

Introduction to Machine Learning

Lecture 3 - Total Derivative and Non-parametric estimations

1 Derivatives of Multivariate Functions

1.1 Scalar function

Consider the scalar function $f : \mathbb{R} \rightarrow \mathbb{R}$:

$$f(x), \quad x \in \mathbb{R}$$

The derivative of f is given by:

$$f' = \frac{df}{dx}$$

That is, we can write:

$$\boxed{df = f' \cdot dx}$$

1.2 Multivariate function - simple example

Consider the multivariate function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$:

$$f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^2$$

The gradient of f is defined by the vector of partial derivatives:

$$\nabla f \triangleq \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix}$$

The total derivative of f is given by:

$$\boxed{df \triangleq \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2}$$

Using inner product notation, we have:

$$\Rightarrow df = \left\langle \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix}, \begin{bmatrix} dx_1 \\ dx_2 \end{bmatrix} \right\rangle = \langle \nabla f, d\mathbf{x} \rangle, \quad d\mathbf{x} \triangleq \begin{bmatrix} dx_1 \\ dx_2 \end{bmatrix}$$

1.2.1 Example

Compute the gradient of the following function using both the direct and the total derivative methods:

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = \sum_{i=1}^2 a_i x_i \quad \mathbf{a}, \mathbf{x} \in \mathbb{R}^2$$

Method I - by definition	Method II - total derivative
$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} \sum_{i=1}^2 a_i x_i \\ \frac{\partial}{\partial x_2} \sum_{i=1}^2 a_i x_i \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$ $\Rightarrow \boxed{\nabla f(\mathbf{x}) = \mathbf{a}}$	$df = \mathbf{a}^T d\mathbf{x} = \langle \mathbf{a}, d\mathbf{x} \rangle$ $\Rightarrow \boxed{\nabla f(\mathbf{x}) = \mathbf{a}}$

1.3 The general case

Consider the differential multivariate function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,
If exist some $\mathbf{g} \in \mathbb{R}^n$ such that:

$$df = \langle \mathbf{g}, d\mathbf{x} \rangle$$

then, \mathbf{g} is the gradient of f , namely, $\mathbf{g} = \nabla f$.

1.3.1 Quadratic example:

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$$

1. Show that without loss of generality, one can assume $\mathbf{A} = \mathbf{A}^T$ (\mathbf{A} is symmetric matrix).
2. Find ∇f .

Solution:

1. We can write \mathbf{A} as a sum of symmetric matrix and anti-symmetric matrix:

$$\mathbf{A} = \underbrace{\frac{\mathbf{A} + \mathbf{A}^T}{2}}_{\mathbf{S}} + \underbrace{\frac{\mathbf{A} - \mathbf{A}^T}{2}}_{\tilde{\mathbf{S}}} = \mathbf{S} + \tilde{\mathbf{S}}$$

now:

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T (\mathbf{S} + \tilde{\mathbf{S}}) \mathbf{x} = \mathbf{x}^T \mathbf{S} \mathbf{x} + \mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x}$$

Notice that $(\tilde{\mathbf{S}} = -\tilde{\mathbf{S}}^T)$:

$$\begin{aligned} \mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x} &= (\mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x})^T, & \mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x} &\in \mathbb{R} \\ \mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x} &= -\mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x} \end{aligned}$$

$$\begin{aligned} &\Rightarrow \mathbf{x}^T \tilde{\mathbf{S}} \mathbf{x} = 0 \\ &\Rightarrow \boxed{\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{S} \mathbf{x}} \end{aligned}$$

Since this is true for any \mathbf{x} we can consider without loss of generality only the symmetric part of \mathbf{A} .

2. Using the product rule

$$\begin{aligned} df &= d\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A} d\mathbf{x} \\ &= \mathbf{x}^T \mathbf{A}^T d\mathbf{x} + \mathbf{x}^T \mathbf{A} d\mathbf{x} \\ &= \mathbf{x}^T (\mathbf{A}^T + \mathbf{A}) d\mathbf{x} \\ &= \left\langle (\mathbf{A} + \mathbf{A}^T) \mathbf{x}, d\mathbf{x} \right\rangle \\ &\Rightarrow \boxed{\nabla f = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}} \end{aligned}$$

if we also consider $\mathbf{A} = \mathbf{A}^T$:

$$\nabla f(\mathbf{x}) = 2\mathbf{A} \mathbf{x}$$

which is vary similar to the scalar case $\frac{d}{dx}(ax^2) = 2ax$.

