

Machine Learning: Optimization 1

Hiroshi Shimodaira and Hao Tang

2025

Ver. 1.0.2

Optimisation problems we saw so far

- Squared error in Linear regression

$$\min_{\mathbf{w}} L$$
$$L = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

Optimisation problems we saw so far

- Squared error in Linear regression

$$\min_{\mathbf{w}} L$$

$$L = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

- Log-loss in logistic regression

$$\min_{\mathbf{w}} \text{NLL}$$

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

Optimisation problems we saw so far

- Squared error in Linear regression

$$\min_{\mathbf{w}} L$$

$$L = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

- Log-loss in logistic regression

$$\min_{\mathbf{w}} \text{NLL}$$

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

- How do we find the optimal solution?

Optimisation problems we saw so far

- Squared error in Linear regression

$$\min_{\mathbf{w}} L$$

$$L = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

- Log-loss in logistic regression

$$\min_{\mathbf{w}} \text{NLL}$$

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

- How do we find the optimal solution?
- How do we know the solution is optimal?

Optimisation problems we saw so far

- Squared error in Linear regression

$$\min_{\mathbf{w}} L$$

$$L = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

- Log-loss in logistic regression

$$\min_{\mathbf{w}} \text{NLL}$$

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

- How do we find the optimal solution?
- How do we know the solution is optimal?
- How do we know the solution is unique?

A toy example of constrained optimisation problem

Question: Given a rope of length L meters, find the rectangle with the largest area that can be formed using the rope.

A toy example of constrained optimisation problem

Question: Given a rope of length L meters, find the rectangle with the largest area that can be formed using the rope.

Optimisation problem: Letting the area $S(x, y) = xy$,

$$\begin{aligned} & \max_{x,y} S(x,y) \\ & \text{subject to: } 2(x+y) = L \end{aligned}$$

A toy example of constrained optimisation problem

Question: Given a rope of length L meters, find the rectangle with the largest area that can be formed using the rope.

Optimisation problem: Letting the area $S(x, y) = xy$,

$$\begin{aligned} & \max_{x,y} S(x, y) \\ & \text{subject to: } 2(x + y) = L \end{aligned}$$

$$\begin{aligned} S &= xy = x\left(\frac{L}{2} - x\right) \\ &= -x^2 + \frac{L}{2}x \\ &= -(x - \frac{L}{4})^2 + \frac{L^2}{16} \end{aligned}$$

A toy example of constrained optimisation problem

Question: Given a rope of length L meters, find the rectangle with the largest area that can be formed using the rope.

Optimisation problem: Letting the area $S(x, y) = xy$,

$$\begin{aligned} & \max_{x,y} S(x, y) \\ & \text{subject to: } 2(x + y) = L \end{aligned}$$

$$\begin{aligned} S &= xy = x\left(\frac{L}{2} - x\right) \\ &= -x^2 + \frac{L}{2}x \\ &= -(x - \frac{L}{4})^2 + \frac{L^2}{16} \end{aligned}$$

$$\frac{dS}{dx} = -2x + \frac{L}{2}, \quad \frac{d^2S}{dx^2} = -2$$

Notation

- Vector $\mathbf{x} \in \mathbb{R}^d$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = [x_1 \ x_2 \ \cdots \ x_d]^\top$$

Notation

- Vector $\mathbf{x} \in \mathbb{R}^d$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = [x_1 \ x_2 \ \cdots \ x_d]^\top$$

- Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$
$$= [a_1 \ a_2 \ \cdots \ a_n]$$

$$= \begin{bmatrix} \tilde{\mathbf{a}}_1^\top \\ \tilde{\mathbf{a}}_2^\top \\ \vdots \\ \tilde{\mathbf{a}}_m^\top \end{bmatrix}$$

Notation (*cont.*)

- Data matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}$$

Notation (*cont.*)

- Data matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nd} \end{bmatrix}$$

Vector matrix multiplication

$$\mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{x} \in \mathbb{R}^d$$

$$\begin{aligned}\mathbf{Ax} &= \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1d} \\ a_{21} & a_{22} & \cdots & a_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{md} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \\ &= \begin{bmatrix} a_{11}x_1 + \cdots + a_{1d}x_d \\ \vdots \\ a_{m1}x_1 + \cdots + a_{md}x_d \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{a}}_1^\top \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_m^\top \mathbf{x} \end{bmatrix}\end{aligned}$$

Gradients and Hessians

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

- Gradient of $f(\mathbf{x})$

$$\frac{\partial f}{\partial \mathbf{x}} = \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \quad (1)$$

NB: See [Wikipedia: Matrix Calculus](#) for layout conventions

Gradients and Hessians

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

- Gradient of $f(\mathbf{x})$

$$\frac{\partial f}{\partial \mathbf{x}} = \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \quad (1)$$

