

# MML lecture extra notes, week Oct 25 - 29, 2021

Linear regression considers solving the following task:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}), \quad L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2. \quad (1)$$

## Arithmetico-geometric sequence

In linear regression, gradient descent returns an update rule as

$$\boldsymbol{\theta}_{t+1} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y}. \quad (2)$$

The solution of this iterative update is related to an arithmetico-geometric sequence. Writing  $\boldsymbol{\theta}_{t+1} + \boldsymbol{\beta} = \mathbf{A}(\boldsymbol{\theta}_t + \boldsymbol{\beta})$  with  $\mathbf{A} := \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}$ , we would like to solve for  $\boldsymbol{\beta}$  such that:

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \mathbf{A}(\boldsymbol{\theta}_t + \boldsymbol{\beta}) - \boldsymbol{\beta} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} \\ \Leftrightarrow -\frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} &= \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} \\ \Leftrightarrow \boldsymbol{\beta} &= -(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = -\boldsymbol{\theta}^*. \end{aligned} \quad (3)$$

So this immediately implies

$$\boldsymbol{\theta}_t - \boldsymbol{\theta}^* = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \Rightarrow \boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*, \quad (4)$$

which means  $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}^*$  if  $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \rightarrow \mathbf{0}$ .

## Rayleigh quotient

Assume  $\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^{-1}$  is symmetric (so that also  $\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top$ ). Consider the following *Rayleigh quotient*

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|_2^2}, \quad \|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x}. \quad (5)$$

Using the fact that  $\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$ , we can define  $\mathbf{z} = \mathbf{Q}^\top \mathbf{x}$  and rewrite the Rayleigh quotient as:

$$R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top \mathbf{x}}{\mathbf{x}^\top \mathbf{Q}\mathbf{Q}^\top \mathbf{x}} = \frac{\mathbf{z}^\top \boldsymbol{\Lambda} \mathbf{z}}{\mathbf{z}^\top \mathbf{z}}. \quad (6)$$

As  $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_D)$  is a diagonal matrix, we have (writing  $\mathbf{z} = (z_1, \dots, z_D)^\top$ )

$$\mathbf{z}^\top \boldsymbol{\Lambda} \mathbf{z} = \sum_{d=1}^D \lambda_d z_d^2. \quad (7)$$

Therefore the Rayleigh quotient can be written as the following weighted average of the eigenvalues

$$R(\mathbf{A}, \mathbf{x}) = \sum_{d=1}^D \frac{z_d^2}{\|\mathbf{z}\|_2^2} \lambda_d, \quad \text{with } \sum_{d=1}^D \frac{z_d^2}{\|\mathbf{z}\|_2^2} = 1. \quad (8)$$

In summary, these derivation indicate that the Rayleigh quotient is bounded as

$$\begin{aligned} \lambda_{\min}(\mathbf{A}) &\leq R(\mathbf{A}, \mathbf{x}) \leq \lambda_{\max}(\mathbf{A}) \\ \Rightarrow \lambda_{\min}(\mathbf{A}) \|\mathbf{x}\|_2^2 &\leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A}) \|\mathbf{x}\|_2^2, \end{aligned} \quad (9)$$

where  $\lambda_{\min}(\mathbf{A})$  and  $\lambda_{\max}(\mathbf{A})$  are the smallest and largest eigenvalues of  $\mathbf{A}$ , respectively.

### An extra exercise

Show that solving linear regression using gradient descent with momentum, if converges, converges to  $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

Hint: consider the simpler case with fixed step size  $\gamma$  and momentum factor  $\alpha$ . Follow the below steps and practice your linear algebra skills :)

1. Write down the update equations for the parameters  $\boldsymbol{\theta}_t$  and the momentum  $\Delta \boldsymbol{\theta}_t$ ;
2. Collect both terms as a long vector  $(\boldsymbol{\theta}_t^\top, \Delta \boldsymbol{\theta}_t^\top)^\top$ , and merge the two linear update equations in step 1 into one “joint” linear equation using block matrices;
3. Apply the analysis techniques in GD for linear regression to show the converged solution (if converges).

**Solution** First note that  $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) = \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta}_t - \mathbf{y})$ . The update equations for both the parameter and the momentum are

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \gamma \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) + \alpha \Delta \boldsymbol{\theta}_t = \boldsymbol{\theta}_t - \frac{\gamma}{\sigma^2} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta}_t - \mathbf{y}) + \alpha \Delta \boldsymbol{\theta}_t \\ \Delta \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t = \alpha \Delta \boldsymbol{\theta}_t - \frac{\gamma}{\sigma^2} \mathbf{X}^\top (\mathbf{X} \boldsymbol{\theta}_t - \mathbf{y}).\end{aligned}\tag{10}$$

Now collecting both equations together into a “joint” linear equation:

$$\begin{bmatrix} \boldsymbol{\theta}_{t+1} \\ \Delta \boldsymbol{\theta}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} & \alpha \mathbf{I} \\ -\frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} & \alpha \mathbf{I} \end{bmatrix} \begin{bmatrix} \boldsymbol{\theta}_t \\ \Delta \boldsymbol{\theta}_t \end{bmatrix} + \begin{bmatrix} \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} \\ \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} \end{bmatrix}.\tag{11}$$

Then we can apply the derivation of arithmetico-geometric sequences above, and show that

$$\begin{bmatrix} \boldsymbol{\theta}_t \\ \Delta \boldsymbol{\theta}_t \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} & \alpha \mathbf{I} \\ -\frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} & \alpha \mathbf{I} \end{bmatrix}^t \begin{bmatrix} \boldsymbol{\theta}_0 - \boldsymbol{\theta}^* \\ \Delta \boldsymbol{\theta}_0 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\theta}^* \\ \mathbf{0} \end{bmatrix},\tag{12}$$

with  $\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ . This equation also says if momentum GD converges, the momentum  $\Delta \boldsymbol{\theta}_t$  will vanish to  $\mathbf{0}$ , which is as expected as  $\Delta \boldsymbol{\theta}_t = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1} \rightarrow \mathbf{0}$ .