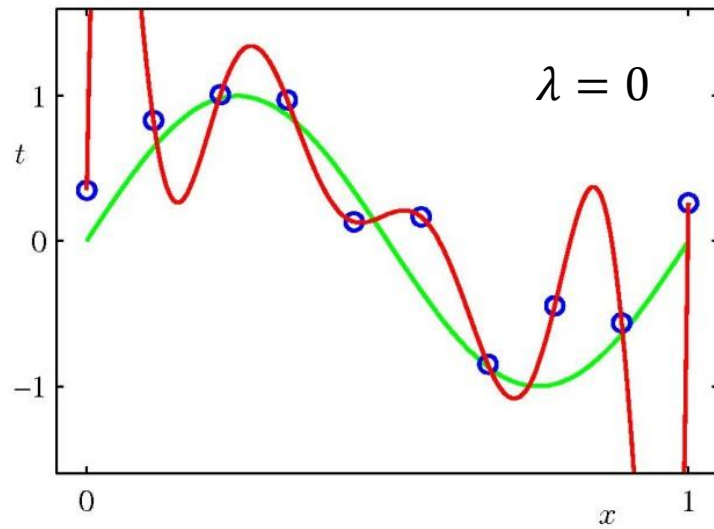
$$J(\mathbf{w}) = \underbrace{\sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(x^{(i)}) \right)^2}_{\text{Approximation: How much model predictions match training data}} + \underbrace{\lambda R(\mathbf{w})}_{\text{Generalization: prefer simple ones; Control the variance of the models}}$$

λ : regularization strength
(hyperparameter)

Generalization: prefer simple ones;
Control the variance of the models

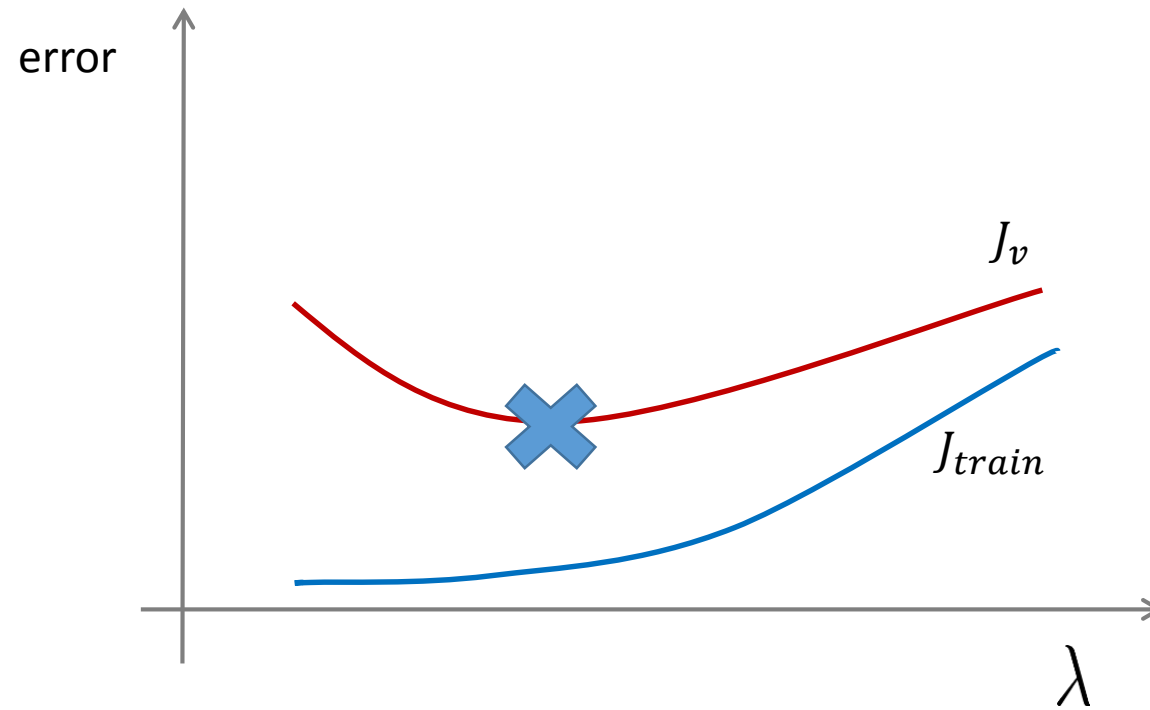
$$\text{e.g. } R(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$$

