

Linear Regression and Classification

Outline

- I. Line fitting and gradient descent
- II. Multivariable linear regression
- III. Linear classifiers
- IV. Logistic regression

I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

I. Linear Regression

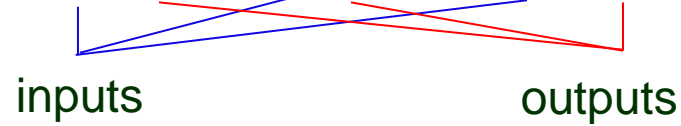
Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

inputs

A blue line originates from the word 'inputs' and extends horizontally to the right, ending under the x_N term of the data points. A short vertical blue line segment drops from the end of this horizontal line to the x_1 term, indicating that the x values are the inputs.

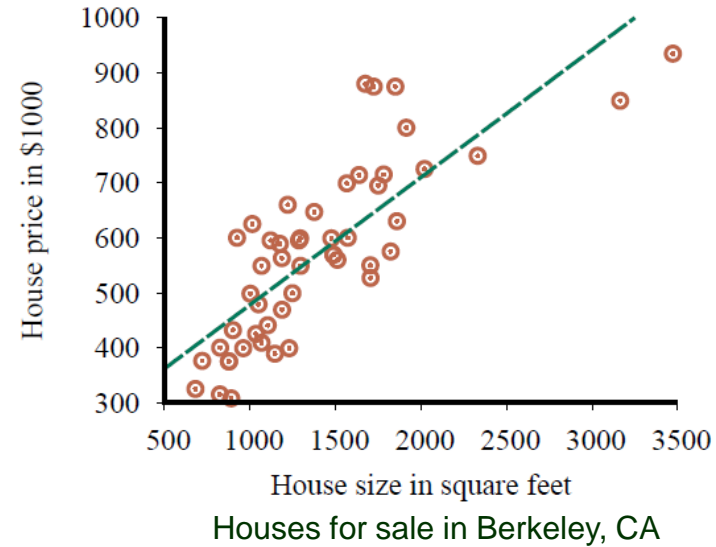
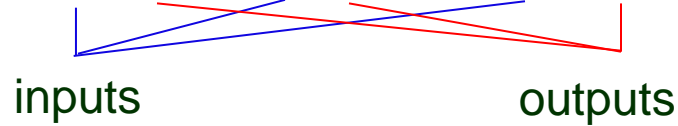
I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



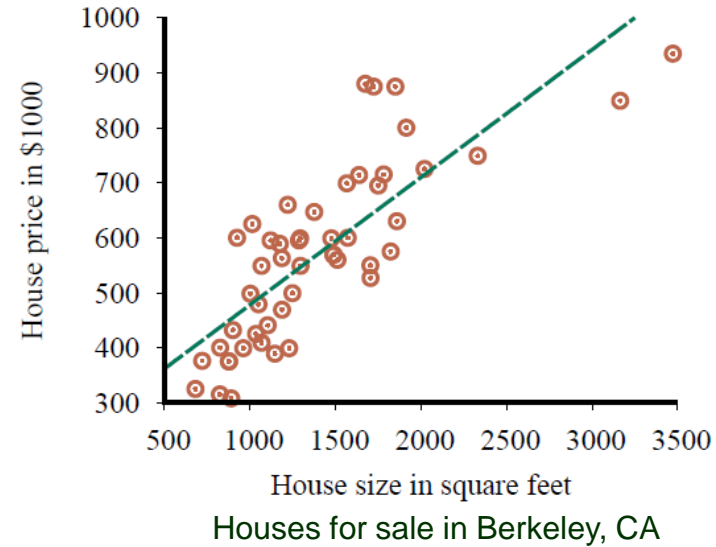
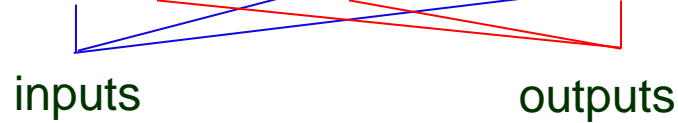
I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

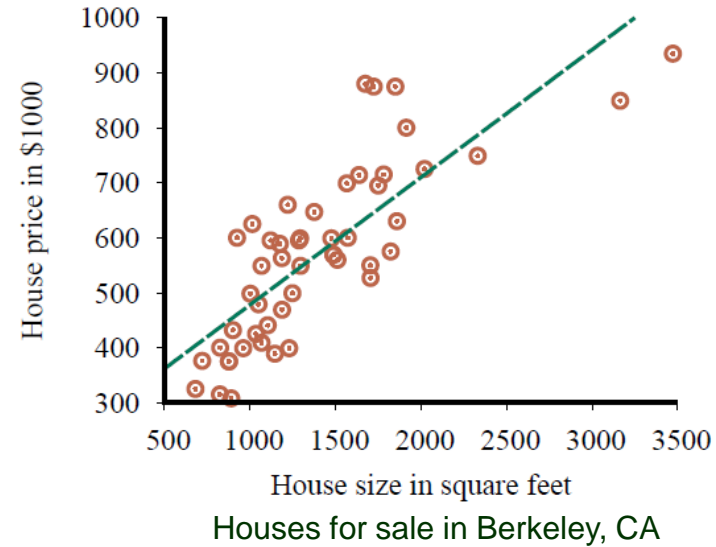
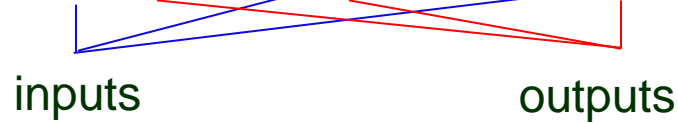


Hypothesis space: univariate linear functions.

$$h_{\mathbf{w}}(x) \equiv w_1 x + w_0$$

I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



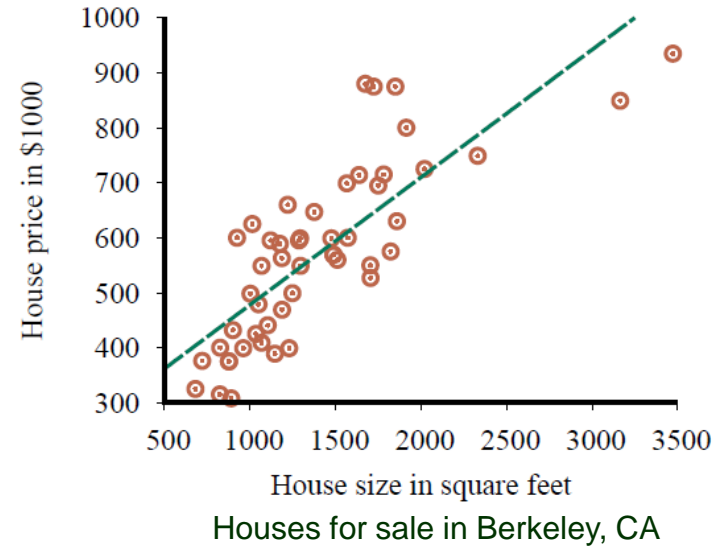
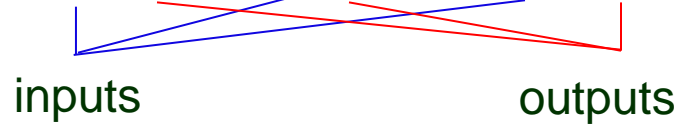
Hypothesis space: univariate linear functions.

$$h_{\mathbf{w}}(x) \equiv w_1 x + w_0$$

(w_0, w_1)

I. Linear Regression

Data points: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$



Hypothesis space: univariate linear functions.

$$h_{\mathbf{w}}(x) \equiv w_1 x + w_0$$

(w_0, w_1)

Linear regression: Find the $h_{\mathbf{w}}$ that best fits the data.

Line Fitting

We find the weights (w_0, w_1) that minimizes the empirical loss.

Use the squared-error loss $L_2(y, h_w) = (y - h_w)^2$, summed over all the points.

$$Loss(h_w) = \sum_{j=1}^N L_2(y_j, h_w(x_j))$$

$$= \sum_{j=1}^N (y_j - h_w(x_j))^2$$

Line Fitting

We find the weights (w_0, w_1) that minimizes the empirical loss.

Use the squared-error loss $L_2(y, h_w) = (y - h_w)^2$, summed over all the points.

$$\begin{aligned} \text{Loss}(h_w) &= \sum_{j=1}^N L_2(y_j, h_w(x_j)) \\ &= \sum_{j=1}^N (y_j - h_w(x_j))^2 \\ &= \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 \end{aligned}$$

Line Fitting

We find the weights (w_0, w_1) that minimizes the empirical loss.

Use the squared-error loss $L_2(y, h_w) = (y - h_w)^2$, summed over all the points.

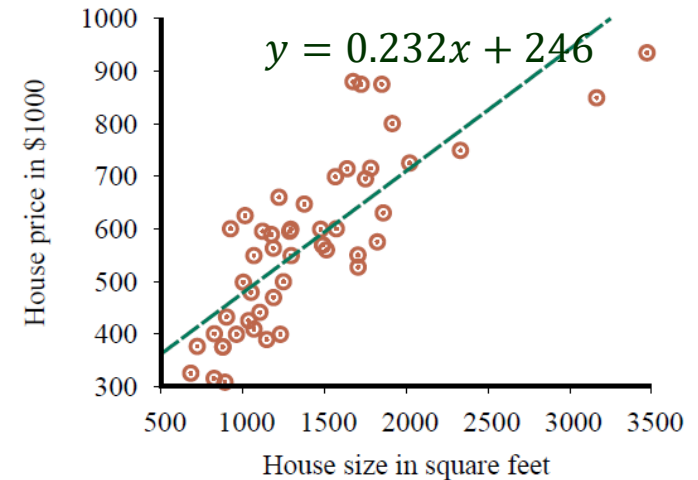
$$\begin{aligned} \text{Loss}(h_w) &= \sum_{j=1}^N L_2(y_j, h_w(x_j)) \\ &= \sum_{j=1}^N (y_j - h_w(x_j))^2 \\ &= \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 \end{aligned}$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \text{Loss}(h_w)$$

Vanishing of Partial Derivatives

At the minimizing \mathbf{w} , the gradient of $Loss(h_{\mathbf{w}})$ must vanish:

$$\nabla Loss(h_{\mathbf{w}}) = \left(\frac{\partial Loss}{\partial w_0}, \frac{\partial Loss}{\partial w_1} \right) = 0$$



Vanishing of Partial Derivatives

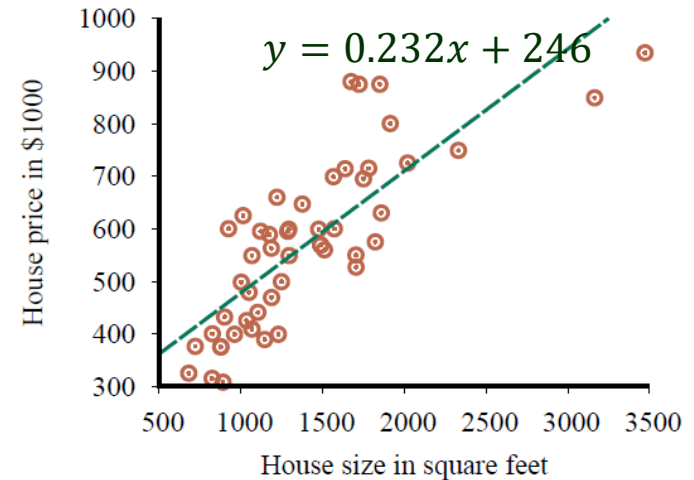
At the minimizing \mathbf{w} , the gradient of $\text{Loss}(h_{\mathbf{w}})$ must vanish:

$$\nabla \text{Loss}(h_{\mathbf{w}}) = \left(\frac{\partial \text{Loss}}{\partial w_0}, \frac{\partial \text{Loss}}{\partial w_1} \right) = 0$$



$$\frac{\partial \text{Loss}}{\partial w_0} = \frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$



Vanishing of Partial Derivatives

At the minimizing \mathbf{w} , the gradient of $\text{Loss}(h_{\mathbf{w}})$ must vanish:

$$\nabla \text{Loss}(h_{\mathbf{w}}) = \left(\frac{\partial \text{Loss}}{\partial w_0}, \frac{\partial \text{Loss}}{\partial w_1} \right) = 0$$



