

Mathematics Primer

Machine Intelligence

Neural Information Processing Group

Summer Term 2018

- 1 Linear Algebra
 - Transpose, Inverse, Rank and Trace
 - Determinant
 - Eigenanalysis
 - Matrix Gradient
- 2 Analysis
 - Metrics
 - Jacobi and Hessian
 - Taylor Series
 - Optimization
- 3 Probability Theory
 - Combinatorics
 - Random Variables and Vectors
 - Conditional Probabilities and Independence
 - Expectations and Moments

Outline

- 1 Linear Algebra
 - Transpose, Inverse, Rank and Trace
 - Determinant
 - Eigenanalysis
 - Matrix Gradient
- 2 Analysis
 - Metrics
 - Jacobi and Hessian
 - Taylor Series
 - Optimization
- 3 Probability Theory
 - Combinatorics
 - Random Variables and Vectors
 - Conditional Probabilities and Independence
 - Expectations and Moments

Matrix Multiplication, Transpose and Inverse

Consider matrices $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{B} \in \mathbb{R}^{m \times p}$ with elements $(\mathbf{A})_{ij} = a_{ij}$, $(\mathbf{B})_{ij} = b_{ij}$.

- The **product** $\mathbf{AB} \in \mathbb{R}^{n \times p}$ has elements $(\mathbf{AB})_{ij} = \sum_{r=1}^m a_{ir} b_{rj}$.
- The **transpose** \mathbf{A}^\top has elements $(\mathbf{A}^\top)_{ij} = a_{ji}$.
- The **inverse** \mathbf{A}^{-1} of a square matrix satisfies $\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$.
- The following identities hold:

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top \quad (1)$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1} \quad (2)$$

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \quad (3)$$

Rank and Trace

Linear independence

A set of vectors $\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$ is **linearly independent**, if $\sum_{i=1}^n \alpha_i \mathbf{a}_i = \mathbf{0}$ holds only if all $\alpha_i = 0$. This means none of the vectors can be expressed as a linear combination of the others.

Rank

The **rank** $\text{rank}(\mathbf{A})$ of a matrix \mathbf{A} is the maximum number of linearly independent rows (or columns).

Trace

The **trace** of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined as $\text{Tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$.

It holds:

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}) \quad (4)$$

Determinant

The **determinant** $\det(\mathbf{A})$ shows certain properties of a square matrix \mathbf{A}

- $\det(\mathbf{A}) = 0$ iff the rows (or columns) are linearly dependent
- $\det(\mathbf{A}) \neq 0$ iff \mathbf{A} is invertible

Note:

- Determinant of the identity matrix: $\det(\mathbf{I}) = 1$
- Determinant of a transposed matrix: $\det(\mathbf{A}) = \det(\mathbf{A}^T)$
- Determinant of a product of two matrices:

$$\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$$

Determinant calculation (general)

Calculation of the determinant of a $n \times n$ -Matrix \mathbf{A} :

$$\det(\mathbf{A}) = \sum_j A_{ij} C_{ij}.$$

Row i can be any row, the result is always the same. The **cofactors** C_{ij} are defined as $C_{ij} = (-1)^{i+j} \det([\mathbf{A}]_{\emptyset ij})$, where $[\mathbf{A}]_{\emptyset ij}$ is the submatrix that remains when the i -th row and j -th column are removed:

$$[\mathbf{A}]_{\emptyset ij} = \begin{pmatrix} A_{11} & A_{12} & \cdots & \emptyset & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & \emptyset & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \emptyset & \ddots & \vdots \\ \emptyset & \emptyset & \emptyset & \emptyset & \emptyset & \emptyset \\ \vdots & \vdots & \ddots & \emptyset & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & \emptyset & \cdots & A_{nn} \end{pmatrix}$$

Determinant calculation (special cases)

$$|\mathbf{A}| = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc$$

$$\begin{aligned} |\mathbf{A}| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\ &= aei + bfg + cdh - ced - bdi - afh \end{aligned}$$

Determinant and Inverse

The **inverse** \mathbf{A}^{-1} of a square matrix \mathbf{A} exists iff $\det(\mathbf{A}) \neq 0$ (matrix not singular).

