

Calculus

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2013

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Functions and Limits

Caution!

The functions f (or \mathbf{f}) discussed here are not required to be linear anymore.

Definition (Limit of a Function)

A function $\mathbf{f} : \mathcal{V} \rightarrow \mathbb{R}^m$, $\mathcal{V} \subseteq \mathbb{R}^n$, has a **limit** $\mathbf{f}^*(\mathbf{a})$ at the point $\mathbf{a} \in \mathcal{V}$ if given any $\varepsilon \in \mathbb{R}$, $\varepsilon > 0$, there exists $\delta \in \mathbb{R}$, $\delta > 0$ such that for all $\mathbf{x} \in \mathcal{V}$, $0 < \|\mathbf{x} - \mathbf{a}\| < \delta$, we have $\|\mathbf{f}(\mathbf{x}) - \mathbf{f}^*(\mathbf{a})\| < \varepsilon$. This is denoted by $\lim_{\mathbf{x} \rightarrow \mathbf{a}} \mathbf{f}(\mathbf{x}) = \mathbf{f}^*(\mathbf{a})$.

Definition (Continuity)

A function $\mathbf{f} : \mathcal{V} \rightarrow \mathbb{R}^m$, $\mathcal{V} \subseteq \mathbb{R}^n$, is **continuous** at \mathbf{a} iff $\mathbf{f}^*(\mathbf{a}) = \mathbf{f}(\mathbf{a})$; that is, given any $\varepsilon > 0$, there exists $\delta > 0$ such that for all $\mathbf{x} \in \mathcal{V}$, $0 < \|\mathbf{x} - \mathbf{a}\| < \delta$, we have $\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{a})\| < \varepsilon$.

Sequences and Convergence (1/2)

- A **sequence** of vectors $\{\mathbf{x}^{(k)}\}_k$ can be think of as the $\mathcal{R}(\mathbf{f})$ for some $\mathbf{f} : \mathbb{N} \rightarrow \mathbb{R}^n$
 - A sequence is **increasing** iff $\mathbf{x}^{(1)} < \mathbf{x}^{(2)} < \dots$, and **nondecreasing** iff $\mathbf{x}^{(1)} \leq \mathbf{x}^{(2)} \leq \dots$
 - Nondecreasing and nonincreasing sequences are called **monotone sequences**

Definition (Limit of a Sequence)

A sequence $\{\mathbf{x}^{(k)}\}_k$ has a **limit** \mathbf{x}^* if given any $\varepsilon \in \mathbb{R}$, $\varepsilon > 0$, there exists $K \in \mathbb{N}$, such that for all $k > K$, we have $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \varepsilon$. This is denoted by $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$.

- A sequence having a limit is called a **convergent** sequence
- Given a sequence $\{\mathbf{x}^{(k)}\}_k$ convergent to \mathbf{a} , we can see that a function $\mathbf{f} : \mathcal{V} \rightarrow \mathbb{R}^m$ is continuous at \mathbf{a} iff $\lim_{k \rightarrow \infty} \mathbf{f}(\mathbf{x}^{(k)}) = \mathbf{f}(\mathbf{a})$ [Proof: Using definitions and the fact that $\lim_{k \rightarrow \infty} \mathbf{f}(\mathbf{x}^{(k)}) = \mathbf{f}(\lim_{k \rightarrow \infty} \mathbf{x}^{(k)})$]

Sequences and Convergence (2/2)

- Given a sequence $\{\mathbf{x}^{(k)}\}_k$ and an increasing sequence of nature numbers $\{m_k\}_k$, we call $\{\mathbf{x}^{(m_k)}\}_k$ the **subsequence** of $\{\mathbf{x}^{(k)}\}_k$
 - A subsequence is obtained by neglecting some elements of a given sequence
- If a sequence converges to \mathbf{x}^* , then all its subsequences converge to \mathbf{x}^* too [Proof]

Extreme Value Theorem

Theorem

Let $f : \Omega \rightarrow \mathbb{R}$ be a continuous function over a compact set $\Omega \subseteq \mathbb{R}^n$. There exist $\mathbf{x}_0, \mathbf{x}_1 \in \Omega$ such that $f(\mathbf{x}_0) \leq f(\mathbf{x}) \leq f(\mathbf{x}_1), \forall \mathbf{x} \in \Omega$; that is, $f(\mathbf{x}_0) = \min_{\mathbf{x} \in \Omega} (f(\mathbf{x}))$ and $f(\mathbf{x}_1) = \max_{\mathbf{x} \in \Omega} f(\mathbf{x})$.

- We say f is **bounded** on Ω iff there exists $l, h \in \mathbb{R}$ such that $l \leq f(\mathbf{x}) \leq h, \forall \mathbf{x} \in \Omega$
 - The above theorem says that f is bounded on Ω if Ω is compact

Min, Max, Inf, and Sup

- Given a subset \mathcal{S} (e.g., $[0, 1)$ or $\{2, 4, 6, \dots\}$) of \mathbb{R} (or any other ordered set where elements can be compared with each other), we have:

Definition (Supremum)

An point $p \in \mathbb{R}$ is called the **supremum**, denoted by $\sup_{s \in \mathcal{S}} s$, iff a) $s \leq p, \forall s \in \mathcal{S}$; b) for any $\varepsilon > 0$, there exist $s \in \mathcal{S}$ such that $s > p - \varepsilon$.

- p is called the maximum iff $p \in \mathcal{S}$

Definition (Infimum)

An point $p \in \mathbb{R}$ is called the **infimum**, denoted by $\inf_{s \in \mathcal{S}} s$, iff a) $s \geq p, \forall s \in \mathcal{S}$; b) for any $\varepsilon > 0$, there exist $s \in \mathcal{S}$ such that $s < p + \varepsilon$.