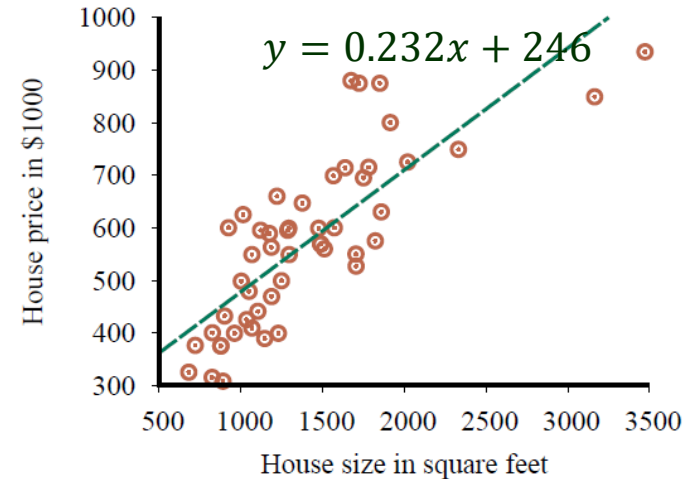
$$\frac{\partial \text{Loss}}{\partial w_0} = \frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$



$$w_1 = \frac{N \sum_{j=1}^N x_j y_j - (\sum_{j=1}^N x_j) \cdot (\sum_{j=1}^N y_j)}{N(\sum_{j=1}^N x_j^2) - (\sum_{j=1}^N x_j)^2}$$

$$w_0 = \frac{1}{N} \left(\sum_{j=1}^N y_j - w_1 \sum_{j=1}^N x_j \right)$$



Vanishing of Partial Derivatives

At the minimizing \mathbf{w} , the gradient of $Loss(h_{\mathbf{w}})$ must vanish:

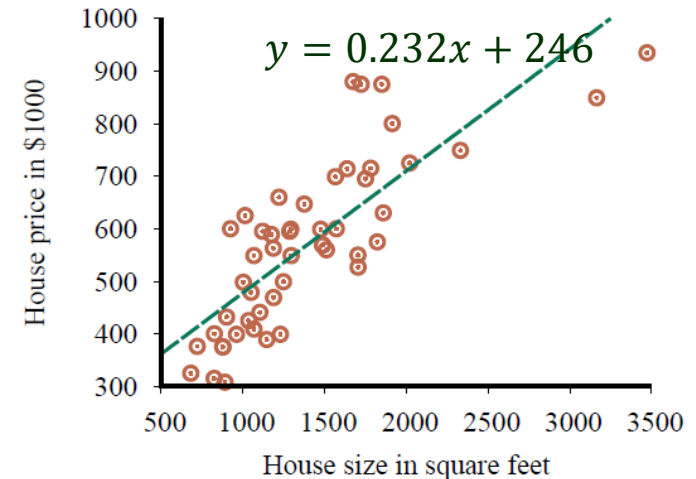
$$\nabla Loss(h_{\mathbf{w}}) = \left(\frac{\partial Loss}{\partial w_0}, \frac{\partial Loss}{\partial w_1} \right) = 0$$

$$\frac{\partial Loss}{\partial w_0} = \frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$

$$\frac{\partial Loss}{\partial w_1} = \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0$$

$$w_1 = \frac{N \sum_{j=1}^N x_j y_j - (\sum_{j=1}^N x_j) \cdot (\sum_{j=1}^N y_j)}{N(\sum_{j=1}^N x_j^2) - (\sum_{j=1}^N x_j)^2}$$

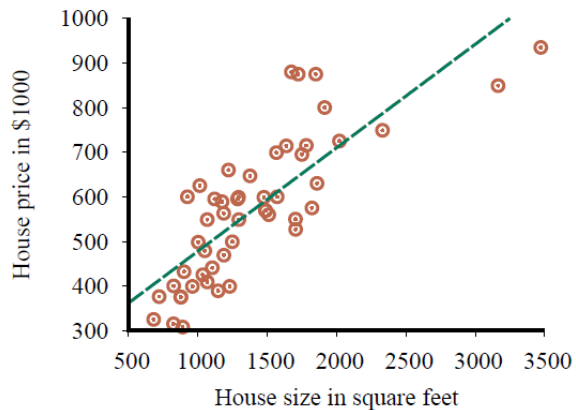
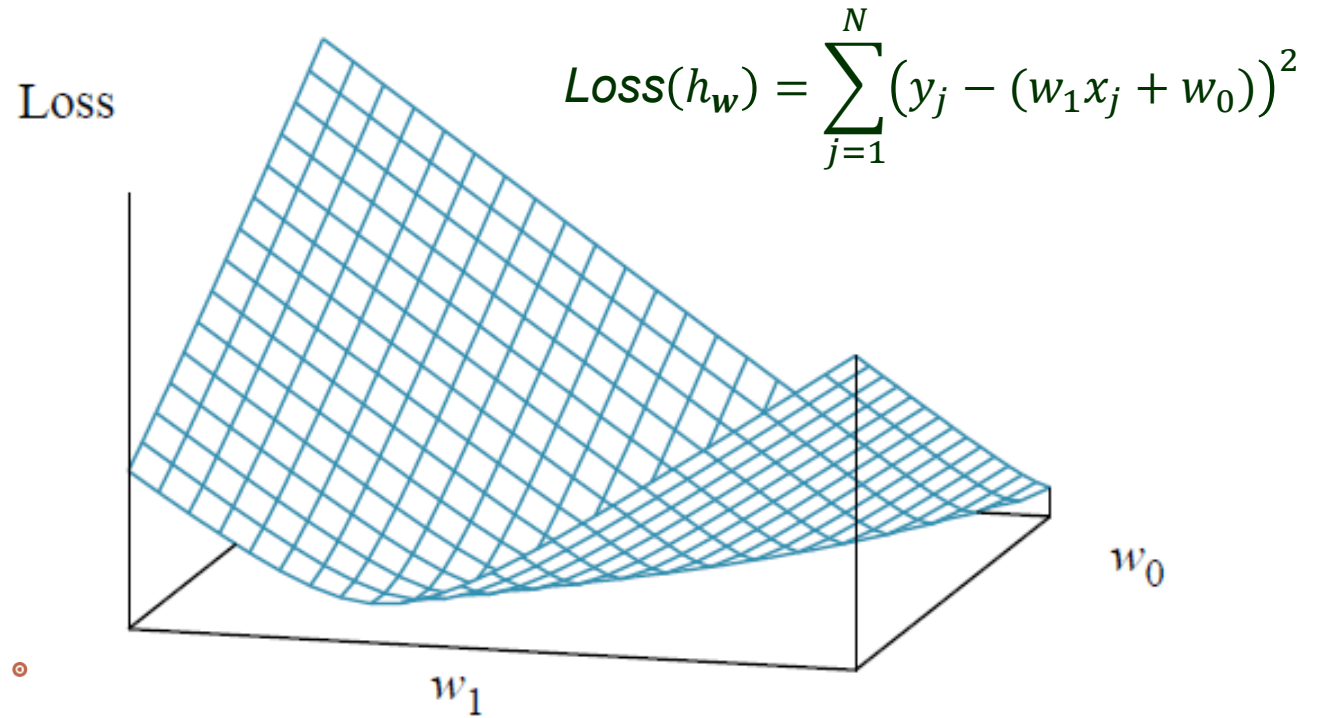
$$w_0 = \frac{1}{N} \left(\sum_{j=1}^N y_j - w_1 \sum_{j=1}^N x_j \right)$$



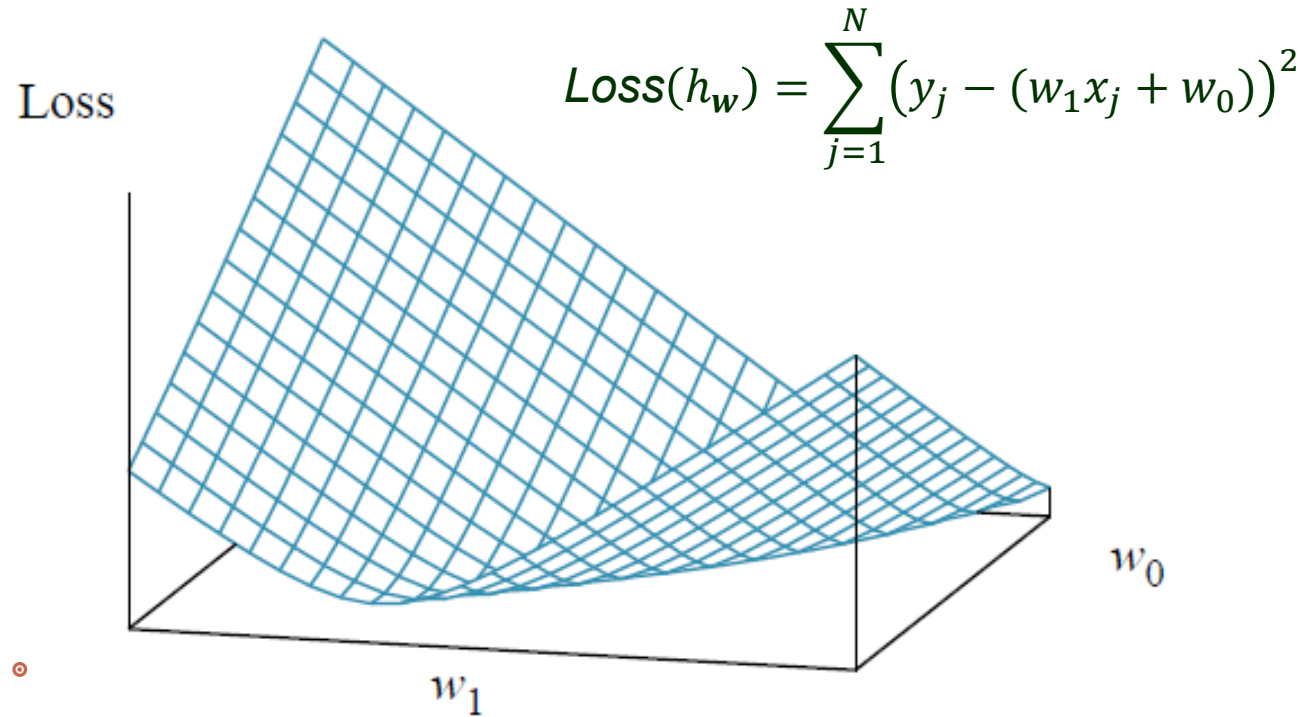
Note: Linear regression (LR) is not appropriate for fitting a line to data points *geometrically* when their x- and y-coordinates have the same physical units. The best-fitting line in LR does **not** minimize the sum of squares of distances of the data points to the line. Here, the distance in the y-direction is used to approximate the point-line distance. Such approximation gets worse the larger the line's slope is, in which case fitting will also become sensitive to small changes in point locations. In addition, the model $h_{\mathbf{w}}(x) \equiv w_1 x + w_0$ cannot represent a vertical line, hence a loss of generality.

The best line fitting method, used in computer vision for extracting straight edges from an image, uses the general line equation $ax + by + c = 0$ subject to the constraint that $a^2 + b^2 = 1$, applying the method of Lagrange multipliers. (See Section 3 of the [COMS 4770/5770 notes on least-squares fitting](#)).

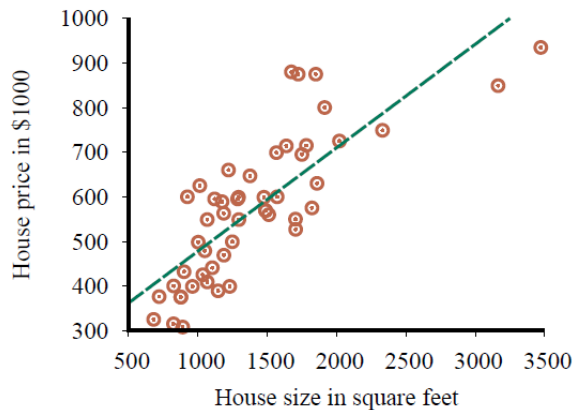
Plot of the Loss Function



Plot of the Loss Function



◆ Convex function with no local minima.



