Machine Learning with Python MTH786U/P 2020/21

Mathematical preliminaries

Mihail Poplavskyi, Queen Mary University of London (QMUL)

LINEAR ALGEBRA

MTH786U/P 2020/2021



MTH786U/P 2020/2021



Matrix-vector multiplication:
$$Xw = y \Leftrightarrow \sum_{i=1}^{d} x_{ij}w_j = y_i \quad \forall i \in \{1,...,s\}$$



Matrix-vector multiplication:
$$Xw = y \Leftrightarrow \sum_{i=1}^{d} x_{ij}w_j = y_i \quad \forall i \in \{1,...,s\}$$

Matrix-matrix multiplication:
$$XW = Y$$
 \Leftrightarrow $\sum_{j=1}^{j=1} x_{ij}w_{jk} = y_{ik}$ for $y \in \mathbb{D}^{s \times d}$

 $X \in \mathbb{R}^{s \times d}, W \in \mathbb{R}^{d \times m}, Y \in \mathbb{R}^{s \times m}$

Matrix:
$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & & \ddots & \vdots \\ x_{s1} & x_{s2} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{s \times d}$$

Matrix-vector multiplication:
$$Xw = y \Leftrightarrow \sum_{i=1}^{a} x_{ij}w_j = y_i \quad \forall i \in \{1,...,s\}$$

Matrix-matrix multiplication:
$$XW = Y$$
 \Leftrightarrow $\sum_{i=1}^{a} x_{ij} w_{jk} = y_{ik}$ for

pduct:
$$X \in \mathbb{R}^{s \times d}, W \in \mathbb{R}^{d \times m}, Y \in \mathbb{R}^{s \times m}$$

Inner / Dot product:
$$\langle x,y\rangle := \sum_{j=1}^{d} x_j y_j = x^{\mathsf{T}} y = y^{\mathsf{T}} x = x \cdot y \quad \text{for } x,y \in \mathbb{R}^{d \times 1}$$

Matrix:
$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & & \ddots & \vdots \\ x_{s1} & x_{s2} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{s \times d}$$

Matrix-vector multiplication:
$$Xw = y \Leftrightarrow \sum_{i=1}^{a} x_{ij}w_j = y_i \quad \forall i \in \{1,...,s\}$$

Matrix-matrix multiplication:
$$XW = Y$$
 \Leftrightarrow $\sum_{i=1}^{a} x_{ij} w_{jk} = y_{ik}$ for

$$j=1$$
 $X \in \mathbb{R}^{s \times d}, W \in \mathbb{R}^{d \times m}, Y \in \mathbb{R}^{s \times m}$

Inner / Dot product:
$$\langle x, y \rangle := \sum_{j=1}^{d} x_j y_j = x^\top y = y^\top x = x \cdot y \quad \text{for } x, y \in \mathbb{R}^{d \times 1}$$

for
$$x, y \in \mathbb{R}^{d \times 1}$$

Norm:
$$||x|| := \sqrt{\langle x, x \rangle}$$

The transpose
$$X^{\top}$$
 of a matrix X is defined as $X^{\top} = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{s1} \\ x_{12} & x_{22} & \dots & x_{s2} \\ \vdots & & \ddots & \vdots \\ x_{1d} & x_{2d} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{d \times s}$

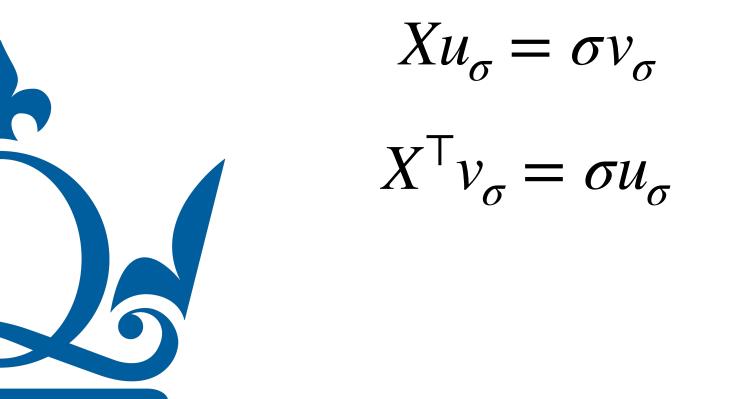


The transpose
$$X^{\top}$$
 of a matrix X is defined as $X^{\top} = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{s1} \\ x_{12} & x_{22} & \dots & x_{s2} \\ \vdots & & \ddots & \vdots \\ x_{1d} & x_{2d} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{d \times s}$

