

Foundations of Machine Learning

Linear Regression and Perceptron

Maheesan Niranjana

School of Electronics and Computer Science
University of Southampton

July 2021

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Linear Regression & Perceptron

- Data: $\{\mathbf{x}_n, f_n\}_{n=1}^N$
Input: $\mathbf{x}_n \in \mathcal{R}^p$; target / output f_n real valued
- Model: $f = \mathbf{w}^t \mathbf{x} + w_0$
Output linear function of input (including a constant w_0)
- Work in $(p + 1)$ dimensional space to avoid treating w_0 separately

$$\mathbf{y} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} \quad \mathbf{a} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- Data: $\{\mathbf{y}_n, f_n\}_{n=1}^N$
- Model: $f = \mathbf{y}^t \mathbf{a}$
- $p + 1$ unknowns held in vector \mathbf{a}

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

$$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Error and Minimization

- $E = \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - f_n\}^2$
- $E = \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\}^2$
- To find the best \mathbf{a} we minimize E – differentiate with respect to each of the unknowns in \mathbf{a} and set to zero.

- $$\frac{\partial E}{\partial a_i} = 2 \sum_{n=1}^N \left\{ \left(\sum_{j=1}^{(p+1)} a_j y_{nj} \right) - f_n \right\} (y_{ni})$$

- There are $(p + 1)$ derivatives (with respect to each a_i)
- Equating them to zero gives $(p + 1)$ equations in $(p + 1)$ unknowns

Solution to Regression

- $(p + 1)$ simultaneous equations to solve:
 i^{th} row, j^{th} column shown

$$\begin{pmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \vdots & \sum_{n=1}^N y_{ni} y_{nj} & \dots \\ \vdots & \dots & \vdots \\ \dots & \dots & \dots \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_{(p+1)} \end{pmatrix} = \begin{pmatrix} \vdots \\ \sum_{n=1}^N f_n y_{ni} \\ \vdots \end{pmatrix}$$

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = \|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- **Homework:** Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- **Homework:** With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Derivation in vector/matrix form

- \mathbf{Y} : $N \times (p + 1)$ matrix n^{th} row is \mathbf{y}_n^t
- \mathbf{f} : $N \times 1$ vector of outputs
- Error $E = ||\mathbf{Y}\mathbf{a} - \mathbf{f}||^2$
- Homework: Verify the error written like this is the same as the one we wrote out in lengthy algebra.
- Gradient

$$\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$$

- Equating the gradient to zero gives

$$\begin{aligned}\mathbf{Y}^t \mathbf{Y} \mathbf{a} &= \mathbf{Y}^t \mathbf{f} \\ \mathbf{a} &= (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}\end{aligned}$$

- Homework: With three data points in one dimensional input space (x_1, f_1) , (x_2, f_2) and (x_3, f_3) and two unknowns, slope (m) and intercept (c) of fitting a straight line, write out all the expressions seen so far.

Solution by Gradient Descent

- Gradient vector: $\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{f})$
- Steepest descent algorithm:

Initialize \mathbf{a} at random
Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \nabla_{\mathbf{a}} E$
Until Convergence

- Second order (Newton's) method

Initialize \mathbf{a} at random
Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \mathbf{H}^{-1} \nabla_{\mathbf{a}} E$
Until Convergence

- Rapid convergence with second order method, but cost of computing and inverting \mathbf{H} can be high (more on this under Neural Networks)

Solution by Gradient Descent

- Gradient vector: $\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{f})$

- Steepest descent algorithm:

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \nabla_{\mathbf{a}} E$

Until Convergence

- Second order (Newton's) method

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \mathbf{H}^{-1} \nabla_{\mathbf{a}} E$

Until Convergence

- Rapid convergence with second order method, but cost of computing and inverting \mathbf{H} can be high (more on this under Neural Networks)

Solution by Gradient Descent

- Gradient vector: $\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{f})$

- Steepest descent algorithm:

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \nabla_{\mathbf{a}} E$

Until Convergence

- Second order (Newton's) method

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \mathbf{H}^{-1} \nabla_{\mathbf{a}} E$

Until Convergence

- Rapid convergence with second order method, but cost of computing and inverting \mathbf{H} can be high (more on this under Neural Networks)