- p is called the minimum iff $p \in \mathcal{S}$

Convergence of Functions (1/2)

- Given a set of data $\{\mathbf{x}^{(t)}\}_{t=1}^N$, suppose we use an ML algorithm to train a model, say $\mathbf{f}^{(N)}$
 - Usually, we want to know how the ML algorithm works when $N \rightarrow \infty$
 - We can think of $\{\mathbf{f}^{(N)}\}_N$ as a sequence, and then investigate the properties of its limit \mathbf{f}^*

Definition (Pointwise Convergence)

A sequence of functions $\{\mathbf{f}^{(N)}\}_N$, where $\mathbf{f}^{(N)} : \mathcal{V} \rightarrow \mathbb{R}^m$ and $\mathcal{V} \subseteq \mathbb{R}^n$, **converges pointwise** to a function $\mathbf{f}^* : \mathcal{V} \rightarrow \mathbb{R}^m$ iff for any $\mathbf{x} \in \mathcal{V}$, we have $\lim_{N \rightarrow \infty} \mathbf{f}^{(N)}(\mathbf{x}) = \mathbf{f}^*(\mathbf{x})$.

- Unfortunately, pointwise convergence is not strong enough to guarantee “reasonable” relations between $\mathbf{f}^{(N)}$ and \mathbf{f}^*
 - E.g., for all $x \in [0, 1]$, a sequence of continuous function $f^{(N)}(x) = x^N$ converges pointwise to $f^*(x) = \begin{cases} 0, & 0 \leq x < 1 \\ 1, & x = 1 \end{cases}$, which is obviously not continuous

Convergence of Functions (2/2)

Definition (Uniform Convergence)

A sequence of functions $\{\mathbf{f}^{(N)}\}_N$, where $\mathbf{f}^{(N)} : \mathcal{V} \rightarrow \mathbb{R}^m$ and $\mathcal{V} \subseteq \mathbb{R}^n$, **converges uniformly** to a function $\mathbf{f}^* : \mathcal{V} \rightarrow \mathbb{R}^m$ iff given any $\varepsilon > 0$, there exists $K \in \mathbb{N}$ such that for all $N \geq K$, we have $\|\mathbf{f}^{(N)} - \mathbf{f}^*\| < \varepsilon$ for all $\mathbf{x} \in \mathcal{V}$.

- Intuitively, $\mathbf{f}^{(N)}$ can be fitted into any given “ ε -tube” around \mathbf{f}^* as long as N is large enough

Theorem

If a sequence of continuous functions $\{\mathbf{f}^{(N)}\}_N$ converges uniformly to \mathbf{f}^ , then \mathbf{f}^* will be continuous.*

- Can be proved by either the “ $\varepsilon/3$ trick” or the “ ε -tube” intuition

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Derivative (1/2)

Definition (Derivative)

A function $f : [s, t] \rightarrow \mathbb{R}$, $s, t \in \mathbb{R}$, is **differentiable** at $a \in (s, t)$ iff $\lim_{\delta \rightarrow 0} \frac{f(a+\delta) - f(a)}{\delta}$ (or equivalently, $\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$) exists. The limit is called the **derivative** of f at a , and is denoted by $f'(a)$, $f^{(1)}(a)$, or $\frac{df}{dx}(a)$.

- “ d ” means the infinitesimal difference, and $f'(a)$ is the slope of a tangent line to f at $f(a)$
- If a function f is differentiable at a , then it is continuous at a (converse is not true, as evidenced by $f(x) = |x|$ and $a = 0$)
- f is said to be differentiable iff it is differentiable at any point of its domain
- f is said to be **continuously differentiable** iff f is differentiable and f' is continuous

Derivative (2/2)

- If f is differentiable, we can think of f' as a function too (although may not be continuous/differentiable)

- E.g., given $f(x) = e^x$, we have

$$\begin{aligned} f'(x) &= \lim_{\delta \rightarrow 0} \frac{e^{(x+\delta)} - e^x}{\delta} = \lim_{\delta \rightarrow 0} \frac{e^x(e^\delta - 1)}{\delta}. \text{ Let } t = e^\delta - 1, \text{ then} \\ f'(x) &= \lim_{t \rightarrow 0} \frac{e^x t}{\ln(1+t)} = e^x \lim_{t \rightarrow 0} \frac{1}{\ln(1+t)^{1/t}} = e^x \frac{1}{\ln(\lim_{t \rightarrow 0} (1+t)^{1/t})} = \\ &e^x \frac{1}{\ln e} = e^x \end{aligned}$$

- $f \in \mathcal{C}^n$ denotes that f is n -times continuously differentiable

Rolle's and Mean Value Theorem

Theorem (Rolle's Theorem)

Given a function $f : [s, t] \rightarrow \mathbb{R}$, where $s, t \in \mathbb{R}$, $f \in \mathcal{C}^1$, and $f(s) = f(t)$. There exists some $u \in (s, t)$ such that $f'(u) = 0$.

- Starting from $f(s)$, f must change its direction at some point when it is getting to $f(t)$
- We can “rotate” the above theorem to get a new one:

Theorem (Mean Value Theorem)

Given a function $f : [s, t] \rightarrow \mathbb{R}$, where $s, t \in \mathbb{R}$ and $f \in \mathcal{C}^1$. There exists some $u \in (s, t)$ such that $f'(u) = \frac{f(t) - f(s)}{t - s}$.

