

pixel.inria.fr

Least squares for programmers

— with color plates —

Dmitry Sokolov

November 15, 2020

Table of Contents

- 1 Maximum likelihood through examples**
- 2 Introduction to systems of linear equations
- 3 Minimization of quadratic functions

Coin toss experiment

We conduct n experiments, two events can happen in each one (“success” or “failure”): one happens with probability p , the other one with probability $1 - p$.

Coin toss experiment

We conduct n experiments, two events can happen in each one (“success” or “failure”): one happens with probability p , the other one with probability $1 - p$.

The probability of getting exactly k successes in these n experiments

$$P(k; n, p) = C_n^k p^k (1 - p)^{n-k}$$

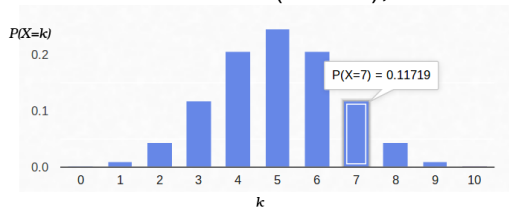
Coin toss experiment

We conduct n experiments, two events can happen in each one (“success” or “failure”): one happens with probability p , the other one with probability $1 - p$.

The probability of getting exactly k successes in these n experiments

$$P(k; n, p) = C_n^k p^k (1 - p)^{n-k}$$

Toss a coin ten times ($n = 10$), count the number of tails:



an ordinary coin ($p = 1/2$)

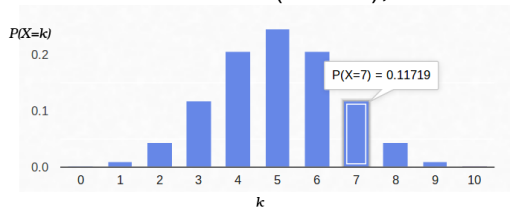
Coin toss experiment

We conduct n experiments, two events can happen in each one (“success” or “failure”): one happens with probability p , the other one with probability $1 - p$.

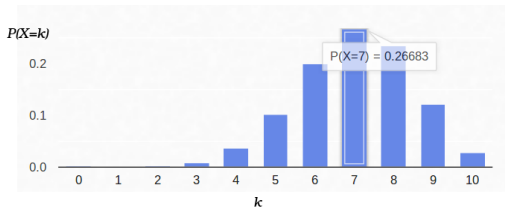
The probability of getting exactly k successes in these n experiments

$$P(k; n, p) = C_n^k p^k (1 - p)^{n-k}$$

Toss a coin ten times ($n = 10$), count the number of tails:



an ordinary coin ($p = 1/2$)



a biased coin ($p = 7/10$)

Coin toss: the likelihood function

Suppose we have a real coin, but we do not know p . However, we can toss it ten times. For example, we have counted seven tails. Would it help us to evaluate p ?

Coin toss: the likelihood function

Suppose we have a real coin, but we do not know p . However, we can toss it ten times. For example, we have counted seven tails. Would it help us to evaluate p ?

Fix $n = 10$ and $k = 7$ in the Bernoulli's formula, leaving p as a free parameter:

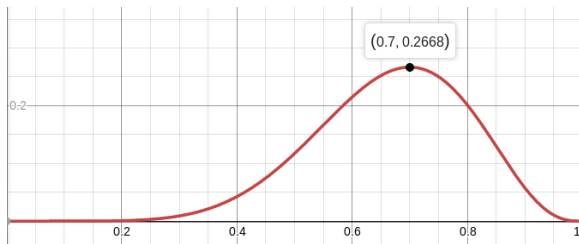
$$\mathcal{L}(p) = C_{10}^7 p^7 (1 - p)^3$$

Coin toss: the likelihood function

Suppose we have a real coin, but we do not know p . However, we can toss it ten times. For example, we have counted seven tails. Would it help us to evaluate p ?

Fix $n = 10$ and $k = 7$ in the Bernoulli's formula, leaving p as a free parameter:

$$\mathcal{L}(p) = C_{10}^7 p^7 (1 - p)^3$$



N.B. the function is continuous!

Coin toss: the maximum likelihood

Let us solve for $\arg \min_p \mathcal{L}(p) = \arg \min_p \log \mathcal{L}(p)$:

Coin toss: the maximum likelihood

Let us solve for $\arg \min_p \mathcal{L}(p) = \arg \min_p \log \mathcal{L}(p)$:

$$\log \mathcal{L}(p) = \log C_{10}^7 + 7 \log p + 3 \log(1 - p)$$

Coin toss: the maximum likelihood

Let us solve for $\arg \min_p \mathcal{L}(p) = \arg \min_p \log \mathcal{L}(p)$:

$$\log \mathcal{L}(p) = \log C_{10}^7 + 7 \log p + 3 \log(1 - p)$$

$$\frac{d \log \mathcal{L}}{dp} = \frac{7}{p} - \frac{3}{1 - p} = 0$$

Coin toss: the maximum likelihood

Let us solve for $\arg \min_p \mathcal{L}(p) = \arg \min_p \log \mathcal{L}(p)$:

$$\log \mathcal{L}(p) = \log C_{10}^7 + 7 \log p + 3 \log(1 - p)$$

$$\frac{d \log \mathcal{L}}{dp} = \frac{7}{p} - \frac{3}{1 - p} = 0$$

That is, the maximum likelihood (about 27%) is reached at the point $p = 7/10$.

