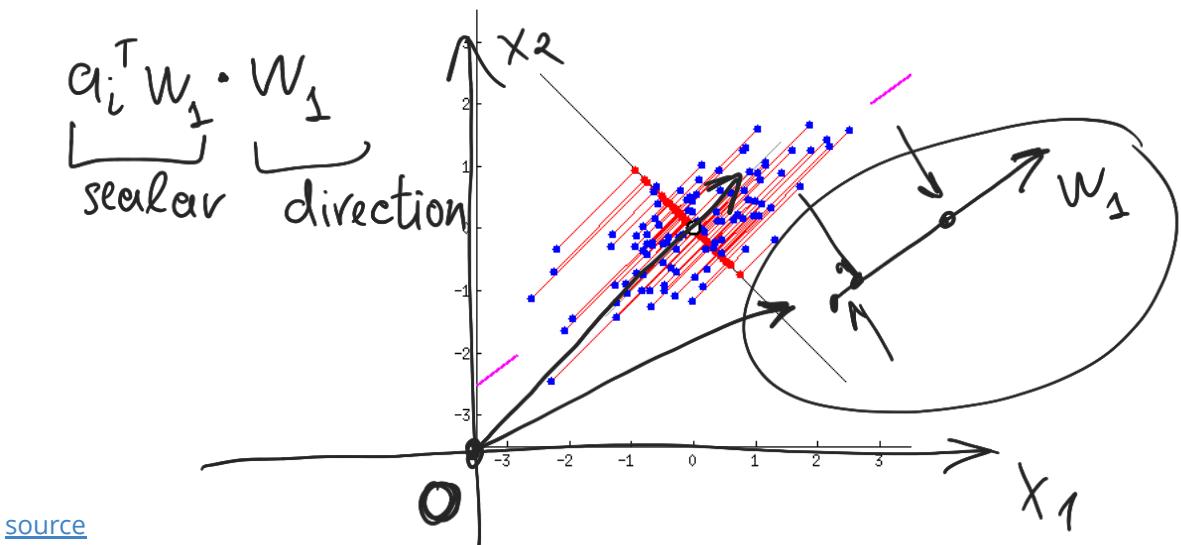


Principal component analysis

Intuition

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.



This idea could be used in a variety 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 two principal components (PC's), we can clearly see, that only two vectors (these PC's) would be enough to differ Finnish people from Italian in particular dataset (celiac disease (Dubois et al. 2010))