Gradient Descent

- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ♦ Instead, the method of gradient descent is used:

Gradient Descent

- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ◆ Instead, the method of gradient descent is used:
 - Start at a point \mathbf{w} in the weight space.
 - Compute an estimate of the gradient of the loss function.
 - Move a small amount in the direction of the negative gradient, i.e., the steepest downhill direction.
 - Repeat until convergence on a point with (local) minimum loss.

Gradient Descent

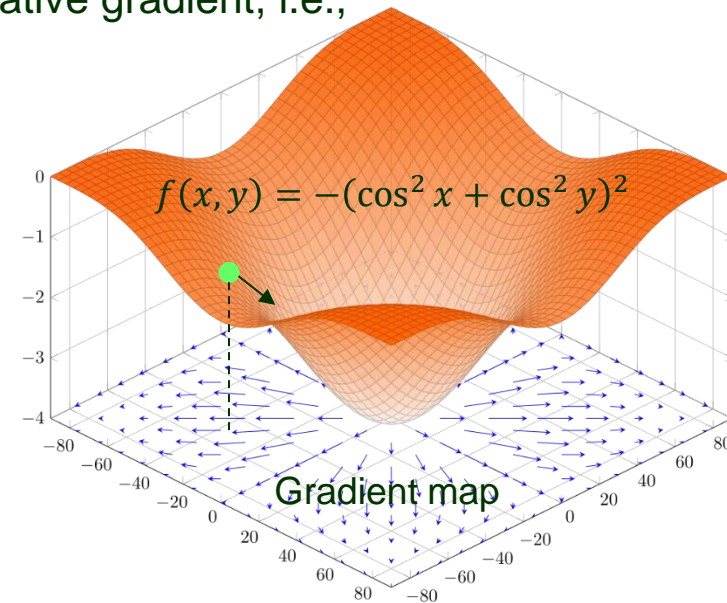
- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ◆ Instead, the method of gradient descent is used:
 - Start at a point \mathbf{w} in the weight space.
 - Compute an estimate of the gradient of the loss function.
 - Move a small amount in the direction of the negative gradient, i.e., the steepest downhill direction.
 - Repeat until convergence on a point with (local) minimum loss.

$\mathbf{w} \leftarrow$ any point in the parameter space

while not converged do

 for each w_i in \mathbf{w} do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w})$$



Gradient Descent

- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ◆ Instead, the method of gradient descent is used:
 - Start at a point \mathbf{w} in the weight space.
 - Compute an estimate of the gradient of the loss function.
 - Move a small amount in the direction of the negative gradient, i.e., the steepest downhill direction.
 - Repeat until convergence on a point with (local) minimum loss.

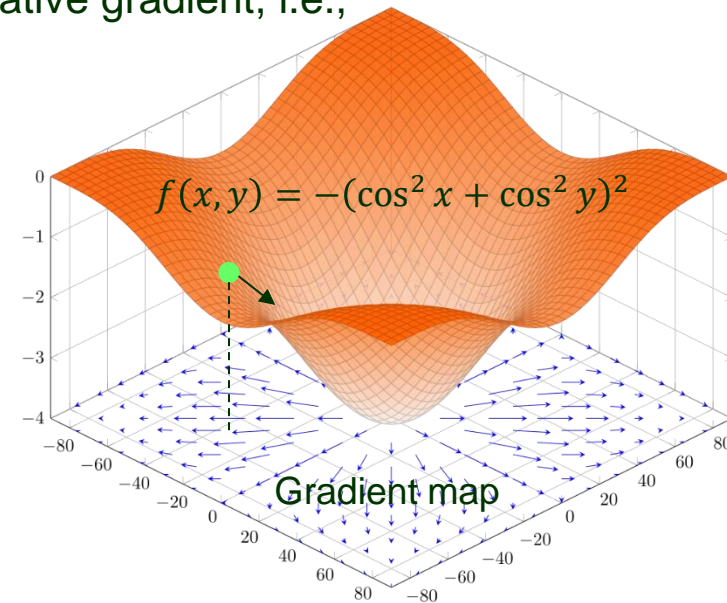
$\mathbf{w} \leftarrow$ any point in the parameter space

while not converged do

for each w_i in \mathbf{w} do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w})$$

step size or *learning rate*



Gradient Descent

- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ♦ Instead, the method of gradient descent is used:
 - Start at a point \mathbf{w} in the weight space.
 - Compute an estimate of the gradient of the loss function.
 - Move a small amount in the direction of the negative gradient, i.e., the steepest downhill direction.
 - Repeat until convergence on a point with (local) minimum loss.

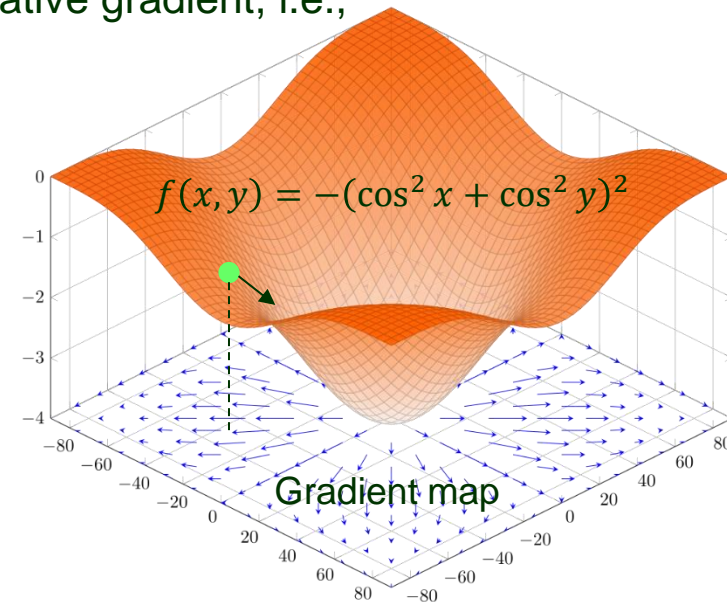
$\mathbf{w} \leftarrow$ any point in the parameter space

while not converged do

for each w_i in \mathbf{w} do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w})$$

step size or *learning rate*



* Section 19.6.2 applies gradient descent to a quadratic loss function, which defeats the purpose since the gradient $\nabla \text{Loss}(h_{\mathbf{w}})$ is linear in \mathbf{w} whose values can be easily determined from solving the linear system $\nabla \text{Loss}(h_{\mathbf{w}}) = 0$.

Gradient Descent

- ♠ For a complex loss function, vanishing of its gradient often results in a system of nonlinear equations in \mathbf{w} that does not have a closed-form solution.
- ♦ Instead, the method of gradient descent is used:
 - Start at a point \mathbf{w} in the weight space.
 - Compute an estimate of the gradient of the loss function.
 - Move a small amount in the direction of the negative gradient, i.e., the steepest downhill direction.
 - Repeat until convergence on a point with (local) minimum loss.

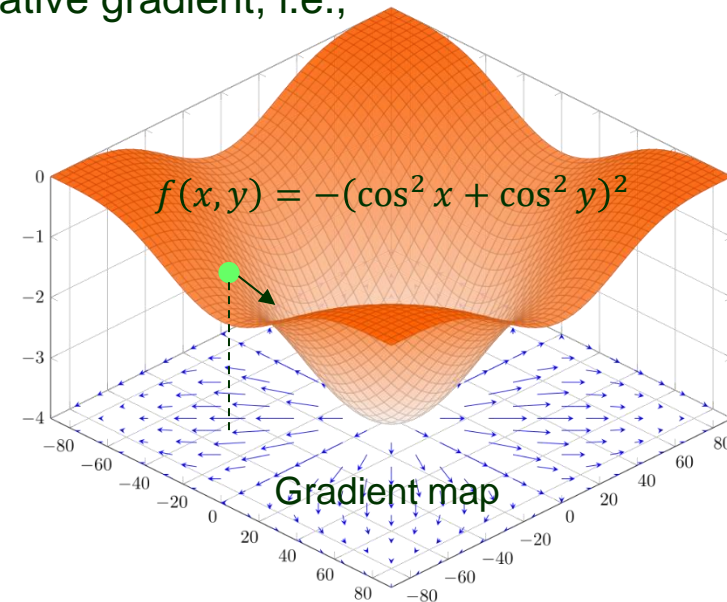
$\mathbf{w} \leftarrow$ any point in the parameter space

while not converged do

for each w_i in \mathbf{w} do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w})$$

step size or *learning rate*



* Section 19.6.2 applies gradient descent to a quadratic loss function, which defeats the purpose since the gradient $\nabla \text{Loss}(h_{\mathbf{w}})$ is linear in \mathbf{w} whose values can be easily determined from solving the linear system $\nabla \text{Loss}(h_{\mathbf{w}}) = 0$.

** To see how gradient descent works, see Section 4 of the [COMS 4770/5770 notes on nonlinear optimization](#).

Multivariable Linear Regression

- ♣ An example is represented by an n -vector $\mathbf{x}_j = (x_{j,1}, \dots, x_{j,n})$.
- ♣ Hypothesis space:

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_nx_n = w_0 + \sum_{i=1}^n w_ix_i$$

Multivariable Linear Regression

- ♣ An example is represented by an n -vector $\mathbf{x}_j = (x_{j,1}, \dots, x_{j,n})$.
- ♣ Hypothesis space:

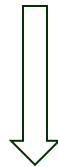
$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_nx_n = w_0 + \sum_{i=1}^n w_ix_i$$

For convenience, we extend \mathbf{x} by adding $x_0 = 1$ such that $\mathbf{x} = (1, x_1, \dots, x_n)$.

Multivariable Linear Regression

- ♣ An example is represented by an n -vector $\mathbf{x}_j = (x_{j,1}, \dots, x_{j,n})$.
- ♣ Hypothesis space:

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_nx_n = w_0 + \sum_{i=1}^n w_ix_i$$



For convenience, we extend \mathbf{x} by adding $x_0 = 1$ such that $\mathbf{x} = (1, x_1, \dots, x_n)$.

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

Multivariable Linear Regression

- ♣ An example is represented by an n -vector $\mathbf{x}_j = (x_{j,1}, \dots, x_{j,n})$.
- ♣ Hypothesis space:

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_nx_n = w_0 + \sum_{i=1}^n w_ix_i$$



For convenience, we extend \mathbf{x} by adding $x_0 = 1$ such that $\mathbf{x} = (1, x_1, \dots, x_n)$.

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

- ♣ Best weight vector:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \sum_j L_2(y_j, \mathbf{w} \cdot \mathbf{x}_j)$$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

$$0 = \nabla L(\mathbf{w}) = \nabla((\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}))$$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

$$0 = \nabla L(\mathbf{w}) = \nabla((\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}))$$



$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0$$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

$$0 = \nabla L(\mathbf{w}) = \nabla((\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}))$$



$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0$$

\mathbf{X} almost always has full rank since $m \gg n$

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

$$0 = \nabla L(\mathbf{w}) = \nabla((\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}))$$



$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0$$



\mathbf{X} almost always has full rank since $m \gg n$

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

Optimal Weights

- Write \mathbf{w} as a column vector, i.e., $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$.
- Vector of m outputs: $\mathbf{y} = (y_1, y_2, \dots, y_m)^T$.
- *Data matrix* ($m \times n$): $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix}$
- Predicted outputs: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Loss over all the training data: $L(\mathbf{w}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$

$$0 = \nabla L(\mathbf{w}) = \nabla((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}))$$



$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0$$



\mathbf{X} almost always has full rank since $m \gg n$

$$\mathbf{w}^* = \mathbf{w} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}_{\text{pseudoinverse of } \mathbf{X}}$$

pseudoinverse of \mathbf{X}

Regularization

Commonly applied on multivariable linear function to avoid overfitting.

$$\text{Cost}(h_{\mathbf{w}}) = \text{EmpLoss}(h_{\mathbf{w}}) + \lambda \text{Complexity}(h_{\mathbf{w}})$$

where

$$\text{Complexity}(h_{\mathbf{w}}) = L_q(\mathbf{w}) = \sum_{i=1}^n |w_i|^q$$

Regularization

Commonly applied on multivariable linear function to avoid overfitting.

$$\text{Cost}(h_{\mathbf{w}}) = \text{EmpLoss}(h_{\mathbf{w}}) + \lambda \text{Complexity}(h_{\mathbf{w}})$$

where

$$\text{Complexity}(h_{\mathbf{w}}) = L_q(\mathbf{w}) = \sum_{i=1}^n |w_i|^q$$

- ◆ L_1 (with $q = 1$) regularization tends to produce a sparse model (in which many weights are set to zero) because it takes the w_0, w_1, \dots, w_n axes seriously.