Solution by Gradient Descent

- Gradient vector: $\nabla_{\mathbf{a}} E = 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{f})$

- Steepest descent algorithm:

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \eta \nabla_{\mathbf{a}} E$

Until Convergence

- Second order (Newton's) method

Initialize \mathbf{a} at random

Update $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} - \mathbf{H}^{-1} \nabla_{\mathbf{a}} E$

Until Convergence

- Rapid convergence with second order method, but cost of computing and inverting \mathbf{H} can be high (more on this under Neural Networks)

Gradient and Stochastic Gradient Descent

- Error $E = \sum_{n=1}^N e_n^2$
- True gradient:

$$\nabla_{\mathbf{a}} E = 2 \sum_{n=1}^N \{\mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n\} (\mathbf{y}_n)$$

- Gradient computed on n^{th} data:

$$\nabla_{\mathbf{a}} e_n = 2 \{\mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n\} (\mathbf{y}_n)$$

Gradient and Stochastic Gradient Descent

- Error $E = \sum_{n=1}^N e_n^2$
- True gradient:

$$\nabla_{\mathbf{a}} E = 2 \sum_{n=1}^N \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

- Gradient computed on n^{th} data:

$$\nabla_{\mathbf{a}} e_n = 2 \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

Gradient and Stochastic Gradient Descent

- Error $E = \sum_{n=1}^N e_n^2$
- True gradient:

$$\nabla_{\mathbf{a}} E = 2 \sum_{n=1}^N \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

- Gradient computed on n^{th} data:

$$\nabla_{\mathbf{a}} e_n = 2 \{ \mathbf{y}_n^t \mathbf{a} - \mathbf{f}_n \} (\mathbf{y}_n)$$

Regularization

- Pseudo inverse solution: $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where γ is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2 + \gamma \|\mathbf{a}\|^2$$

Here a quadratic penalty term has been included

- Homework: Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an l_1 norm (sum of absolute values of \mathbf{a} , i.e. $\sum_{j=1}^p |a_j|$); See **Lab 4**.

Regularization

- Pseudo inverse solution: $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where γ is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2 + \gamma \|\mathbf{a}\|^2$$

Here a quadratic penalty term has been included

- Homework: Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an l_1 norm (sum of absolute values of \mathbf{a} , i.e. $\sum_{j=1}^p |a_j|$); See **Lab 4**.

Regularization

- Pseudo inverse solution: $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where γ is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2 + \gamma \|\mathbf{a}\|^2$$

Here a quadratic penalty term has been included

- Homework: Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an l_1 norm (sum of absolute values of \mathbf{a} , i.e. $\sum_{j=1}^p |a_j|$); See **Lab 4**.

Regularization

- Pseudo inverse solution: $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where γ is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2 + \gamma \|\mathbf{a}\|^2$$

Here a quadratic penalty term has been included

- **Homework:** Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an l_1 norm (sum of absolute values of \mathbf{a} , i.e. $\sum_{j=1}^p |a_j|$); See **Lab 4**.

Regularization

- Pseudo inverse solution: $\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{f}$
- This can be ill conditioned, so we could *regularize* by

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y} + \gamma \mathbf{I})^{-1} \mathbf{Y}^t \mathbf{f}$$

where γ is a small constant.

- We achieve precisely this by minimizing an error of the form

$$\|\mathbf{Y}\mathbf{a} - \mathbf{f}\|^2 + \gamma \|\mathbf{a}\|^2$$

Here a quadratic penalty term has been included

- Homework: Differentiate this error and derive the regularized solution
- Sparse solutions are obtained by regularizing with an l_1 norm (sum of absolute values of \mathbf{a} , i.e. $\sum_{j=1}^p |a_j|$); See **Lab 4**.

Perceptron

A suitable performance measure

- Number of misclassified examples as measure of error
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum \mathbf{y}_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with $\mathbf{y}_n^t \mathbf{a} > 0$ for positive class and $\mathbf{y}_n^t \mathbf{a} < 0$ for the negative class; we then switch the signs of negative class examples and required $\mathbf{y}_n^t \mathbf{a} > 0$ for all the training data; so for the misclassified examples $-\sum \mathbf{y}_n^t \mathbf{a}$ should be as small as possible.