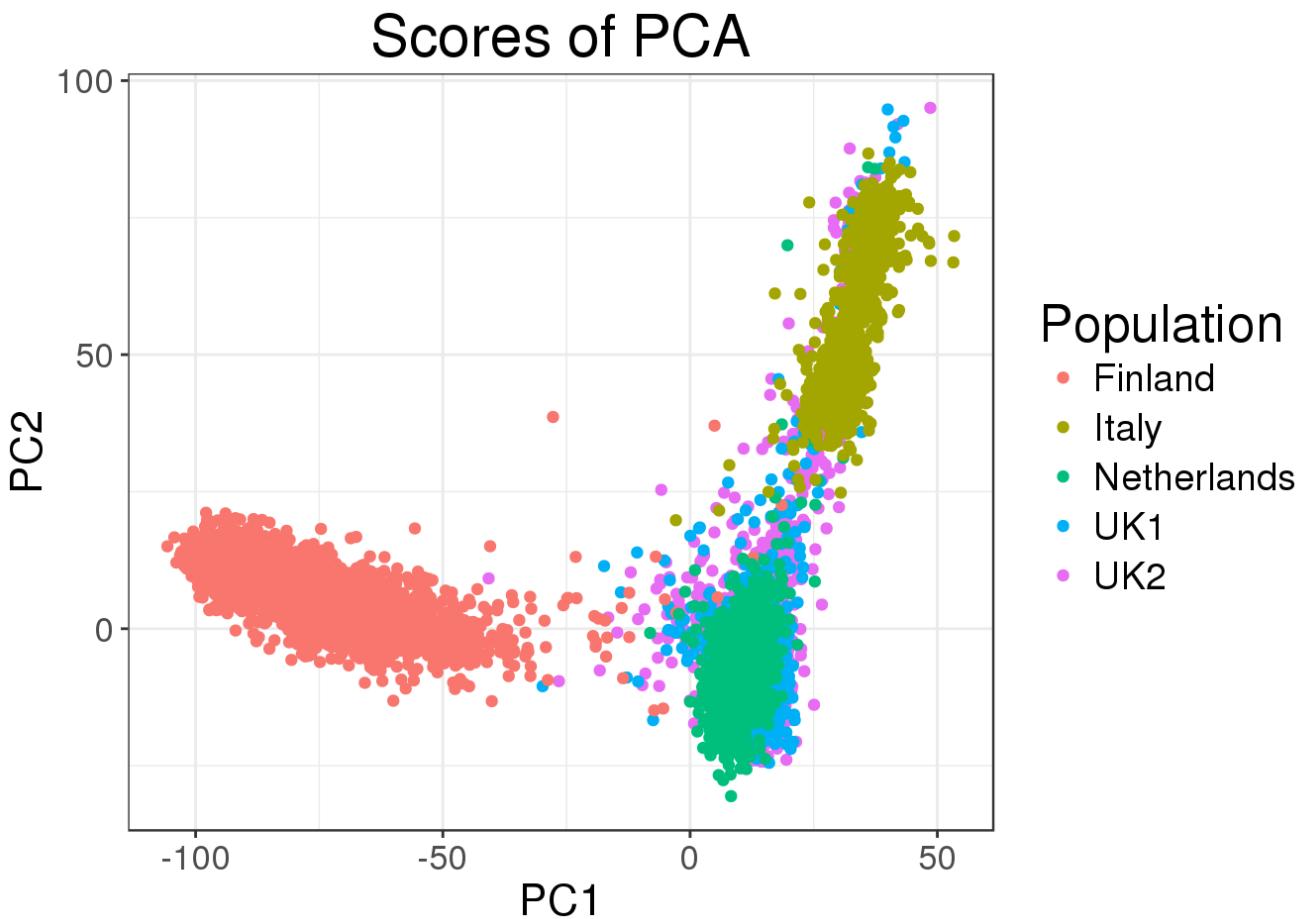
$$f(w) = w^T A^T A w \rightarrow \max$$

$$L(w, \lambda) = w^T A^* A w + \lambda w^T w$$

$$\frac{dL}{dw} = 2 A^* A w + 2\lambda w = 0$$

$$A^* A w = -\lambda w$$

eigen vector
of
 $A^* A$



[source](#)

Problem

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

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{a}_{(i)} \cdot \mathbf{w})^2 \right\}$$

or

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \{ \|\mathbf{Aw}\|^2 \} = \arg \max_{\|\mathbf{w}\|=1} \{ \mathbf{w}^\top \mathbf{A}^\top \mathbf{Aw} \}$$

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

$$\mathbf{w}_{(1)} = \arg \max \left\{ \frac{\mathbf{w}^\top \mathbf{A}^\top \mathbf{Aw}}{\mathbf{w}^\top \mathbf{w}} \right\}$$

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

So, we can conclude, that the following mapping:

$$\underset{n \times k}{\Pi} = \underset{n \times d}{A} \cdot \underset{d \times k}{W}$$

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

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

Firstly, we write down SVD decomposition of our matrix:

$$A = U\Sigma W^\top$$

and to its transpose:

$$\begin{aligned} A^\top &= (U\Sigma W^\top)^\top \\ &= (W^\top)^\top \Sigma^\top U^\top \\ &= W\Sigma^\top U^\top \\ &= W\Sigma U^\top \end{aligned}$$

