

Directional derivatives for matrix calculus

Benjamin Chu

Graduate program in Biomathematics
Department of Computational Medicine
University of California, Los Angeles

January 28, 2021

Introduction

- In statistics, we often want to maximize things, requiring gradients and Hessians
- One way to obtain gradients and Hessians is through directional derivatives
- Compared to traditional matrix calculus, its primary advantages are:
 1. Gradient and Hessian are implicit in the 1st and 2nd directional derivatives
 2. Variables remain intact until gradients/Hessians are explicitly needed (e.g. \mathbf{X} , \mathbf{Y} stay as \mathbf{X} , \mathbf{Y} instead of $\text{vec}(\mathbf{X})$...etc)
 3. Scalar/vector/matrix valued functions have scalar/vector/matrix valued directional derivatives

Definition

The (Hadamard semi-) directional derivative of $f(\mathbf{x})$ in the direction \mathbf{v} is the limit

$$d_{\mathbf{v}}f(\mathbf{x}) = \lim_{\substack{h \rightarrow 0 \\ \mathbf{w} \rightarrow \mathbf{v}}} \frac{f(\mathbf{x} + h\mathbf{w}) - f(\mathbf{x})}{h}.$$

Directional derivatives enjoys the chain, sum, and product rules (Proposition 3.2.4 of MM optimization)

$$\begin{aligned}d_{\mathbf{v}}[f \circ g(\mathbf{x})] &= d_{d_{\mathbf{v}}g(\mathbf{x})}f[g(\mathbf{x})] \\d_{\mathbf{v}}[f(\mathbf{x}) + g(\mathbf{x})] &= d_{\mathbf{v}}f(\mathbf{x}) + d_{\mathbf{v}}g(\mathbf{x}) \\d_{\mathbf{v}}[f(\mathbf{x})g(\mathbf{x})] &= d_{\mathbf{v}}f(\mathbf{x})g(\mathbf{x}) + f(\mathbf{x})d_{\mathbf{v}}g(\mathbf{x})\end{aligned}$$

Furthermore, for a differentiable function $f(\mathbf{x})$, we have $d_{\mathbf{v}}f(\mathbf{x}) = df(\mathbf{x})\mathbf{v}$, where $df(\mathbf{x}) = \nabla f(\mathbf{x})^t$ is the first differential of $f(\mathbf{x})$. For a differentiable function $f(\mathbf{x}, \mathbf{y})$ of two variables, we also have $d_{(\mathbf{u}, \mathbf{v})}f(\mathbf{x}, \mathbf{y}) = d_{(\mathbf{u}, 0)}f(\mathbf{x}, \mathbf{y}) + d_{(0, \mathbf{v})}f(\mathbf{x}, \mathbf{y})$.

Second Differentials

The second directional derivative is similarly defined

$$d_{\mathbf{u}}[d_{\mathbf{v}}f(\mathbf{x})] = \lim_{\substack{h \rightarrow 0 \\ \tilde{\mathbf{w}} \rightarrow \mathbf{u}}} \frac{d_{\mathbf{v}}f(\mathbf{x} + h\tilde{\mathbf{w}}) - d_{\mathbf{v}}f(\mathbf{x})}{h}.$$

For $f(\mathbf{x})$ twice differentiable, the second differential $d_{\mathbf{u}}[d_{\mathbf{v}}f(\mathbf{x})]$ is linear in \mathbf{u} for \mathbf{v} fixed and linear in \mathbf{v} for \mathbf{u} fixed. It is accordingly a bilinear form. When we set $\mathbf{u} = \mathbf{v}$, we get a quadratic form $d_{\mathbf{v}}^2f(\mathbf{x})$. At this juncture we have not yet associated a matrix to the quadratic form. In fact, if \mathbf{x} is a matrix, then the matrix is replaced by a tensor. We can identify a matrix if we vectorize \mathbf{x} by stacking its columns into a matrix.

Extracting gradients from directional derivatives

For vector \mathbf{x} , we can extract gradients and Hessians for twice differentiable functions as

$$d_{\mathbf{v}}f(\mathbf{x}) = \nabla f(\mathbf{x}) \cdot \mathbf{v} = (\nabla f(\mathbf{x}))^t \mathbf{v}$$

$$d_{\mathbf{v}}^2 f(\mathbf{x}) = \mathbf{v}^t \mathbf{H} \mathbf{v}$$

For matrix \mathbf{X} , we extract these as

$$d_{\mathbf{V}}f(\mathbf{X}) = \text{vec}(\nabla f(\mathbf{X}))^t \text{vec}(\mathbf{V}) = \text{tr}((\nabla f(\mathbf{X}))^t \mathbf{V})$$

$$d_{\mathbf{V}}^2 f(\mathbf{X}) = \text{vec}(\mathbf{V})^t \mathbf{H} \text{vec}(\mathbf{V})$$

Proof.

