

Matrix Algebra and Optimization for Statistics and Machine Learning

Yiyuan She

Department of Statistics, Florida State University

- ▶ Matrix differentiation

Kronecker product and the vec operator

- ▶ Given $A \in \mathbb{R}^{n \times p}$ and $B \in \mathbb{R}^{m \times q}$, $A \otimes B$ is an $nm \times pq$ matrix defined by

$$\begin{bmatrix} a_{11}B & \dots & a_{1p}B \\ \vdots & & \vdots \\ a_{n1}B & \dots & a_{np}B \end{bmatrix}$$

- ▶ $\text{vec}(A)$ is obtained by stacking the columns of A :

$$[a_{11}, \dots, a_{n1}, \dots, a_{1p}, \dots, a_{np}]^T$$

Properties

- ▶ Some basics: $(A \otimes B) \otimes C = A \otimes (B \otimes C)$,
 $(A + B) \otimes (C + D) = A \otimes C + A \otimes D + B \otimes C + B \otimes D$,
 $a \otimes A = A \otimes a = aA$
- ▶ Vector **outer-product** can be viewed as a special case:

$$\alpha \otimes \beta^T = \alpha \beta^T = \beta^T \otimes \alpha$$

- ▶ Although $A \otimes B$ and $B \otimes A$ do not equal in general, there exist permutations so that $P(A \otimes B)Q = B \otimes A$

- ▶ **Mixed-product** property:

$$(A \otimes B)(C \otimes D) = AC \otimes BD$$

where conformability is assumed

- The RHS is often much more efficient to compute
- Assume all are square: $2n^4 + n^6$ vs. $2n^3 + n^4$
- ▶ Let $A \in \mathbb{R}^{m \times m}$, $B \in \mathbb{R}^{n \times n}$. From the above property, $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$.
 - It actually holds generally for the MP inverse

- ▶ In addition, $(A \otimes B)^T = A^T \otimes B^T$
 - **No change in order**
 - A special case helps to remember: $A = I$
 - [Check the 2nd argument first when seeing \otimes]
- ▶ $|A \otimes B| = |A|^n |B|^m$ (a nice result)
- ▶ $tr(A \otimes B) = tr(A)tr(B)$

- ▶ $\text{vec}(\alpha\beta^T) = \beta \otimes \alpha$ for any vectors α, β
- ▶ More generally,

$$\text{vec}(AXB^T) = (B \otimes A) \text{vec}(X)$$

- ▶ Let $A \in \mathbb{R}^{m \times n}$. Then

$$K_{m,n} \text{vec}(A) = \text{vec}(A^T)$$

where $K_{m,n}$ is the **commutation** matrix.

- ▶ Let $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$. Then

$$K_{p,m}(A \otimes B)K_{n,q} = B \otimes A$$

- ▶ In fact, for any $X \in \mathbb{R}^{q \times n}$, $K_{p,m}(A \otimes B) \operatorname{vec} X = K_{p,m} \operatorname{vec}(BXA^T) = \operatorname{vec}(AX^TB^T) = (B \otimes A) \operatorname{vec}(X^T) = (B \otimes A)K_{q,n} \operatorname{vec} X$, and $K_{n,m}^{-1} = K_{n,m}^T = K_{m,n}$.

Partial derivatives and the Jacobian matrix

- ▶ Given a (smooth) vector function $f : S \rightarrow \mathbb{R}^m$ with $S \subset \mathbb{R}^n$, the Jacobian matrix of f at x is

$$\mathbf{D}f(x) = [d_{i,j}] = [\mathbf{D}_j f_i(x)] \in \mathbb{R}^{m \times n},$$

where $\mathbf{D}_j f_i(x)$ is the partial derivative $\frac{\partial f_i(x)}{\partial x_j}$

- ▶ The gradient $\nabla f(x) = (\mathbf{D}f(x))^T$ (but we will often use a *non-vectorized* version in matrix optimization)
- ▶ When f is **real**-valued, $\nabla f(x)$ has the **same** size as x ; here we also denote $\nabla f(x)$ by $\frac{\partial f}{\partial x}$ (derivative: $\frac{\partial f(x)}{\partial \mathbf{x}^T}$)