1.3.2 ℓ_2 - norm**Question 1** Find ∇f of:

$$f(\mathbf{x}) = \|\mathbf{x} - \mathbf{b}\|_2^2$$

Solution:

$$f(\mathbf{x}) = \|\mathbf{x} - \mathbf{b}\|_2^2 = (\mathbf{x} - \mathbf{b})^T (\mathbf{x} - \mathbf{b})$$

$$\begin{aligned} \Rightarrow df &= d\mathbf{x}^T (\mathbf{x} - \mathbf{b}) + (\mathbf{x} - \mathbf{b})^T d\mathbf{x} \\ &= (\mathbf{x} - \mathbf{b})^T d\mathbf{x} + (\mathbf{x} - \mathbf{b})^T d\mathbf{x} \\ &= 2(\mathbf{x} - \mathbf{b})^T d\mathbf{x} \\ &= \langle 2(\mathbf{x} - \mathbf{b}), d\mathbf{x} \rangle \end{aligned}$$

$$\Rightarrow \boxed{\nabla f(\mathbf{x}) = 2(\mathbf{x} - \mathbf{b})}$$

Question 2 Find ∇h of:

$$h(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

Solution I:

$$h(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2 = (\mathbf{Ax} - \mathbf{b})^T (\mathbf{Ax} - \mathbf{b})$$

$$\begin{aligned} \Rightarrow dh &= (\mathbf{A}d\mathbf{x})^T (\mathbf{Ax} - \mathbf{b}) + (\mathbf{Ax} - \mathbf{b})^T \mathbf{A}d\mathbf{x} \\ &= 2(\mathbf{Ax} - \mathbf{b})^T \mathbf{A}d\mathbf{x} \\ &= \langle 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{b}), d\mathbf{x} \rangle \end{aligned}$$

$$\Rightarrow \boxed{\nabla h(\mathbf{x}) = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{b})}$$

Solution II - Using the chain rule:

$$h(\mathbf{x}) = f(\mathbf{Ax}) = f(\mathbf{u}), \quad \mathbf{u} \triangleq \mathbf{Ax}$$

$$\begin{aligned} \Rightarrow dh &= \nabla^T f(\mathbf{u}) d\mathbf{u} \\ &= \nabla^T f(\mathbf{u}) \mathbf{A}d\mathbf{x} \\ &= \langle \mathbf{A}^T \nabla f(\mathbf{u}), d\mathbf{x} \rangle \end{aligned}$$

$$\boxed{\nabla h(\mathbf{x}) = \mathbf{A}^T \nabla f(\mathbf{u}) = \mathbf{A}^T \nabla f(\mathbf{Ax}) = 2\mathbf{A}^T (\mathbf{Ax} - \mathbf{b})}$$

1.3.3 Chain rule

Find ∇h of:

$$h(\mathbf{x}) = \varphi(f(\mathbf{x}))$$

where $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar function and $f, \nabla f$ are known.

Solution:

let $u \triangleq f(\mathbf{x})$

$$dh = \varphi'(u) du = \varphi'(f(\mathbf{x})) \langle \nabla f(\mathbf{x}), d\mathbf{x} \rangle = \langle \varphi'(f) \nabla f, d\mathbf{x} \rangle$$

$$\Rightarrow \nabla h(\mathbf{x}) = \varphi'(f(\mathbf{x})) \nabla f(\mathbf{x})$$

1.4 Matrix derivative example

Find $\nabla_A f$ of:

$$f(\mathbf{A}) = \mathbf{x}^T \mathbf{A} \mathbf{y}$$

Solution:

$$\begin{aligned} df &= \mathbf{x}^T d\mathbf{A} \mathbf{y} \\ &= \text{Tr} \{ \mathbf{x}^T d\mathbf{A} \mathbf{y} \} \\ &= \text{Tr} \{ \mathbf{y} \mathbf{x}^T d\mathbf{A} \} \\ &= \langle \mathbf{x} \mathbf{y}^T, d\mathbf{A} \rangle \end{aligned}$$

$$\Rightarrow \boxed{\nabla_A f = \mathbf{x} \mathbf{y}^T}$$

2 Non-parametric Estimation

Consider a random variable X , with some unknown probability function of p_X .