Partial Derivative for Multivariate Functions

Definition (Partial Derivative)

The **partial derivative** of a function $f: \mathcal{V} \rightarrow \mathbb{R}$, $\mathcal{V} \subseteq \mathbb{R}^n$, in the direction along the i th component at point \mathbf{a} is defined as

$$\lim_{\delta \rightarrow 0} \frac{f(a_1, \dots, a_i + \delta, \dots, a_n) - f(a_1, \dots, a_n)}{\delta}, \text{ denoted by } \frac{\partial f}{\partial x_i}(\mathbf{a}).$$

- We can only look at one component a time
- Given $f(x_1, x_2) = (x_1 + x_2)^2$, we have $\frac{\partial f}{\partial x_1}(x_1, x_2) =$
 $\lim_{\delta \rightarrow 0} \frac{(x_1 + \delta + x_2)^2 - (x_1 + x_2)^2}{\delta} = \lim_{\delta \rightarrow 0} \frac{\delta(2x_1 + 2x_2)}{\delta} = 2x_1 + 2x_2$
 - Simply treat x_2 as constant here

Integral

- Given a function $f: \mathcal{V} \rightarrow \mathbb{R}$, $\mathcal{V} \subseteq \mathbb{R}^n$, the set $\{[\mathbf{x}^\top, f(\mathbf{x})]^\top : \mathbf{x} \in \mathcal{V}\}$ is called the **graph** of f
- Now consider a function $f: [s, t] \rightarrow \mathbb{R}$, $s, t \in \mathbb{R}$, how do you approximate the area between the curve $y = f(x)$ and the x -axis in the graph of f ?
 - Partition $[s, t]$ evenly into n segments of width $\frac{t-s}{n}$, and let h_i (or l_i), $1 \leq i \leq n$, be the highest (or lowest) value of f in each segment
 - We can approximate the area by $H(n) = \sum_{i=1}^n h_i(\frac{t-s}{n})$ (or $L(n) = \sum_{i=1}^n l_i(\frac{t-s}{n})$)
- The larger the value of n , the more precise the approximation

Definition (Integral)

A function $f: [s, t] \rightarrow \mathbb{R}$, $s, t \in \mathbb{R}$, is **integrable** iff both $\lim_{n \rightarrow \infty} H(n)$ and $\lim_{n \rightarrow \infty} L(n)$ exist and are equal to each other. The limit is called the **integral** of f , denoted by $\int_s^t f(x) dx$.

Fundamental Theorem of Calculus

Theorem (Fundamental Theorem of Calculus)

Given a function $f : [s, t] \rightarrow \mathbb{R}$, $s, t \in \mathbb{R}$. We have:

a) The function $F : [s, t] \rightarrow \mathbb{R}$, $F(x) = \int_s^x f(z) dz$, is differentiable and

$$\frac{dF}{dx}(x) = \frac{d}{dx} \int_s^x f(z) dz = f(x);$$

b) If there exists a differentiable function $G : [s', t'] \rightarrow \mathbb{R}$, $[s, t] \subseteq [s', t']$, such that $\frac{dG}{dx}(x) = f(x)$ for every $x \in [s, t]$, then $\int_s^t f(x) dx = G(t) - G(s)$.

- F is called the **indefinite integral** (or **antiderivative**) of f and is a function of “accumulated area” from s
 - F can be “reversed” by differentiation and f is the “rate of accumulation”
- $\int_s^t f(x) dx$ is called the **definite integral** formally and represents an area
 - Definite integral can be computed by using indefinite integrals (which are usually easier to get)

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- **Derivative of Vector-Valued Functions**
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Isn't It Straightforward?

- Recall that the derivative of a real-valued function f at a is defined as $f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$
 - This is not applicable to $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$, as we cannot divide vectors
- We need a more general definition where the vectors can be fitted in with
 - Note that $f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$ iff $0 = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} - f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a) - f'(a)(x - a)}{x - a}$
 - Since the limit equals 0, the sign of numerator and denominator at the right hand side does not matter; that is, the above equation is equivalent to $0 = \lim_{x \rightarrow a} \frac{|f(x) - f(a) - f'(a)(x - a)|}{|x - a|} = \lim_{x \rightarrow a} \frac{|f(x) - (f'(a)(x - a) + f(a))|}{|x - a|}$
 - Now we can replace $|\cdot|$ with a vector norm $\|\cdot\|$
- But what does it mean?

Isn't It Straightforward?

- Recall that the derivative of a real-valued function f at a is defined as $f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$
 - This is not applicable to $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$, as we cannot divide vectors
- We need a more general definition where the vectors can be fitted in with
 - Note that $f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$ iff $0 = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a} - f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a) - f'(a)(x - a)}{x - a}$
 - Since the limit equals 0, the sign of numerator and denominator at the right hand side does not matter; that is, the above equation is equivalent to $0 = \lim_{x \rightarrow a} \frac{|f(x) - f(a) - f'(a)(x - a)|}{|x - a|} = \lim_{x \rightarrow a} \frac{|f(x) - (f'(a)(x - a) + f(a))|}{|x - a|}$
 - Now we can replace $|\cdot|$ with a vector norm $\|\cdot\|$
- But what does it mean? In the graph of f , $f'(a)(x - a) + f(a)$ is a tangent line to f at $f(a)$

Derivative of Vector-Valued Functions

- The notion of a tangent line can be generalized into an **affine function** $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathcal{A}(\mathbf{x}) = \mathcal{L}(\mathbf{x}) + \mathbf{c}$, where $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear transformation and $\mathbf{c} \in \mathbb{R}^m$
 - An affine function is a “point” in an affine space

Theorem (Derivative)

A function $\mathbf{f} : \mathcal{V} \rightarrow \mathbb{R}^m$, $\mathcal{V} \subseteq \mathbb{R}^n$, is **differentiable** at $\mathbf{a} \in \mathcal{V}$ iff there exists a linear transformation $\mathcal{L}(\mathbf{a}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that

$$\lim_{\delta \rightarrow 0} \frac{\|\mathbf{f}(\mathbf{a} + \delta) - (\mathcal{L}(\mathbf{a})(\delta) + \mathbf{f}(\mathbf{a}))\|}{\|\delta\|} = 0 \text{ (or equivalently,}$$

$$\lim_{\mathbf{x} \rightarrow \mathbf{a}} \frac{\|\mathbf{f}(\mathbf{x}) - (\mathcal{L}(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \mathbf{f}(\mathbf{a}))\|}{\|\mathbf{x} - \mathbf{a}\|} = 0). \mathcal{L}(\mathbf{a}) \text{ is called the } \mathbf{derivative} \text{ of } \mathbf{f} \text{ at } \mathbf{a}, \text{ denoted by } \mathbf{f}'(\mathbf{a}).$$

- Since $\mathbf{f}'(\mathbf{a})$ is linear, it can be represented by a matrix \mathbf{J}_a
 - How does \mathbf{J}_a look like?

