

inp. $linear\ space \rightarrow u, s, v^T$

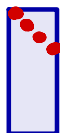
Applications / Principal component analysis

собств. вектора/числа

сингулярные вектора/числа

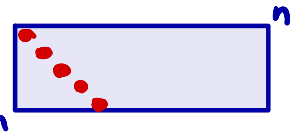
Intuition

$$U^T U = I \\ V^T V = I$$



$$A = U \Sigma V^T$$

$m \times n$ $m \times m$ $m \times n$ $n \times n$



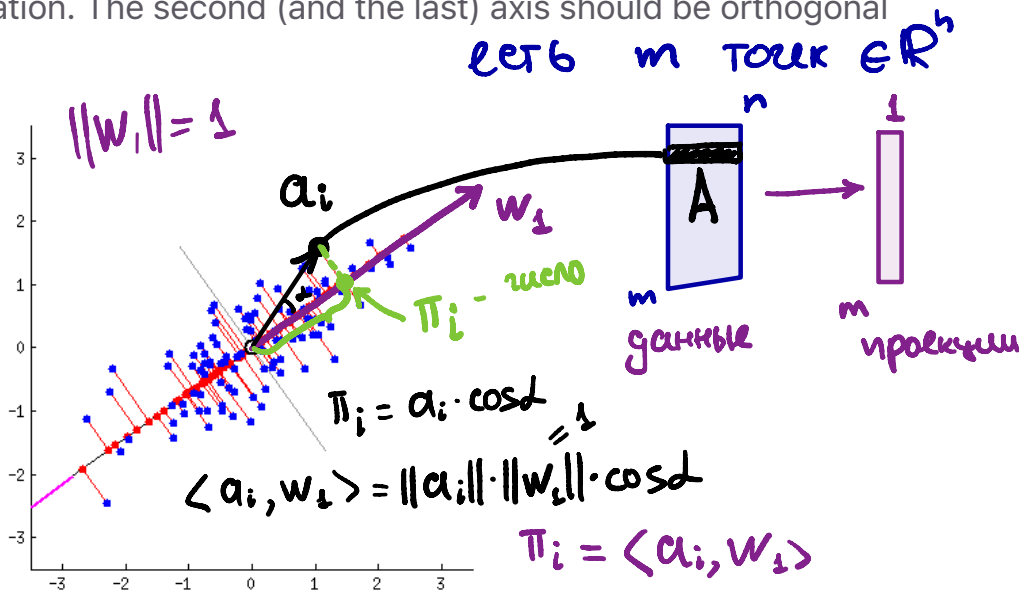
Imagine, that you have a dataset of points. Your goal is to choose orthogonal axes, that describe your data the most informative way. To be precise, we choose first axis in such a way, that maximize the variance (expressiveness) of the projected data. All the following axes have to be orthogonal to the previously chosen ones, while satisfy largest possible variance of the projections.

Let's take a look at the simple 2d data. We have a set of blue points on the plane. We can easily see that the projections on the first axis (red dots) have maximum variance at the final position of the animation. The second (and the last) axis should be orthogonal to the previous one.

задача: понизить размерность
(в некотором смысле оптимально)

выбрав такие оси, проецируя на которые

MAX дисперсию (вариативность)



есть m точек $\in \mathbb{R}^n$

$$\|w_1\| = 1$$

a_i

w_1

π_i - значение

$$\pi_i = a_i \cdot \cos \alpha$$

$$\langle a_i, w_1 \rangle = \|a_i\| \cdot \|w_1\| \cdot \cos \alpha$$

$$\pi_i = \langle a_i, w_1 \rangle$$

source

This idea could be used in a variety of ways. For example, it might happen, that projection of complex data on the principal plane (only 2 components) bring you enough intuition for clustering. The picture below plots projection of the labeled dataset onto the first to principal components (PCs), we can clearly see, that only two vectors (these PCs) would be enough to differ Finnish people from Italian in particular dataset (celiac disease (Dubois et al. 2010))