Regularization

Commonly applied on multivariable linear function to avoid overfitting.

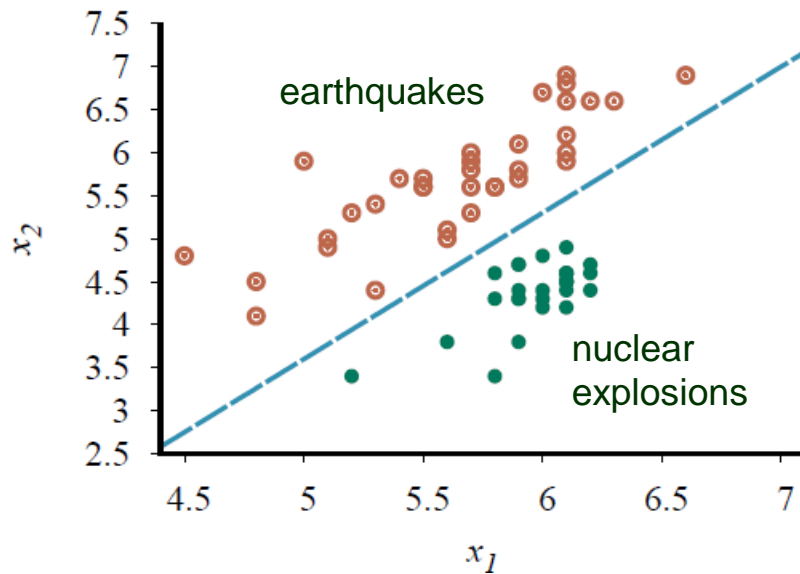
$$\text{Cost}(h_{\mathbf{w}}) = \text{EmpLoss}(h_{\mathbf{w}}) + \lambda \text{Complexity}(h_{\mathbf{w}})$$

where

$$\text{Complexity}(h_{\mathbf{w}}) = L_q(\mathbf{w}) = \sum_{i=1}^n |w_i|^q$$

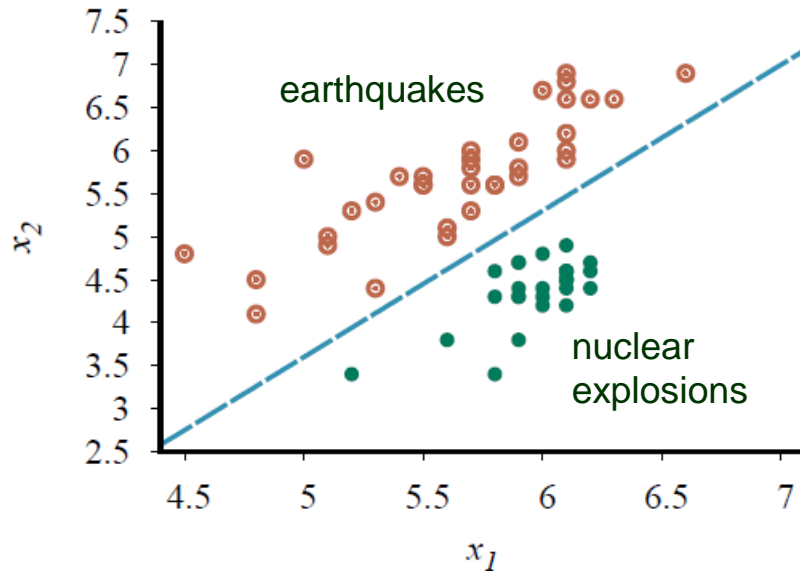
- ◆ L_1 (with $q = 1$) regularization tends to produce a sparse model (in which many weights are set to zero) because it takes the w_0, w_1, \dots, w_n axes seriously.
- ♠ L_2 (with $q = 2$) regularization takes the dimension axes arbitrarily.

III. Linear Classifiers

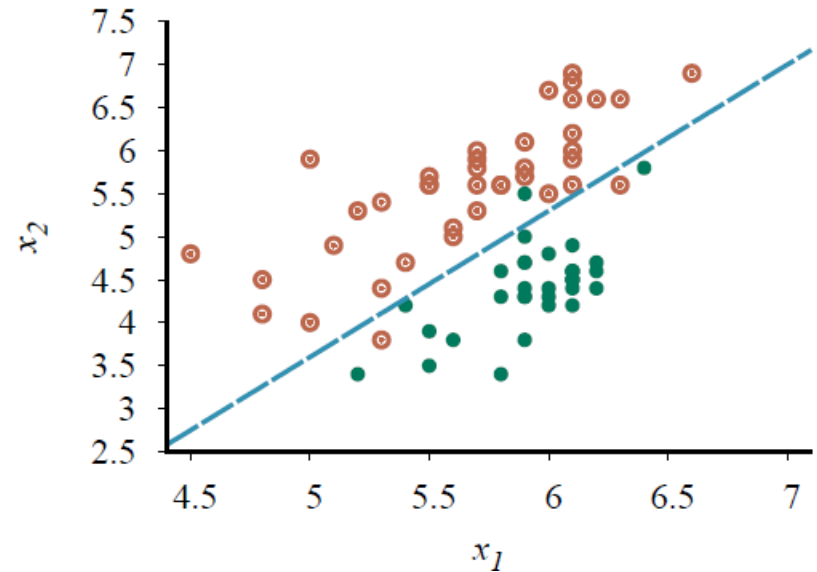


Seismic data for earthquakes and nuclear explosions:
 x_1 and x_2 respectively refer to body and surface wave magnitudes computed from the seismic signal.

III. Linear Classifiers

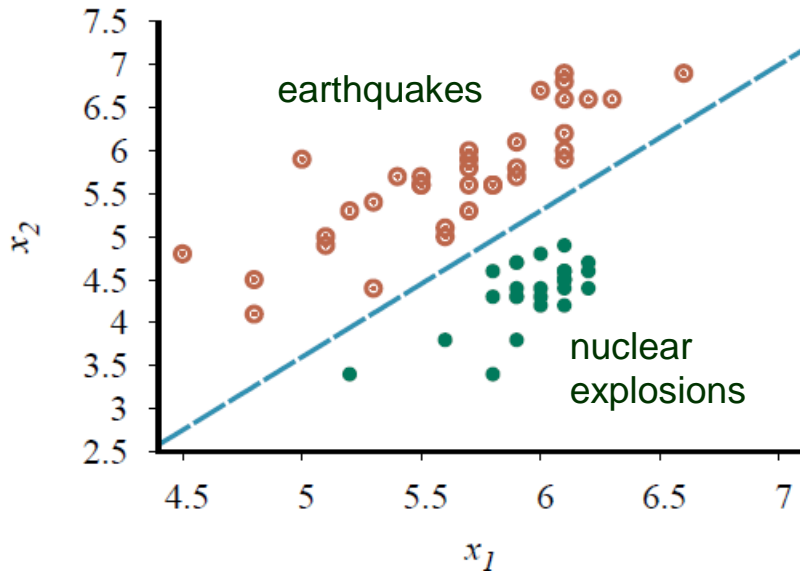


Seismic data for earthquakes and nuclear explosions: x_1 and x_2 respectively refer to body and surface wave magnitudes computed from the seismic signal.

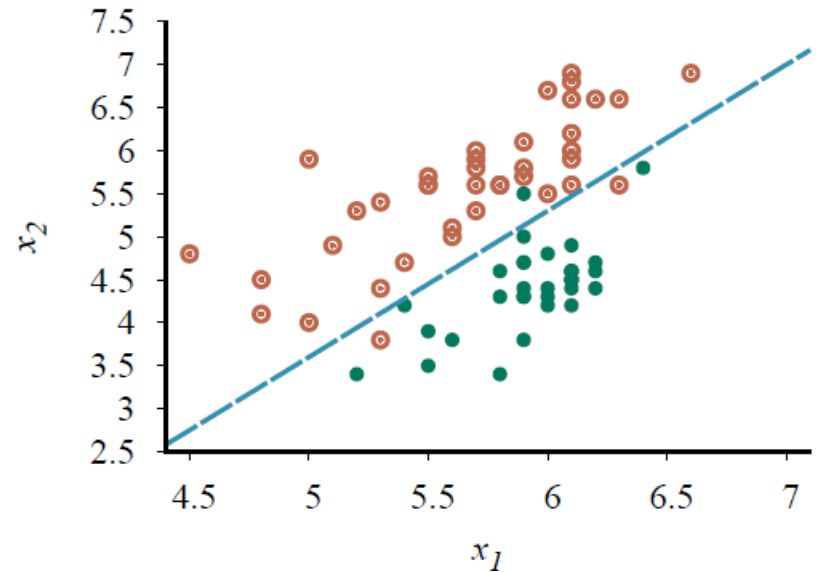


Same domain with more data points.

III. Linear Classifiers



Seismic data for earthquakes and nuclear explosions: x_1 and x_2 respectively refer to body and surface wave magnitudes computed from the seismic signal.



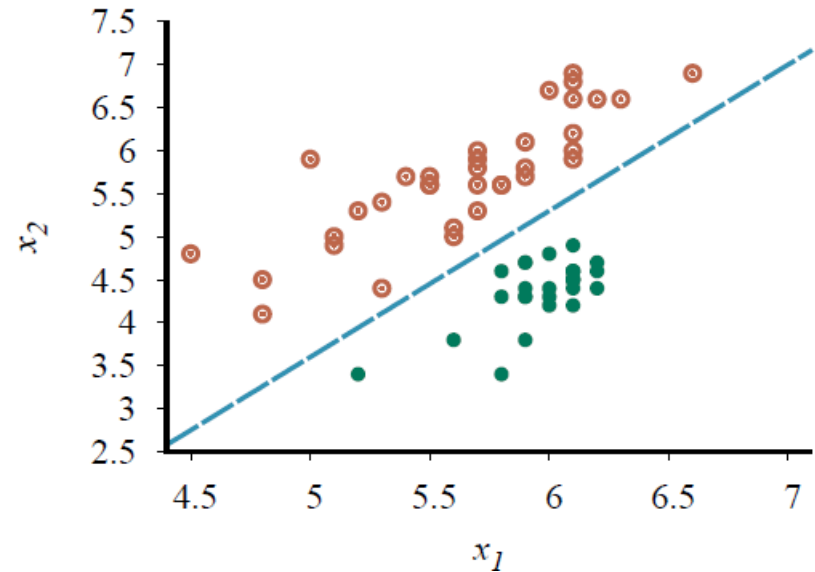
Same domain with more data points.

Task Learn a hypothesis that will take new (x_1, x_2) points and return 0 for earthquakes and 1 for explosions.

Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

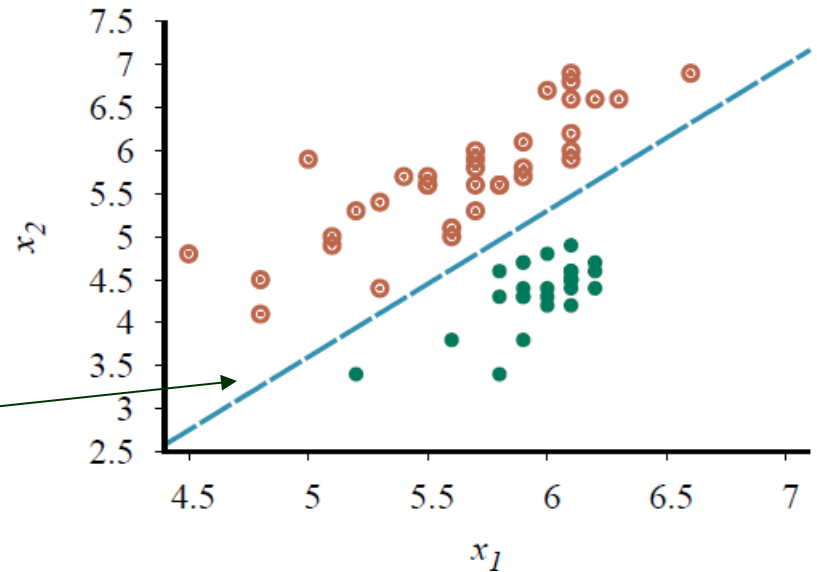


Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

e.g., $-4.9 + 1.7x_1 - x_2 = 0$

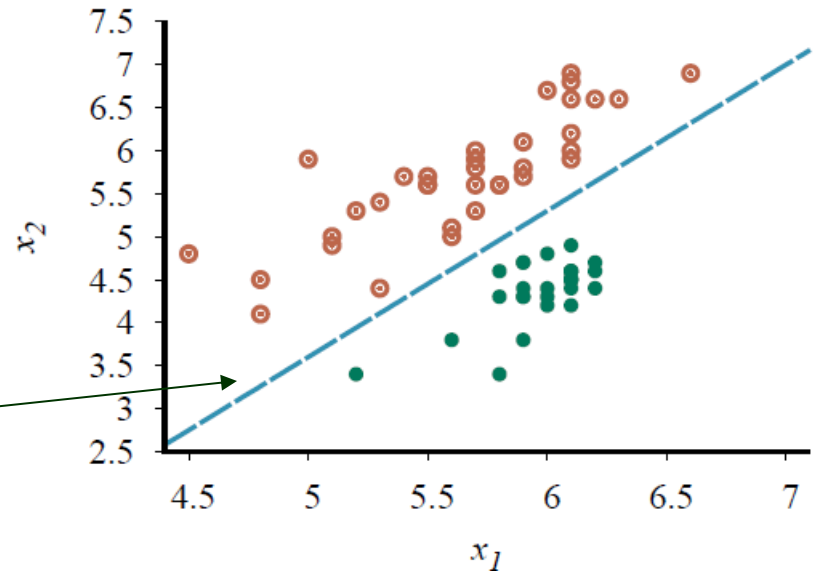


Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

e.g., $-4.9 + 1.7x_1 - x_2 = 0$



- $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$

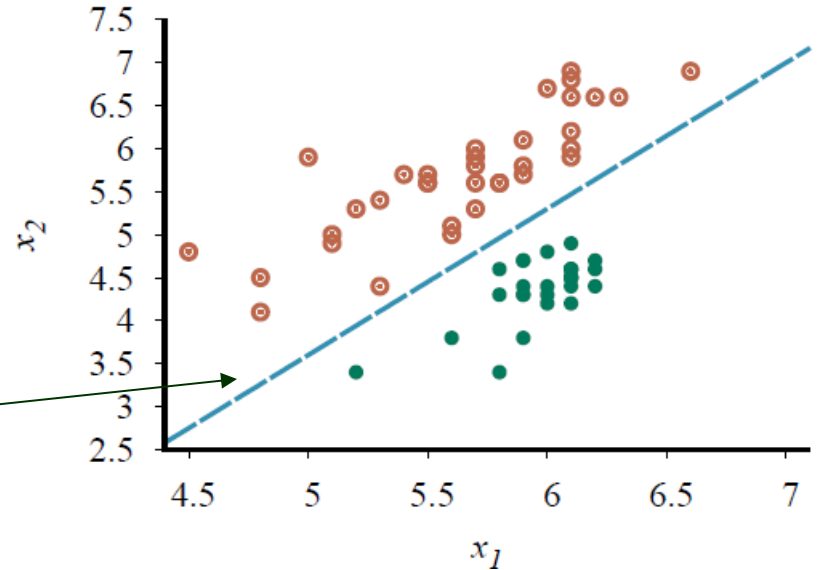
- $\mathbf{x} = (x_0, x_1, \dots, x_n)$
|
= 1

Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

e.g., $-4.9 + 1.7x_1 - x_2 = 0$



- $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$

- $\mathbf{x} = (x_0, x_1, \dots, x_n)$
|
= 1

- *Classification hypothesis:*

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

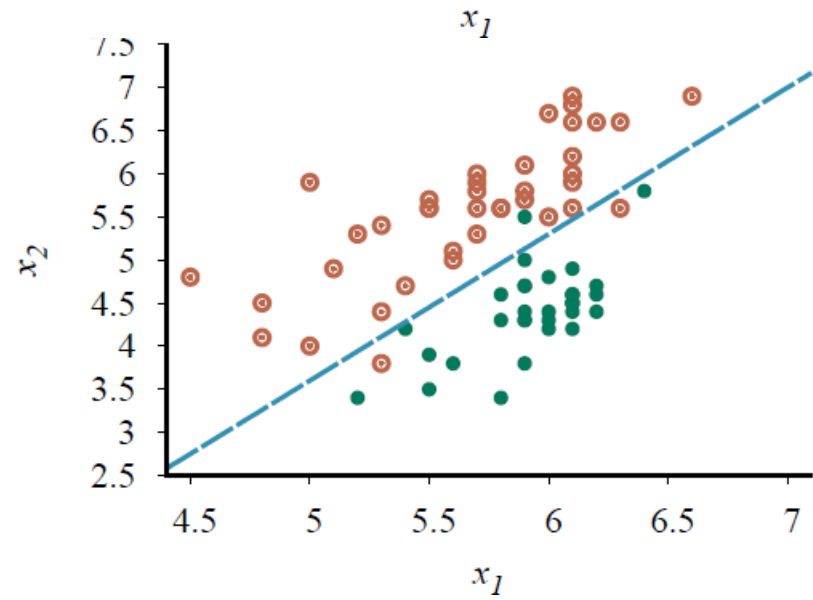
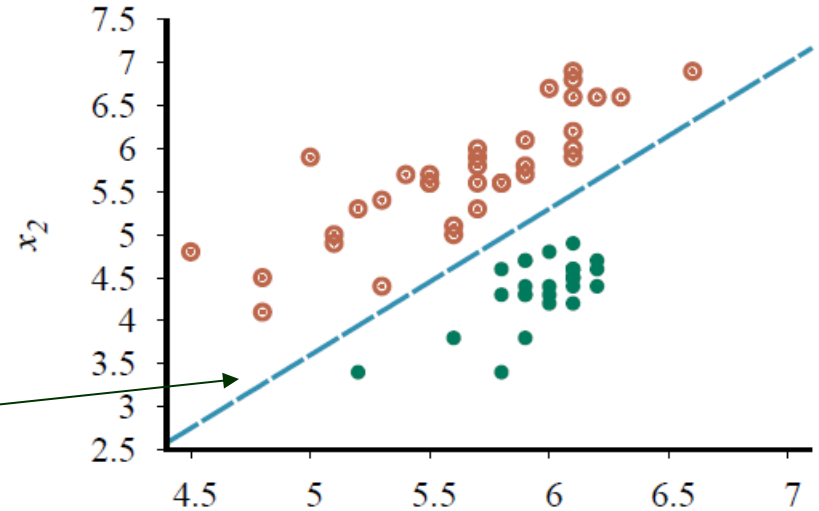
e.g., $-4.9 + 1.7x_1 - x_2 = 0$

- $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$

- $\mathbf{x} = (x_0, x_1, \dots, x_n)$
|
= 1

- *Classification hypothesis:*

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$



Linear Separator

A *decision boundary* is a line that separates two classes.

A *linear separator* is a linear decision boundary.

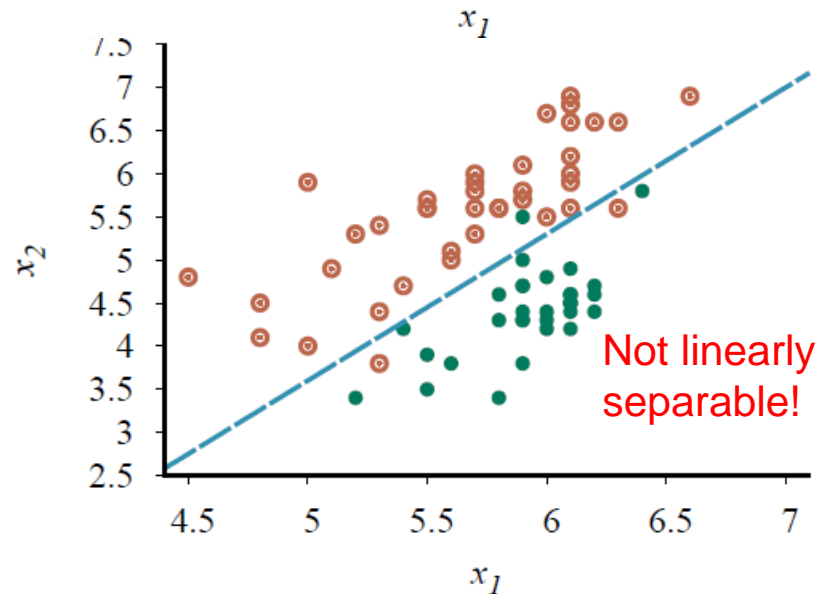
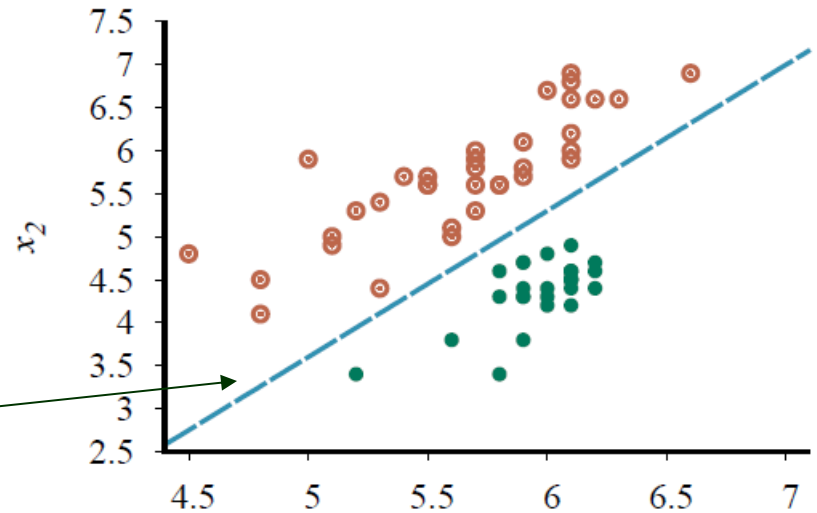
e.g., $-4.9 + 1.7x_1 - x_2 = 0$

- $\mathbf{w} = (w_0, w_1, \dots, w_n)^T$

- $\mathbf{x} = (x_0, x_1, \dots, x_n)$
|
= 1

- *Classification hypothesis:*

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$



Learning Rule

♠ Gradient $\nabla h_{\mathbf{w}}$ either vanishes or is undefined.

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

Learning Rule

♠ Gradient $\nabla h_{\mathbf{w}}$ either vanishes or is undefined.

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

♦ Use the *perceptron learning rule* (essentially borrowed from gradient descent):

$$w_i \leftarrow w_i + \alpha(y_j - h_{\mathbf{w}}(\mathbf{x}_j))x_{j,i} \quad \text{on a single example } (\mathbf{x}_j, y)$$

Learning Rule

♠ Gradient $\nabla h_{\mathbf{w}}$ either vanishes or is undefined.

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

♦ Use the *perceptron learning rule* (essentially borrowed from gradient descent):

$$w_i \leftarrow w_i + \alpha(y_j - h_{\mathbf{w}}(\mathbf{x}_j))x_{j,i} \quad \text{on a single example } (\mathbf{x}_j, y)$$

- $y_j = h_{\mathbf{w}}(\mathbf{x}_j)$. The output is correct, so no change of weights.

Learning Rule

♠ Gradient $\nabla h_{\mathbf{w}}$ either vanishes or is undefined.

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

◆ Use the *perceptron learning rule* (essentially borrowed from gradient descent):

$$w_i \leftarrow w_i + \alpha(y_j - h_{\mathbf{w}}(\mathbf{x}_j))x_{j,i} \quad \text{on a single example } (\mathbf{x}_j, y)$$

- $y_j = h_{\mathbf{w}}(\mathbf{x}_j)$. The output is correct, so no change of weights.
- $y_j = 1$ but $h_{\mathbf{w}}(\mathbf{x}_j) = 0$. w_i is increased if $x_{j,i} > 0$ and decreased if $x_{j,i} < 0$. In both situations, $\mathbf{x}\mathbf{w}$ increases with the intention to output 1.

