Machine Learning

Gradients of scalar functions w.r.t. matrices

Emanuele Rodolà rodola@di.uniroma1.it



Linear regression: Matrix notation

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} a \\ b \end{pmatrix}}_{\boldsymbol{\theta}}$$

This expresses all the equations $y_i = ax_i + b$ at once and makes the linearity w.r.t. a,b evident.

The MSE is simply:

$$\ell(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 = \mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Linear regression: Matrix notation

$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} a \\ b \end{pmatrix}}_{\boldsymbol{\theta}}$$

This expresses all the equations $y_i = ax_i + b$ at once and makes the linearity w.r.t. a, b evident.

The MSE is simply:

$$\ell(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 = \mathbf{y}^{\mathsf{T}}\mathbf{y} - 2\mathbf{y}^{\mathsf{T}}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\boldsymbol{\theta}$$

Setting the gradient w.r.t. θ to zero:

$$-2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} = \mathbf{0}$$

$$\mathbf{y}^{\mathsf{T}}\mathbf{y} - 2\mathbf{y}^{\mathsf{T}}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\boldsymbol{\theta} \quad \stackrel{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\mathsf{T}}\mathbf{y} + 2\mathbf{X}^{\mathsf{T}}\mathbf{X}\boldsymbol{\theta}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \overset{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \begin{pmatrix} \theta_1 & \cdots & \theta_n \end{pmatrix} \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_n \end{pmatrix}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \overset{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{n} \sum_{j=1}^{n} a_{ij} \theta_{i} \theta_{j}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \overset{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \sum_{i=1}^n \sum_{j=1}^n a_{ij} \theta_i \theta_j \\ \vdots \\ \frac{\partial}{\partial \theta_n} \sum_{i=1}^n \sum_{j=1}^n a_{ij} \theta_i \theta_j \end{pmatrix}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \overset{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \begin{pmatrix} \sum_{j} a_{1j} \theta_{j} + \sum_{i} a_{i1} \theta_{i} \\ \vdots \\ \sum_{j} a_{nj} \theta_{j} + \sum_{i} a_{in} \theta_{i} \end{pmatrix}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \overset{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \begin{pmatrix} \sum_{i} (a_{1i} + a_{i1}) \theta_{i} \\ \vdots \\ \sum_{i} (a_{ni} + a_{in}) \theta_{i} \end{pmatrix}$$

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \stackrel{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = (\mathbf{A} + \mathbf{A}^{\top})\boldsymbol{\theta}$$

How did we compute this gradient?

$$\mathbf{y}^{\top}\mathbf{y} - 2\mathbf{y}^{\top}\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta} \quad \stackrel{\nabla_{\boldsymbol{\theta}}}{\Longrightarrow} \quad -2\mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\theta}$$

Example:
$$f(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{A} \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = (\mathbf{A} + \mathbf{A}^{\top})\boldsymbol{\theta}$$

If A is symmetric (e.g., $A = X^{T}X$), then:

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = 2\mathbf{A}\boldsymbol{\theta}$$

Gradient w.r.t. a vector or a matrix

Solve the following exercises:

- Compute $\nabla_{\boldsymbol{\theta}} \left(\mathbf{y}^{\top} \mathbf{y} 2 \mathbf{y}^{\top} \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\theta}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\theta} \right)$
- Compute $\nabla_{\mathbf{\Theta}} \| \mathbf{Y}^{\top} \mathbf{X}^{\top} \mathbf{\Theta} \|_F^2$
- Set $\nabla_{\Theta} \| \mathbf{Y}^{\top} \mathbf{X}^{\top} \mathbf{\Theta} \|_F^2 = \mathbf{0}$ and solve for $\mathbf{\Theta}$ (this gives the solution reported in "Linear regression: Higher dimensions" in the main deck of slides)

In the first two exercises, you can either compute all the derivatives by yourself, or you can use pre-computed formulas from a book.

For example, the gradient for the term $\nabla_{\Theta} \mathrm{tr}(\Theta^{\top} X X^{\top} \Theta)$ (which arises in the second exercise) is found in Equation (108) of the Matrix Cookbook.

Suggested reading

For several pre-computed matrix derivatives, refer to the book:

K.B. Petersen & M.S. Pedersen, "The Matrix Cookbook". Technical University of Denmark, 2012

Public download link: https://www2.imm.dtu.dk/pubdb/edoc/imm3274.pdf