Нормировка: Убедиться, что $\sum_i a_i = 0$
если это не так

$$\sum_i a_i = 0$$

Problem

$$A_{std} \rightarrow \begin{matrix} \mu \\ A \end{matrix} - \mu$$

$$\mu = \frac{1}{n} \sum_{i=1}^n a_i \in \mathbb{R}^n$$

The first component should be defined in order to maximize variance. Suppose, we've already normalized the data, i.e. $\sum_i a_i = 0$, then sample variance will become the sum

т.к. главные компоненты, то $\mathbb{E} \Pi_i = 0$

of all squared projections of data points to our vector $w_{(1)}$, which implies the following optimization problem:

Выб. луч. $\Pi_i = \mathbb{E}(\Pi_i)^2 - (\mathbb{E} \Pi_i)^2 \rightarrow 0$

$a_i^T \cdot w = (Aw)_i$
 $m \times n \cdot n \times 1 = m \times 1$

R - CB

$Var(R) = \mathbb{E}R^2 - (\mathbb{E}R)^2$

$w_{(1)} = \arg \max_{\|w\|=1} \left\{ \sum_i (a_{(i)}^T \cdot w)^2 \right\}$

or

$w_{(1)} = \arg \max_{\|w\|=1} \{ \|Aw\|^2 \} = \arg \max_{\|w\|=1} \{ w^T A^T A w \}$

$\|Aw\|^2 = \langle Aw, Aw \rangle = (Aw)^T \cdot Aw = w^T \underbrace{A^T A}_{n \times n} w$

since we are looking for the unit vector, we can reformulate the problem:

ОТВЕТ:

$w_1 =$ соев. вектор $A^T A$

$w_{(1)} = \arg \max \left\{ \frac{w^T A^T A w}{w^T w} \right\}$

$\lambda_{\min} \leq \frac{w^T X w}{w^T w} \leq \lambda_{\max}$

It is known, that for positive semidefinite matrix $A^T A$ such vector is nothing else, but eigenvector of $A^T A$, which corresponds to the largest eigenvalue. The following components will give you the same results (eigenvectors).

So, we can conclude, that the following mapping:

$\lambda(A^T A) = \sigma^2(A)$

проекция

$\Pi = A \cdot W$

PC

соев. век
 $Xw = \lambda w / w^T$
 $w^T X w = \lambda \cdot w^T w$
 $\lambda = \frac{w^T X w}{w^T w}$

describes the projection of data onto the k principal components, where W contains first (by the size of eigenvalues) k eigenvectors of $A^T A$.

Now we'll briefly derive how SVD decomposition could lead us to the PCA.

Firstly, we write down SVD decomposition of our matrix:

$A = U \Sigma W^T$

$m \times n$ $n \times n$ $n \times n$
 данные Σ U

and to its transpose:

$A^T = (U \Sigma W^T)^T$
 $= (W^T)^T \Sigma^T U^T$
 $= W \Sigma^T U^T$
 $= W \Sigma U^T$

$A^T A = W \Sigma U^T \cdot U \Sigma W^T = W \underbrace{\Sigma^T \Sigma}_I W^T =$ соев. вектор
 $= W \Sigma^2 W^T$ спектральное разложение

Then, consider matrix AA^T :

PCA: соев. векторы матрицы

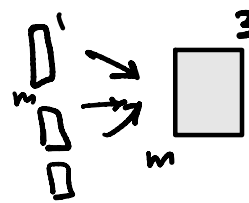
$A^T A$ - оптимальные главные компоненты

w_1, w_2, w_3

$A \cdot w_1$

$A \cdot w_2$

$A \cdot w_3$



$$\begin{aligned}
 A^T A &= (W \Sigma U^T)(U \Sigma V^T) \\
 &= W \Sigma I \Sigma W^T \\
 &= W \Sigma \Sigma W^T \\
 &= W \Sigma^2 W^T
 \end{aligned}$$

Which corresponds to the eigendecomposition of matrix $A^T A$, where W stands for the matrix of eigenvectors of $A^T A$, while Σ^2 contains eigenvalues of $A^T A$.