Differential of a vector function

- ▶ Again assume $f : S \rightarrow \mathbb{R}^m$ with $S \subset \mathbb{R}^n$ and let x be an interior point of S . f is differentiable at x if

$$f(x + dx) - f(x) = df(x; dx) + o(\|dx\|), \forall dx \in S \cap B(0; r)$$

with $df(x; dx) = A(x) dx$, $\lim_{dx \rightarrow 0} o(\|dx\|)/\|dx\| = 0$

- ▶ $A(x) : \text{derivative}$ of f at x
- ▶ If $D_j f_i(\cdot)$ exist in some neighborhood of c and are continuous at $x = c$, then f is differentiable at c and

$$A(c) = Df(c) =: \left. \frac{\partial f(x)}{\partial x^T} \right|_{x=c}$$

Matrix functions

- ▶ For a matrix function $F : S \rightarrow \mathbb{R}^{m \times p}$ with $S \subset \mathbb{R}^{n \times q}$, we define its Jacobian matrix at X as

$$DF(X) = \frac{\partial \text{vec} F(X)}{\partial (\text{vec}(X))^T}$$

- ▶ That is, if we introduce $f(\text{vec} X) = \text{vec} F(X)$,
 $DF(X) = Df(\text{vec} X)$

- ▶ Similarly, we apply vec to define differentiability:

$$F(X + dX) - F(X) = dF(X; dX) + o(dX), \text{ or} \\ \text{vec} F(X + dX) - \text{vec} F(X) = \text{vec} dF(X; dX) + o(dX),$$

where $\text{vec} dF(X; dX) = A(X) \text{vec}(dX)$

- ▶ $dF(X; dX) \in \mathbb{R}^{m \times p}$ is called the **differential** of F at X with increment dX
- ▶ $A(X) \in \mathbb{R}^{mp \times nq}$ is called the derivative of F at X

The chain rule and Cauchy invariance

- ▶ Consider a smooth (vector) function $h(t) = g \circ f(t) = g(f(t))$ at t and $x = f(t)$
- ▶ Then $Dh(t) = Dg(x) Df(t)$ (matrix product)
- ▶ It means that $dh(t; dt) = dg(x; df(t; dt))$ for any dt
(LHS = $Dg(x) Df(t) dt = Dg(x) df(t; dt) =$ RHS)
- ▶ We will use the shorthand notation dh from now on

A practical differentiation scheme

- ▶ Let $F(X)$ be a matrix function of X that is smooth
 1. Get the differential dF as a linear function of dX
 2. Vectorize dF as a linear transformation of $d \operatorname{vec} X$
- ▶ The transformation matrix gives all partial derivatives or the Jacobian matrix $DF(X)$

An illustration

- ▶ Let $F(X) = AXB$.
- ▶ Then $dF(X) = A(dX)B$
- ▶ Vectorization: $d \operatorname{vec} F(X) = (B^T \otimes A) d \operatorname{vec} X$
- ▶ Hence $DF(X) = B^T \otimes A$

Scalar functions

- ▶ In optimization, we typically deal with a real-valued objective function ϕ .
- ▶ Knowing its gradient often suffices (cf. the 1st step)
- ▶ Here, it is more convenient to use $\frac{\partial \phi(X)}{\partial X} := [\frac{\partial \phi(X)}{\partial X_{i,j}}]$
- ▶ We still call it the gradient $\nabla \phi(X)$ (abuse of notation!)
- ▶ [If you want to use a Newton-type algorithm, applying the vectorization step could be helpful]

Inner-product form of differential

- ▶ Concretely, when ϕ is a scalar function defined on a matrix X , we often try to express the differential as

$$d\phi(X) = \langle A, dX \rangle$$

- ▶ Then $\nabla\phi(X)(= \frac{\partial\phi(X)}{\partial X}) = A$, while $D\phi(X) = (\text{vec}A)^T$ (as $d\phi(X) = (\text{vec}A)^T \text{vec}X$) gives $\nabla\phi(X) = \text{vec}A$!
- ▶ Note the advantage of representing a problem in terms of a matrix in algorithm design and in implementation