Perceptron

A suitable performance measure

- Number of misclassified examples as measure of error
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum \mathbf{y}_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with $\mathbf{y}_n^t \mathbf{a} > 0$ for positive class and $\mathbf{y}_n^t \mathbf{a} < 0$ for the negative class; we then switch the signs of negative class examples and required $\mathbf{y}_n^t \mathbf{a} > 0$ for all the training data; so for the misclassified examples $-\sum \mathbf{y}_n^t \mathbf{a}$ should be as small as possible.

Perceptron

A suitable performance measure

- Number of misclassified examples as measure of error
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum \mathbf{y}_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with $\mathbf{y}_n^t \mathbf{a} > 0$ for positive class and $\mathbf{y}_n^t \mathbf{a} < 0$ for the negative class; we then switch the signs of negative class examples and required $\mathbf{y}_n^t \mathbf{a} > 0$ for all the training data; so for the misclassified examples $-\sum \mathbf{y}_n^t \mathbf{a}$ should be as small as possible.

Perceptron

A suitable performance measure

- Number of misclassified examples as measure of error
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum \mathbf{y}_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with $\mathbf{y}_n^t \mathbf{a} > 0$ for positive class and $\mathbf{y}_n^t \mathbf{a} < 0$ for the negative class; we then switch the signs of negative class examples and required $\mathbf{y}_n^t \mathbf{a} > 0$ for all the training data; so for the misclassified examples $-\sum \mathbf{y}_n^t \mathbf{a}$ should be as small as possible.

Perceptron

A suitable performance measure

- Number of misclassified examples as measure of error
Piecewise constant (cannot differentiate)
- Suitable error measure:

$$E_P = - \sum \mathbf{y}_n^t \mathbf{a}$$

- Summation taken over misclassified examples
- We started with $\mathbf{y}_n^t \mathbf{a} > 0$ for positive class and $\mathbf{y}_n^t \mathbf{a} < 0$ for the negative class; we then switch the signs of negative class examples and required $\mathbf{y}_n^t \mathbf{a} > 0$ for all the training data; so for the misclassified examples $-\sum \mathbf{y}_n^t \mathbf{a}$ should be as small as possible.

Perceptron

Learning rule

- Gradient:

$$\frac{\partial E}{\partial \mathbf{a}} = - \sum \mathbf{y}_n$$

- Gradient algorithm: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \sum \mathbf{y}_n$
- Stochastic gradient algorithm:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_n$$

- Note what \mathbf{y}_n is. It is an item of data that is taken at random and happens to be misclassified by the current value of \mathbf{a} at iteration k .

Perceptron

Learning rule

- Gradient:

$$\frac{\partial E}{\partial \mathbf{a}} = - \sum \mathbf{y}_n$$

- Gradient algorithm: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \sum \mathbf{y}_n$

- Stochastic gradient algorithm:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_n$$

- Note what \mathbf{y}_n is. It is an item of data that is taken at random and happens to be misclassified by the current value of \mathbf{a} at iteration k .

Perceptron

Learning rule

- Gradient:

$$\frac{\partial E}{\partial \mathbf{a}} = - \sum \mathbf{y}_n$$

- Gradient algorithm: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \sum \mathbf{y}_n$
- Stochastic gradient algorithm:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_n$$

- Note what \mathbf{y}_n is. It is an item of data that is taken at random and happens to be misclassified by the current value of \mathbf{a} at iteration k .

Perceptron

Learning rule

- Gradient:

$$\frac{\partial E}{\partial \mathbf{a}} = - \sum \mathbf{y}_n$$

- Gradient algorithm: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \sum \mathbf{y}_n$
- Stochastic gradient algorithm:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_n$$

- Note what \mathbf{y}_n is. It is an item of data that is taken at random and happens to be misclassified by the current value of \mathbf{a} at iteration k .

Perceptron

Convergence of the learning rule

- Learning Rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$
where $\mathbf{y}(k)$ is a misclassified input.
- Training criterion
 - We start with requiring $\mathbf{a}^t \mathbf{y}(k) \leq 0$, depending on the example belonging to class 1 or class 2.
 - If we switch the signs of examples of class 2, we require $\mathbf{a}^t \mathbf{y}(k) > 0$ for all k .
- On misclassified data $\mathbf{a}^t \mathbf{y}(k) < 0$
- If $\hat{\mathbf{a}}$ is a solution (separable data), for all k , $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$ for this update rule. *i.e.* the learning rule brings the guess closer to a valid solution.