At the end:

$$\begin{aligned}
 &U \Sigma W^T \\
 \Pi &= A \cdot W = \\
 &= \underline{U \Sigma W^T} W = U \Sigma
 \end{aligned}$$

- ① $A = U \Sigma W^T$
- ② проекции $\Pi_r = U_r \cdot \Sigma_r$
- ③

The latter formula provide us with easy way to compute PCA via SVD with any number of principal components:

$$\Pi_r = U_r \Sigma_r$$

↓ проекции по r PC

Examples

🌻 Iris dataset

Consider the classical Iris dataset

матрица 10000×10

$A \in \mathbb{R}$

хотим взять $r = 2$

PC 1, PC 2

w_1, w_2

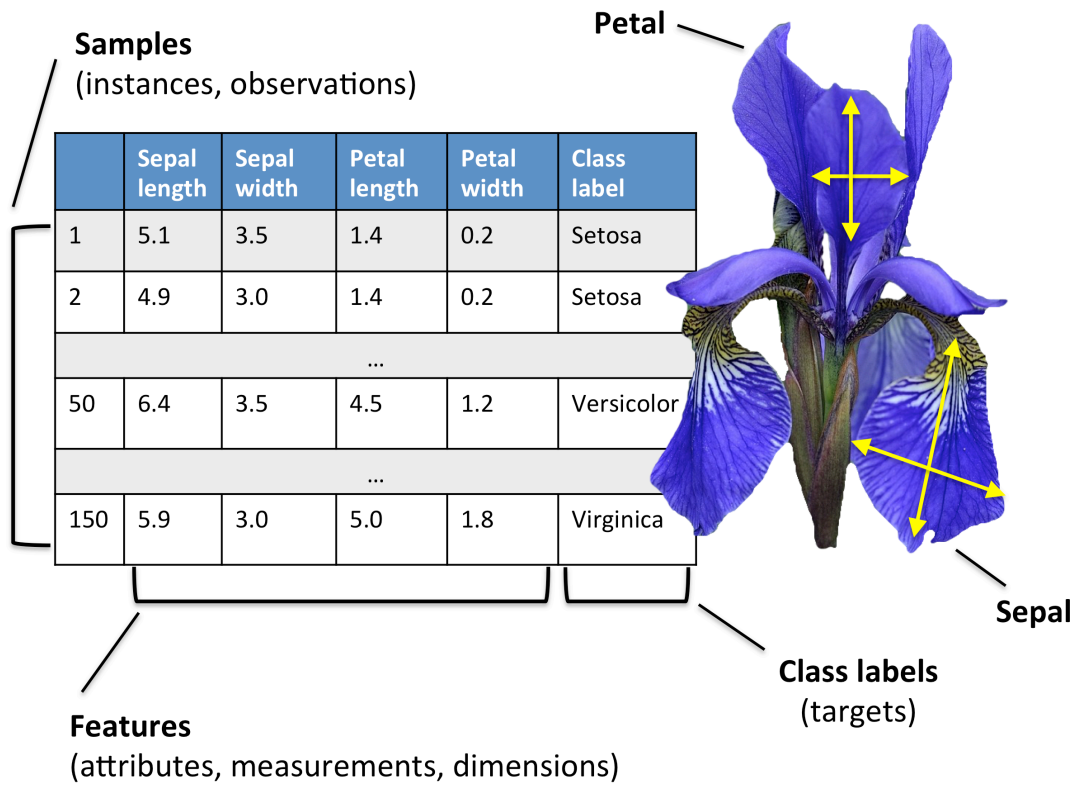
Π_1, Π_2 для каждого C_i

сжатое представление $\begin{matrix} \Pi_1, \Pi_2 \\ \vdots \\ \vdots \end{matrix}$