Calculation of the inverse matrix:

$$\mathbf{A}^{-1} = \frac{\text{adj}[\mathbf{A}]}{\det(\mathbf{A})}$$

where the **adjoint** $\text{adj}[\mathbf{A}]$ of \mathbf{A} is the matrix whose elements are the cofactors:

$$(\text{adj}[\mathbf{A}])_{ij} = C_{ji}$$

The determinant of an inverse matrix is given by

$$\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$$

Eigendecomposition of a Matrix

Problem: Find the Eigenvectors and Eigenvalues of a $N \times N$ matrix \mathbf{A} .

- Consider the system of linear equations:

$$\begin{aligned}\mathbf{A}\mathbf{x} &= \lambda\mathbf{x} \\ (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} &= \mathbf{0}\end{aligned}$$

- Solutions: N Eigenvectors $\mathbf{x} = \mathbf{v}_i$ and corresponding Eigenvalues $\lambda = \lambda_i$
- $\mathbf{B}\mathbf{x} = \mathbf{0}$ has non-trivial solutions iff $\det(\mathbf{B}) = 0$
- Therefore, non-trivial λ are the roots of the **characteristic polynomial**:

$$p(\lambda) \equiv \det(\mathbf{A} - \lambda\mathbf{I}) = 0$$

Eigenvalues and Eigenvectors

Characteristic Equation:

$$p(\lambda) \equiv \det(\mathbf{A} - \lambda \mathbf{I}) = 0$$

- Polynomial of order N
- N (not necessarily distinct) solutions
- Number of non-zero Eigenvalues: $\text{rank}(\mathbf{A})$
- In general: Eigenvalues are complex
- For symmetric matrices ($\mathbf{A} = \mathbf{A}^\top$): Eigenvalues are real
- Determinant: $\det(\mathbf{A}) = \prod_{i=1}^M \lambda_i$
- Trace: $\text{Tr}(\mathbf{A}) = \sum_{i=1}^M \lambda_i$

Matrix Gradient

The **gradient** of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is given by

$$\nabla f \equiv \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^\top$$

Examples:

- linear $f : \mathbf{x} \mapsto \mathbf{a}^\top \mathbf{x}$ $\nabla f(\mathbf{x}) = \mathbf{a}$
- quadratic $f : \mathbf{x} \mapsto \mathbf{x}^\top \mathbf{A} \mathbf{x}$ $\nabla f(\mathbf{x}) = (\mathbf{A}^\top + \mathbf{A})\mathbf{x}$

Consider a scalar-valued function f of the elements of an $n \times m$ matrix \mathbf{W} , $f : \mathbf{W} \mapsto \mathbb{R}$, $f(\mathbf{W}) = f(w_{11}, \dots, w_{nm})$.

The **matrix gradient** of f w.r.t. \mathbf{W} is defined as

$$\frac{\partial f}{\partial \mathbf{W}} = \begin{pmatrix} \frac{\partial f}{\partial w_{11}} & \dots & \frac{\partial f}{\partial w_{n1}} \\ \vdots & & \vdots \\ \frac{\partial f}{\partial w_{1m}} & \dots & \frac{\partial f}{\partial w_{nm}} \end{pmatrix}$$

Outline

- 1 Linear Algebra
 - Transpose, Inverse, Rank and Trace
 - Determinant
 - Eigenanalysis
 - Matrix Gradient
- 2 Analysis
 - Metrics
 - Jacobi and Hessian
 - Taylor Series
 - Optimization
- 3 Probability Theory
 - Combinatorics
 - Random Variables and Vectors
 - Conditional Probabilities and Independence
 - Expectations and Moments

Definitions from Functional Analysis

Functions, Functionals and Operators

Two sets \mathcal{M} and \mathcal{N} are connected by a **functional dependency**, if to each $x \in \mathcal{M}$ there corresponds a unique element $y \in \mathcal{N}$. This functional dependency is called

- a **function** if \mathcal{M} and \mathcal{N} are sets of numbers
- a **functional** if \mathcal{M} is a set of functions and \mathcal{N} a set of numbers
- an **operator** if both sets are sets of functions

Example: Linear integral operator T with kernel $K(t, x)$:

$$Tf(x) = \int_a^b K(t, x)f(t)dt$$

Infimum and Supremum

Infimum, Supremum

Let D be a subset of \mathbb{R} . A number K is called **supremum (infimum)** of D , if K is the smallest upper bound (largest lower bound) of D :

$$x \leq K \text{ (} x \geq K \text{), } \forall x \in D$$

We write: $\sup D = K$ ($\inf D = K$).

Examples:

- For the closed interval $D = [a, b]$, $a \leq b$: $\sup D = b$, $\inf D = a$.
- For $D = \{\frac{n}{n+1}, n \in \mathbb{N}\}$: $\sup D = 1$.