The gradients are covered in Prop 3.2.1 of MM Optimization. For matrix gradients, the trace function induces an inner product on matrices by analogy with the ordinary inner product for vectors. The formula $\text{tr}(\mathbf{A}^t \mathbf{B}) = \text{vec}(\mathbf{A})^t \text{vec}(\mathbf{B})$ converts the trace inner product to the vector inner product. □

Warm up 1: find $d_{\mathbf{V}}[f(\mathbf{X})^{-1}]$ for \mathbf{X} and $f(\mathbf{X})$ matrices.

$$\begin{aligned}d_{\mathbf{V}}f(\mathbf{X})^{-1} &= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{f(\mathbf{X} + t\mathbf{W})^{-1} - f(\mathbf{X})^{-1}}{t} \\&= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} f(\mathbf{X} + t\mathbf{W})^{-1} \frac{1 - f(\mathbf{X} + t\mathbf{W})f(\mathbf{X})^{-1}}{t} \\&= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} f(\mathbf{X} + t\mathbf{W})^{-1} \frac{f(\mathbf{X}) - f(\mathbf{X} + t\mathbf{W})}{t} f(\mathbf{X})^{-1} \\&= -f(\mathbf{X})^{-1}d_{\mathbf{V}}f(\mathbf{X})f(\mathbf{X})^{-1}.\end{aligned}$$

Warm up 2: directional derivative of log determinants

For a symmetric positive definite \mathbf{X} and a symmetric direction \mathbf{V} ,

$$\begin{aligned}d_{\mathbf{V}} \ln \det(\mathbf{X}) &= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\ln \det(\mathbf{X} + t\mathbf{W}) - \ln \det(\mathbf{X})}{t} \\&= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\ln \left[\det(\mathbf{X}^{1/2}) \det(\mathbf{I} + t\mathbf{X}^{-1/2}\mathbf{W}\mathbf{X}^{-1/2})\mathbf{X}^{1/2} \right] - \ln \det(\mathbf{X})}{t} \\&= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\ln \det(\mathbf{I} + t\mathbf{X}^{-1/2}\mathbf{W}\mathbf{X}^{-1/2})}{t} \\&= \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\sum \ln(1 + t\lambda_i)}{t} \quad \left(\lambda = \text{eigenvalues of } \mathbf{X}^{-1/2}\mathbf{W}\mathbf{X}^{-1/2} \right) \\&\approx \lim_{\substack{h \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\sum t\lambda_i}{t} \quad (\ln(1+x) \approx x, \text{ for small } x) \\&= \sum \lambda_i = \text{tr} \left(\mathbf{X}^{-1/2}\mathbf{V}\mathbf{X}^{-1/2} \right) = \text{tr}(\mathbf{X}^{-1}\mathbf{V}).\end{aligned}$$

Detailed example: evaluate gradient of LMM

The loglikelihood for the i th sample in a linear mixed model is

$$\ell_i(\boldsymbol{\beta}, \mathbf{L}, \sigma^2) = -\frac{1}{2} \log \det \Omega_i - \frac{1}{2} (\mathbf{y} - \mathbf{X}_i \boldsymbol{\beta})^T \Omega_i^{-1} (\mathbf{y} - \mathbf{X}_i \boldsymbol{\beta}).$$

Prove that the gradient is

$$\begin{aligned}\nabla_{\boldsymbol{\beta}} \ell_i(\boldsymbol{\beta}, \mathbf{L}, \sigma^2) &= \mathbf{X}_i^T \Omega_i^{-1} \mathbf{r}_i, \\ \nabla_{\sigma^2} \ell_i(\boldsymbol{\beta}, \mathbf{L}, \sigma^2) &= -\frac{1}{2} \text{tr}(\Omega_i^{-1}) + \frac{1}{2} \mathbf{r}_i^T \Omega_i^{-2} \mathbf{r}_i, \\ \nabla_{\mathbf{L}} \ell_i(\boldsymbol{\beta}, \mathbf{L}, \sigma^2) &= -\mathbf{Z}_i^T \Omega_i^{-1} \mathbf{Z}_i \mathbf{L} + \mathbf{Z}_i^T \Omega_i^{-1} \mathbf{r}_i \mathbf{r}_i^T \Omega_i^{-1} \mathbf{Z}_i \mathbf{L},\end{aligned}$$

where $\Omega_i = \sigma^2 \mathbf{I} + \mathbf{Z}_i \Sigma \mathbf{Z}_i^t = \sigma^2 \mathbf{I} + \mathbf{Z}_i \mathbf{L} \mathbf{L}^t \mathbf{Z}_i^t$ and $\mathbf{r} = \mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}$. We optimize over the Cholesky factor \mathbf{L} instead of Σ directly because the later needs to be positive semidefinite. The diagonal entries of \mathbf{L} must be positive, simpler constraint.