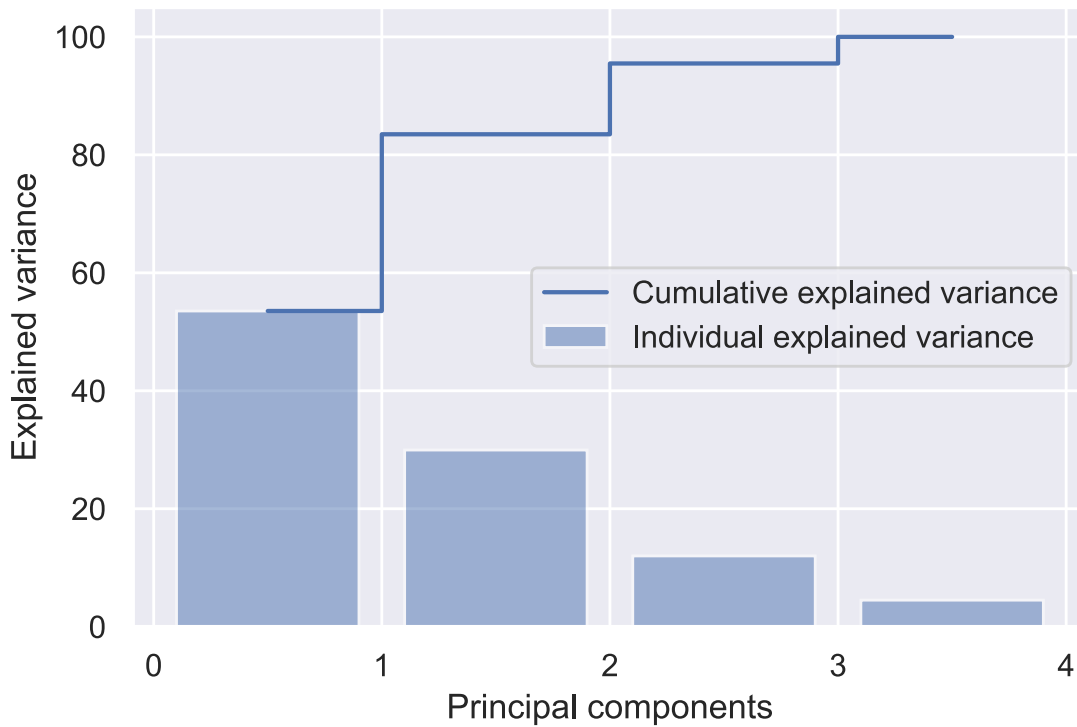
$\tilde{A} = \Pi \cdot W$
 $m \times 2 \times n$

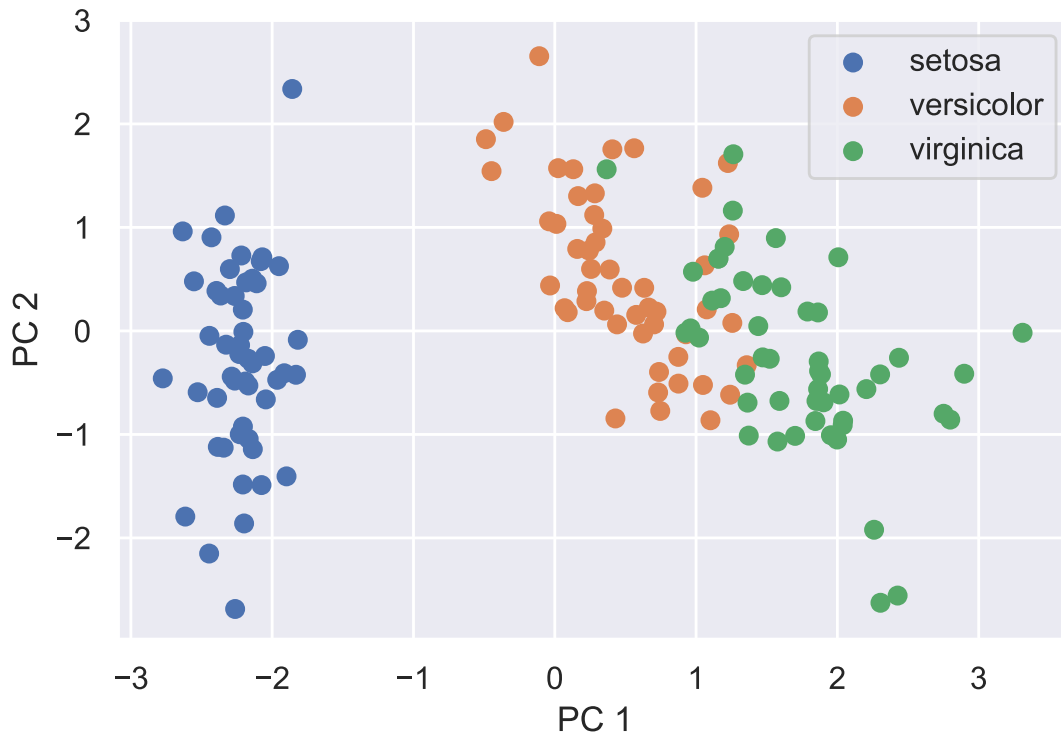
$\tilde{\alpha}_j = \Pi_1^j \cdot w_1 + \Pi_2^j \cdot w_2$