We can find so-called singular vectors u_{σ}, v_{σ} and singular values σ that satisfy

$$Xu_{\sigma} = \sigma v_{\sigma}$$

$$X^{\mathsf{T}} v_{\sigma} = \sigma u_{\sigma}$$



The transpose
$$X^{\top}$$
 of a matrix X is defined as $X^{\top} = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{s1} \\ x_{12} & x_{22} & \dots & x_{s2} \\ \vdots & & \ddots & \vdots \\ x_{1d} & x_{2d} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{d \times s}$

We can find so-called singular vectors u_{σ}, v_{σ} and singular values σ that satisfy



$$Xu_{\sigma} = \sigma v_{\sigma} \implies X^{\mathsf{T}} X u_{\sigma} = \sigma^{2} u_{\sigma}$$

$$X^{\mathsf{T}} v_{\sigma} = \sigma u_{\sigma} \implies XX^{\mathsf{T}} v_{\sigma} = \sigma^{2} v_{\sigma}$$

The transpose
$$X^{\top}$$
 of a matrix X is defined as $X^{\top} = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{s1} \\ x_{12} & x_{22} & \dots & x_{s2} \\ \vdots & & \ddots & \vdots \\ x_{1d} & x_{2d} & \dots & x_{sd} \end{pmatrix} \in \mathbb{R}^{d \times s}$

We can find so-called singular vectors u_{σ} , v_{σ} and singular values σ that satisfy



$$Xu_{\sigma} = \sigma v_{\sigma} \\ X^{\mathsf{T}} v_{\sigma} = \sigma u_{\sigma} \implies X^{\mathsf{T}} X u_{\sigma} = \sigma^{2} u_{\sigma} \\ XX^{\mathsf{T}} v_{\sigma} = \sigma^{2} v_{\sigma} \implies \sigma = \frac{\|X u_{\sigma}\|}{\|u_{\sigma}\|} = \frac{\|X^{\mathsf{T}} v_{\sigma}\|}{\|v_{\sigma}\|}$$

We can find so-called singular vectors u_{σ} , v_{σ} and singular values σ that satisfy

$$Xu_{\sigma} = \sigma v_{\sigma} \\ X^{\top} Xu_{\sigma} = \sigma^{2} u_{\sigma} \\ X^{\top} v_{\sigma} = \sigma u_{\sigma}$$
 \Longrightarrow $\sigma = \frac{\|Xu_{\sigma}\|}{\|u_{\sigma}\|} = \frac{\|X^{\top}v_{\sigma}\|}{\|v_{\sigma}\|}$

It can be shown that for every matrix $X \in \mathbb{R}^{s \times d}$ there exist $\{\sigma_i\}_{i=1}^{\min(s,d)}$ with $\sigma_1 \ge \sigma_2 \ge \dots \ge \sigma_{\min(s,d)}$ and vectors $\{u_{\sigma_i}\}_{i=1}^{\min(s,d)}$ and $\{v_{\sigma_i}\}_{i=1}^{\min(s,d)}$ such that



$$Xw = \sum_{j=1}^{\min(s,d)} \sigma_j \langle w, u_j \rangle v_j$$
 and $X^{\top}y = \sum_{j=1}^{\min(s,d)} \sigma_j \langle y, v_j \rangle u_j$

$$X^{\top} y = \sum_{j=1}^{\min(s,d)} \sigma_j \langle y, v_j \rangle u_j$$

for all $w \in \mathbb{R}^d$ and $y \in \mathbb{R}^s$



MTH786U/P 2020/2021



MTH786U/P 2020/2021

Assume we have a function $f(x_1, x_2, ..., x_d)$

and let
$$x = (x_1, x_2, ..., x_d)$$



Assume we have a function $f(x_1, x_2, ..., x_d)$

and let $x = (x_1, x_2, ..., x_d)$. Then the gradient

of f is the vector of partial derivatives:

$$\nabla f(x)^{\top} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \vdots \\ \frac{\partial}{\partial x_d} f(x) \end{bmatrix}$$