Regularization



[Bishop]

Choosing the regularization parameter



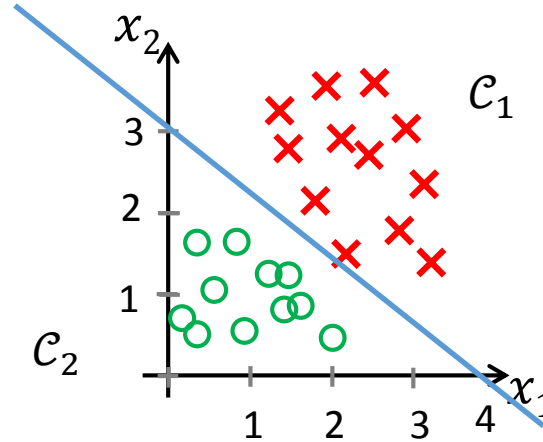
Classification problem

- Given: Training set
 - labeled set of N input-output pairs $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - $y \in \{1, \dots, K\}$
- Goal: Given an input \mathbf{x} , assign it to one of K classes
- Examples:
 - Image classification
 - Speech recognition
 - ...

Linear Classifier example

- Two class example:

$$-\frac{3}{4}x_1 - x_2 + 3 = 0$$



if $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$ then \mathcal{C}_1
else \mathcal{C}_2

$$\mathbf{w} = \begin{bmatrix} -\frac{3}{4} & -1 \end{bmatrix}$$
$$w_0 = 3$$

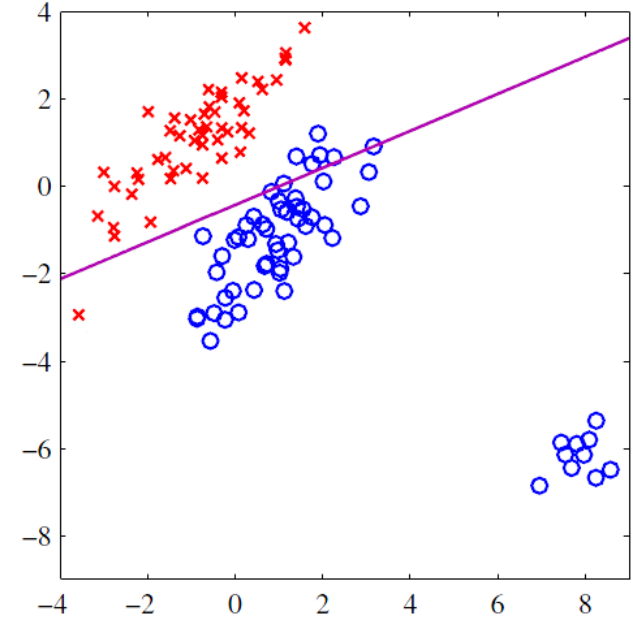
Square error loss function for classification!

$K = 2$

Square error loss is not suitable for classification:

- Least square loss penalizes ‘too correct’ predictions (that they lie a long way on the correct side of the decision)
- Least square loss also lack robustness to noise

$$J(\mathbf{w}) = \sum_{i=1}^N (wx^{(i)} + w_0 - y^{(i)})^2$$



Parametric classifier: Multiclass


- $f(\mathbf{x}; \mathbf{W}) = [f_1(\mathbf{x}, \mathbf{W}), \dots, f_K(\mathbf{x}, \mathbf{W})]^T$
- $\mathbf{W} = [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_K]$ contains one vector of parameters for each class

Parametric classifier: Linear

- $f(\mathbf{x}; \mathbf{W}) = [f_1(\mathbf{x}, \mathbf{W}), \dots, f_K(\mathbf{x}, \mathbf{W})]$
- $\mathbf{W} = [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_K]$ contains one vector of parameters for each class
 - In linear classifiers, \mathbf{W} is $d \times K$ where d shows number of features
 - $\mathbf{W}^T \mathbf{x}$ provides us a vector
- $f(\mathbf{x}; \mathbf{W})$ contains K numbers giving class scores for the input \mathbf{x}