In Lecture 2, we assumed some model for p_X (X is Gaussian, X is uniform, etc') and we only estimated the model's parameters. In this section, we will estimate p_X without any model assumptions.

2.1 Cumulative Distribution Function (CDF) estimation

Reminder The CDF of the random variable X is given by:

$$F_X(x) \triangleq \Pr\{X \leq x\}$$

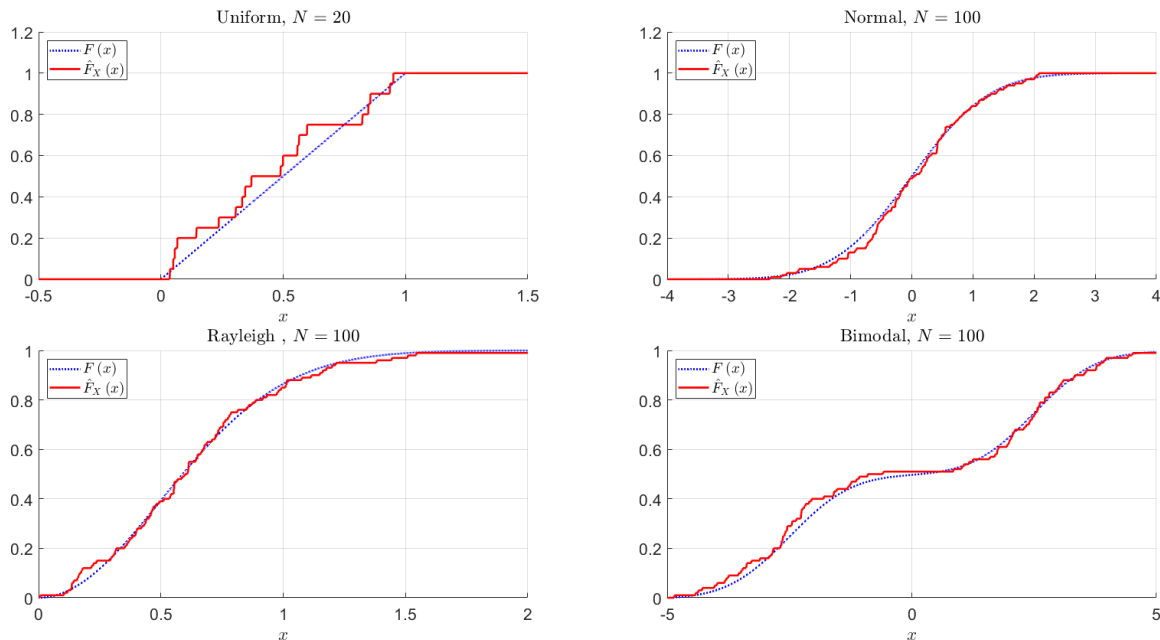
and the probability density p_X (assuming X is continuous) is given by:

$$p_X(x) = F'(x) = \frac{d}{dx}F(x)$$

CDF non-parametric estimation Consider $\{x_i\}_{i=1}^N$, N i.i.d realizations of X and the following estimation for F_X :

$$\hat{F}_X(x) = \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \leq x\}$$

Example We generate N points $\{x_i\}_{i=1}^N$ from different distribution and plot the CDF estimation.



Exercise

1. Compute the bias of $\hat{F}_X(x_0)$ (for some x_0).
2. Compute the MSE ($F_X(x_0)$) (for some x_0).
3. What is the MSE for $N \rightarrow \infty$?