Assume we have a function $f(x_1, x_2, ..., x_d)$

and let $x = (x_1, x_2, ..., x_d)$. Then the gradient

of f is the vector of partial derivatives:

$$\nabla f(x)^{\top} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \vdots \\ \frac{\partial}{\partial x_d} f(x) \end{bmatrix}$$

Example:
$$f(x_1, x_2) = (x_1x_2 - y)^2$$



Assume we have a function $f(x_1, x_2, ..., x_d)$

and let $x = (x_1, x_2, ..., x_d)$. Then the gradient

 $\nabla f(x)^{\top} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \vdots \\ \frac{\partial}{\partial x_d} f(x) \end{bmatrix}$

of f is the vector of partial derivatives:

Example:
$$f(x_1, x_2) = (x_1x_2 - y)^2$$



Then

$$\frac{\partial}{\partial x_1} f(x) = 2(x_1 x_2 - y)x_2$$

$$\frac{\partial}{\partial x_2} f(x) = 2x_1(x_1 x_2 - y)$$

Assume we have a function $f(x_1, x_2, ..., x_d)$

and let $x = (x_1, x_2, ..., x_d)$. Then the gradient

$$\nabla f(x)^{\top} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \\ \vdots \\ \frac{\partial}{\partial x_d} f(x) \end{bmatrix}$$

of f is the vector of partial derivatives:

Example:
$$f(x_1, x_2) = (x_1x_2 - y)^2$$



Then

$$\frac{\partial}{\partial x_1} f(x) = 2(x_1 x_2 - y)x_2$$
$$\frac{\partial}{\partial x_2} f(x) = 2x_1(x_1 x_2 - y)$$
$$\frac{\partial}{\partial x_2} f(x) = 2x_1(x_1 x_2 - y)$$

$$\Rightarrow \nabla f(x)^{\mathsf{T}} = 2 \begin{pmatrix} x_1 x_2^2 - y x_2 \\ x_1^2 x_2 - x_1 y \end{pmatrix}$$

We can extend this to functions $f: \mathbb{R}^d \to \mathbb{R}^s$ with multiple outputs via the Jacobian matrix $J_f: \mathbb{R}^d \to \mathbb{R}^{s \times d}$ defined as

$$J_{f}(x) := \begin{bmatrix} \frac{\partial f_{1}}{\partial x_{1}} & \cdots & \frac{\partial f_{1}}{\partial x_{d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{s}}{\partial x_{1}} & \cdots & \frac{\partial f_{s}}{\partial x_{d}} \end{bmatrix}$$



We can extend this to functions $f: \mathbb{R}^d \to \mathbb{R}^s$ with multiple outputs via the Jacobian matrix $J_f: \mathbb{R}^d \to \mathbb{R}^{s \times d}$ defined as

$$J_{f}(x) := \begin{bmatrix} \frac{\partial f_{1}}{\partial x_{1}} & \cdots & \frac{\partial f_{1}}{\partial x_{d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{s}}{\partial x_{1}} & \cdots & \frac{\partial f_{s}}{\partial x_{d}} \end{bmatrix}$$

And we can even define a second-order derivative matrix (known as Hessian) via

$$H_f(x) := J_{\nabla f}(x)$$

PROBABILITY & STATISTICS

MTH786U/P 2020/2021

Assume we have a random variable X with a finite no. of outcomes $x_1, x_2, ..., x_s$ and probabilities $\rho_1 = P(X = x_1), \rho_2 = P(X = x_2), ..., \rho_s = P(X = x_s)$. The expectation of X is defined as

$$\mathbb{E}_{x}[x_{i}] := \sum_{i=1}^{S} x_{i} \rho_{i}$$



Assume we have a random variable X with a finite no. of outcomes $x_1, x_2, ..., x_s$ and probabilities $\rho_1 = P(X = x_1), \rho_2 = P(X = x_2), ..., \rho_s = P(X = x_s)$. The expectation of X is defined as

$$\mathbb{E}_{x}[x_{i}] := \sum_{i=1}^{S} x_{i} \rho_{i}$$

Example: s = 3, $x_1 = 1$, $x_2 = 11/10$, $x_3 = 1/2$ and $\rho_1 = 1/2$, $\rho_2 = 1/3$, $\rho_3 = 1/6$



