Optimization in Machine Learning

Mathematical Concepts Differentiation and Derivatives





Learning goals

- Definition of smoothness
- Uni- & multivariate differentiation
- Gradient, partial derivatives
- Jacobian matrix
- Hessian matrix
- Lipschitz continuity

UNIVARIATE DIFFERENTIABILITY

Definition: A function $f : S \subseteq \mathbb{R} \to \mathbb{R}$ is said to be **differentiable** for each inner point $x \in S$ if the following limit exists:

$$f'(x) := \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

Intuitively: *f* can be approved locally by a lin. fun. with slope m = f'(x).



Left: Function is differentiable everywhere. Right: Not differentiable at the red point.

× 0 0 × ×

SMOOTH VS. NON-SMOOTH

- **Smoothness** of a function $f : S \to \mathbb{R}$ is measured by the number of its continuous derivatives
- C^k is class of k-times continuously differentiable functions (f ∈ C^k means f^(k) exists and is continuous)
- In this lecture, we call *f* "smooth", if at least $f \in C^1$



 f_1 is smooth, f_2 is continuous but not differentiable, and f_3 is non-continuous.

Optimization in Machine Learning - 2 / 13

 $\times \times$

MULTIVARIATE DIFFERENTIABILITY

Definition: $f : S \subseteq \mathbb{R}^d \to \mathbb{R}$ is **differentiable** in $\mathbf{x} \in S$ if there exists a (continuous) linear map $\nabla f(\mathbf{x}) : S \subseteq \mathbb{R}^d \to \mathbb{R}^d$ with

$$\lim_{\mathbf{h}\to 0} \frac{f(\mathbf{x}+\mathbf{h}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T \cdot \mathbf{h}}{||\mathbf{h}||} = 0$$

× × ×

Geometrically: The function can be locally approximated by a tangent hyperplane. Source: https://github.com/jermwatt/machine_learning_refined.

GRADIENT

• Linear approximation is given by the gradient:

$$\nabla f = \frac{\partial f}{\partial x_1} \boldsymbol{e}_1 + \dots + \frac{\partial f}{\partial x_d} \boldsymbol{e}_d = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_d}\right)^T$$

- Elements of the gradient are called partial derivatives.
- To compute $\partial f / \partial x_i$, regard *f* as function of x_i only (others fixed)

Example:
$$f(\mathbf{x}) = x_1^2/2 + x_1x_2 + x_2^2 \Rightarrow \nabla f(\mathbf{x}) = (x_1 + x_2, x_1 + 2x_2)^T$$





× 0 0 × 0 × ×

DIRECTIONAL DERIVATIVE

The **directional derivative** tells how fast $f : S \to \mathbb{R}$ is changing w.r.t. an arbitrary direction \mathbf{v} :

$$D_{\mathbf{v}}f(\mathbf{x}) := \lim_{h \to 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h} = \nabla f(\mathbf{x})^T \cdot \mathbf{v}.$$

Example: The directional derivative for $\mathbf{v} = (1, 1)$ is:

$$D_{\mathbf{v}}f(\mathbf{x}) = \nabla f(\mathbf{x})^{T} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{\partial f}{\partial x_{1}} + \frac{\partial f}{\partial x_{2}}$$

NB: Some people require that $||\mathbf{v}|| = 1$. Then, we can identify $D_{\mathbf{v}}f(\mathbf{x})$ with the instantaneous rate of change in direction \mathbf{v} – and in our example we would have to divide by $\sqrt{2}$.



PROPERTIES OF THE GRADIENT

- Orthogonal to level curves/surfaces of a function
- Points in direction of greatest increase of f





Proof: Let **v** be a vector with $||\mathbf{v}|| = 1$ and θ the angle between **v** and $\nabla f(\mathbf{x})$.

$$D_{\mathbf{v}}f(\mathbf{x}) = \nabla f(\mathbf{x})^{\mathsf{T}}\mathbf{v} = \|\nabla f(\mathbf{x})\| \|\mathbf{v}\| \cos(\theta) = \|\nabla f(\mathbf{x})\| \cos(\theta)$$

by the cosine formula for dot products and $\|\mathbf{v}\| = 1$. $\cos(\theta)$ is maximal if $\theta = 0$, hence if \mathbf{v} and $\nabla f(\mathbf{x})$ point in the same direction. (Alternative proof: Apply Cauchy-Schwarz to $\nabla f(\mathbf{x})^T \mathbf{v}$ and look for equality.) Analogous: Negative gradient $-\nabla f(\mathbf{x})$ points in direction of greatest *de*crease

PROPERTIES OF THE GRADIENT / 2

Mod. Branin function with neg. grads.