$\tilde{A} \xrightarrow{r \rightarrow \min(m,n)} A$



source We have the dataset matrix $A \in \mathbb{R}^{150 \times 4}$





Code

 [Open in Colab](#)

Related materials

- [Wikipedia](#)
- [Blog post](#)
- [Blog post](#)

Useful definitions and notations

We will treat all vectors as column vectors by default. The space of real vectors of length n is denoted by \mathbb{R}^n , while the space of real-valued $m \times n$ matrices is denoted by $\mathbb{R}^{m \times n}$.

Basic linear algebra background

The standard **inner product** between vectors x and y from \mathbb{R}^n is given by

$$\langle x, y \rangle = x^\top y = \sum_{i=1}^n x_i y_i = y^\top x = \langle y, x \rangle$$

Here x_i and y_i are the scalar i -th components of corresponding vectors.

The standard **inner product** between matrices X and Y from $\mathbb{R}^{m \times n}$ is given by

$$\langle X, Y \rangle = \text{tr}(X^\top Y) = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij} = \text{tr}(Y^\top X) = \langle Y, X \rangle$$

The determinant and trace can be expressed in terms of the eigenvalues

$$\det A = \prod_{i=1}^n \lambda_i, \quad \text{tr} A = \sum_{i=1}^n \lambda_i$$

Don't forget about the cyclic property of a trace for a square matrices A, B, C, D :

$$\text{tr}(ABCD) = \text{tr}(DABC) = \text{tr}(CDAB) = \text{tr}(BCDA)$$

The largest and smallest eigenvalues satisfy

$$\lambda_{\min}(A) = \inf_{x \neq 0} \frac{x^\top A x}{x^\top x}, \quad \lambda_{\max}(A) = \sup_{x \neq 0} \frac{x^\top A x}{x^\top x}$$

and consequently $\forall x \in \mathbb{R}^n$ (Rayleigh quotient):

$$\lambda_{\min}(A) x^\top x \leq x^\top A x \leq \lambda_{\max}(A) x^\top x$$

A matrix $A \in \mathbb{S}^n$ (set of square symmetric matrices of dimension n) is called **positive (semi)definite** if for all $x \neq 0$ (for all x): $x^\top A x > (\geq) 0$. We denote this as

$$A \succ (\succeq) 0.$$

$$f(x) = x^T A x = 100x_1^2 + x_2^2$$

$$A_2 = \begin{pmatrix} 100 & 0 \\ 0 & 1 \end{pmatrix} \kappa(A_2) = \frac{100}{1}$$

$$A_2 = \begin{pmatrix} 1.5 & 0 \\ 0 & 1 \end{pmatrix} \kappa(A_2) = 1.5$$

The **condition number** of a nonsingular matrix is defined as

$$A_3 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \kappa(A_3) = 1$$

$$\kappa(A) = \|A\| \|A^{-1}\|$$

$$\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

Matrix and vector multiplication

$$A \succ 0$$

$$\Rightarrow \kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$$

Let A be a matrix of size $m \times n$, and B be a matrix of size $n \times p$, and let the product AB be:

$$C = AB$$

then C is a $m \times p$ matrix, with element (i, j) given by:

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

Let A be a matrix of shape $m \times n$, and x be $n \times 1$ vector, then the i -th component of the product:

$$z = Ax$$

is given by:

$$z_i = \sum_{k=1}^n a_{ik} x_k$$

Finally, just to remind:

- $C = AB \quad C^T = B^T A^T$
- $AB \neq BA$
- $e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k$
- $e^{A+B} \neq e^A e^B$ (but if A and B are commuting matrices, which means that $AB = BA$, $e^{A+B} = e^A e^B$)
- $\langle x, Ay \rangle = \langle A^T x, y \rangle$

Gradient

Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, then vector, which contains all first order partial derivatives:

$$\nabla f(x) = \frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

named gradient of $f(x)$. This vector indicates the direction of steepest ascent. Thus, vector $-\nabla f(x)$ means the direction of the steepest descent of the function in the point. Moreover, the gradient vector is always orthogonal to the contour line in the point.

Hessian

Let $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, then matrix, containing all the second order partial derivatives:

$$f''(x) = \frac{\partial^2 f}{\partial x_i \partial x_j} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

In fact, Hessian could be a tensor in such a way: $(f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m)$ is just 3d tensor, every slice is just hessian of corresponding scalar function $(H(f_1(x)), H(f_2(x)), \dots, H(f_m(x)))$.

Jacobian

The extension of the gradient of multidimensional $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is the following matrix:

$$f'(x) = \frac{df}{dx^T} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

Summary

$$f(x) : X \rightarrow Y; \quad \frac{\partial f(x)}{\partial x} \in G$$

X	Y	G	Name
\mathbb{R}	\mathbb{R}	\mathbb{R}	$f'(x)$ (derivative)
\mathbb{R}^n	\mathbb{R}	\mathbb{R}^n	$\frac{\partial f}{\partial x_i}$ (gradient)
\mathbb{R}^n	\mathbb{R}^m	$\mathbb{R}^{m \times n}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$	\mathbb{R}	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$

General concept

Naive approach

The basic idea of naive approach is to reduce matrix/vector derivatives to the well-known scalar derivatives.

Matrix notation of a function

$$f(x) = c^\top x$$



Scalar notation of a function

$$f(x) = \sum_{i=1}^n c_i x_i$$

Matrix notation of a gradient

$$\nabla f(x) = c$$



$$\frac{\partial f(x)}{\partial x_k} = c_k$$

Simple derivative

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial (\sum_{i=1}^n c_i x_i)}{\partial x_k}$$

One of the most important practical tricks here is to separate indices of sum (i) and

partial derivatives (k). Ignoring this simple rule tends to produce mistakes.

Differential approach

The guru approach implies formulating a set of simple rules, which allows you to calculate derivatives just like in a scalar case. It might be convenient to use the differential notation here.

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Differentials $df = f(x+dx) - f(x) \quad \|dx\| \rightarrow 0$

After obtaining the differential notation of df we can retrieve the gradient using the following formula:

$$df(x) = \langle \nabla f(x), dx \rangle$$

Задача: получить ∇f по df
 1) Построить df
 2) Прислать к df
 3) Это ответ

Then, if we have differential of the above form and we need to calculate the second derivative of the matrix/vector function, we treat "old" dx as the constant dx_1 , then calculate $d(df) = d^2 f(x)$

$$d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx \rangle = \langle H_f(x) dx_1, dx \rangle$$

$df = \langle \dots, dx \rangle$
 3) Это ответ

Properties

Let A and B be the constant matrices, while X and Y are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$

Пример 1 $f(x) = \|x\|^2$ $\nabla f = ?$ $\nabla f \in \mathbb{R}^n$
 $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Решение:

$$\begin{aligned} 1) \quad df &= d(\|x\|^2) = d(\langle x, x \rangle) = \\ &= \langle dx, x \rangle + \langle x, dx \rangle = \langle x, dx \rangle + \langle x, dx \rangle = \\ &= \langle 2x, dx \rangle \\ &\rightarrow \nabla f = 2x \in \mathbb{R}^n \end{aligned}$$