Metric Space

Metric

A metric (or distance function) on a set X is a non-negative mapping

$$d : X \times X \rightarrow \mathbb{R}^+$$

$$(x, y) \mapsto d(x, y)$$

with the following characteristics

- ① Positive definiteness: $d(x, y) = 0$ iff $x = y$, $d(x, y) > 0$ otherwise
- ② Symmetry: $d(x, y) = d(y, x)$, $\forall x, y \in X$
- ③ Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$, $\forall x, y, z \in X$

- The pair (X, d) forms a **metric space**
- $d(x, y)$ is called the distance between x and y .

Jacobi and Hessian

- The matrix of the partial derivatives of a vector-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is known as **Jacobi matrix** and given by

$$Jf \equiv \frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

- The square matrix of second-order partial derivatives of a scalar-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called **Hessian matrix** and given by

$$Hf \equiv \frac{\partial^2 f}{\partial x^2} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Taylor Series

Taylor Series in \mathbb{R}

Let $f : I \rightarrow \mathbb{R}$ be an infinitely often differentiable function, and $x_0 \in I$. Then the Taylor series around x_0 is defined as

$$\begin{aligned} f(x) &= \sum_{n=0}^{\infty} \frac{1}{n!} \left. \frac{d^n f(x)}{dx^n} \right|_{x_0} (x - x_0)^n \\ &= f(x_0) + f'(x_0) \cdot (x - x_0) + \frac{1}{2} f''(x_0) \cdot (x - x_0)^2 + \dots \end{aligned}$$

Taylor Series in \mathbb{R}^n

Let f be an infinitely smooth scalar-valued function with domain in \mathbb{R}^n :

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f_{(\mathbf{x}_0)}^\top (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H} f_{(\mathbf{x}_0)} (\mathbf{x} - \mathbf{x}_0) + \dots$$

Local Extrema

Let f be a scalar-valued function $\mathbb{R}^n \rightarrow \mathbb{R}$.

Critical Points

A point x_0 , where $\nabla f(x_0) = 0$ is called a critical point of f .

Local Extrema

A critical point x_0 of f is

- a minimum of f , if all Eigenvalues of $(Hf)(x_0)$ are positive (the Hessian is **positive definite**)
- a maximum of f , if all Eigenvalues of $(Hf)(x_0)$ are negative (the Hessian is **negative definite**)
- no extremum of f , in all other cases (the Hessian is **indefinite**)

Convexity

Convex Functions

Let $U \subset \mathbb{R}^N$ be open and convex. A function $f : U \rightarrow \mathbb{R}$ is called (strictly) convex, if for all $x_1, x_2 \in U$ with $x_1 \neq x_2$ and all $0 < \lambda < 1$

$$f(\lambda x_1 + (1 - \lambda)x_2)(<) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

Concave Functions

f is called concave, if $-f$ is convex.

The Lagrange Method (Equality Constraints)

Problem: Maximization of a function $f(\mathbf{w}): \mathbb{R}^n \rightarrow \mathbb{R}$ under some **equality** constraints.

$$\max f(\mathbf{w}), \quad \text{s.t.} \quad g_i(\mathbf{w}) = 0, \quad \forall i \in \{1, \dots, k\}$$

Solution: Form the **Lagrangian**

$$\mathcal{L}(\mathbf{w}, \lambda_1, \dots, \lambda_k) = f(\mathbf{w}) + \sum_{i=1}^k \lambda_i g_i(\mathbf{w}),$$

where $\lambda_1, \dots, \lambda_k$ are called Lagrange multipliers. Find the stationary points (saddle points) of the Lagrangian w.r.t. both \mathbf{w} and all the λ_i :

$$\frac{\partial \mathcal{L}(\mathbf{w}, \lambda_1, \dots, \lambda_k)}{\partial \mathbf{w}} = \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} + \sum_{i=1}^k \lambda_i \frac{\partial g_i(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{0}$$

and

$$\frac{\partial \mathcal{L}(\mathbf{w}, \lambda_1, \dots, \lambda_k)}{\partial \lambda_i} = g_i(\mathbf{w}) = 0, \forall i.$$

The Lagrange Method (Inequality Constraints)

Now: Maximization of a function $f(\mathbf{w})$ under some **inequality** constraints.

$$\max f(\mathbf{w}), \quad \text{s.t.} \quad h_i(\mathbf{w}) \leq 0, \quad \forall i \in \{1, \dots, k\}$$

Solution: Find the stationary points of the Lagrangian

$$\mathcal{L}(\mathbf{w}, \lambda_1, \dots, \lambda_k) = f(\mathbf{w}) + \sum_{i=1}^k \lambda_i h_i(\mathbf{w}),$$

w.r.t. \mathbf{w} under the constraints

$$h_i(\mathbf{w}) \leq 0, \forall i$$

$$\lambda_i \geq 0, \forall i$$

$$\lambda_i \cdot h_i(\mathbf{w}) = 0, \forall i,$$

which are known as the **Karush-Kuhn-Tucker (KKT) conditions**.