- Hessian of $f(\mathbf{x})$

$$\mathbf{H}_f = \frac{\partial^2 f}{\partial \mathbf{x}^2} = \nabla^2 f = \begin{bmatrix} \frac{\partial f}{\partial x_1^2} & \cdots & \frac{\partial f}{\partial x_1 x_d} \\ & \vdots & \\ \frac{\partial f}{\partial x_d x_1} & \cdots & \frac{\partial f}{\partial x_d^2} \end{bmatrix} \quad (2)$$

NB: See [Wikipedia: Matrix Calculus](#) for layout conventions

Useful formula 1

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (3)$$

Useful formula 1

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (3)$$

$$\mathbf{a}^\top \mathbf{x} = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d$$

Useful formula 1

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (3)$$

$$\mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial x_1} \\ \vdots \\ \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial x_d} \end{bmatrix}$$

Useful formula 1

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (3)$$

$$\mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial x_1} \\ \vdots \\ \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial x_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_d \end{bmatrix} = \mathbf{a}$$

Useful formula 2

$$\frac{\mathbf{b}^\top \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{b}$$

Useful formula 2

$$\frac{\mathbf{b}^\top \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{b}$$

Let $\mathbf{a}^\top = \mathbf{b}^\top \mathbf{A}$

$$\frac{\partial \mathbf{b}^\top \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}}$$

Useful formula 2

$$\frac{\mathbf{b}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{b}$$

Let $\mathbf{a}^\top = \mathbf{b}^\top \mathbf{A}$

$$\frac{\partial \mathbf{b}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} = (\mathbf{b}^\top \mathbf{A})^\top = \mathbf{A}^\top \mathbf{b}$$

Useful formula 3

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

Useful formula 3

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}}_1^\top \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_d^\top \mathbf{x} \end{bmatrix}$$

Useful formula 3

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

$$\begin{aligned}\mathbf{x}^T \mathbf{A} \mathbf{x} &= \mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}}_1^T \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_d^T \mathbf{x} \end{bmatrix} \\ &= x_1 \tilde{\mathbf{a}}_1^T \mathbf{x} + x_2 \tilde{\mathbf{a}}_2^T \mathbf{x} + \cdots + x_d \tilde{\mathbf{a}}_d^T \mathbf{x}\end{aligned}$$

Useful formula 3

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

$$\begin{aligned}\mathbf{x}^T \mathbf{A} \mathbf{x} &= \mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}}_1^\top \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_d^\top \mathbf{x} \end{bmatrix} \\ &= x_1 \tilde{\mathbf{a}}_1^\top \mathbf{x} + x_2 \tilde{\mathbf{a}}_2^\top \mathbf{x} + \cdots + x_d \tilde{\mathbf{a}}_d^\top \mathbf{x}\end{aligned}$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial x_1} \\ \vdots \\ \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial x_d} \end{bmatrix}$$

Useful formula 3

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

$$\begin{aligned}\mathbf{x}^T \mathbf{A} \mathbf{x} &= \mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}}_1^T \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_d^T \mathbf{x} \end{bmatrix} \\ &= x_1 \tilde{\mathbf{a}}_1^T \mathbf{x} + x_2 \tilde{\mathbf{a}}_2^T \mathbf{x} + \cdots + x_d \tilde{\mathbf{a}}_d^T \mathbf{x}\end{aligned}$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial x_1} \\ \vdots \\ \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial x_d} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{a}}_1^T \mathbf{x} + \mathbf{a}_1^T \mathbf{x} \\ \vdots \\ \tilde{\mathbf{a}}_d^T \mathbf{x} + \mathbf{a}_d^T \mathbf{x} \end{bmatrix} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

MSE in linear regression

$$L = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

MSE in linear regression

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{N} \left[\begin{array}{ccc} (\mathbf{w}^\top \mathbf{x}_1 - y_1) & \cdots & (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{array} \right] \begin{bmatrix} (\mathbf{w}^\top \mathbf{x}_1 - y_1) \\ \vdots \\ (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{bmatrix} \end{aligned}$$

MSE in linear regression

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{N} \left[(\mathbf{w}^\top \mathbf{x}_1 - y_1) \quad \cdots \quad (\mathbf{w}^\top \mathbf{x}_N - y_N) \right] \begin{bmatrix} (\mathbf{w}^\top \mathbf{x}_1 - y_1) \\ \vdots \\ (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{bmatrix} \\ &= \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \end{aligned}$$

MSE in linear regression

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{N} \left[\begin{array}{ccc} (\mathbf{w}^\top \mathbf{x}_1 - y_1) & \cdots & (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{array} \right] \begin{bmatrix} (\mathbf{w}^\top \mathbf{x}_1 - y_1) \\ \vdots \\ (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{bmatrix} \\ &= \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