Differential calculus

► Some properties:

- $d(U + V) = dU + dV$
- $d(UV) = (dU)V + U(dV)$
- $d(U^T) = (dU)^T$, $d \operatorname{vec} U = \operatorname{vec} dU$, $d \operatorname{tr}(U) = \operatorname{tr} dU$
- $d(U \otimes V) = (dU) \otimes V + U \otimes (dV)$

► A few useful facts:

- $dAx = A dx$, $\nabla(Ax) = A^T$
- $d|X| = \langle |X|(X^T)^{-1}, dX \rangle$ (use the adjoint matrix)
- $dX^{-1} = -X^{-1}(dX)X^{-1}$

- ▶ We show the last result assuming X^{-1} is differentiable
- ▶ First we show the conclusion by definition

$$\begin{aligned}
 & (X + dX)^{-1} - X^{-1} \\
 &= (UDV^T + dX)^{-1} - VD^{-1}U^T \\
 &= VD^{-1/2}\{(I + D^{-1/2}U^T(dX)VD^{-1/2})^{-1} - I\}D^{-1/2}U^T \\
 &= VD^{-1/2}\{I - D^{-1/2}U^T(dX)VD^{-1/2} - I\}D^{-1/2}U^T + o(dX) \\
 &= -X^{-1}(dX)X^{-1} + o(dX)
 \end{aligned}$$

- ▶ We can also use $XX^{-1} = I \Rightarrow d(XX^{-1}) = 0 \Rightarrow$
 $(dX)X^{-1} + X(dX^{-1}) = 0 \Rightarrow dX^{-1} = -X^{-1}(dX)X^{-1}$

Examples

- ▶ $\phi(x) = x^T A x$:

$$\begin{aligned} d\phi(x) &= d(x^T A x) = (dx)^T A x + x^T A dx \\ &= x^T (A^T + A) dx \Rightarrow D\phi(x) = x^T (\textcolor{red}{A} + \textcolor{red}{A}^T) \end{aligned}$$

- ▶ $\phi(x) = \|y - Ax\|_2^2/2$: Let $e = y - Ax$, $\phi(x) = e^T e/2$.
 $d\phi(x) = (y - Ax)^T d(y - Ax) = (Ax - y)^T A dx$.
Therefore, $\textcolor{blue}{\nabla}\phi(x) = x^T (Ax - y)$

- ▶ $\phi(x) = \exp(x^T x)$:

$$\begin{aligned} d\phi(x) &= \exp(x^T x) d(x^T x) = \exp(x^T x)((dx)^T x + x^T dx) \\ &= 2 \exp(x^T x) x^T dx \end{aligned}$$

- ▶ GLM: $\min_{\beta} l(\beta) = -\langle y, X\beta \rangle + \langle 1, b(X\beta) \rangle$

$$\begin{aligned} dl(\beta) &= -\langle y, X d\beta \rangle + \langle 1, b'(X\beta) \circ (X d\beta) \rangle \\ &= -\langle X^T y, d\beta \rangle + \langle b'(X\beta), X d\beta \rangle \\ &= \langle X^T (b'(X\beta) - y), d\beta \rangle \end{aligned}$$

► $\phi(X) = \text{tr}(X)$: $\nabla \phi(X) = I$

► $\phi(X) = \log |X|$:

$$\begin{aligned} d\phi(X) &= |X|^{-1} \langle |X| (X^T)^{-1}, dX \rangle = \langle (X^T)^{-1}, dX \rangle \\ &\Rightarrow \nabla(\log |X|) = (X^T)^{-1} \end{aligned}$$

► $F(x) = xx^T$ (a matrix function defined on a vector)

$$\begin{aligned} d(xx^T) &= (dx)x^T + x(dx)^T \text{ (combine? vec!)} \\ &= (x \otimes I) d \text{vec} x + (I \otimes x) d \text{vec}(x^T) \\ &= (x \otimes I + I \otimes x) d \text{vec} x \end{aligned}$$

(What if x is replaced by X ?)