упр. 1 $f(x) = \|x\|_F^2$, $x \in \mathbb{R}^{m \times n}$ $\nabla f = ?$
 $\nabla f \in \mathbb{R}^{m \times n}$

Решение:

$$\begin{aligned} 1) \quad df &= d(\|x\|_F^2) = d(\langle x, x \rangle) = \\ &= \langle dx, x \rangle + \langle x, dx \rangle = 2 \langle x, dx \rangle = \\ &= \langle 2x, dx \rangle \Rightarrow \nabla f = 2 \cdot x \end{aligned}$$

Пример 2 $f(x) = \frac{1}{2} x^T A x - b^T x + c$ $f: \mathbb{R}^n \rightarrow \mathbb{R}$
 $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$
 $\nabla f = ?$

Решение:

$$\begin{aligned} 1) \quad df &= d\left(\frac{1}{2} x^T A x - b^T x + c\right) = \\ &= d\left(\frac{1}{2} \langle x, Ax \rangle - \langle b, x \rangle + c\right) = \\ &= d\left(\frac{1}{2} \langle x, Ax \rangle\right) - d(\langle b, x \rangle) + d(c) = \end{aligned}$$

$x^T y = \langle x, y \rangle$

$$= \frac{1}{2} (\langle dx, Ax \rangle + \langle x, d(Ax) \rangle) - \langle db, x \rangle - \langle dx, b \rangle + 0 =$$

$$= \frac{1}{2} (\langle Ax, dx \rangle + \langle x, Adx \rangle) - \langle b, dx \rangle =$$

$$= \frac{1}{2} (\langle Ax, dx \rangle + \langle A^T x, dx \rangle) - \langle b, dx \rangle =$$

$$= \frac{1}{2} \langle Ax + A^T x, dx \rangle - \langle b, dx \rangle =$$

$$= \langle \frac{1}{2}(A+A^T)x - b, dx \rangle$$



∇f

$$\Rightarrow \nabla f = \frac{1}{2}(A+A^T)x - b$$

$n \times 1 \quad n \times n \quad n \times 1 \quad n \times 1$

если $A \geq 0 \Rightarrow A = A^T$

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}^T = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\nabla f = Ax - b$$

задача мин. регрессии в ℓ_2 пространстве

Упр 2

$$f(x) = \|Ax - b\|^2 + \frac{\lambda}{2} \|x\|^2 \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$\nabla f = ?$

Решение:

$$1) df = d\left(\|Ax - b\|^2 + \frac{\lambda}{2} \|x\|^2\right) =$$

$$= d(\|Ax - b\|^2) + d\left(\frac{\lambda}{2} \|x\|^2\right) =$$

$$= d(\langle y, y \rangle) =$$

$$= 2\langle y, dy \rangle =$$

$$y = Ax - b$$

$$\downarrow$$

$$dy = d(Ax - b) =$$

$$= Adx$$

$$= \frac{\lambda}{2} d(\|x\|^2) =$$

$$= \frac{\lambda}{2} d(\langle x, x \rangle) =$$

$$= \frac{\lambda}{2} \cdot 2\langle x, dx \rangle =$$

$$= \langle \lambda x, dx \rangle$$

$$= \langle 2(Ax-b), Adx \rangle =$$

$$= \langle 2A^T(Ax-b), dx \rangle$$

$$\boxed{\nabla f = \lambda x + 2A^T(Ax-b)} \in \mathbb{R}^n$$

(n) (n) 1 (n) (n) (n) (n) 1 (n) 1

Ynp.

$$f(x) = \text{tr}(X)$$

$$\nabla f = ? \in \mathbb{R}^{n \times n}$$

$$1) df = d(\text{tr}(X)) =$$

$$= d(\text{tr}(I \cdot X)) =$$

$$= d(\langle X, \underbrace{I^T}_I \rangle) = \langle dx, I \rangle + \langle X, dI \rangle =$$

$$= \langle I, dx \rangle \Rightarrow \boxed{\nabla f = I}$$

$$\langle X, Y \rangle = \text{tr}(X^T Y) \\ = \text{tr}(Y^T X)$$

Мы не умеем считать ∇f . Как считать f'' ?