MSE in linear regression

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \\ &= \frac{1}{N} \left[\begin{array}{ccc} (\mathbf{w}^\top \mathbf{x}_1 - y_1) & \cdots & (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{array} \right] \begin{bmatrix} (\mathbf{w}^\top \mathbf{x}_1 - y_1) \\ \vdots \\ (\mathbf{w}^\top \mathbf{x}_N - y_N) \end{bmatrix} \\ &= \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \\ &= \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} - 2\mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{y}^\top \mathbf{y} \right) \end{aligned}$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{1}{N} \left(\frac{\partial \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} - 2 \frac{\partial \mathbf{y}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \right)$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= \frac{1}{N} \left(\frac{\partial \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} - 2 \frac{\partial \mathbf{y}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \right) \\ &= \frac{1}{N} \left((\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} \right)\end{aligned}$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= \frac{1}{N} \left(\frac{\partial \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} - 2 \frac{\partial \mathbf{y}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \right) \\ &= \frac{1}{N} \left((\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} \right) \\ &= \frac{2}{N} (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y})\end{aligned}$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= \frac{1}{N} \left(\frac{\partial \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} - 2 \frac{\partial \mathbf{y}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \right) \\ &= \frac{1}{N} \left((\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} \right) \\ &= \frac{2}{N} (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y}) \\ &= \frac{2}{N} \mathbf{X}^\top (\mathbf{X} \mathbf{w} - \mathbf{y})\end{aligned}$$

Gradient of MSE in linear regression

$$L = \frac{1}{N} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right)$$

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}} &= \frac{1}{N} \left(\frac{\partial \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} - 2 \frac{\partial \mathbf{y}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \right) \\ &= \frac{1}{N} \left((\mathbf{X}^\top \mathbf{X} + (\mathbf{X}^\top \mathbf{X})^\top) \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} \right) \\ &= \frac{2}{N} (\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y}) \\ &= \frac{2}{N} \mathbf{X}^\top (\mathbf{X} \mathbf{w} - \mathbf{y})\end{aligned}$$

Letting $\frac{\partial L}{\partial \mathbf{w}} = \mathbf{0}$ yields

$$\mathbf{X}^\top (\mathbf{X} \mathbf{w} - \mathbf{y}) = \mathbf{0}$$

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

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

Log loss in logistic regression training

Letting $(1, \mathbf{x})$ and (b, \mathbf{w}) denoted as \mathbf{x} and \mathbf{w} , respectively,

$$\begin{aligned} \text{LL} &= \log \prod_{i=1}^N p(y_i | \mathbf{x}_i, \theta) = \sum_{i=1}^N \log \frac{1}{1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)} \\ &= \sum_{i=1}^N -\log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right) \end{aligned}$$

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

Gradient of log loss

$$\text{NLL} = \sum_{i=1}^N \log \left(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)$$

$$\begin{aligned}\frac{\partial \text{NLL}}{\partial \mathbf{w}} &= \sum_{i=1}^N \frac{\exp(-y_i \mathbf{w}^\top \mathbf{x}_i)}{1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)} (-y_i \mathbf{x}_i) \\ &= \sum_{i=1}^N \left(1 - \frac{1}{1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)} \right) (-y_i \mathbf{x}_i)\end{aligned}$$

What if we use 0/1 labels instead of -1/+1?

$$\text{NLL} = -\frac{1}{N} \sum_{i=1}^N y_i \log s_i + (1 - y_i) \log(1 - s_i)$$

where $s_i = \sigma(\mathbf{w}^\top \mathbf{x}_i) = \frac{1}{1 + \exp(-(\mathbf{w}^\top \mathbf{x}_i))}$.

$$\frac{\partial \text{NLL}}{\partial \mathbf{w}} = -\frac{1}{N} \sum_{i=1}^N y_i \frac{1}{s_i} \frac{\partial s_i}{\partial \mathbf{w}} + (1 - y_i) \frac{-1}{1 - s_i} \frac{\partial s_i}{\partial \mathbf{w}}$$

$$\text{Since } \frac{\partial s_i}{\partial \mathbf{w}} = \sigma(\mathbf{w}^\top \mathbf{x}_i)(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i))\mathbf{x}_i = s_i(1 - s_i)\mathbf{x}_i$$

$$\frac{\partial \text{NLL}}{\partial \mathbf{w}} = -\frac{1}{N} \sum_{i=1}^N (y_i(1 - s_i) - (1 - y_i)s_i) \mathbf{x}_i$$

$$= \frac{1}{N} \sum_{i=1}^N (s_i - y_i)\mathbf{x}_i$$

Vectorisation in Python

- Consider a linear regression model $\mathbf{w}^\top \mathbf{x} + b$ that has been trained already. How are the outputs calculated for a test samples $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$?
- Consider a logistic regression model $p(y=1 | \mathbf{x}) = \frac{1}{1+\exp(\mathbf{w}^\top \mathbf{x} + b)}$. How are the outputs calculated for the test samples?

Quizzes

- Show $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$
- Find $\frac{\partial \mathbf{x}^T \mathbf{B}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}}$
- Letting $\mathbf{X} \in \mathbb{R}^{N \times d}$ and $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$,
 - show that \mathbf{S} is a symmetric matrix,
 - show that $\mathbf{x}^\top \mathbf{S} \mathbf{x} \geq 0$ for any $\mathbf{x} \in \mathbb{R}^d$.