Learning Rule

♠ Gradient $\nabla h_{\mathbf{w}}$ either vanishes or is undefined.

$$h_{\mathbf{w}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0 & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

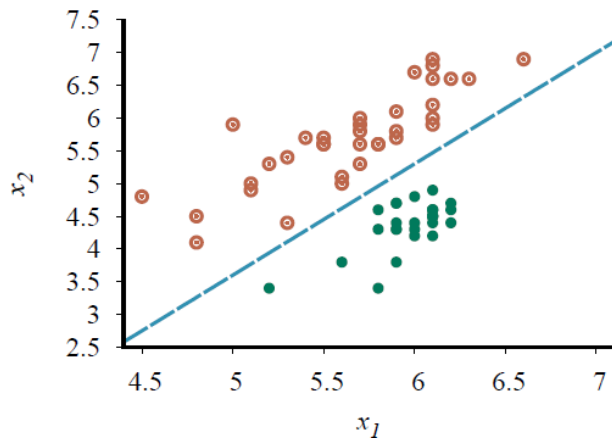
◆ Use the *perceptron learning rule* (essentially borrowed from gradient descent):

$$w_i \leftarrow w_i + \alpha(y_j - h_{\mathbf{w}}(\mathbf{x}_j))x_{j,i} \quad \text{on a single example } (\mathbf{x}_j, y)$$

- $y_j = h_{\mathbf{w}}(\mathbf{x}_j)$. The output is correct, so no change of weights.
- $y_j = 1$ but $h_{\mathbf{w}}(\mathbf{x}_j) = 0$. w_i is increased if $x_{j,i} > 0$ and decreased if $x_{j,i} < 0$. In both situations, $\mathbf{x}\mathbf{w}$ increases with the intention to output 1.
- $y_j = 0$ but $h_{\mathbf{w}}(\mathbf{x}_j) = 1$. w_i is decreased if $x_{j,i} > 0$ and increased if $x_{j,i} < 0$. In both situations, $\mathbf{x}\mathbf{w}$ decreases with the intention to output 0.

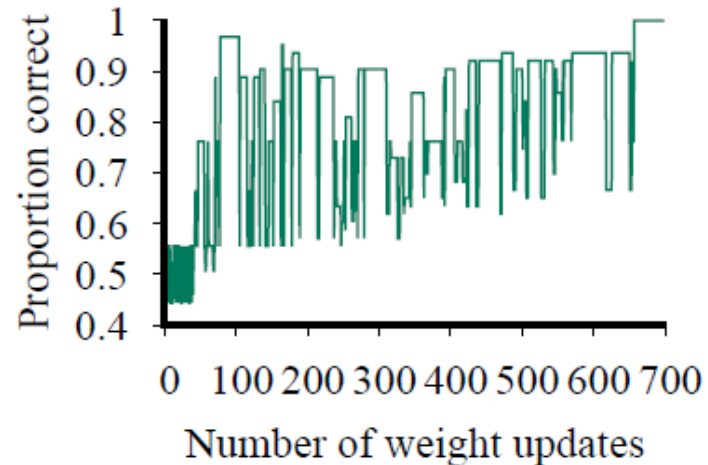
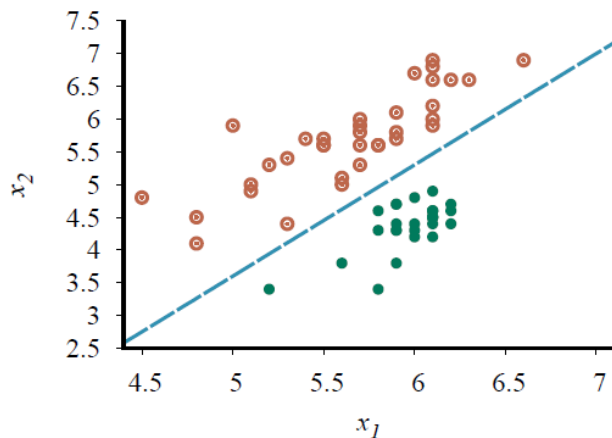
Training Curves for Perceptron Learning

- The learning rule is applied one example at a time.
- A *training curve* measures the classifier performance on a fixed training set as learning proceeds one example at a time on the same set.



Training Curves for Perceptron Learning

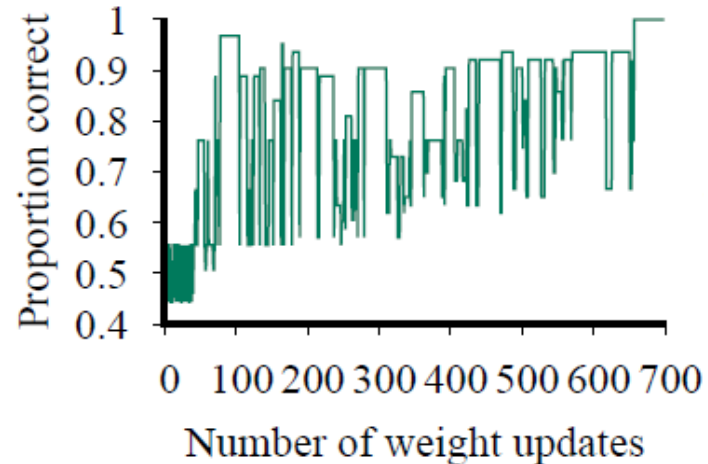
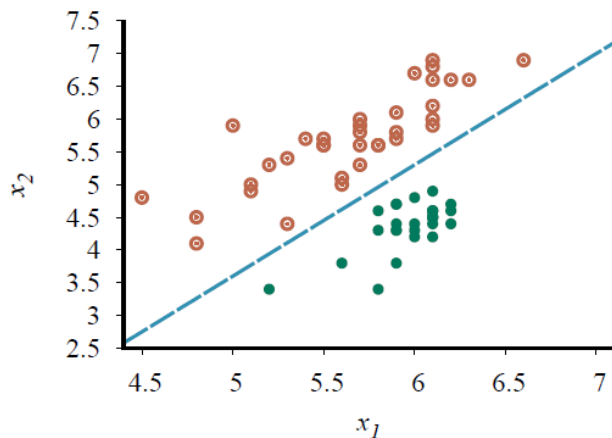
- The learning rule is applied one example at a time.
- A *training curve* measures the classifier performance on a fixed training set as learning proceeds one example at a time on the same set.



$$\alpha = 1$$

Training Curves for Perceptron Learning

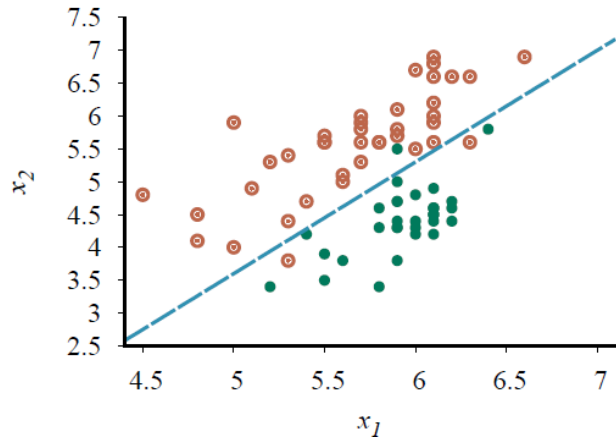
- The learning rule is applied one example at a time.
- A *training curve* measures the classifier performance on a fixed training set as learning proceeds one example at a time on the same set.



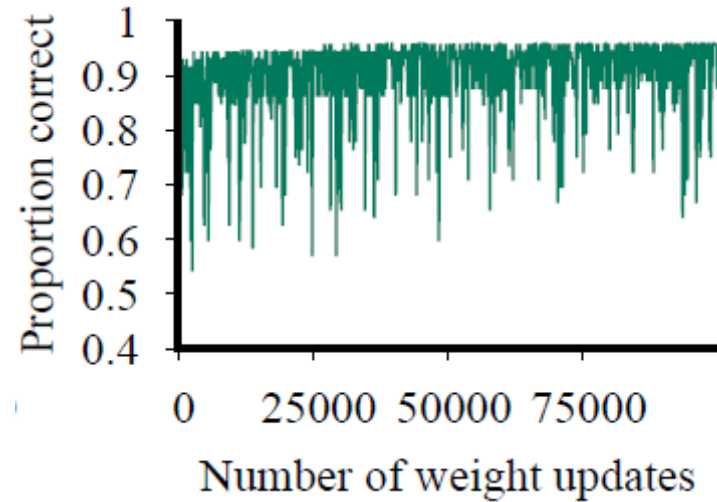
$\alpha = 1$

- 657 steps before convergence
- 63 examples, each used 10 times on average

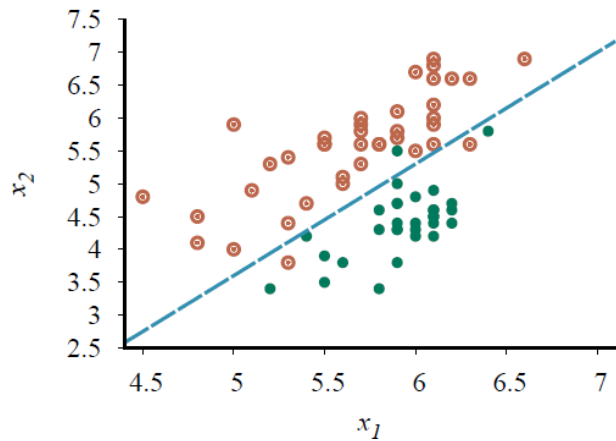
Training Curves (cont'd)



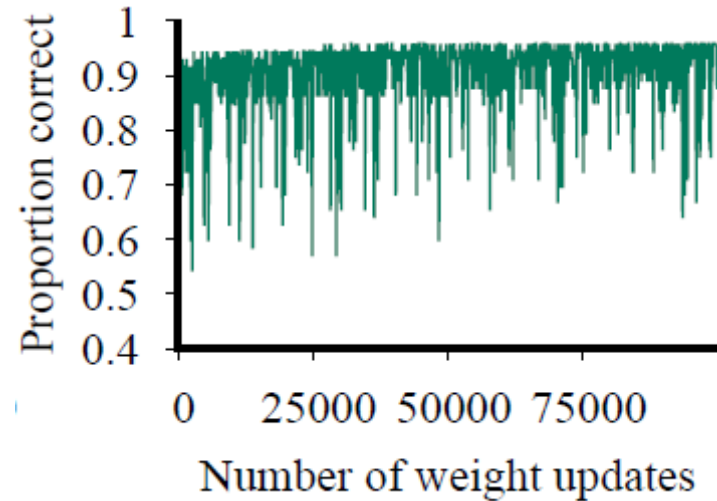
Data not linearly separable.



Training Curves (cont'd)

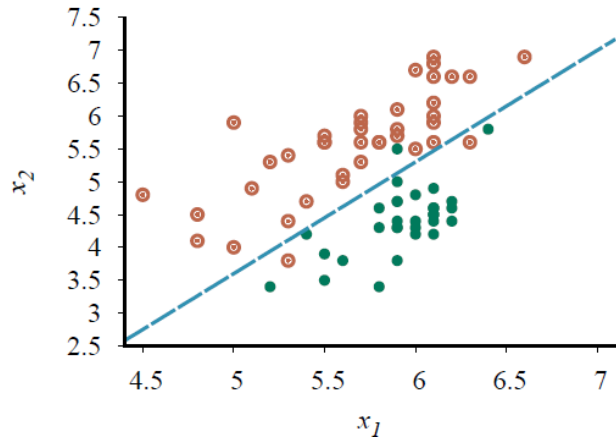


Data not linearly separable.

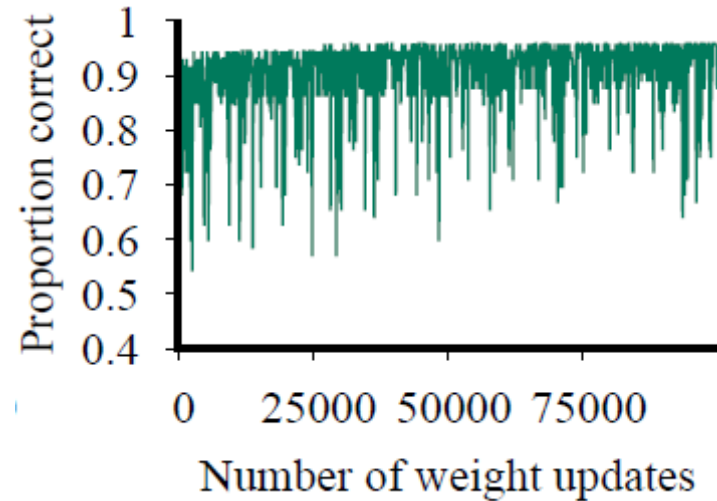


- Fails to converge after 10,000 steps.