Perceptron

Convergence of the learning rule

- Learning Rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$
where $\mathbf{y}(k)$ is a misclassified input.
- Training criterion
 - We start with requiring $\mathbf{a}^t \mathbf{y}(k) \leq 0$, depending on the example belonging to class 1 or class 2.
 - If we switch the signs of examples of class 2, we require $\mathbf{a}^t \mathbf{y}(k) > 0$ for all k .
- On misclassified data $\mathbf{a}^t \mathbf{y}(k) < 0$
- If $\hat{\mathbf{a}}$ is a solution (separable data), for all k , $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$ for this update rule. *i.e.* the learning rule brings the guess closer to a valid solution.

Perceptron

Convergence of the learning rule

- Learning Rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$
where $\mathbf{y}(k)$ is a misclassified input.
- Training criterion
 - We start with requiring $\mathbf{a}^t \mathbf{y}(k) \leq 0$, depending on the example belonging to class 1 or class 2.
 - If we switch the signs of examples of class 2, we require $\mathbf{a}^t \mathbf{y}(k) > 0$ for all k .
- On misclassified data $\mathbf{a}^t \mathbf{y}(k) < 0$
- If $\hat{\mathbf{a}}$ is a solution (separable data), for all k , $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$ for this update rule. *i.e.* the learning rule brings the guess closer to a valid solution.

Perceptron

Convergence of the learning rule

- Learning Rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$
where $\mathbf{y}(k)$ is a misclassified input.
- Training criterion
 - We start with requiring $\mathbf{a}^t \mathbf{y}(k) \leq 0$, depending on the example belonging to class 1 or class 2.
 - If we switch the signs of examples of class 2, we require $\mathbf{a}^t \mathbf{y}(k) > 0$ for all k .
- On misclassified data $\mathbf{a}^t \mathbf{y}(k) < 0$
- If $\hat{\mathbf{a}}$ is a solution (separable data), for all k , $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$ for this update rule. *i.e.* the learning rule brings the guess closer to a valid solution.

Perceptron

Convergence of the learning rule

- Learning Rule: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}(k)$
where $\mathbf{y}(k)$ is a misclassified input.
- Training criterion
 - We start with requiring $\mathbf{a}^t \mathbf{y}(k) \leq 0$, depending on the example belonging to class 1 or class 2.
 - If we switch the signs of examples of class 2, we require $\mathbf{a}^t \mathbf{y}(k) > 0$ for all k .
- On misclassified data $\mathbf{a}^t \mathbf{y}(k) < 0$
- If $\hat{\mathbf{a}}$ is a solution (separable data), for all k , $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$
- We prove convergence by showing:
 $\|\mathbf{a}^{(k+1)} - \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \hat{\mathbf{a}}\|^2$ for this update rule. *i.e.* the learning rule brings the guess closer to a valid solution.

Perceptron

Convergence of the learning rule (cont'd)

- For perceptron criterion, the magnitude of \mathbf{a} is not relevant (only the direction is). Hence for some scalar α , we wish to show

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2$$

- From the update formula

$$\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}} = \mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}} + \mathbf{y}(k)$$

- Taking magnitudes

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}})^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

- If we drop the negative term $\mathbf{a}^{(k)t} \mathbf{y}(k)$ from RHS, the equality becomes an inequality

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

Perceptron

Convergence of the learning rule (cont'd)

- For perceptron criterion, the magnitude of \mathbf{a} is not relevant (only the direction is). Hence for some scalar α , we wish to show

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2$$

- From the update formula

$$\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}} = \mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}} + \mathbf{y}(k)$$

- Taking magnitudes

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}})^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

- If we drop the negative term $\mathbf{a}^{(k)t} \mathbf{y}(k)$ from RHS, the equality becomes an inequality

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

Perceptron

Convergence of the learning rule (cont'd)

- For perceptron criterion, the magnitude of \mathbf{a} is not relevant (only the direction is). Hence for some scalar α , we wish to show

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2$$

- From the update formula

$$\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}} = \mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}} + \mathbf{y}(k)$$

- Taking magnitudes

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}})^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

- If we drop the negative term $\mathbf{a}^{(k)t} \mathbf{y}(k)$ from RHS, the equality becomes an inequality