Then, consider matrix $A^\top A$:

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

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

At the end:

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

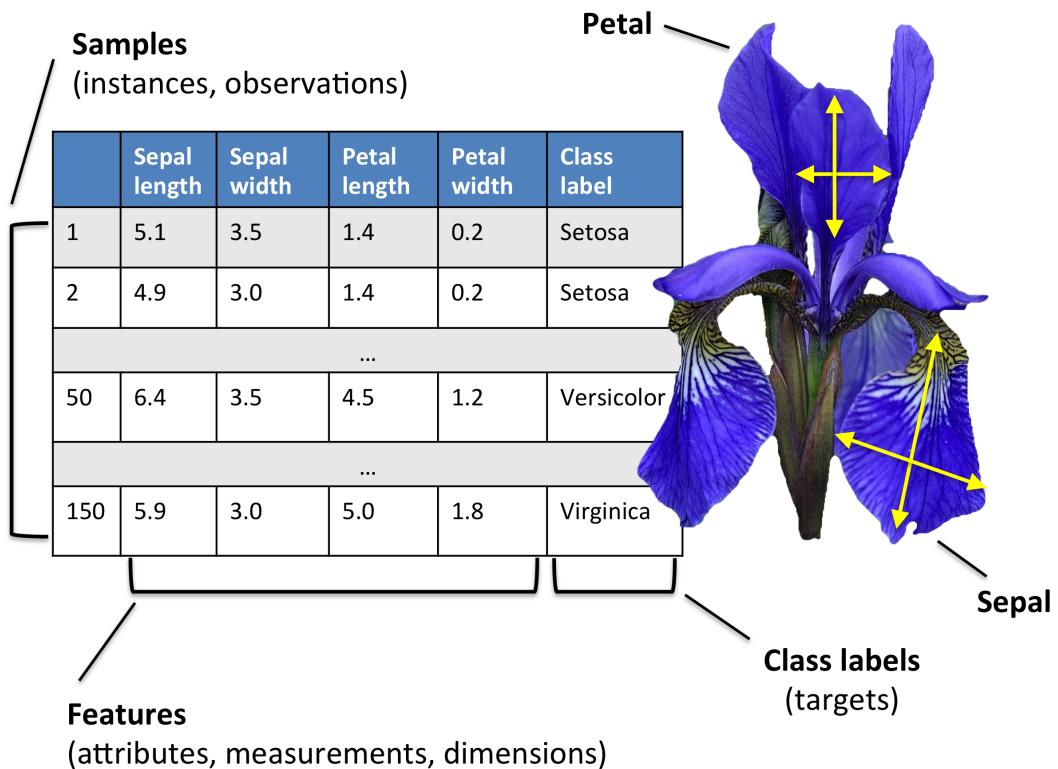
The latter formula provide us with easy way to compute PCA via SVD with any number of principal components:

$$\Pi_r = U_r \Sigma_r$$

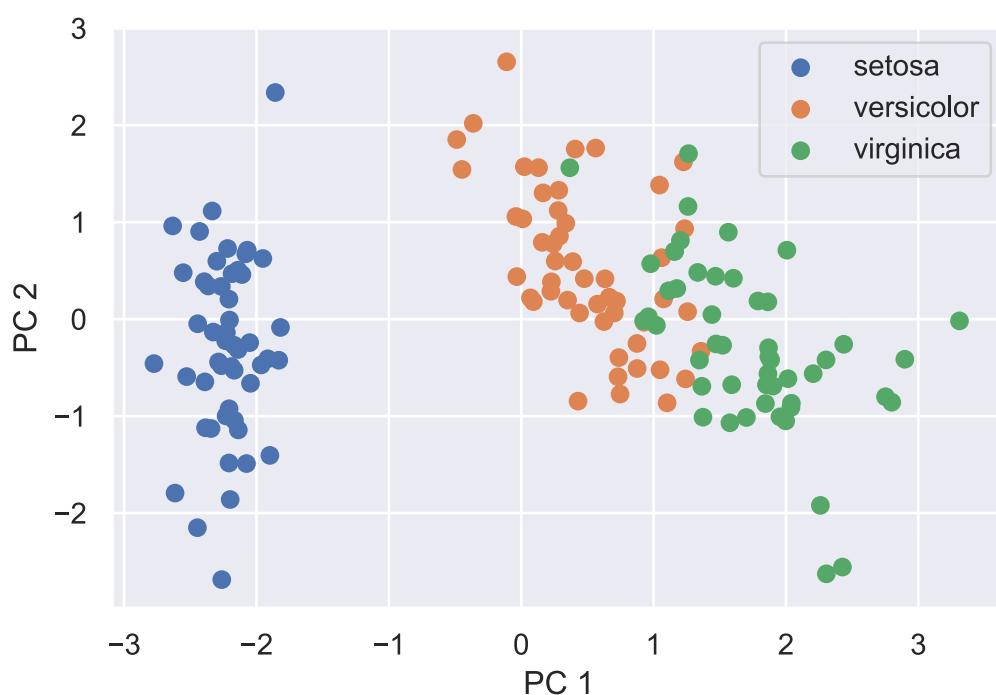
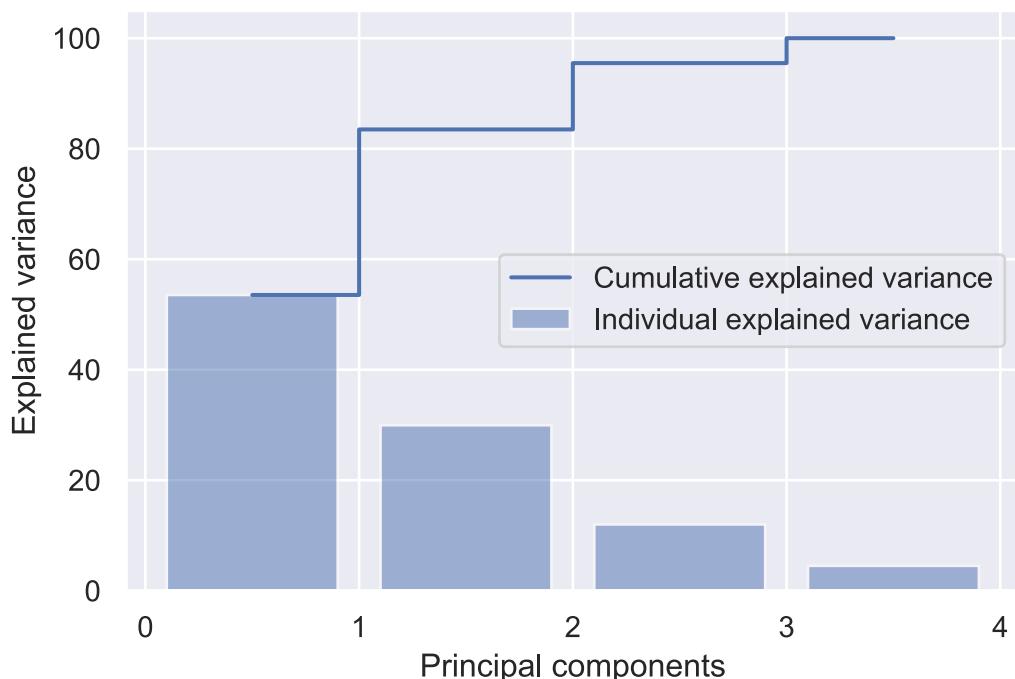
Examples

Iris dataset

Consider the classical Iris dataset



[source](#) We have the dataset matrix $A \in \mathbb{R}^{150 \times 4}$



Code

[Open in Colab](#)

Related materials

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

Matrix calculus

Useful definitions and notations

We will treat all vectors as column vectors by default.