Solution:

1.

$$\begin{aligned}
 \mathbb{E}[\hat{F}_X(x_0)] &= \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \leq x_0\}\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{I}\{x_i \leq x_0\}] \\
 &= \frac{1}{N} \sum_{i=1}^N \Pr\{x_i \leq x_0\} = \Pr\{x_1 \leq x_0\} \\
 &= F_X(x_0) \\
 \Rightarrow b(\hat{F}_X(x_0)) &= \mathbb{E}[\hat{F}_X(x_0)] - F_X(x_0) = 0
 \end{aligned}$$

Hence, \hat{F}_X is unbiased.

2. In Lecture 1 we proved that:

$$\text{MSE}(\hat{\theta}) \triangleq \mathbb{E}\left[(\hat{\theta} - \theta)^2\right] = b^2(\hat{\theta}) + V(\hat{\theta})$$

$$\begin{aligned}
 \Rightarrow \text{MSE}(\hat{F}_X(x_0)) &= \underbrace{b^2(\hat{F}_X(x_0))}_{=0} + \text{Var}(\hat{F}_X(x_0)) = \text{Var}\left(\frac{1}{N} \sum_{i=1}^N \underbrace{\mathbf{I}\{x_i \leq x_0\}}_{\triangleq Y_i}\right) \\
 &= \frac{1}{N^2} \text{Var}\left(\sum_{i=1}^N Y_i\right), \quad \{Y_i\}_i \text{ are i.i.d} \\
 &= \frac{1}{N^2} N \text{Var}(Y_1) = \frac{\mathbb{E}[Y_1^2] - \mathbb{E}^2[Y_1]}{N} \\
 &= \frac{F_X(x_0) - F_X^2(x_0)}{N}
 \end{aligned}$$

3. Using the previous result, we have:

$$\text{MSE}(\hat{F}_X(x_0)) \xrightarrow{N \rightarrow \infty} 0$$

2.2 Histogram

Let $\{x_i\}_{i=1}^N$ be N i.i.d realizations of $X \in \mathcal{X}$.

We can split the domain \mathcal{X} into K **disjoint** intervals $\{R_k\}_{k=1}^K$ (see figure below) such that:

$$\mathcal{X} = \bigsqcup_{k=1}^K R_k, \quad (\text{disjoint union})$$

Then, for any $x \in R_k$ we estimate the PDF by:

$$\hat{p}_X(x) = \frac{1}{|R_k|} \frac{1}{N} \underbrace{\sum_{i=1}^N \mathbf{I}\{x_i \in R_k\}}_{(*)}, \quad x \in R_k$$

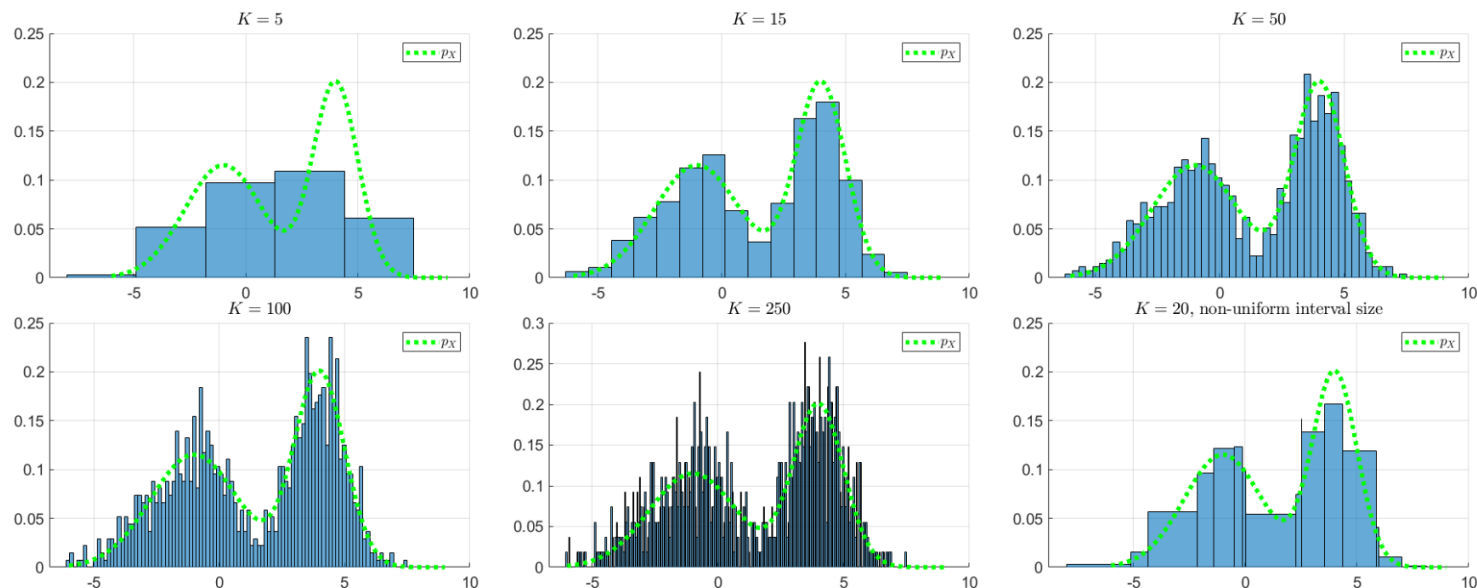
where $|R_k|$ is the length (or volume) of R_k .