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

Perceptron

Convergence of the learning rule (cont'd)

- For perceptron criterion, the magnitude of \mathbf{a} is not relevant (only the direction is). Hence for some scalar α , we wish to show

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2$$

- From the update formula

$$\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}} = \mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}} + \mathbf{y}(k)$$

- Taking magnitudes

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 = \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 + 2(\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}})^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

- If we drop the negative term $\mathbf{a}^{(k)t} \mathbf{y}(k)$ from RHS, the equality becomes an inequality

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha \hat{\mathbf{a}}^t \mathbf{y}(k) + \|\mathbf{y}(k)\|^2$$

Perceptron

Convergence of the learning rule (cont'd)

- Of the three terms on the right hand side, we know $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$, because $\hat{\mathbf{a}}$ is assumed to be a solution.
- If we select

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2$$

$$\gamma = \min_i \hat{\mathbf{a}}^t \mathbf{y}_i$$

i.e. largest of the positive term and smallest of the negative term,
then for $\alpha = \beta^2 / \gamma$,

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - \beta^2$$

- (Note the inequality remains true when the right hand side is replaced by a quantity larger than what it previously was.)
- Every correction takes the guess closer to a true solution.
- From an initialization $\mathbf{a}^{(1)}$, we will find a solution in *at most*
 $k_0 = \frac{\|\mathbf{a}^{(1)} - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}$ updates.

Perceptron

Convergence of the learning rule (cont'd)

- Of the three terms on the right hand side, we know $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$, because $\hat{\mathbf{a}}$ is assumed to be a solution.
- If we select

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2$$

$$\gamma = \min_i \hat{\mathbf{a}}^t \mathbf{y}_i$$

i.e. largest of the positive term and smallest of the negative term,
then for $\alpha = \beta^2 / \gamma$,

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - \beta^2$$

- (Note the inequality remains true when the right hand side is replaced by a quantity larger than what it previously was.)
- Every correction takes the guess closer to a true solution.
- From an initialization $\mathbf{a}^{(1)}$, we will find a solution in *at most*
 $k_0 = \frac{\|\mathbf{a}^{(1)} - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}$ updates.

Perceptron

Convergence of the learning rule (cont'd)

- Of the three terms on the right hand side, we know $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$, because $\hat{\mathbf{a}}$ is assumed to be a solution.
- If we select

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2$$

$$\gamma = \min_i \hat{\mathbf{a}}^t \mathbf{y}_i$$

i.e. largest of the positive term and smallest of the negative term,
then for $\alpha = \beta^2 / \gamma$,

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - \beta^2$$

- (Note the inequality remains true when the right hand side is replaced by a quantity larger than what it previously was.)
- Every correction takes the guess closer to a true solution.
- From an initialization $\mathbf{a}^{(1)}$, we will find a solution in *at most*
 $k_0 = \frac{\|\mathbf{a}^{(1)} - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}$ updates.

Perceptron

Convergence of the learning rule (cont'd)

- Of the three terms on the right hand side, we know $\hat{\mathbf{a}}^t \mathbf{y}(k) > 0$, because $\hat{\mathbf{a}}$ is assumed to be a solution.
- If we select

$$\beta^2 = \max_i \|\mathbf{y}_i\|^2$$

$$\gamma = \min_i \hat{\mathbf{a}}^t \mathbf{y}_i$$

i.e. largest of the positive term and smallest of the negative term, then for $\alpha = \beta^2 / \gamma$,

$$\|\mathbf{a}^{(k+1)} - \alpha \hat{\mathbf{a}}\|^2 < \|\mathbf{a}^{(k)} - \alpha \hat{\mathbf{a}}\|^2 - \beta^2$$

- (Note the inequality remains true when the right hand side is replaced by a quantity larger than what it previously was.)
- Every correction takes the guess closer to a true solution.
- From an initialization $\mathbf{a}^{(1)}$, we will find a solution in *at most* $k_0 = \frac{\|\mathbf{a}^{(1)} - \alpha \hat{\mathbf{a}}\|^2}{\beta^2}$ updates.

- Linear regression
 - Solution as pseudo inverse
 - Solution by gradient descent
 - Regularization
- Perceptron
 - Setting up a suitable error function
 - Convergence of the algorithm