- ▶ $F(X) = X$ ($X \in \mathbb{R}^{n \times q}$): $dF(X) = dX \Rightarrow$
 $d \operatorname{vec} F(X) = I_{nq} d \operatorname{vec} X \Rightarrow DF(X) = I_{nq}$
- ▶ $F(X) = X^T$ ($X \in \mathbb{R}^{n \times q}$): $dF(X) = dX^T \Rightarrow$
 $d \operatorname{vec} F(X) = d \operatorname{vec}(X^T) = K_{n,q} d \operatorname{vec} X \Rightarrow$
 $DF(X) = K_{n,q}$
- ▶ $F(X) = X^{-1}$: $dF(X) = -X^{-1}(dX)X^{-1} \Rightarrow$
 $d \operatorname{vec} F(X) = -(X^T)^{-1} \otimes X^{-1} d \operatorname{vec} X$

Multivariate meta-analysis (regression)

- ▶ Assume a random effects meta-regression model with multiple outcomes: $y_i = X_i\beta + \epsilon_i + \delta_i$, $i = 1, \dots, m$
- ▶ $y_i \in \mathbb{R}^n$: n outcomes from m **studies**. $X_i \in \mathbb{R}^{n \times p}$: design matrices containing study-level predictors
- ▶ $\epsilon_i \sim \mathcal{N}(0, \Sigma_i)$, $\delta_i \sim \mathcal{N}(0, \Sigma)$, $\epsilon_i \perp \delta_j$
- ▶ Σ_i : **within-study** covariances (known)
- ▶ Goal: Estimate β (fixed) as well as the **between-study** covariance Σ from the m studies

- ▶ MLE: $\min_{\beta, \Sigma \succeq 0} \sum_{i=1}^m \log \det(\Sigma_i + \Sigma) + \sum_{i=1}^m (y_i - X_i \beta)^T (\Sigma_i + \Sigma)^{-1} (y_i - X_i \beta)$ (Q: regularization?)
- ▶ We will assume β is **known** for simplicity (o/w: BCD)
- ▶ Let $R_i = (y_i - X_i \beta)(y_i - X_i \beta)^T$.

$$\min_{\Sigma \succeq 0} f(\Sigma) = \sum_{i=1}^m \log \det(\Sigma_i + \Sigma) + \sum_{i=1}^m \text{tr}((\Sigma_i + \Sigma)^{-1} R_i)$$

- ▶ A special case: $\Sigma_i = 0$ and $\Theta = \Sigma^{-1}$
- ▶ The problem is not convex and has a psd constraint

- ▶ We do not consider the **constrained** optimization for now. Let's use it to practice our differential skills.
- ▶ The gradient equation is

$$\sum (\Sigma + \Sigma_i)^{-1} - \sum (\Sigma + \Sigma_i)^{-1} R_i (\Sigma + \Sigma_i)^{-1} = 0, \text{ or}$$

$$\sum (\Sigma + \Sigma_i)^{-1} (\Sigma + \Sigma_i - R_i) (\Sigma + \Sigma_i)^{-1} = 0.$$

- ▶ Assuming Σ_i are not very different, $\sum \Sigma + \Sigma_i - R_i \approx 0$
- ▶ **GLS** (Berkey et al 96): $\Sigma = \sum (R_i - \Sigma_i)/m$ (given β)
- ▶ [But is the assumption reasonable in general? Does the estimate satisfy the psd constraint?]