Outline

- 1 Linear Algebra
 - Transpose, Inverse, Rank and Trace
 - Determinant
 - Eigenanalysis
 - Matrix Gradient
- 2 Analysis
 - Metrics
 - Jacobi and Hessian
 - Taylor Series
 - Optimization
- 3 Probability Theory
 - Combinatorics
 - Random Variables and Vectors
 - Conditional Probabilities and Independence
 - Expectations and Moments

Combinatorics

Consider a set consisting of n elements. The **power set** is the set of all subsets, its cardinality is 2^n .

- **Permutation:** arrangement of n elements in a certain order
 - # **without** repetitions: $P_n = n!$
 - # **with** repetitions ($k \leq n$ repeated elements): $P_n^{(k)} = \frac{n!}{k!}$
- **Combination:** choice of k out of n elements regardless of order
 - # **without** repetitions: $C_n^{(k)} = \binom{n}{k} = \frac{n!}{k!(n-k)!}$
 - # **with** repetitions: $C_n^{(k)} = \binom{n+k-1}{k}$
- **Variation:** choice of k out of n elements taking their order into account
 - # **without** repetitions $V_n^{(k)} = k! \binom{n}{k}$
 - # **with** repetitions: $V_n^{(k)} = n^k$

Random Variable

Consider a set Ω of elementary events w , e.g. all possible outcomes of an experiment. The mapping

$$\Omega \rightarrow R \subset \mathbb{R}$$

$$w \rightarrow X(w) \equiv X$$

is called a **random variable**.

- If R consists of a finite or countable infinite number of elements, then X is called a **discrete** random variable.
- If $R = \mathbb{R}$ or R consists of intervals from \mathbb{R} , then X is called a **continuous** random variable.

Example: Roll dice

w_1 : 1 comes up $\rightarrow X(w_1) = 1, \dots, w_6$: 6 comes up $\rightarrow X(w_6) = 6$

Distribution of a Random Variable

The **cumulative distribution function (cdf)** or simply **distribution function** of a random variable X at point z is defined as the probability that $X \leq z$:

$$F_X(z) = P(X \leq z)$$

- Allowing z to vary in $(-\infty, \infty)$ defines the cdf for all values of X .
- $0 \leq F_X \leq 1$, a nondecreasing and continuous function for continuous X .

Example: Roll ideal dice, where $P(X = i) = \frac{1}{6} \forall i$

$$F_X(z) = \begin{cases} 0 & \text{for } z < 1 \\ 1/6 & \text{for } 1 \leq z < 2 \\ 2/6 & \text{for } 2 \leq z < 3 \\ \dots & \\ 1 & \text{for } z \geq 6 \end{cases}$$

Probability Density of a Continuous Variable

The **probability density function (pdf)** p_X of a continuous X is obtained as the derivative of its cdf:

$$p_X(z) = \left. \frac{dF_X(x)}{dx} \right|_{x=z}$$

In practice, the cdf is computed from the known pdf using the inverse relationship

$$F_X(z) = \int_{-\infty}^z p_X(t) dt$$

Example: Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$

■ **cdf** $F(z) \equiv P(X \leq z) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$

■ **pdf** $p(z) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$

Distribution of a Random Vector

The **distribution function** of a **random vector** \mathbf{X} :

$$\Omega \rightarrow \mathbb{R}^N \subset \mathbb{R}^N$$

$$w \rightarrow \mathbf{X}(w) \equiv \mathbf{X}$$

at a point \mathbf{z} is given by

$$F_{\mathbf{X}}(\mathbf{z}) = P(\mathbf{X} \leq \mathbf{z})$$

Distribution of a Random Vector

Example

Toss a German 2 Euro and a German 20 Cent coin.