Assume we have a random variable X with a finite no. of outcomes $x_1, x_2, ..., x_s$ and probabilities $\rho_1 = P(X = x_1), \rho_2 = P(X = x_2), ..., \rho_s = P(X = x_s)$. The expectation of X is defined as

$$\mathbb{E}_{x}[x_{i}] := \sum_{i=1}^{S} x_{i} \rho_{i}$$

Example: s=3, $x_1=1$, $x_2=11/10$, $x_3=1/2$ and $\rho_1=1/2$, $\rho_2=1/3$, $\rho_3=1/6$



$$\Longrightarrow \mathbb{E}_{x}[x_{i}] = \sum_{i=1}^{3} x_{i} \rho_{i} = \frac{1}{2} + \frac{11}{30} + \frac{1}{12} = \frac{19}{20} = 0.95$$

(weighted average)

Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$



Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$

Example: uniform random variable X in [a,b] with $\rho(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$



Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$

Example: uniform random variable X in [a,b] with $\rho(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$

$$\Longrightarrow \quad \mathbb{E}_{x}[x] = \int_{\mathbb{R}} x \, \rho(x) \, dx = \frac{1}{b-a} \int_{a}^{b} x \, dx$$

Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$

Example: uniform random variable X in [a,b] with $\rho(x)=\begin{cases} \frac{1}{b-a} & x\in[a,b]\\ 0 & \text{otherwise} \end{cases}$

$$\Longrightarrow \quad \mathbb{E}_{x}[x] = \int_{\mathbb{R}} x \, \rho(x) \, dx = \frac{1}{b-a} \int_{a}^{b} x \, dx = \frac{b^2 - a^2}{2(b-a)}$$

Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$

Example: uniform random variable X in [a,b] with $\rho(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$

$$\implies \mathbb{E}_{x}[x] = \int_{\mathbb{R}} x \, \rho(x) \, dx = \frac{1}{b-a} \int_{a}^{b} x \, dx = \frac{b^2 - a^2}{2(b-a)} = \frac{(b-a)(b+a)}{2(b-a)}$$

Assume we have an absolutely continuous random variable X with probability density function ρ . The expectation of X is defined as

$$\mathbb{E}_{x}[x] := \int_{\mathbb{R}} x \, \rho(x) \, dx$$

Example: uniform random variable X in [a,b] with $\rho(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$

$$\implies \mathbb{E}_{x}[x] = \int_{\mathbb{R}} x \, \rho(x) \, dx = \frac{1}{b-a} \int_{a}^{b} x \, dx = \frac{b^2 - a^2}{2(b-a)} = \frac{(b-a)(b+a)}{2(b-a)} = \frac{a+b}{2}$$

The variance of a random variable X is defined as

$$\operatorname{Var}_{x}[x] := \mathbb{E}_{x} \left[\left(x - \mathbb{E}_{x}[x] \right)^{2} \right]$$
$$= \mathbb{E}_{x}[x^{2}] - \mathbb{E}_{x}[x]^{2}$$



The variance of a random variable X is defined as

$$Var_x[x] := \mathbb{E}_x \left[\left(x - \mathbb{E}_x[x] \right)^2 \right]$$
$$= \mathbb{E}_x[x^2] - \mathbb{E}_x[x]^2$$

Its square-root



Two random variables X and Y are independent if their joint PDF factors, i.e.

$$\rho(x, y) = \rho_X(x) \rho_Y(y)$$



Two random variables X and Y are independent if their joint PDF factors, i.e.

$$\rho(x, y) = \rho_X(x) \rho_Y(y)$$

An arbitrary no. of n random variables $\{X_i\}_{i=1}^n$ is independent if

$$\rho(x_1, ..., x_n) = \prod_{i=1}^n \rho_{X_i}(x_i)$$



Two random variables X and Y are independent if their joint PDF factors, i.e.

$$\rho(x, y) = \rho_X(x) \rho_Y(y)$$

An arbitrary no. of n random variables $\{X_i\}_{i=1}^n$ is independent if

$$\rho(x_1, ..., x_n) = \prod_{i=1}^n \rho_{X_i}(x_i)$$

The collection of random variables is independent and identically distributed (i.i.d.) if in addition we have

$$\rho_{X_1} = \rho_{X_2} = \cdots = \rho_{X_n}$$