- ▶ We could use Newton's method, which requires the gradient information of

$$g(\Sigma) = \sum (\Sigma + \Sigma_i)^{-1} - \sum (\Sigma + \Sigma_i)^{-1} R_i (\Sigma + \Sigma_i)^{-1}$$

- ▶ What is $H^{-1}g$ (no PSD constraint, no stepsize)? We did not introduce vech and Hessian..
- ▶ Using the differentiation technique, $dg = -\sum (\Sigma + \Sigma_i)^{-1} d\Sigma (\Sigma + \Sigma_i)^{-1} + \sum (\Sigma + \Sigma_i)^{-1} d\Sigma (\Sigma + \Sigma_i)^{-1} R_i (\Sigma + \Sigma_i)^{-1} + \sum (\Sigma + \Sigma_i)^{-1} R_i (\Sigma + \Sigma_i)^{-1} d\Sigma (\Sigma + \Sigma_i)^{-1}$

- ▶ So the second-order approximation of $f(\Sigma + \Delta)$ is

$$\begin{aligned} & \sum tr\{(\Sigma + \Sigma_i)^{-1}\Delta\} - tr\{(\Sigma + \Sigma_i)^{-1}R_i(\Sigma + \Sigma_i)^{-1}\Delta\} \\ & - \frac{1}{2}tr\{(\Sigma + \Sigma_i)^{-1}\Delta\Sigma(\Sigma + \Sigma_i)^{-1}\Delta\} \\ & + tr\{(\Sigma + \Sigma_i)^{-1}\Delta(\Sigma + \Sigma_i)^{-1}R_i(\Sigma + \Sigma_i)^{-1}\Delta\} \end{aligned}$$

- ▶ The solution for the update Δ satisfies

$$\begin{aligned} & \sum (\Sigma + \Sigma_i)^{-1}[\Delta - R_i(\Sigma + \Sigma_i)^{-1}\Delta - \Delta(\Sigma + \Sigma_i)^{-1}R_i](\Sigma + \Sigma_i)^{-1} \\ & = \sum (\Sigma + \Sigma_i)^{-1}[\Sigma + \Sigma_i - R_i](\Sigma + \Sigma_i)^{-1}, \end{aligned}$$

- ▶ Now we can vectorize it to solve a linear system

- ▶ When Σ_i are not very different, we have

$$\sum (\Delta - R_i(\Sigma + \Sigma_i)^{-1}\Delta - \Delta(\Sigma + \Sigma_i)^{-1}R_i) \approx \sum (\Sigma + \Sigma_i - R_i)$$

- ▶ It is an example of the **Lyapunov** equation

$$A\Delta + \Delta A^T + Q = 0$$

with $A = \sum (R_i(\Sigma + \Sigma_i)^{-1} - \frac{1}{2}I)$, $Q = \sum (\Sigma + \Sigma_i - R_i)$

- ▶ There are standard algorithms (often treating it as a special Sylvester equation). Matlab: $\Delta = \text{lyap}(A, Q)$.

Subgradients and subdifferential

- ▶ Another very useful concept is subdifferential
- ▶ Let f be a real-valued convex function. $g(x)$ is a subgradient of f at x if

$$f(y) \geq f(x) + \langle g(x), y - x \rangle, \forall y$$

- ▶ Subdifferential $\partial f(x)$: the set of all subgradients at x
- ▶ Why studying subgradients?
 - Nicely, if f is convex, $\partial f(x)$ is never null
 - If x is a minimizer of a convex f , then $0 \in \partial f(x)$
 - Finding just one subgradient is often easy

Two examples

- ▶ $f(x) = |x|$:

$$\partial f(0) = \begin{cases} 1, & x > 0 \\ [-1, 1], & x = 0 \\ -1, & x < 0 \end{cases}$$

- ▶ $f(X) = \|X\|_*$. Let $X = UDV^T$ (compact form)

$$\begin{aligned} \partial \|X\|_* &= \{UV^T + W : U^T W = 0, W V = 0, \|W\|_2 \leq 1\} \\ &= \{\textcolor{red}{U}V^T + U_\perp Z V_\perp^T : \|Z\|_2 \leq 1\} \end{aligned}$$

- ▶ Quite useful in computation and in theory

Soft-thresholding

- ▶ For example, let's consider a simple lasso problem $\min_{\beta} \frac{1}{2} \|y - \beta\|_F^2 + \lambda \|\beta\|_1$ ($\lambda \geq 0$, $X = I$)
- ▶ Let β^o be the optimal solution. Then it satisfies

$$\beta - y + \lambda \partial \|\beta\|_1 \ni 0, \text{ or}$$

$$\beta - y + \lambda s(\beta) = 0 \text{ for some } s(\beta) \in \partial \|\beta\|_1$$