Training Curves (cont'd)

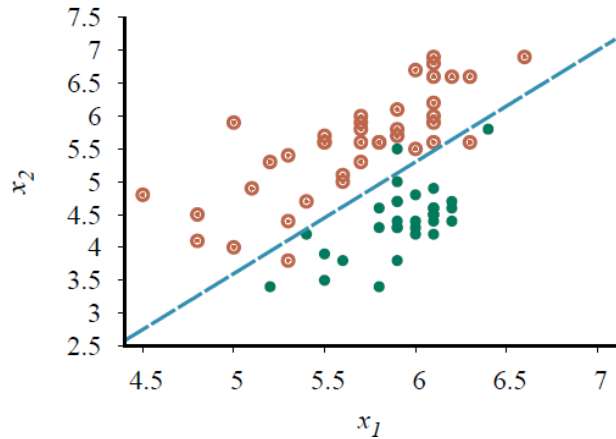


Data not linearly separable.

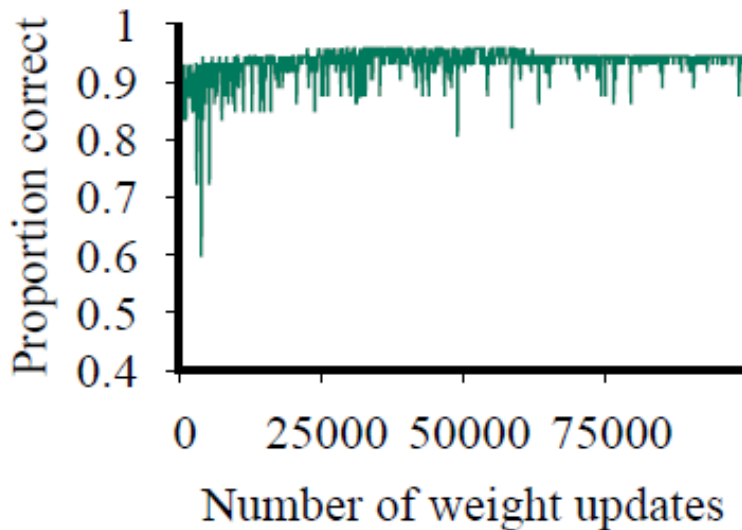
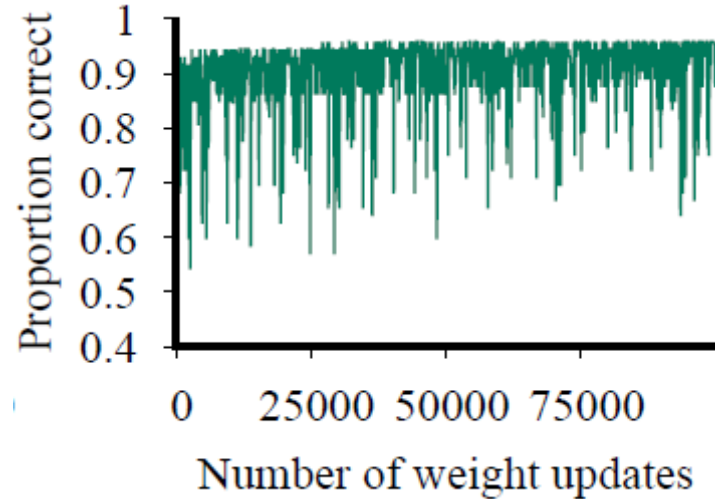


- Fails to converge after 10,000 steps.
- Let α decay as $O(1/t)$ where $t = \#$ iterations.
$$w_i \leftarrow w_i + \alpha(y_j - h_w(\mathbf{x}_j))x_{j,i}$$

Training Curves (cont'd)



Data not linearly separable.



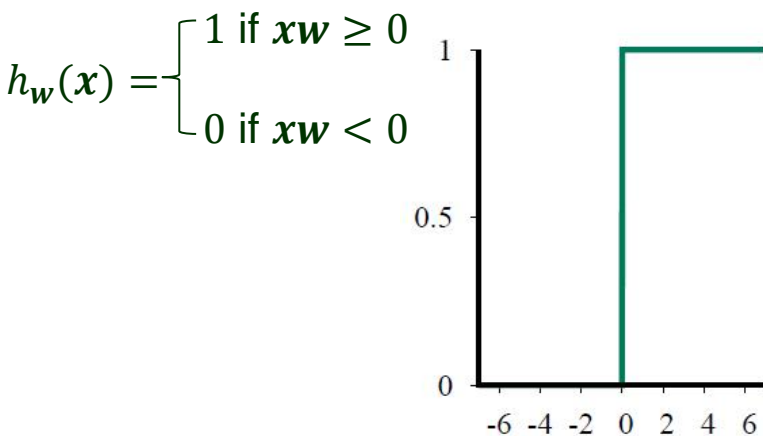
- Fails to converge after 10,000 steps.
- Let α decay as $O(1/t)$ where $t = \#$ iterations.

$$w_i \leftarrow w_i + \alpha(y_j - h_w(\mathbf{x}_j))x_{j,i}$$

← e.g., $\alpha(t) = 1000/(1000 + t)$

IV. Logistic Function

- ♠ Current hypothesis function is not continuous, let alone differentiable.
- ♠ This makes learning with the perceptron rule very unpredictable.
- ♠ It would be better if some examples could be classified as unclear borderline cases.



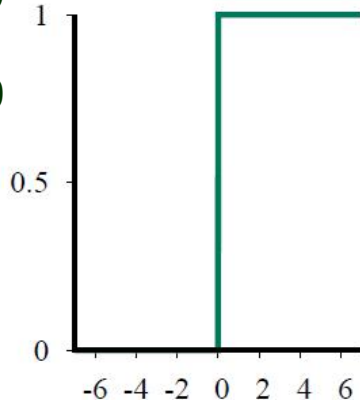
IV. Logistic Function

- ♠ Current hypothesis function is not continuous, let alone differentiable.
- ♠ This makes learning with the perceptron rule very unpredictable.
- ♠ It would be better if some examples could be classified as unclear borderline cases.
- ◆ Use a continuous, differential function to soften the threshold

Logistic function:

$$\text{Logistics}(z) = g(z) = \frac{1}{1+e^{-z}}$$

$$h_w(x) = \begin{cases} 1 & \text{if } xw \geq 0 \\ 0 & \text{if } xw < 0 \end{cases}$$



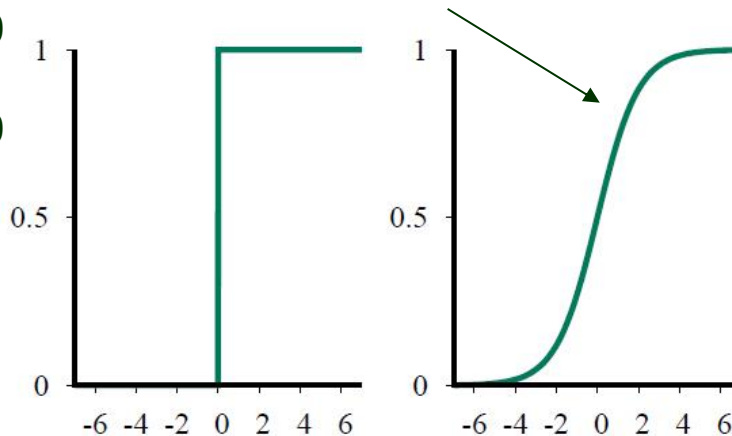
IV. Logistic Function

- ♠ Current hypothesis function is not continuous, let alone differentiable.
- ♠ This makes learning with the perceptron rule very unpredictable.
- ♠ It would be better if some examples could be classified as unclear borderline cases.
- ◆ Use a continuous, differential function to soften the threshold

Logistic function:

$$\text{Logistics}(z) = g(z) = \frac{1}{1+e^{-z}}$$

$$h_w(x) = \begin{cases} 1 & \text{if } xw \geq 0 \\ 0 & \text{if } xw < 0 \end{cases}$$



IV. Logistic Function

- ♠ Current hypothesis function is not continuous, let alone differentiable.
- ♠ This makes learning with the perceptron rule very unpredictable.
- ♠ It would be better if some examples could be classified as unclear borderline cases.
- ◆ Use a continuous, differential function to soften the threshold

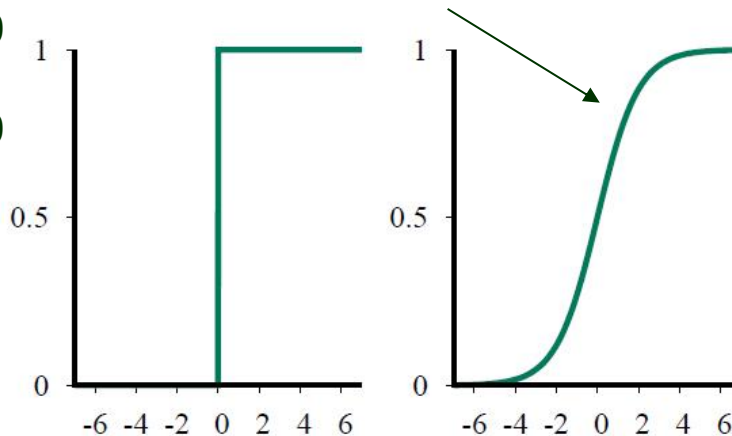
Logistic function:

$$\text{Logistics}(z) = g(z) = \frac{1}{1+e^{-z}}$$

Hypothesis function:

$$h_w(x) = \text{Logistics}(xw) = \frac{1}{1 + e^{-xw}}$$

$$h_w(x) = \begin{cases} 1 & \text{if } xw \geq 0 \\ 0 & \text{if } xw < 0 \end{cases}$$



IV. Logistic Function

- ♠ Current hypothesis function is not continuous, let alone differentiable.
- ♠ This makes learning with the perceptron rule very unpredictable.
- ♠ It would be better if some examples could be classified as unclear borderline cases.
- ◆ Use a continuous, differential function to soften the threshold

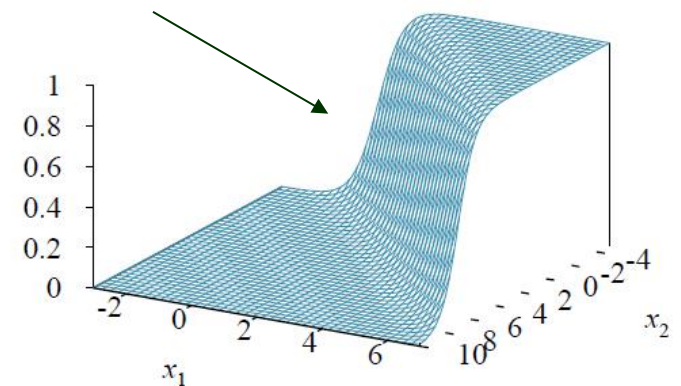
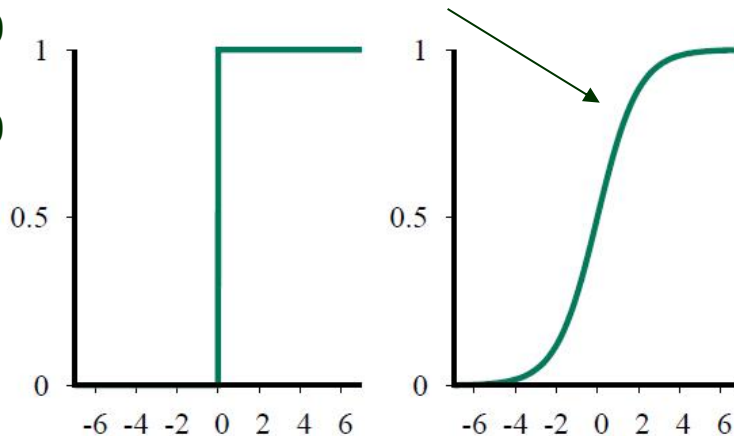
Logistic function:

$$\text{Logistics}(z) = g(z) = \frac{1}{1+e^{-z}}$$

Hypothesis function:

$$h_w(x) = \text{Logistics}(xw) = \frac{1}{1+e^{-xw}}$$

$$h_w(x) = \begin{cases} 1 & \text{if } xw \geq 0 \\ 0 & \text{if } xw < 0 \end{cases}$$



Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

Still apply gradient descent.