Jacobian Matrices (1/2)

- Any function $\mathbf{f} : \mathcal{V} \rightarrow \mathbb{R}^m$, $\mathcal{V} \subseteq \mathbb{R}^n$ and $\mathbf{f}(\mathbf{v}) = \mathbf{w}$ can be rewritten as:

$$\begin{pmatrix} f_1(v_1, \dots, v_n) = w_1 \\ \vdots \\ f_m(v_1, \dots, v_n) = w_m \end{pmatrix}$$

- If \mathbf{f} is linear, each real-valued f_i , $1 \leq i \leq m$, can be represented by a row vector $\mathbf{j}_i \in \mathbb{R}^n$ such that $\mathbf{j}_i^\top \mathbf{v} = w_i$ (Remember the system of linear equations?)

- Let $\mathbf{J}_a = \begin{bmatrix} \mathbf{j}_1^\top \\ \vdots \\ \mathbf{j}_m^\top \end{bmatrix}$ be the matrix representation of $\mathbf{f}'(\mathbf{a})$

Jacobian Matrices (2/2)

- Let $\delta = \delta \mathbf{e}_j$, where $1 \leq j \leq n$ and $\delta \in \mathbb{R}$
- Looking at the definition $\lim_{\delta \rightarrow 0} \frac{\|f(\mathbf{a} + \delta) - (\mathcal{L}(\mathbf{a})(\delta) + f(\mathbf{a}))\|}{\|\delta\|} = 0$
row-by-row, we have $\lim_{\delta \rightarrow 0} \frac{f_i(\mathbf{a} + \delta \mathbf{e}_j) - (\delta \mathbf{j}_i^\top \mathbf{e}_j + f_i(\mathbf{a}))}{\delta} = 0$ for each i and j
 - This implies that $\lim_{\delta \rightarrow 0} \frac{f_i(\mathbf{a} + \delta \mathbf{e}_j) - f_i(\mathbf{a})}{\delta} = \mathbf{j}_i^\top \mathbf{e}_j$
 - The right hand side denotes the element of \mathbf{J}_a at the i th row and the j th column
 - The left hand side is the partial derivative $\frac{\partial f_i}{\partial x_j}(\mathbf{a})$ by definition

- $\mathbf{J}_a = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{a}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{a}) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{a}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{a}) \end{bmatrix}$ is called the **Jacobian matrix** (or **derivative matrix**) of \mathbf{f} at \mathbf{a}

Gradient (1/2)

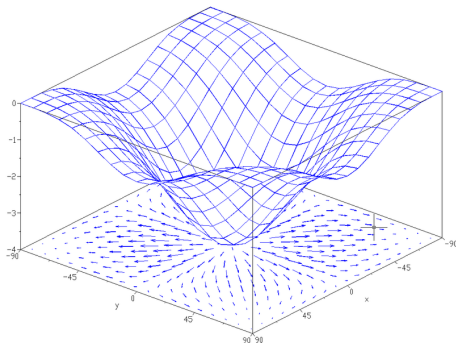
Definition

If a function $f: \mathcal{V} \rightarrow \mathbb{R}$, $\mathcal{V} \subseteq \mathbb{R}^n$, is differentiable, then the *gradient* of f is defined by $\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right]^\top = f'(\mathbf{x})^\top$.

- The Jacobian matrix of \mathbf{f} at \mathbf{a} can be rewritten as $\mathbf{J}_\mathbf{a} = \begin{bmatrix} \nabla f_1(\mathbf{a})^\top \\ \vdots \\ \nabla f_m(\mathbf{a})^\top \end{bmatrix}$
- From the previous page, we can see that in the graph of f_i , $\{\mathbf{x} : \nabla f_i(\mathbf{a})^\top (\mathbf{x} - \mathbf{a}) + f_i(\mathbf{a})\}$ is a tangent hyperplane to f_i at $f_i(\mathbf{a})$
 - The norm of $\nabla f_i(\mathbf{a})$ is the “slope” of the tangent hyperplane
 - $\nabla f_i(\mathbf{a})$ also acts as the direction that for a given small displacement from \mathbf{a} , f_i increase more in the direction of $\nabla f_i(\mathbf{a})$ than in any other direction [Proof]

Gradient (2/2)

- The gradient $\nabla f(\mathbf{x})$ is a function from \mathbb{R}^n to \mathbb{R}^n and can be pictured as a **vector field**
 - Each vector in the field is the direction of maximum rate of increase of f
 - E.g., consider $f(x_1, x_2) = -(\cos^2 x_1 + \cos^2 x_2)^2$:



Hessian Matrices (1/2)

Definition (Hessian Matrix)

Given a differentiable function $f: \mathcal{V} \rightarrow \mathbb{R}$, $\mathcal{V} \subseteq \mathbb{R}^n$. If $\nabla f(\mathbf{x})$ is differentiable, we have the derivative of $\nabla f(\mathbf{x})$ as:

$$H(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix},$$

which is called the **Hessian matrix** of f at \mathbf{x} .

- $\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x})$ means taking the partial derivative of f in the direction x_j first, and then x_i

Hessian Matrices (2/2)

Theorem (Clairaut's/Schwarz's Theorem)

If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuously differentiable at x , then its Hessian matrix at x is symmetric.