Linear classifier

- Output obtained from $\mathbf{W}^T \mathbf{x} + \mathbf{b}$



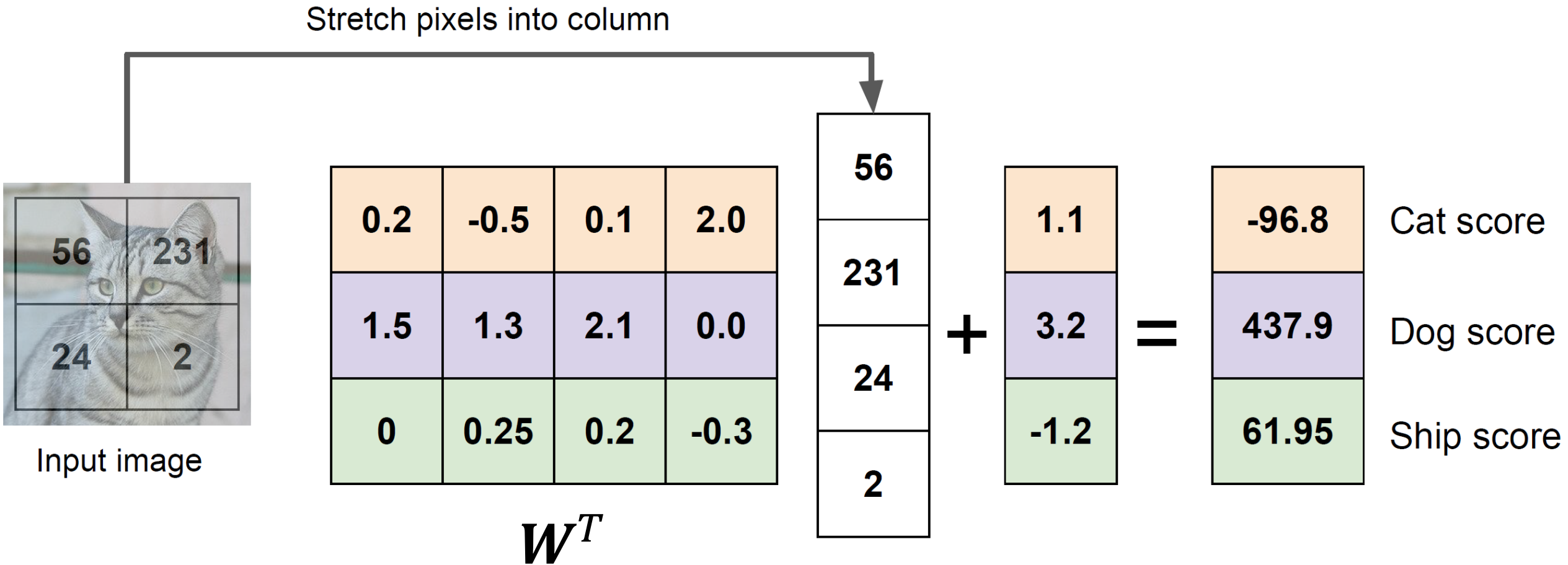
28 × 28

→ $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_{784} \end{bmatrix}$

$\mathbf{W}^T = \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_{10} \end{bmatrix}_{10 \times 784}$