- $w_1 = \{2 \text{ Euro: eagle, 20 Cent: gate}\} \rightarrow \mathbf{X}(w_1) = (1, 1)^\top$
- $w_2 = \{2 \text{ Euro: eagle, 20 Cent: number}\} \rightarrow \mathbf{X}(w_2) = (1, 2)^\top$
- $w_3 = \{2 \text{ Euro: number, 20 Cent: gate}\} \rightarrow \mathbf{X}(w_3) = (2, 1)^\top$
- $w_4 = \{2 \text{ Euro: number, 20 Cent: number}\} \rightarrow \mathbf{X}(w_4) = (2, 2)^\top$

$$F_{\mathbf{X}}(\mathbf{z}) = \begin{cases} 0 & \text{for } (z_1 < 1) \vee (z_2 < 1) \\ 1/4 & \text{for } (1 \leq z_1 < 2) \wedge (1 \leq z_2 < 2) \\ 1/2 & \text{for } (1 \leq z_1 < 2) \wedge (2 \leq z_2) \\ 1/2 & \text{for } (2 \leq z_1) \wedge (1 \leq z_2 < 2) \\ 1 & \text{for } (2 \leq z_1) \wedge (2 \leq z_2) \end{cases}$$

Conditional Probabilities

Conditional Probabilities

Consider two discrete random variables X and Y . The conditional probability of Y given X :

$$P(Y = y|X = x) = \frac{P(X = x, Y = y)}{P(X = x)}, \quad P(X = x) \neq 0$$

Conditional Probability Densities

Consider two continuous random vectors \mathbf{X} , \mathbf{Y} and their joint probability density. The conditional probability density of \mathbf{Y} given \mathbf{X} : Probability for finding $\mathbf{Y} \in [\mathbf{y}, \mathbf{y} + d\mathbf{y}]$ if we already know that $\mathbf{X} \in [\mathbf{x}, \mathbf{x} + d\mathbf{x}]$.

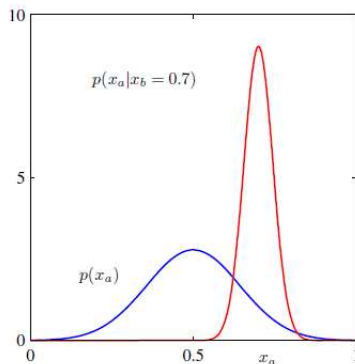
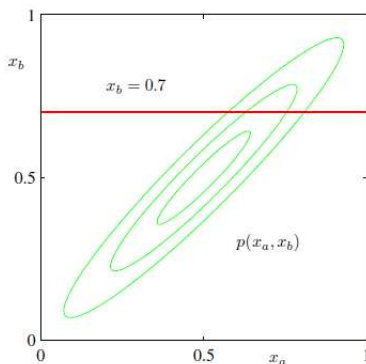
$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})} \quad \text{almost everywhere in } \mathbf{X}$$

Independence

Statistical Independence of Continuous Random Vectors

The random vectors \mathbf{X} and \mathbf{Y} are statistically independent iff

$$p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}) \quad \text{or equivalently} \quad p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$



Marginals

Law of Total Probability (Discrete Random Variables)

Marginalisation over Y :

$$P(X = x) = \sum_k P(X = x, Y = y_k)$$

Marginal Densities (Continuous Random Vectors)

Given the joint density $p_{X,Y}(\mathbf{x}, \mathbf{y})$ of two random vectors \mathbf{X} and \mathbf{Y} , the marginal density $p_X(\mathbf{x})$ is obtained by integrating over the other random vector:

$$p_X(\mathbf{x}) = \int_{-\infty}^{\infty} p_{X,Y}(\mathbf{x}, \tilde{\mathbf{y}}) d\tilde{\mathbf{y}}$$

Bayes' Theorem

Bayes' Theorem (Discrete Random Variables)

$$\begin{aligned}P(Y = y|X = x) &= \frac{P(X = x|Y = y)P(Y = y)}{P(X = x)} \\&= \frac{P(X = x|Y = y)P(Y = y)}{\sum_k P(X = x|Y = y_k)P(Y = y_k)}\end{aligned}$$

Bayes' Theorem (Continuous Random Vectors)

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{\int p(\mathbf{x}|\tilde{\mathbf{y}})p(\tilde{\mathbf{y}})d\tilde{\mathbf{y}}}$$