- If the second partial derivatives of f is not continuous, then there is no such a guarantee
- Here is an example:

$$f(x_1, x_2) = \begin{cases} \frac{x_1 x_2 (x_1^2 - x_2^2)}{(x_1^2 + x_2^2)}, & [x_1, x_2]^T \neq [0, 0]^T \\ 0, & [x_1, x_2]^T = [0, 0]^T. \end{cases}$$

The Hessian matrix of f at the point $[0, 0]^T$ is not symmetric
[Homework]

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Differentiation Rules

Theorem (Chain Rule)

Let $\mathbf{f} : (s, t) \rightarrow \mathcal{D}$ and $g : \mathcal{D} \rightarrow \mathbb{R}$ be differentiable functions, where $\mathcal{D} \subseteq \mathbb{R}^n$ is an open set and $s, t \in \mathbb{R}$. The composite function $g \circ \mathbf{f} : (s, t) \rightarrow \mathbb{R}$ is

differentiable and $(g \circ \mathbf{f})'(x) = g'(\mathbf{f}(x))\mathbf{f}'(x) = \nabla g(\mathbf{f}(x))^\top \begin{bmatrix} f_1'(x) \\ \vdots \\ f_n'(x) \end{bmatrix}$.

- g' and \mathbf{f}' are derivatives but with respect to different variables

Theorem (Product Rule)

Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be differentiable functions. Then the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $h(\mathbf{x}) = \mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$, is differentiable and $h'(\mathbf{x}) = \mathbf{f}(\mathbf{x})^\top \mathbf{g}'(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top \mathbf{f}'(\mathbf{x})^\top$.

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Definition (Level Set)

The **level set** of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at level $c \in \mathbb{R}$ is the set of points $S = \{\mathbf{x} : f(\mathbf{x}) = c\}$.

- When $n = 2$, S is a plane curve
- When $n = 3$, S is a surface

Level Sets and Gradients

Theorem

The vector $\nabla f(\mathbf{a})$ is orthogonal to the tangent vector to an arbitrary curve passing through \mathbf{a} on a level set at level $f(\mathbf{a})$.

Proof.

Let S be the level set at level $f(\mathbf{a})$ and $\gamma: \mathbb{R} \rightarrow \mathbb{R}^n$ be a curve lying on S passing through \mathbf{a} ; that is, there exists c such that $\gamma(c) = \mathbf{a}$. Suppose $\gamma'(c) = \mathbf{v} \neq \mathbf{0}$ so \mathbf{v} is a tangent vector to γ at \mathbf{a} . We have $(f \circ \gamma)'(c) = f'(\gamma(c))\gamma'(c) = f'(\mathbf{a})\mathbf{v}$. Since γ lies on S , we have $f \circ \gamma(t) = f(\mathbf{a})$ for all $t \in \mathbb{R}$. The function $f \circ \gamma$ is a constant, implying that $(f \circ \gamma)'(c) = 0$. So $f'(\mathbf{a})\mathbf{v} = \nabla f(\mathbf{a})^\top \mathbf{v} = 0$. □

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Taylor's Theorem (1/2)

Theorem

Given a function $f : [a, b] \rightarrow \mathbb{R}$ and $f \in \mathcal{C}^m$. We have

$$f(b) = f(a) + \frac{f^{(1)}(a)}{1!}(b-a) + \frac{f^{(2)}(a)}{2!}(b-a)^2 + \cdots + \frac{f^{(m-1)}(a)}{(m-1)!}(b-a)^{m-1} + R_m,$$

where $R_m = \frac{f^{(m)}(c)}{m!}(b-a)^m$ for some $c \in (a, b)$.

Proof.

Define $R(x) =$

$$f(b) - f(x) - \frac{f^{(1)}(x)}{1!}(b-x) - \frac{f^{(2)}(x)}{2!}(b-x)^2 - \cdots - \frac{f^{(m-1)}(x)}{(m-1)!}(b-x)^{m-1}.$$

We show that $R(a) = \frac{f^{(m)}(c)}{m!}(b-a)^m$ for some $c \in (a, b)$. Note that

$$R^{(1)}(x) =$$

$$-f^{(1)}(x) + \left[f^{(1)}(x) - \frac{f^{(2)}(x)}{1!}(b-x) \right] + \left[\frac{f^{(2)}(x)}{1!}(b-x) - \frac{f^{(3)}(x)}{2!}(b-x)^2 \right] + \cdots + \left[\frac{f^{(m-1)}(x)}{(m-2)!}(b-x)^{m-2} - \frac{f^{(m)}(x)}{(m-1)!}(b-x)^{m-1} \right] = -\frac{f^{(m)}(x)}{(m-1)!}(b-x)^{m-1}.$$

Define $g(x) = R(x) - \left(\frac{b-x}{b-a}\right)^m R(a)$. It's easy to check that $g(a) = g(b) = 0$.