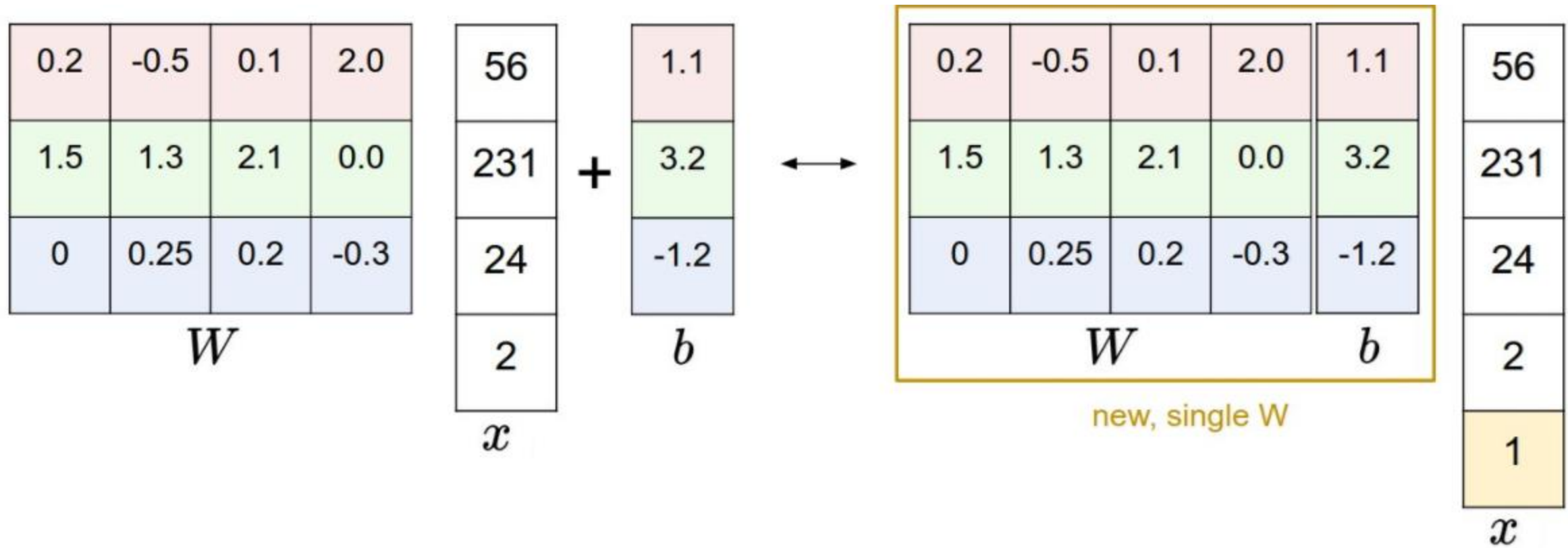
$\mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_{10} \end{bmatrix}$

Example



How can we tell whether this W and b is good or bad?

Bias can also be included in the W matrix



Multi-class SVM

$$J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N L^{(i)} + \lambda R(\mathbf{W})$$

Hinge loss:

$$L^{(i)} = \sum_{j \neq y^{(i)}} \max\left(0, 1 + s_j - s_{y^{(i)}}\right)$$

$$\begin{aligned} s_j &\equiv f_j(\mathbf{x}^{(i)}; \mathbf{W}) \\ &= \mathbf{w}_j^T \mathbf{x}^{(i)} \end{aligned}$$

$$= \sum_{j \neq y^{(i)}} \max\left(0, 1 + \mathbf{w}_j^T \mathbf{x}^{(i)} - \mathbf{w}_{y^{(i)}}^T \mathbf{x}^{(i)}\right)$$

L2 regularization:

$$R(\mathbf{W}) = \sum_{k=1}^K \sum_{l=1}^d w_{lk}^2$$

Multi-class SVM loss: Example



3 training examples, 3 classes.
With some W the scores are $W^T x$

$$s_j = \mathbf{w}_j^T \mathbf{x}^{(i)}$$

$$L^{(i)} = \sum_{j \neq y^{(i)}} \max(0, 1 + s_j - s_{y^{(i)}})$$

$$\frac{1}{N} \sum_{i=1}^N L^{(i)} = \frac{1}{3} (2.9 + 0 + 12.9) = 5.7$$

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

$$\begin{aligned} L^{(1)} &= \max(0, 1 + 5.1 - 3.2) \\ &\quad + \max(0, 1 - 1.7 - 3.2) \\ &= \max(0, 2.9) + \max(0, -3.9) \\ &= 2.9 + 0 \end{aligned}$$

$$\begin{aligned} L^{(2)} &= \max(0, 1 + 1.3 - 4.9) \\ &\quad + \max(0, 1 + 2 - 4.9) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \end{aligned}$$

$$\begin{aligned} L^{(3)} &= \max(0, 2.2 - (-3.1) + 1) \\ &\quad + \max(0, 2.5 - (-3.1) + 1) \\ &= \max(0, 6.3) + \max(0, 6.6) \\ &= 6.3 + 6.6 = 12.9 \end{aligned}$$

Some questions?

$$L^{(i)} = \sum_{j \neq y^{(i)}} \max(0, 1 + s_j - s_{y^{(i)}})$$

- Q1: What if the sum was over all classes? (including $j = y_i$)
- Q2: What if we used mean instead of sum?
- Q3: What if we used $L^{(i)} = \sum_{j \neq y^{(i)}} \max(0, 1 + s_j - s_{y^{(i)}})^2$?
- Q4: what is the min/max possible?
- Q5: why do we use regularization term?