In words, $(*)$ is the number of realizations inside the interval R_k .

Example:

Let $\{x_i\}_{i=1}^N$ be $N = 1,000$ i.i.d realizations from an unknown p_X .

We plot the histogram for different values of K :



Notes:

- Remember the bias variance tradeoff.

$$\text{MSE}(\hat{\theta}) = b^2(\hat{\theta}) + V(\hat{\theta})$$

For small values of K the variance is small but the bias is large,
whereas for large values of K the bias is small but the variance is large.

- A reasonable choice is $K = \sqrt{N}$.

Exercise:

Show that \hat{p}_X is a valid density function, namely show that:

1. $\hat{p}_X(x) \geq 0, \forall x$
2. $\int_{\mathcal{X}} \hat{p}_X(x) dx = 1$

Solution:

1. $\hat{p}_X(x) \geq 0$ is immediate from the definition

$$\hat{p}_X(x) = \frac{1}{|R_k|} \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \in R_k\} \geq 0$$

2.

$$\begin{aligned} \int_{\mathcal{X}} \hat{p}_X(x) dx &= \sum_{k=1}^K \int_{R_k} \hat{p}_X(x) dx \\ &= \sum_{k=1}^K \int_{R_k} \frac{1}{|R_k|} \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \in R_k\} dx \\ &= \sum_{k=1}^K \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \in R_k\} \underbrace{\int_{R_k} \frac{1}{|R_k|} dx}_{=1} \\ &= \frac{1}{N} \underbrace{\sum_{k=1}^K \sum_{i=1}^N \mathbf{I}\{x_i \in R_k\}}_{=N} \\ &= 1 \end{aligned}$$

2.3 Kernel Density Estimation

When there are not enough data points to properly estimate p_X , one should consider using KDE.

For a given kernel h , the Kernel Density Estimation (KDE) is given by:

$$\hat{p}_h(x) = \frac{1}{N} \sum_{i=1}^N h(x - x_i)$$

where:

1. $h \geq 0$
2. $\int_{\mathcal{X}} h(x) dx = 1$

Common choices for h are:

1. Rectangular window:

$$h(x) = \frac{1}{\alpha} \begin{cases} 1 & |x| \leq \frac{\alpha}{2} \\ 0 & \text{else} \end{cases}, \quad \alpha > 0$$

2. Gaussian window:

$$h(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}, \quad \sigma^2 > 0$$

- Note: We can write $\hat{p}_h(x)$ as a convolution with h (δ is the Dirac delta function):

$$\hat{p}_h(x) = \frac{1}{N} \sum_{i=1}^N h(x - x_i) = \left(\frac{1}{N} \sum_{i=1}^N \delta(x - x_i) \right) * h$$

Example:

Let $\{x_i\}_{i=1}^N$ be $N = 100$ i.i.d realizations from the unknown p_X .