By Rolle's theorem there exists some $c \in (a, b)$ such that $0 = g^{(1)}(c) =$

$$R^{(1)}(c) + \frac{m(b-c)^{m-1}}{(b-a)^m} R(a) = -\frac{f^{(m)}(c)}{(m-1)!}(b-c)^{m-1} + \frac{m(b-c)^{m-1}}{(b-a)^m} R(a),$$

implying $R(a) = \frac{f^{(m)}(c)}{m!}(b-a)^m$. □

Taylor's Theorem (2/2)

- An well-known application is Taylor series:
 - $e^x = 1 + x + \frac{x^2}{2!} + \dots = \sum_{n=0}^{\infty} \frac{x^n}{n!}$ for all x (expanding e^x at $a=0$)
 - $\ln(1+x) = \sum_{n=1}^{\infty} (-1)^{n+1} \frac{x^n}{n}$ for $|x| < 1$, which implies $\ln(1+x) \approx x$ for $|x| \ll 1$
- A function $f: \mathcal{D} \rightarrow \mathbb{R}$, where \mathcal{D} is an open interval, is said to be **analytic** iff for any $x, a \in \mathcal{D}$, f can be written as $f(x) = \sum_{n=0}^{\infty} c_n (x-a)^n$ for some $c_n \in \mathbb{R}$
 - An analytic f is easy to analyze, e.g., $c_n = \frac{f^{(n)}(a)}{n!}$
 - f is analytic iff give any x , the Taylor series at a converges to $f(x)$

Order of Convergence

- Consider $\mathbf{f}(\mathbf{x}) : \mathcal{V} \rightarrow \mathbb{R}^m$ and $\mathbf{g}(\mathbf{x}) : \mathcal{V} \rightarrow \mathbb{R}^m$, where $\mathcal{V} \subseteq \mathbb{R}^n$ includes $\mathbf{0}$
- We denote $\mathbf{f}(\mathbf{x}) = o(\mathbf{g}(\mathbf{x}))$ iff $\mathbf{f}(\mathbf{x})$ goes to $\mathbf{0}$ faster than $\mathbf{g}(\mathbf{x})$ does
 - Specifically, $\lim_{\mathbf{x} \rightarrow \mathbf{0}} \frac{\|\mathbf{f}(\mathbf{x})\|}{\|\mathbf{g}(\mathbf{x})\|} = 0$
- We denote $\mathbf{f}(\mathbf{x}) = O(\mathbf{g}(\mathbf{x}))$ iff $\mathbf{f}(\mathbf{x})$ goes to $\mathbf{0}$ faster than or equal to $\mathbf{g}(\mathbf{x})$ does
 - Specifically, for a sufficiently small $\delta \in \mathbb{R}$, there exists $c \in \mathbb{R}$ such that if $\|\mathbf{x}\| < \delta$, then $\frac{\|\mathbf{f}(\mathbf{x})\|}{\|\mathbf{g}(\mathbf{x})\|} \leq c$
- E.g., $x^2 = o(x)$, $x = O(x)$, $[x^3, 2x^2]^\top = o([x, 0]^\top)$
- Don't mix this up with the order of growth

Taylor Theorem for Multivariate Functions (1/2)

- Recall that a function $f : [a, b] \rightarrow \mathbb{R}$, $f \in \mathcal{C}^m$, can be written as
$$f(b) = f(a) + \frac{f^{(1)}(a)}{1!}(b-a) + \frac{f^{(2)}(a)}{2!}(b-a)^2 + \cdots + \frac{f^{(m-1)}(a)}{(m-1)!}(b-a)^{m-1} + \frac{f^{(m)}(a+\delta(b-a))}{m!}(b-a)^m$$
, where $\delta \in (0, 1)$ is a constant
- By the continuity of $f^{(m)}$, we have $\lim_{(b-a) \rightarrow 0} f^{(m)}(a + \delta(b-a)) = f^{(m)}(a) \Rightarrow \lim_{(b-a) \rightarrow 0} \frac{f^{(m)}(a+\delta(b-a)) - f^{(m)}(a)}{1} = 0$; that is, $f^{(m)}(a + \delta(b-a)) - f^{(m)}(a) = o(1) \Rightarrow f^{(m)}(a + \delta(b-a)) = f^{(m)}(a) + o(1)$
- We can rewrite f as $f(b) = f(a) + \frac{f^{(1)}(a)}{1!}(b-a) + \frac{f^{(2)}(a)}{2!}(b-a)^2 + \cdots + \frac{f^{(m)}(a)}{m!}(b-a)^m + o((b-a)^m)$, since $o(1) \frac{(b-a)^m}{m!} = o((b-a)^m)$
- If $f \in \mathcal{C}^{m+1}$, we can further rewrite f as $f(b) = f(a) + \frac{f^{(1)}(a)}{1!}(b-a) + \frac{f^{(2)}(a)}{2!}(b-a)^2 + \cdots + \frac{f^{(m)}(a)}{m!}(b-a)^m + O((b-a)^{m+1})$
 - $R_{m+1} = \frac{f^{(m+1)}(c)}{(m+1)!}(b-a)^{m+1} = O((b-a)^{m+1})$, as $f^{(m+1)}$ is bound on the compact set $[a, b]$ and therefore can be regarded as a constant

Taylor Theorem for Multivariate Functions (2/2)

Theorem

Given $f: \mathcal{V} \rightarrow \mathbb{R}$, where $\mathcal{V} \subseteq \mathbb{R}^n$ is an open set and $f \in \mathcal{C}^2$. For any $\mathbf{x}, \mathbf{a} \in \mathcal{V}$, there exists $\mathbf{c} = \mathbf{a} + c(\mathbf{x} - \mathbf{a})/\|\mathbf{x} - \mathbf{a}\|$ for some $c \in (0, \|\mathbf{x} - \mathbf{a}\|)$ such that $f(\mathbf{x}) = f(\mathbf{a}) + \frac{1}{1!} \nabla f(\mathbf{a})^\top (\mathbf{x} - \mathbf{a}) + R_2$, where $R_2 = \frac{1}{2!} (\mathbf{x} - \mathbf{a})^\top \mathbf{H}(\mathbf{c})(\mathbf{x} - \mathbf{a})$.

Proof.