$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) = \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2$$

Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

Still apply gradient descent.

$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2 \\ &\vdots \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot g'(\mathbf{x}\mathbf{w}) \cdot x_i \end{aligned}$$

Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

Still apply gradient descent.

$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2 \\ &\vdots \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot g'(\mathbf{x}\mathbf{w}) \cdot x_i \end{aligned}$$

$$\begin{aligned} g'(\mathbf{x}\mathbf{w}) &= g(\mathbf{x}\mathbf{w})(1 - g(\mathbf{x}\mathbf{w})) \\ &= h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \end{aligned}$$

Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

Still apply gradient descent.

$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2 \\ &\quad \vdots \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot g'(\mathbf{x}\mathbf{w}) \cdot x_i \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \cdot x_i \end{aligned} \quad \begin{aligned} g'(\mathbf{x}\mathbf{w}) &= g(\mathbf{x}\mathbf{w})(1 - g(\mathbf{x}\mathbf{w})) \\ &= h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \end{aligned}$$

Logistic Regression

Fit the model $h_{\mathbf{w}}(\mathbf{x}) = \text{Logistics}(\mathbf{x}\mathbf{w})$ to minimize loss on a data set.

Still apply gradient descent.

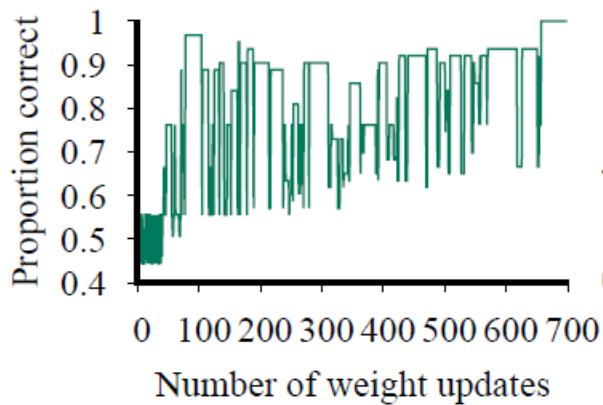
$$\text{Logistics}(z) = g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \frac{\partial}{\partial w_i} \text{Loss}(\mathbf{w}) &= \frac{\partial}{\partial w_i} (y - h_{\mathbf{w}}(\mathbf{x}))^2 \\ &\vdots \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot g'(\mathbf{x}\mathbf{w}) \cdot x_i \\ &= -2(y - h_{\mathbf{w}}(\mathbf{x})) \cdot h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \cdot x_i \end{aligned} \quad \begin{aligned} g'(\mathbf{x}\mathbf{w}) &= g(\mathbf{x}\mathbf{w})(1 - g(\mathbf{x}\mathbf{w})) \\ &= h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \end{aligned}$$

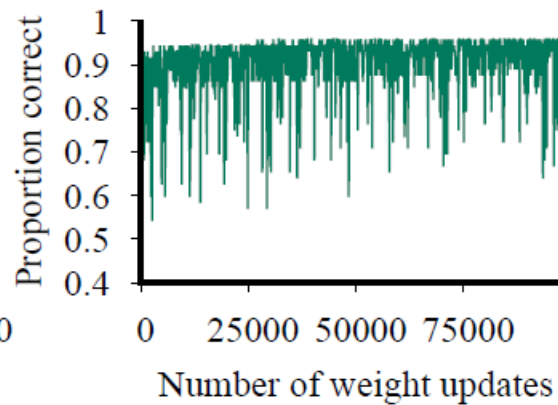
Weight update:

$$w_i \leftarrow w_i + \alpha(y - h_{\mathbf{w}}(\mathbf{x})) \cdot h_{\mathbf{w}}(\mathbf{x})(1 - h_{\mathbf{w}}(\mathbf{x})) \cdot x_i$$

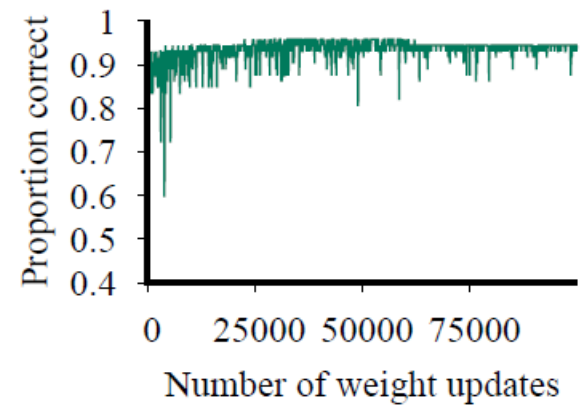
Improvements on Training Results



$\alpha = 1$

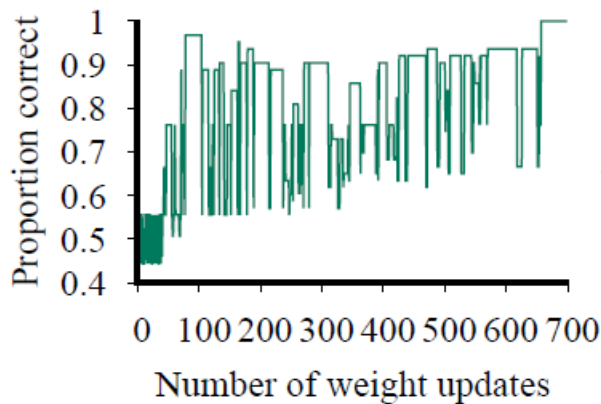
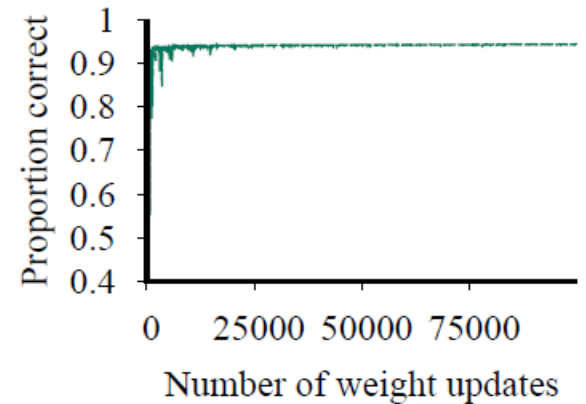
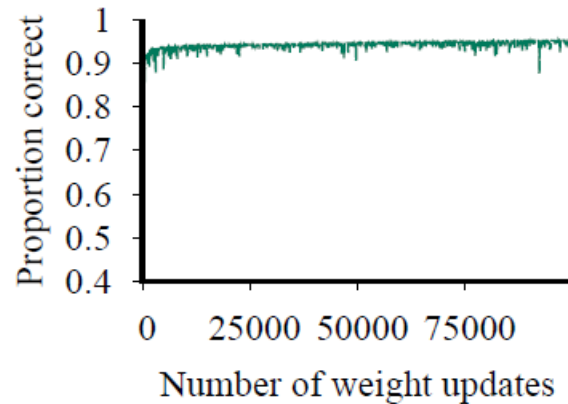
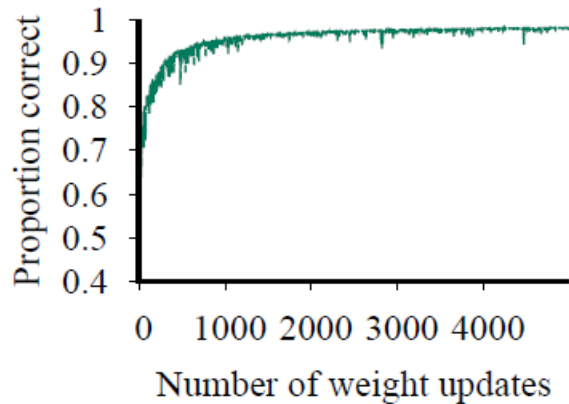


$\alpha = 1$

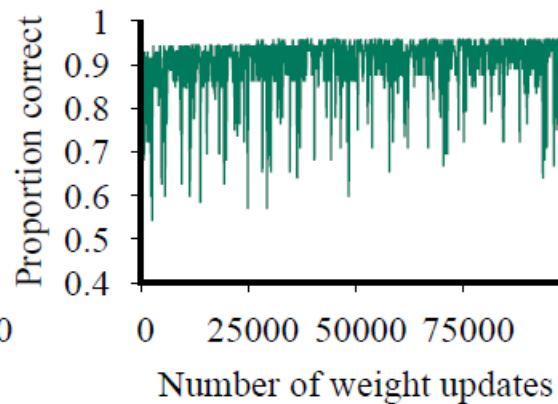


$\alpha(t) = 1000/(1000 + t)$

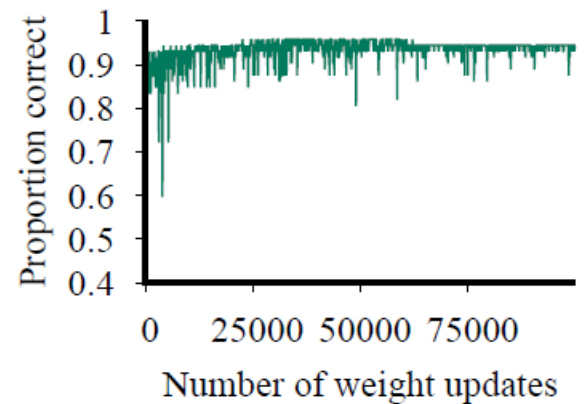
Improvements on Training Results



$\alpha = 1$

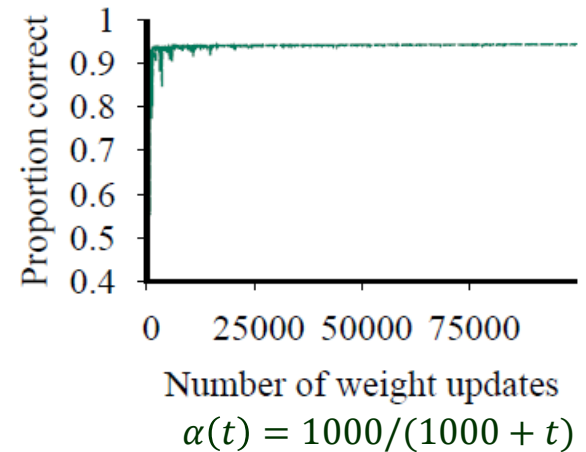
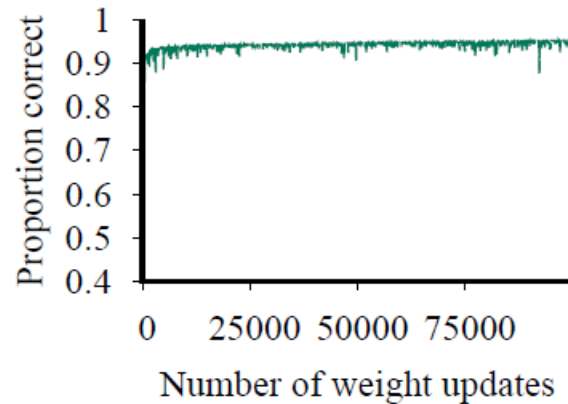
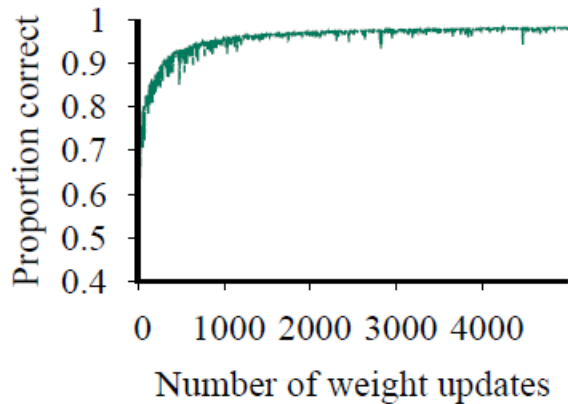


$\alpha = 1$



$\alpha(t) = 1000/(1000 + t)$

Improvements on Training Results



Logistic regression converges far more quickly and reliably.