× 0 0 × × ×

JACOBIAN MATRIX

For vector-valued function $f = (f_1, ..., f_m)^T$, $f_j : S \to \mathbb{R}$, the **Jacobian** matrix $J_f : S \to \mathbb{R}^{m \times d}$ generalizes gradient by placing all ∇f_j in its rows:

$$J_{f}(\mathbf{x}) = \begin{pmatrix} \nabla f_{1}(\mathbf{x})^{T} \\ \vdots \\ \nabla f_{m}(\mathbf{x})^{T} \end{pmatrix} = \begin{pmatrix} \frac{\partial f_{1}(\mathbf{x})}{\partial x_{1}} & \cdots & \frac{\partial f_{1}(\mathbf{x})}{\partial x_{d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{m}(\mathbf{x})}{\partial x_{1}} & \cdots & \frac{\partial f_{m}(\mathbf{x})}{\partial x_{d}} \end{pmatrix}$$

• Jacobian gives best linear approximation of distorted volumes



Source: Wikipedia

× × 0 × × ×

JACOBIAN DETERMINANT

Let $f \in C^1$ and $\mathbf{x}_0 \in S$.

Inverse function theorem: Let $\mathbf{y}_0 = f(\mathbf{x}_0)$. If $\det(J_f(\mathbf{x}_0)) \neq 0$, then

• *f* is invertible in a neighborhood of \mathbf{x}_0 ,

2
$$f^{-1} \in C^1$$
 with $J_{f^{-1}}(\mathbf{y}_0) = J_f(\mathbf{x}_0)^{-1}$

- $|\det(J_f(\mathbf{x}_0))|$: factor by which *f* expands/shrinks volumes near \mathbf{x}_0
- If $det(J_f(\mathbf{x}_0)) > 0$, *f* preserves orientation near \mathbf{x}_0
- If $det(J_f(\mathbf{x}_0)) < 0$, *f* reverses orientation near \mathbf{x}_0



HESSIAN MATRIX

For real-valued function $f : S \to \mathbb{R}$, the **Hessian** matrix $H : S \to \mathbb{R}^{d \times d}$ contains all their second derivatives (if they exist):

$$H(\mathbf{x}) = \nabla^2 f(\mathbf{x}) = \left(\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}\right)_{i,j=1,\dots,a}$$

Note: Hessian of *f* is Jacobian of ∇f

Example: Let $f(\mathbf{x}) = \sin(x_1) \cdot \cos(2x_2)$. Then:

$$H(\mathbf{x}) = \begin{pmatrix} -\cos(2x_2) \cdot \sin(x_1) & -2\cos(x_1) \cdot \sin(2x_2) \\ -2\cos(x_1) \cdot \sin(2x_2) & -4\cos(2x_2) \cdot \sin(x_1) \end{pmatrix}$$

- If $f \in C^2$, then *H* is symmetric
- Many local properties (geometry, convexity, critical points) are encoded by the Hessian and its spectrum (→ later)



LOCAL CURVATURE BY HESSIAN

Eigenvector corresponding to largest (resp. smallest) **eigenvalue** of Hessian points in direction of largest (resp. smallest) **curvature**

Example (previous slide): For $\boldsymbol{a} = (-\pi/2, 0)^T$, we have

$$H(\boldsymbol{a}) = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$$

and thus $\lambda_1 = 4, \lambda_2 = 1$, $\mathbf{v}_1 = (0, 1)^T$, and $\mathbf{v}_2 = (1, 0)^T$.



× 0 0 × 0 × × **LIPSCHITZ CONTINUITY** Function $h : S \to \mathbb{R}^m$ is Lipschitz continuous if slopes are bounded:

 $\|h(\mathbf{x}) - h(\mathbf{y})\| \le L \|\mathbf{x} - \mathbf{y}\|$ for each $\mathbf{x}, \mathbf{y} \in S$ and some L > 0

- **Examples** (d = m = 1): sin(x), |x|
- Not examples: 1/x (but *locally* Lipschitz continuous), \sqrt{x}
- If m = d and h differentiable:

h Lipschitz continuous with constant $L \iff J_h \preccurlyeq L \cdot \mathbf{I}_d$

Note: $\mathbf{A} \preccurlyeq \mathbf{B} : \iff \mathbf{B} - \mathbf{A}$ is positive semidefinite, i.e., $\mathbf{v}^{T}(\mathbf{B} - \mathbf{A})\mathbf{v} \ge 0 \quad \forall \mathbf{v} \neq 0$

Proof of " \Rightarrow " for d = m = 1:

$$h'(x) = \lim_{\epsilon \to 0} \frac{h(x+\epsilon) - h(x)}{\epsilon} \le \lim_{\epsilon \to 0} \underbrace{\left| \frac{h(x+\epsilon) - h(x)}{\epsilon} \right|}_{\le L} \le \lim_{\epsilon \to 0} L = L$$

[**Proof** of " \Leftarrow " by mean value theorem: Show that $\lambda_{\max}(J_h) \leq L$.]

LIPSCHITZ GRADIENTS

• Let $f \in C^2$. Since $\nabla^2 f$ is Jacobian of $h = \nabla f$ (m = d):

 ∇f Lipschitz continuous with constant $L \iff \nabla^2 f \preccurlyeq L \cdot \mathbf{I}_d$

- Equivalently, eigenvalues of $\nabla^2 f$ are bounded by L
- Interpretation: Curvature in any direction is bounded by L
- Lipschitz gradients occur frequently in machine learning \implies Fairly weak assumption
- Important for analysis of gradient descent optimization
 - \implies Descent lemma (later)