Matrix and vector multiplication

Let A be $m \times n$, and B be $n \times p$, and let the product AB be

$$C = AB$$

$m \times n \times p$

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 $m \times n$, and x be $n \times 1$, then the typical element of the product

$$z = Ax$$

is given by

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

Finally, just to remind:

- $C = AB \quad C^\top = B^\top A^\top$
- $AB \neq BA$
- $e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k \quad \leftarrow$
- $e^{A+B} \neq e^A e^B \quad \leftarrow$

Gradient

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

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

Hessian

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

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

But actually, 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 $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

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

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}^m | \mathbb{R}^n | \mathbb{R}^n | $\frac{\partial f}{\partial x_i}$ (gradient) |
| \mathbb{R}^n | \mathbb{R}^m | $\mathbb{R}^{n \times m}$ | $\frac{\partial f_i}{\partial x_j}$ (jacobian) |
| $\mathbb{R}^{m \times n}$ | \mathbb{R} | $\mathbb{R}^{m \times n}$ | $\frac{\partial f}{\partial x_{ij}}$ |

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

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~~

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

Simple derivative

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

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

One of the most important practical trick here is to separate indicies of sum (i) and partial derivatives (k). Ignoring this simple rule tends to produce mistakes.

Guru 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 convinient to use the differential notation here.

Differentials

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

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

Than, 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 , than calculate $d(df)$

$$d^2 f(x) = \langle \nabla^2 f(x) dx_1, dx_2 \rangle = \langle H_f(x) dx_1, dx_2 \rangle$$

Properties

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

- $dA = 0$
- $d(\alpha X) = \alpha(dx)$
- $d(AXB) = A(dx)B$
- $d(X + Y) = dX + dY$
- $d(X^\top) = (dX)^\top$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$

- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-\top}, dX \rangle$
- $d\text{tr } X = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$

References

- [Good introduction](#)
- [The Matrix Cookbook](#)
- [MSU seminars](#) (Rus.)
- [Online tool](#) for analytic expression of a derivative.

Example 1

Find $\nabla f(x)$, if $f(x) = \frac{1}{2}x^T Ax + b^T x + c$.

$$\frac{\partial f}{\partial x_k} = \frac{\partial \frac{1}{2} \sum_i x_i (Ax)_i}{\partial x_k} + \frac{\partial (\sum_i b_i x_i)}{\partial x_k} + \frac{\partial c}{\partial x_k} =$$

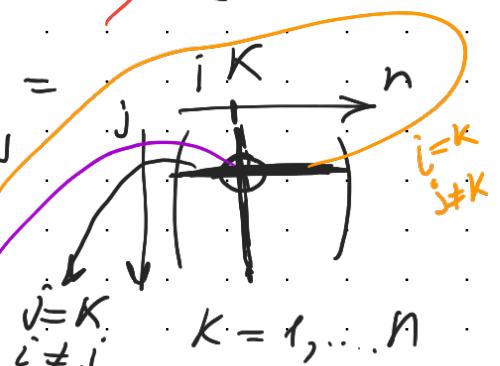
$$= \frac{1}{2} \sum_{ij} \frac{\partial (x_i a_{ij} x_j)}{\partial x_k} + \sum_i b_i \frac{\partial x_i}{\partial x_k} =$$

$$= \frac{1}{2} \sum_t (x_t a_{tk} + x_t a_{kt}) + b_k$$

1) Naive approach:

$$f(x) = \frac{1}{2} \sum_i x_i (Ax)_i + \sum_i b_i x_i + c$$

$$= \frac{1}{2} \sum_{ij} \frac{\partial (x_i a_{ij} x_j)}{\partial x_k} + \sum_i b_i \frac{\partial x_i}{\partial x_k} =$$