- ▶ If $\beta_j^o > 0$, then $y_j - \lambda = \beta_j^o$ (and so $y_j > \lambda$); if $\beta_j^o < 0$, $y_j + \lambda = \beta_j^o$ (and $y_j < -\lambda$); if $\beta_j^o = 0$, we know $|y_j| \leq \lambda$
- ▶ Hence $\beta^o = \Theta(y; \lambda)$ where $\Theta(t; \lambda) = 1_{|t| > \lambda}(t - \text{sgn}(t)\lambda)$.

Singular-value (soft)thresholding

- ▶ Now consider the problem $\min_B \frac{1}{2}\|Y - B\|_F^2 + \lambda\|B\|_*$
- ▶ Let B_o be the optimal solution with $B_o = UD_oV$. Then $B_o - Y + \lambda UV^T + \lambda U_\perp ZV_\perp^T = 0$ for some $\|Z\|_2 \leq 1$, or

$$\begin{aligned} Y &= U(D_o + \lambda I)V^T + U_\perp \lambda ZV_\perp^T \\ &= [U \ U_\perp U_z] \begin{bmatrix} D_o + \lambda I & \\ & \lambda D_z \end{bmatrix} [V \ V_\perp V_z]^T \end{aligned}$$

- ▶ It implies that $B_o = \Theta^\sigma(Y; \lambda) = U_y \Theta(D_y; \lambda) V_y^T$ (well defined!), by (soft-)thresholding Y 's singular values
- ▶ These are examples of **proximity** operators

Subgradient update

- ▶ If $f(x^+) < f(x)$ (convex), then $0 \leq f(x^+) - f(x) - \langle g(x), x^+ - x \rangle < -\langle g(x), x^+ - x \rangle$ or $\langle g(x), x^+ - x \rangle < 0$, $\forall g(x) \in \partial f(x)$
- ▶ How about a subgradient update $x^+ = x - \alpha g(x)$?
 - $f(x^+) < f(x)$ may not hold; even $\lim_{\alpha \downarrow 0} \frac{1}{\alpha} [f(x^+(\alpha)) - f(x)] \leq 0$ may not be true (f may not be in $\mathcal{C}^{(1)}$)
- ▶ We often apply subgradient algorithms on the dual.
- ▶ Subgradient methods can be slow, but when $\partial f(x)$ is a **singleton**, we get faster convergence (why?)

- ▶ In handling the dual problem, we often need to know the first-order information of $d(\lambda) = \sup_{x \in A} L(x; \lambda)$
- ▶ Then $\partial f(\lambda)$ is the closure of

$$\text{conv} \cup \{\partial L(x; \cdot) : L(x; \lambda) = d(\lambda), x \in A\}$$

- ▶ When $\sup_{x \in A} L(x; \lambda)$ has a **unique** solution, d is differentiable at λ !
 - ADMM uses the fact (actually, a variant) to update the dual variable efficiently

Example

- ▶ We know $f(X) = \lambda_{\max}(X)$ ($X \in \mathbf{S}^n$) is convex
- ▶ What is its subdifferential?
 1. $f(X) = \sup_{\|\alpha\|_2^2=1} \alpha^T X \alpha$
 2. $d(\alpha^T X \alpha) = d\text{tr}(\alpha^T X \alpha) = \text{tr } d(\alpha \alpha^T X) = \langle \alpha \alpha^T, dX \rangle$
 3. $\partial f(X) = \mathbf{conv}\{\alpha \alpha^T : \alpha^T X \alpha = \|X\|_2, \|\alpha\|_2^2 = 1\}$
- ▶ So for **any** eigenvector α (with norm 1) associated with the largest eigenvalue of X , $\alpha \alpha^T$ gives a subgradient
- ▶ What if the geometric multiplicity of $\lambda_{\max}(X)$ is 1?
- ▶ What about $f(X) = \|X\|_2$ for any X ($X \neq 0$, $X = 0$)?