Define $z: \mathbb{R} \rightarrow \mathbb{R}^n$ by $z(\delta) = \mathbf{a} + \delta(\mathbf{x} - \mathbf{a})/\|\mathbf{x} - \mathbf{a}\|$ and $\phi: \mathbb{R} \rightarrow \mathbb{R}$ by $\phi(\delta) = f \circ z(\delta) = f(\mathbf{a} + \delta(\mathbf{x} - \mathbf{a})/\|\mathbf{x} - \mathbf{a}\|)$, we can see that $f(\mathbf{x}) = \phi(\|\mathbf{x} - \mathbf{a}\|)$ and by Taylor's theorem,

$f(\mathbf{x}) = \phi(0) + \frac{\phi^{(1)}(0)}{1!} \|\mathbf{x} - \mathbf{a}\| + \frac{\phi^{(2)}(c)}{2!} \|\mathbf{x} - \mathbf{a}\|^2$. Note

$\phi^{(1)}(\delta) = f^{(1)}(z(\delta))z^{(1)}(\delta) = \nabla f(z(\delta))^\top \left(\frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|} \right) = \left(\frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|} \right)^\top \nabla f(z(\delta))$

and $\phi^{(2)}(\delta) = \left(\frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|} \right)^\top \mathbf{H}(z(\delta))z^{(1)}(\delta) = \left(\frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|} \right)^\top \mathbf{H}(z(\delta)) \left(\frac{\mathbf{x} - \mathbf{a}}{\|\mathbf{x} - \mathbf{a}\|} \right)$.

Substituting $\phi^{(1)}(0)$ and $\phi^{(2)}(c)$ in $f(\mathbf{x})$ we have the proof. \square

- We can also write $R_2 = \frac{1}{2!} (\mathbf{x} - \mathbf{a})^\top \mathbf{H}(\mathbf{a})(\mathbf{x} - \mathbf{a}) + o(\|\mathbf{x} - \mathbf{a}\|^2)$ [Proof]
- Or $R_2 = \frac{1}{2!} (\mathbf{x} - \mathbf{a})^\top \mathbf{H}(\mathbf{a})(\mathbf{x} - \mathbf{a}) + O(\|\mathbf{x} - \mathbf{a}\|^3)$ if $f \in \mathcal{C}^3$ [Proof]

Mean Value Theorem Revisited

Theorem (Mean Value Theorem)

Given a function $f : \mathcal{V} \rightarrow \mathbb{R}^m$, where $\mathcal{V} \subseteq \mathbb{R}^n$ is open and $f \in \mathcal{C}^1$. For any $\mathbf{b}, \mathbf{a} \in \mathcal{V}$, there exists $\mathbf{M} = \begin{bmatrix} \nabla f_1(\mathbf{c}^{(1)})^\top \\ \vdots \\ \nabla f_m(\mathbf{c}^{(m)})^\top \end{bmatrix}$ for some $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(m)} \in \mathcal{V}$ such that $f(\mathbf{b}) - f(\mathbf{a}) = \mathbf{M}(\mathbf{b} - \mathbf{a})$.

- Can be easily proved by the Taylor's Theorem [Proof]

Outline

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Vector Derivatives

- Recall that the derivative of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ at a point $\mathbf{x} \in \mathbb{R}^n$ can be written as a Jacobian matrix
$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$
- Given $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in \mathbb{R}^m$, define $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{m \times n}$ such that $\left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)_{i,j} = \frac{\partial y_i}{\partial x_j}$
 - We can express the above Jacobian matrix succinctly as $\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}$
 - $\frac{\partial}{\partial \mathbf{x}}(\mathbf{A}\mathbf{x}) = \mathbf{A}$ and $\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \mathbf{I}$
- $\frac{\partial \mathbf{y}}{\partial x} \in \mathbb{R}^{m \times 1}$ for $x \in \mathbb{R}$; and $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times n}$ for $y \in \mathbb{R}$
 - $\frac{\partial}{\partial \mathbf{x}}(\mathbf{a}^\top \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{a}) = \mathbf{a}^\top$ for any $\mathbf{a} \in \mathbb{R}^n$
 - $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}^\top$
- Differentiation rules are applicable
 - $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^\top \mathbf{A}\mathbf{x}) = \mathbf{x}^\top \frac{\partial}{\partial \mathbf{x}}(\mathbf{A}\mathbf{x}) + (\mathbf{A}\mathbf{x})^\top \left(\frac{\partial \mathbf{x}}{\partial \mathbf{x}}\right)^\top = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$ for any $\mathbf{A} \in \mathbb{R}^{n \times n}$ [Proof]

Matrix Derivatives

- Given $x \in \mathbb{R}$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$, define
 - $\frac{\partial \mathbf{A}}{\partial x} \in \mathbb{R}^{m \times n}$ such that $\left(\frac{\partial \mathbf{A}}{\partial x}\right)_{i,j} = \frac{\partial a_{i,j}}{\partial x}$
 - $\frac{\partial x}{\partial \mathbf{A}} \in \mathbb{R}^{m \times n}$ such that $\left(\frac{\partial x}{\partial \mathbf{A}}\right)_{i,j} = \frac{\partial x}{\partial a_{i,j}}$
- x should be related to \mathbf{A} (e.g., $a_{i,j}$, $tr(\mathbf{A})$, or $det(\mathbf{A})$, etc.)
 - $\frac{\partial \mathbf{A}}{\partial a_{i,j}}$ is a matrix whose element at the i th row and j th column equals 1, and others 0
- Although looked similar to vector derivatives, matrix derivatives have no obvious geometric implications and are used mainly to simplify the calculation of partial derivatives
- $\frac{\partial}{\partial x}(\mathbf{A}\mathbf{B}) = \mathbf{A} \frac{\partial \mathbf{B}}{\partial x} + \frac{\partial \mathbf{A}}{\partial x} \mathbf{B}$ [Proof]
 - $\frac{\partial}{\partial x}(\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}$ [Proof: $\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$ and apply the above]

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Derivatives of Traces

- $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{AB}) = \frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{BA}) = \mathbf{B}^\top$, as
$$\frac{\partial}{\partial a_{ij}} \text{tr}(\mathbf{AB}) = \frac{\partial}{\partial a_{ij}} \sum_{r=1}^n \sum_{s=1}^n a_{r,s} b_{s,r} = b_{j,i}$$
 - $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}) = \frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{AI}) = \mathbf{I}$
- $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}^\top \mathbf{B}) = \mathbf{B}$ [Proof]
- $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{ABA}^\top) = \mathbf{A}(\mathbf{B} + \mathbf{B}^\top)$ [Proof]

Derivatives of Determinants (1/2)

Theorem

Given an invertible matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}$, we have

$$\frac{\partial}{\partial x} \ln(\det(\mathbf{A})) = \text{tr}(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}).$$

Proof.