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

Example 2

$$\frac{\partial (x_i a_{ik} x_k)}{\partial x_k} = x_i a_{ik}$$

$$\frac{\partial (x_k a_{kj} x_j)}{\partial x_k} = a_{kj} x_j$$

$$\frac{\partial a_{kk} x_k^2}{\partial x_k} = 2 x_k a_{kk}$$

2) Guru's approach

$$f(x) = \frac{1}{2} \langle x, Ax \rangle + \langle b, x \rangle + c$$

$$df = \frac{1}{2} d(\langle x, Ax \rangle) + d(\langle b, x \rangle) + dc =$$

$$= \frac{1}{2} (\langle dx, Ax \rangle + \langle x, A^T dx \rangle) + \langle b, dx \rangle + 0 =$$

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

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

$$df(x) = dx \cdot (x^T A + A^T x) =$$

$$= dx^T (A^T x)$$

Find the gradient $\nabla f(x)$ and hessian $f''(x)$, if $f(x) = \frac{1}{2} \|Ax - b\|_2^2$.

$$1) df = \frac{1}{2} d(\langle Ax - b, Ax - b \rangle) = \frac{1}{2} \cdot 2 \cdot \langle Ax - b, d(Ax - b) \rangle$$

$$= \langle Ax - b, A dx \rangle = \langle A^T(Ax - b), dx \rangle \quad \boxed{\nabla f(x) = A^T(Ax - b)}$$

$$2) g(x) = \langle A^T(Ax - b), dx_1 \rangle$$

$$dg = \langle d(A^T(Ax - b)), dx_1 \rangle = \langle A^T d(Ax - b), dx_1 \rangle =$$

$$= \langle A^T A dx, dx_1 \rangle = \langle dx, (A^T A)^T dx_1 \rangle = \langle A^T A dx, dx \rangle$$

$$d(df) = \langle H_f dx_1, dx_2 \rangle$$

Example 3

$$dx = dx_2$$

$$H_f = A^T A$$

Find $\nabla f(x), f''(x)$, if $f(x) = -e^{-x^T x}$.

$$1) df = d(-e^{-x^T x}) = - \left(e^{-x^T x} \cdot d(\langle x, x \rangle) \right) =$$

$$= e^{-x^T x} \cdot 2 \langle x, dx \rangle \Rightarrow \boxed{\nabla f(x) = e^{-x^T x} \cdot 2x}$$

$$2) g(x) = \langle 2e^{-x^T x} x, dx_1 \rangle$$

$$dg = 2 \langle d(e^{-x^T x} \cdot x), dx_1 \rangle = 2 \langle d(e^{-x^T x}) \cdot x + e^{-x^T x} \cdot dx, dx_1 \rangle =$$

$$= 2 \langle \langle -e^{-x^T x} \cdot 2x, dx \rangle \cdot x + e^{-x^T x} \cdot dx, dx_1 \rangle =$$

$$= 2 \langle -2e^{-x^T x} \langle x, dx \rangle \cdot x + e^{-x^T x} dx, dx_1 \rangle = 2e^{-x^T x} \langle -2 \underbrace{x^T dx}_{} \cdot x + dx, dx_1 \rangle =$$

$$= 2e^{-x^T x} \langle -2x x^T dx + dx, dx_1 \rangle =$$

Example 4

Find the gradient $\nabla f(x)$ and hessian $f''(x)$, if

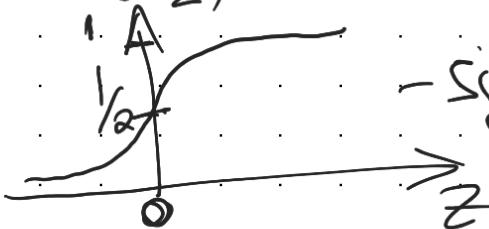
$$f(x) = \ln(1 + \exp(a, x)) = \boxed{2e^{-x^T x} \langle (-2xx^T + I) dx, dx_1 \rangle =}$$