1st term: $d_{\mathbf{v}}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})$ wrt to $\boldsymbol{\beta}$

Using symmetry of Ω_i ,

$$\begin{aligned} & -\frac{1}{2}d_{\mathbf{v}}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta}) \\ &= -\frac{1}{2} \lim_{\substack{t \rightarrow 0 \\ \mathbf{w} \rightarrow \mathbf{v}}} \frac{(\mathbf{y} - \mathbf{X}_i(\boldsymbol{\beta} + t\mathbf{w}))^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i(\boldsymbol{\beta} + t\mathbf{w})) - (\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})}{t} \\ &= -\frac{1}{2} \lim_{\substack{t \rightarrow 0 \\ \mathbf{w} \rightarrow \mathbf{v}}} \frac{-t(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1} \mathbf{X}_i \mathbf{w} - t(\mathbf{X}_i \mathbf{w})^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta}) + O(t^2)}{t} \\ &= (\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1} \mathbf{X}_i \mathbf{v}. \end{aligned}$$

To extract gradients, we set

$$(\nabla_{\boldsymbol{\beta}} \ell)^T \mathbf{v} \equiv (\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta})^T \Omega_i^{-1} \mathbf{X}_i \mathbf{v}$$

So as desired,

$$\nabla_{\boldsymbol{\beta}} \ell = \mathbf{X}_i^T \Omega_i^{-1}(\mathbf{y} - \mathbf{X}_i\boldsymbol{\beta}).$$

2nd term part I: $d_v \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)$ wrt to σ^2

Chain rule says we must compute

$$d_{d_v \sigma^2 \mathbf{I}} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t).$$

First evaluate the direction

$$d_v \sigma^2 \mathbf{I} = \lim_{\substack{t \rightarrow 0 \\ w \rightarrow v}} \frac{(\sigma^2 + tw) \mathbf{I} - \sigma^2 \mathbf{I}}{t} = \lim_{\substack{t \rightarrow 0 \\ w \rightarrow v}} \frac{tw \mathbf{I}}{t} = v \mathbf{I}$$

Thus

$$\begin{aligned} & -\frac{1}{2} d_{d_v \sigma^2 \mathbf{I}} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t) \\ &= -\frac{1}{2} d_{v \mathbf{I}} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t) \\ &= -\frac{1}{2} \operatorname{tr} ((\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} v \mathbf{I}) \quad (\text{warm up 2}) \end{aligned}$$

2nd term part II: $d_v \mathbf{r}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \mathbf{r}$ wrt to σ^2

In our notation $\Omega_i = \sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t$, so

$$\begin{aligned}
 & d_v \left[-\frac{1}{2} \mathbf{r}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \mathbf{r} \right] \\
 &= -\frac{1}{2} \mathbf{r}^t d_v [(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1}] \mathbf{r} \\
 &= \frac{1}{2} \mathbf{r}^t \Omega_i^{-1} d_v (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t) \Omega_i^{-1} \mathbf{r} \quad (\text{warm up 1}) \\
 &= \frac{1}{2} \mathbf{r}^t \Omega_i^{-1} \lim_{\substack{t \rightarrow 0 \\ w \rightarrow v}} \frac{(\sigma^2 + tw) \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t - (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)}{t} \Omega_i^{-1} \mathbf{r} \\
 &= \frac{1}{2} \mathbf{r}^t \Omega_i^{-1} \lim_{\substack{t \rightarrow 0 \\ w \rightarrow v}} \frac{tw \mathbf{I}}{t} \Omega_i^{-1} \mathbf{r} \\
 &= \frac{1}{2} \mathbf{r}^t \Omega_i^{-1} v \mathbf{I} \Omega_i^{-1} \mathbf{r} \\
 &= \frac{1}{2} \mathbf{r}^t \Omega_i^{-2} \mathbf{r} v \quad (v \text{ is constant})
 \end{aligned}$$

2nd term: extracting gradients

In summary, the full 1st directional derivative wrt to σ^2 is

$$d_v \ell = -\frac{1}{2} \text{tr} \left((\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} v \mathbf{I} \right) + \frac{1}{2} \mathbf{r}^t \Omega_i^{-2} \mathbf{r} v$$

To extract gradients, we set

$$v(\nabla_{\sigma^2} \ell) \equiv -\frac{1}{2} \text{tr} \left((\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \right) v + \frac{1}{2} \mathbf{r}^t \Omega_i^{-2} \mathbf{r} v$$

So as desired,

$$\nabla_{\sigma^2} \ell = -\frac{1}{2} \text{tr} \left(\Omega_i^{-1} \right) + \frac{1}{2} \mathbf{r}^t \Omega_i^{-2} \mathbf{r}$$