We only proof the case where $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$ is symmetric here. We have

$\frac{\partial}{\partial x} \ln(\det(\mathbf{A})) = \frac{\partial}{\partial x} \ln(\det(\mathbf{U})\det(\mathbf{D})\det(\mathbf{U})^{-1}) = \frac{\partial}{\partial x} \ln(\det(\mathbf{D})) =$
 $\frac{\partial}{\partial x} \ln(\prod_{i=1}^n \lambda_i) = \sum_{i=1}^n \frac{\partial}{\partial x} \ln \lambda_i = \sum_{i=1}^n \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} = \text{tr}(\mathbf{D}^{-1} \frac{\partial \mathbf{D}}{\partial x})$. Note \mathbf{U} is orthogonal, and diagonalizable, so there exists an antisymmetric matrix $\mathbf{W} = \frac{1}{x} \ln \mathbf{U}$ such that $\mathbf{U} = e^{\mathbf{W}x}$. By the chain rule we have

$$\frac{\partial \mathbf{U}}{\partial x} = e^{\mathbf{W}x} \left(\frac{\partial}{\partial x} \mathbf{W}x \right) = \mathbf{U}\mathbf{W} \text{ and } \frac{\partial \mathbf{U}^\top}{\partial x} = \frac{\partial}{\partial x} e^{\mathbf{W}^\top x} = \frac{\partial}{\partial x} e^{-\mathbf{W}x} = -\mathbf{U}^\top \mathbf{W}.$$

Therefore, $\text{tr}(\mathbf{D}^{-1} \frac{\partial \mathbf{D}}{\partial x}) = \text{tr}((\mathbf{U}^\top \mathbf{A} \mathbf{U})^{-1} \frac{\partial}{\partial x} (\mathbf{U}^\top \mathbf{A} \mathbf{U})) =$

$$\text{tr}((\mathbf{U}^\top \mathbf{A} \mathbf{U})^{-1} (\frac{\partial \mathbf{U}^\top}{\partial x} \mathbf{A} \mathbf{U} + \mathbf{U}^\top \frac{\partial \mathbf{A}}{\partial x} \mathbf{U} + \mathbf{U}^\top \mathbf{A} \frac{\partial \mathbf{U}}{\partial x})) =$$

$$\text{tr}((\mathbf{U}^\top \mathbf{A} \mathbf{U})^{-1} (-\mathbf{U}^\top \mathbf{W} \mathbf{A} \mathbf{U} + \mathbf{U}^\top \frac{\partial \mathbf{A}}{\partial x} \mathbf{U} + \mathbf{U}^\top \mathbf{A} \mathbf{U} \mathbf{W})) =$$

$$\text{tr}(-\mathbf{U}^\top \mathbf{A}^{-1} \mathbf{W} \mathbf{A} \mathbf{U}) + \text{tr}(\mathbf{U}^\top \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{U}) + \text{tr}(\mathbf{W}), \text{ which can be simplified}$$

to $\text{tr}(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x})$ by the cyclic property of trace. □

Derivatives of Determinants (2/2)

- $\frac{\partial}{\partial \mathbf{A}} \ln(\det(\mathbf{A})) = (\mathbf{A}^{-1})^\top$
 - Let $a_{i,j}$ and $b_{i,j}$ be the elements of \mathbf{A} and \mathbf{A}^{-1} respectively, then $\frac{\partial}{\partial a_{i,j}} \ln(\det(\mathbf{A})) = \text{tr}(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial a_{i,j}}) = \sum_{r=1}^n \sum_{s=1}^n b_{r,s} \frac{\partial a_{s,r}}{\partial a_{i,j}} = b_{j,i}$

Outline

1 Calculus, The Basics

- Sequences and Limits
- Derivative and Integral of Real-Valued Functions
- Derivative of Vector-Valued Functions
- Differentiation Rules
- Level Sets and Gradients
- Taylor's Theorem

2 Matrix Calculus

- Vector and Matrix Derivatives
- Derivatives of Traces and Determinants**

3 Calculus of Variations**

- Functionals

Functionals

- Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = ax + b$ (or $f(x|a, b) = ax + b$)
 - x is an **argument** and a and b are parameters
- Let \mathcal{S} be the set of functions $f : \mathcal{V} \rightarrow \mathcal{W}$, we can define a **functional** $F : \mathcal{S} \rightarrow \mathcal{W}$, $F[f]$, with f as the argument
- E.g., value of a function $f : \mathbb{R} \rightarrow \mathbb{R}$ at x : $F : \mathcal{S} \rightarrow \mathbb{R}$, $F[f] = f(x)$
 - x is a **parameter**
 - We can write $F[f]$ as $F[f|x]$
- E.g., definite integral of a function $f : \mathbb{R} \rightarrow \mathbb{R}$: $I : \mathcal{S} \rightarrow \mathbb{R}$,
 $I[f] = \int_a^b f(x) dx$
 - a and b are parameters
- E.g., expectation of $f : \mathbb{R} \rightarrow \mathbb{R}$ defined over the values of a random variable X : $E : \mathcal{S} \rightarrow \mathbb{R}$, $E[f(X)] = \int f(x) p_X(x) dx$