$$\boxed{H_f = 2e^{-x^T x} \cdot (I - 2xx^T)}$$

$$f(x) = \ln(1 + \exp(a, x))$$

$$\cdot df = \frac{d(1 + \exp(a, x))}{1 + \exp(a, x)} = \frac{\exp(a, x) \cdot d(a, x)}{1 + \exp(a, x)}$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\nabla f = \frac{\exp(a, x)}{1 + \exp(a, x)} \cdot a$$

$$\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$$

Example 5

Find $f'(X)$, if $f(X) = \det X$

$$X \in \mathbb{S}^{n \times n}_{++}$$

Note: here under $f'(X)$ assumes first order approximation of $f(X)$ using Taylor series:

$$f(X + \Delta X) \approx f(X) + \text{tr}(f'(X)^T \Delta X) \quad \langle \nabla f(X), \Delta X \rangle$$

$$\begin{aligned} \cdot f(X + dX) &= \det(X + dX) = \det(X(I + X^{-1}dX)) = \det X \cdot \det(I + X^{-1}dX) = \\ &= \det X \cdot \prod_{i=1}^n (1 + \lambda_i^T dX) = \det X \prod_{i=1}^n (1 + \lambda_i) \end{aligned}$$

assume $\|X^{-1}dX\| \rightarrow 0$

$$\begin{aligned} \cdot f(X + dX) - f(X) &= \det X \prod_{i=1}^n (1 + \lambda_i^T) - \det X = \det X \left(\prod_{i=1}^n (1 + \lambda_i^T) - 1 \right) = \\ &\approx \det X \cdot \left(1 + \sum_{i=1}^n \lambda_i^T - 1 \right) = \\ &\approx \det X \sum_{i=1}^n \lambda_i^T = \det X \cdot \text{tr}(X^{-1}dX) \end{aligned}$$

$$f(X + dX) - f(X)$$

$$\langle \nabla f(X), dX \rangle$$

$$\nabla f(X) = \det X \cdot (X^{-1})^T$$

Example 6

$$\langle \det X \cdot X^{-1}, dX \rangle$$

Calculate: $\frac{\partial}{\partial X} \sum \text{eig}(X)$, $\frac{\partial}{\partial X} \prod \text{eig}(X)$, $\frac{\partial}{\partial X} \text{tr}(X)$, $\frac{\partial}{\partial X} \det(X)$

$$\frac{\partial \text{tr}(X)}{\partial X} = \frac{\partial \text{tr}(I \cdot X)}{\partial X} = \frac{\partial \langle X, I \rangle}{\partial X} =$$

$$= I \quad \text{tr}(X) = \sum_i x_{ii}$$

$$\frac{\partial f}{\partial x_{kp}} = \frac{\partial (\sum_i x_{ii})}{\partial x_{kp}} = \begin{cases} 1, & k=p \\ 0, & k \neq p \end{cases}$$

Example 7

Find $\nabla f(X)$, if $f(X) = \langle S, X \rangle - \log \det X$

$(\det X \neq 0)$

$$f(X) = \text{tr}(S^* X) - \log \det X$$

$$df = \langle S, dX \rangle - \frac{d(\det X)}{\det X} = \langle S, dX \rangle - \frac{\langle \det X X^{-1}, dX \rangle}{\det X} =$$

$$\boxed{\nabla f(X) = S - X^{-1}}$$

$$= \langle S, dX \rangle - \langle X^{-1}, dX \rangle =$$

$$= \langle S - X^{-1}, dX \rangle$$

Example 8

$$\langle A, B \rangle = \text{tr}(A^* B) = \text{tr}(B^* A)$$

Calculate the derivatives of the loss function with respect to parameters $\frac{\partial L}{\partial W}, \frac{\partial L}{\partial b}$ for the single object x_i (or, $n = 1$)

Learning