Decomposition

Factorization of a joint pdf (or cdf), as given by the Chain Rule:

$$p(x_1, \dots, x_d) = p(x_1)p(x_2|x_1) \dots p(x_d|x_1, \dots, x_{d-1})$$

- Special case: Statistical Independence

$$p(x_1, \dots, x_d) = p(x_1)p(x_2) \dots p(x_d) = \prod_{k=1}^d p(x_k)$$

- Special case: 1st order Markov chain

$$p(x_1, \dots, x_d) = p(x_d|x_{d-1})p(x_{d-1}|x_{d-2}) \dots p(x_2|x_1)p(x_1)$$

Expectations

- In Practice: Probability density usually unknown
- However: Expectations of functions can be directly estimated from the data

The expectation of a scalar-, vector- or matrix-valued function $\mathbf{g}(\mathbf{X})$ of a random vector \mathbf{X} , as defined below, can be estimated from a dataset of k **i.i.d.** samples $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(k)}$:

$$\langle \mathbf{g}(\mathbf{X}) \rangle \equiv \int_{-\infty}^{\infty} \mathbf{g}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \approx \frac{1}{k} \sum_{j=1}^k \mathbf{g}(\mathbf{x}^{(j)})$$

- Linearity: $\langle a\mathbf{X} + b\mathbf{Y} + c \rangle = a\langle \mathbf{X} \rangle + b\langle \mathbf{Y} \rangle + c$
- $p_{\mathbf{X}}$ known \Rightarrow Expectations of arbitrary function available
- Expectations for all functions f known $\Rightarrow p_{\mathbf{x}}$ can be determined
 \Rightarrow Statistics of \mathbf{X} completely known

Moments

Moments of a random vector $\mathbf{X} = (X_1, \dots, X_n)$ are typical expectations used to characterize it. They are obtained when $\mathbf{g}(\mathbf{X})$ consists of products of components of \mathbf{X} .

Examples:

- 1st order: $\langle X_i \rangle = \int p(x_i) x_i dx_i$... mean value μ_i , $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)$
- 2nd order: $\langle X_i X_j \rangle$... correlation between X_i, X_j
- 3rd order: $\langle X_i X_j X_k \rangle$... e.g. skewness

Correlation Matrix

The correlation matrix of a random vector \mathbf{X} contains all second order moments $\langle X_i X_j \rangle$:

$$\mathbf{R}_X = \langle \mathbf{X} \mathbf{X}^\top \rangle$$

- Symmetry: $\mathbf{R}_X = \mathbf{R}_X^\top$
- Positive semidefinite: $\mathbf{a}^\top \mathbf{R}_X \mathbf{a} \geq 0, \forall \mathbf{a}$
 - \Rightarrow all eigenvalues real and nonnegative
 - \Rightarrow all eigenvectors are mutually orthogonal

Covariance Matrix

The covariance matrix of a random vector \mathbf{X} is given by

$$\mathbf{C}_X \equiv \langle (\mathbf{X} - \boldsymbol{\mu}_X)(\mathbf{X} - \boldsymbol{\mu}_X)^\top \rangle = \langle \mathbf{X} \mathbf{X}^\top \rangle - \boldsymbol{\mu}_X \boldsymbol{\mu}_X^\top = \mathbf{R}_X - \boldsymbol{\mu}_X \boldsymbol{\mu}_X^\top$$

and the components C_{ij} are calculated as

$$C_{ij} = \langle X_i X_j \rangle - \mu_i \mu_j = \iint p(x_i, x_j) x_i x_j dx_i dx_j - \mu_i \mu_j.$$

- $C_{ii} = \sigma_i^2$... variance of X_i
- For zero mean, correlation and covariance matrix are identical

Uncorrelatedness and Independence

Two random vectors \mathbf{X} and \mathbf{Y} are **uncorrelated** iff their cross-covariance matrix $\mathbf{C}_{XY} = \langle \mathbf{X}\mathbf{Y}^\top \rangle - \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top = \mathbf{0}$.

- **Uncorrelatedness** implies that

$$\mathbf{R}_{XY} = \langle \mathbf{X}\mathbf{Y}^\top \rangle = \langle \mathbf{X} \rangle \langle \mathbf{Y}^\top \rangle = \boldsymbol{\mu}_X \boldsymbol{\mu}_Y^\top,$$

while **independence** implies that

$$\langle g(\mathbf{X})h(\mathbf{Y}) \rangle = \langle g(\mathbf{X}) \rangle \langle h(\mathbf{Y}) \rangle \quad \text{for any } g, h$$

\Rightarrow Independence much stronger property than uncorrelatedness

- Special property of **Gaussian distributions**:
uncorrelatedness = independence