$$\boxed{df = \langle \nabla f, dx \rangle}$$

КАК СЧИТАТЬ f'' ? $\nabla^2 f$

$$1) df = \langle \nabla f, dx \rangle$$

$$dx := dx_1 \quad \leftarrow \begin{array}{l} \text{считаем} \\ dx_i = \text{const} \end{array}$$

$$2) d(df) = d^2 f = d(\langle \nabla f, dx_1 \rangle) =$$

$$= \langle d(\nabla f), dx_1 \rangle$$

$$3) \text{Проберем к бугру: } d^2 f = \langle \underbrace{\quad}_{f''} dx_1, dx \rangle$$

Пример:

$$f(x) = \frac{1}{2} x^T A x - b^T x + c$$

1)

$$df = \left\langle \frac{1}{2} (A + A^T) x - b, dx \right\rangle$$

$$dx := dx_1$$

2) Проверим $d^2 f = d \left(\left\langle \frac{1}{2} (A + A^T) x - b, dx_1 \right\rangle \right) =$

$$= \left\langle d \left(\frac{1}{2} (A + A^T) x - b \right), dx_1 \right\rangle =$$

$$db = 0$$

$$= \left\langle \frac{1}{2} d((A + A^T) x), dx_1 \right\rangle =$$

$$= \left\langle \frac{1}{2} (A + A^T) dx, dx_1 \right\rangle =$$

МАТРИЦА

dx, dx_1

$$f'' = \frac{1}{2} (A + A^T)$$

$$= \left\langle dx, \frac{1}{2} (A + A^T) dx_1 \right\rangle =$$

$$(A + A^T)^T = A^T + A = A + A^T$$

$$= \left\langle \frac{1}{2} (A + A^T) dx_1, dx \right\rangle$$

Упр.

$$f(x) = \|Ax - b\|^2 + \frac{\lambda}{2} \|x\|^2$$

$$f'' = ?$$

Решение:

$$df = \left\langle 2A^T(Ax - b) + \lambda x, dx \right\rangle$$

матрица $dx^T dx$

$$d(df) = d^2 f =$$

$$d^2 f = \left\langle \dots dx, dx_1 \right\rangle$$

$$= \left\langle d(2A^T(Ax - b) + \lambda x), dx_1 \right\rangle =$$

матрица f''

\ominus

$$\begin{aligned}
 2A^T \cdot d(Ax - b) &= & d(\lambda x) &= \lambda \cdot dx \\
 &= 2A^T \cdot (d(Ax) - 0) = \\
 &= 2A^T \cdot A dx
 \end{aligned}$$

$$\begin{aligned}
 \Leftrightarrow \langle 2A^T A dx + \lambda dx, dx \rangle &= & (A^T A)^T &= \\
 &= \langle (2A^T A + \lambda I) dx, dx \rangle & I \cdot v = v &= A^T A^T = \\
 & & &= A^T A
 \end{aligned}$$

$\lambda dx = \lambda \cdot I \cdot dx$
 ~~$\lambda = \lambda I$~~

$f'' = 2A^T A + \lambda I$

$$\|dx\| \rightarrow 0$$

$$\begin{aligned}
 d(dx) & \quad dx = 1e-5 \\
 (dx)^2 &= 1e-10
 \end{aligned}$$

$$d(\langle f'(x), dx \rangle) =$$

$$= \langle df', dx \rangle + \langle f', d(dx) \rangle$$

- $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

References

- [Convex Optimization](#) book by S. Boyd and L. Vandenberghe - Appendix A. Mathematical background.
- [Numerical Optimization](#) by J. Nocedal and S. J. Wright. - Background Material.
- [Matrix decompositions Cheat Sheet](#).
- [Good introduction](#)
- [The Matrix Cookbook](#)
- [MSU seminars](#) (Rus.)
- [Online tool](#) for analytic expression of a derivative.
- [Determinant derivative](#)