Other regularization terms

- **L2 regularization**
- L1 regularization
- Elastic net (L1 + L2)

$$\sum_{k=1}^K \sum_{l=1}^d w_{lk}^2$$

$$\sum_{k=1}^K \sum_{l=1}^d |w_{lk}|$$

$$\beta \sum_{k=1}^K \sum_{l=1}^d w_{lk}^2 + \sum_{k=1}^K \sum_{l=1}^d |w_{lk}|$$

Softmax Classifier (Multinomial Logistic Regression)

softmax function

$$P(Y = k | X = \mathbf{x}^{(i)}) = \frac{e^{s_k}}{\sum_{j=1}^K e^{s_j}} \quad s_k = f_k(\mathbf{x}^{(i)}; W) = \mathbf{w}_k^T \mathbf{x}^{(i)}$$

- Maximum log likelihood is equivalent to minimize the negative of log likelihood of the correct class:

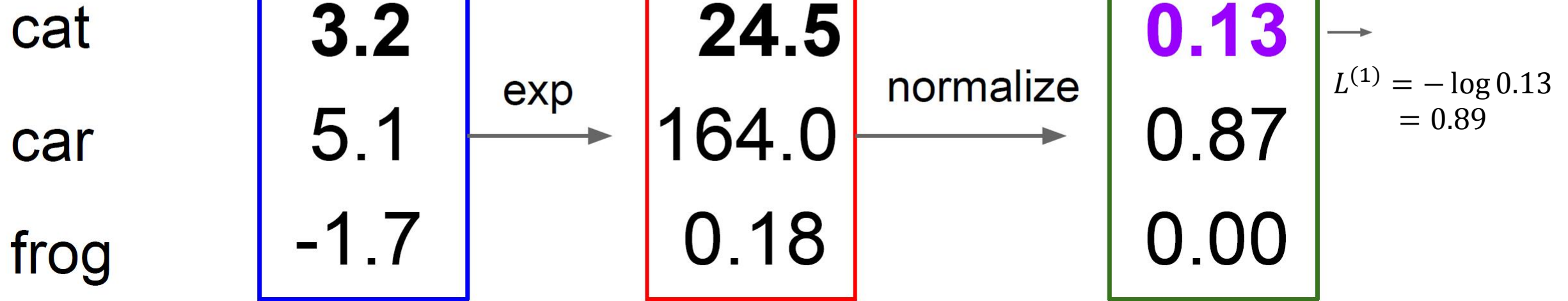
$$\begin{aligned} L^{(i)} &= -\log P(Y = y^{(i)} | X = \mathbf{x}^{(i)}) \\ &= -s_{y^{(i)}} + \log \sum_{j=1}^K e^{s_j} \end{aligned} \quad \text{Cross-entropy loss}$$

Softmax classifier loss: example



$$L^{(i)} = -\log \frac{e^{s_{y^{(i)}}}}{\sum_{j=1}^K e^{s_j}}$$

unnormalized probabilities



Cross entropy

$$H(q, p) = - \sum_x q(x) \log p(x)$$

- For the loss of the softmax classifier:
 - p: estimated class probabilities $\frac{e^{s_k}}{\sum_{j=1}^K e^{s_j}}$
 - q: the true distribution
 - all probability mass is on the correct class $q(Y = y^{(i)}) = 1$ ($q(Y \neq y^{(i)}) = 0$).

Relation to KL divergence

$$H(q, p) = H(q) + D_{KL}(q||p)$$

- Since $H(q)$ for the loss of softmax classifier is zero:
 - Minimizing cross entropy is equivalent to minimizing the KL divergence between the two distributions (a measure of distance).
 - cross-entropy loss wants the predicted distribution to have all of its mass on the correct answer.

Recap

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \quad \text{Softmax}$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W) \quad \text{Full loss}$$

We need $\nabla_W L$ to update weights

Resources

- Deep Learning Book, Chapter 5.
- Please see the following notes:
 - <http://cs231n.github.io/linear-classify/>
 - <http://cs231n.github.io/optimization-1/>