3rd term part I: $d_{\mathbf{V}} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)$ wrt to \mathbf{L}

Chain rule says we must compute

$$d_{d_{\mathbf{V}} \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t).$$

First evaluate the direction

$$\begin{aligned} d_{\mathbf{V}} \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t &= \lim_{\substack{t \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\mathbf{Z}(\mathbf{L} + t\mathbf{W})(\mathbf{L} + t\mathbf{W})^t \mathbf{Z}^t - \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t}{t} \\ &= \lim_{\substack{t \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{t \mathbf{Z} \mathbf{L} \mathbf{W}^t \mathbf{Z}^t + t \mathbf{Z} \mathbf{W} \mathbf{L}^t \mathbf{Z}^t + t^2 \mathbf{Z} \mathbf{W} \mathbf{W}^t \mathbf{Z}^t}{t} \\ &= \mathbf{Z} \mathbf{L} \mathbf{V}^t \mathbf{Z}^t + \mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t = 2 \mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t. \end{aligned}$$

Thus

$$\begin{aligned} &d_{\mathbf{V}} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t) \\ &= d_{2 \mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t} \ln \det(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t) && \text{(chain rule)} \\ &= 2 \operatorname{tr}((\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} (\mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t)) && \text{(warm up 2)} \\ &= 2 \operatorname{tr}([\mathbf{L}^t \mathbf{Z}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \mathbf{Z}] \mathbf{V}) \end{aligned}$$

3rd term part II: $d_{\mathbf{V}} \mathbf{r}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z})^{-1} \mathbf{r}$ wrt \mathbf{L}

In our notation $\Omega = \sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}$, so

$$\begin{aligned}
 & d_{\mathbf{V}} \mathbf{r}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z})^{-1} \mathbf{r} \\
 &= \mathbf{r}^t d_{\mathbf{V}} [(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z})^{-1}] \mathbf{r} \\
 &= -\mathbf{r}^t \Omega_i^{-1} d_{\mathbf{V}} [(\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z})] \Omega_i^{-1} \mathbf{r} \quad (\text{warm up 1}) \\
 &= -\mathbf{r}^t \Omega_i^{-1} \left[\lim_{\substack{t \rightarrow 0 \\ \mathbf{W} \rightarrow \mathbf{V}}} \frac{\sigma^2 \mathbf{I} + \mathbf{Z}(\mathbf{L} + t\mathbf{W})(\mathbf{L} + t\mathbf{W})^t \mathbf{Z}^t - \sigma^2 \mathbf{I} - \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t}{t} \right] \Omega_i^{-1} \mathbf{r} \\
 &= -\mathbf{r}^t \Omega_i^{-1} [\mathbf{Z} \mathbf{L} \mathbf{V} \mathbf{Z}^t + \mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t] \Omega_i^{-1} \mathbf{r} \\
 &= -2 \operatorname{tr} (\mathbf{r}^t \Omega_i^{-1} \mathbf{Z} \mathbf{V} \mathbf{L}^t \mathbf{Z}^t \Omega_i^{-1} \mathbf{r}) \\
 &= -2 \operatorname{tr} (\mathbf{L}^t \mathbf{Z}^t \Omega_i^{-1} \mathbf{r} \mathbf{r}^t \Omega_i^{-1} \mathbf{Z} \mathbf{V})
 \end{aligned}$$

3rd term: extracting gradients

To extract gradients, we set

$$\begin{aligned}\text{tr}((\nabla_{\mathbf{L}}\ell)^t \mathbf{V}) &\equiv -\text{tr}([\mathbf{L}^t \mathbf{Z}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \mathbf{Z}] \mathbf{V}) \\ &\quad + \text{tr}([\mathbf{L}^t \mathbf{Z}^t \Omega_i^{-1} \mathbf{r} \mathbf{r}^t \Omega_i^{-1} \mathbf{Z}] \mathbf{V})\end{aligned}$$

So as desired,

$$\nabla_{\mathbf{L}}\ell = -\mathbf{Z}^t (\sigma^2 \mathbf{I} + \mathbf{Z} \mathbf{L} \mathbf{L}^t \mathbf{Z}^t)^{-1} \mathbf{Z} \mathbf{L} + \mathbf{Z}^t \Omega_i^{-1} \mathbf{r} \mathbf{r}^t \Omega_i^{-1} \mathbf{Z} \mathbf{L}.$$

Learning resources

1. Chapter 3 of MM Optimization by Kenneth Lange
2. Chapter 3 of Introduction to optimization and semidifferential calculus by Michel Delfour
3. Chapters 5, 6, 8, 9 of Matrix differential calculus with applications in statistics and econometrics by Magnus and Neudecker