Coin toss: the maximum likelihood

Let us solve for $\arg \min_p \mathcal{L}(p) = \arg \min_p \log \mathcal{L}(p)$:

$$\log \mathcal{L}(p) = \log C_{10}^7 + 7 \log p + 3 \log(1 - p)$$

$$\frac{d \log \mathcal{L}}{dp} = \frac{7}{p} - \frac{3}{1 - p} = 0$$

That is, the maximum likelihood (about 27%) is reached at the point $p = 7/10$. Just in case, let us check the second derivative:

$$\frac{d^2 \log \mathcal{L}}{dp^2} = -\frac{7}{p^2} - \frac{3}{(1 - p)^2}$$

At the point $p = 7/10$ it is negative, therefore this point is indeed a maximum of the function \mathcal{L} :

$$\frac{d^2 \log \mathcal{L}}{dp^2}(0.7) \approx -48 < 0$$

Least squares through maximum likelihood

Let us measure a constant value; all measurements are inherently noisy.

For example, if we measure the battery voltage N times, we get N different measurements:

$$\{U_j\}_{j=1}^N$$

Suppose that each measurement U_j is i.i.d. and subject to a Gaussian noise, e.g. it is equal to the real value plus the Gaussian noise. The probability density can be expressed as follows:

$$p(U_j) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(U_j - U)^2}{2\sigma^2}\right),$$

where U is the (unknown) value and σ is the noise amplitude (can be unknown).

Least squares through maximum likelihood

$$\log \mathcal{L}(U, \sigma) = \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right)$$

Least squares through maximum likelihood

$$\begin{aligned}\log \mathcal{L}(U, \sigma) &= \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) \\ &= \sum_{j=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) =\end{aligned}$$

Least squares through maximum likelihood

$$\begin{aligned}\log \mathcal{L}(U, \sigma) &= \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) \\ &= \sum_{j=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) = \sum_{j=1}^N \left(\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(U_j - U)^2}{2\sigma^2} \right)\end{aligned}$$

Least squares through maximum likelihood

$$\begin{aligned}\log \mathcal{L}(U, \sigma) &= \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) \\&= \sum_{j=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(U_j - U)^2}{2\sigma^2} \right) \right) = \sum_{j=1}^N \left(\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(U_j - U)^2}{2\sigma^2} \right) \\&= \underbrace{-N \left(\log \sqrt{2\pi} + \log \sigma \right)}_{\text{does not depend on } \{U_j\}_{j=1}^N} - \frac{1}{2\sigma^2} \sum_{j=1}^N (U_j - U)^2\end{aligned}$$

Under Gaussian noise

$$\arg \max_U \log \mathcal{L} = \arg \min_U \sum_{j=1}^N (U_j - U)^2$$

Least squares through maximum likelihood

$$\frac{\partial \log \mathcal{L}}{\partial U} = \frac{1}{\sigma^2} \sum_{j=1}^N (U_j - U) = 0$$

Least squares through maximum likelihood

$$\frac{\partial \log \mathcal{L}}{\partial U} = \frac{1}{\sigma^2} \sum_{j=1}^N (U_j - U) = 0$$

The most plausible estimation of the unknown value U is the simple average of all measurements:

$$U = \frac{\sum_{j=1}^N U_j}{N}$$

Least squares through maximum likelihood

$$\frac{\partial \log \mathcal{L}}{\partial U} = \frac{1}{\sigma^2} \sum_{j=1}^N (U_j - U) = 0$$

The most plausible estimation of the unknown value U is the simple average of all measurements:

$$U = \frac{\sum_{j=1}^N U_j}{N}$$

And the most plausible estimation of σ turns out to be the standard deviation:

$$\frac{\partial \log \mathcal{L}}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{j=1}^N (U_j - U)^2 = 0$$

Least squares through maximum likelihood

$$\frac{\partial \log \mathcal{L}}{\partial U} = \frac{1}{\sigma^2} \sum_{j=1}^N (U_j - U) = 0$$

The most plausible estimation of the unknown value U is the simple average of all measurements:

$$U = \frac{\sum_{j=1}^N U_j}{N}$$

And the most plausible estimation of σ turns out to be the standard deviation:

$$\frac{\partial \log \mathcal{L}}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{j=1}^N (U_j - U)^2 = 0 \qquad \sigma = \sqrt{\frac{\sum_{j=1}^N (U_j - U)^2}{N}}$$

Such a convoluted way to obtain a simple average of all measurements. . .

Linear regression

It is much harder for less trivial examples. Suppose we have N measurements $\{x_j, y_j\}_{j=1}^N$, and we want to fit a straight line onto it.

Linear regression

It is much harder for less trivial examples. Suppose we have N measurements $\{x_j, y_j\}_{j=1}^N$, and we want to fit a straight line onto it.

$$\begin{aligned}\log \mathcal{L}(a, b, \sigma) &= \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y_j - ax_j - b)^2}{2\sigma^2} \right) \right) = \\ &= \underbrace{-N \left(\log \sqrt{2\pi} + \log \sigma \right)}_{\text{does not depend on } a, b} - \frac{1}{2\sigma^2} \underbrace{\sum_{j=1}^N (y_j - ax_j - b)^2}_{:=S(a,b)}\end{aligned}$$

Linear regression

It is much harder for less trivial examples. Suppose we have N measurements $\{x_j, y_j\}_{j=1}^N$, and we want to fit a straight line onto it.

$$\begin{aligned}\log \mathcal{L}(a, b, \sigma) &= \log \left(\prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y_j - ax_j - b)^2}{2\sigma^2} \right) \right) = \\ &= \underbrace{-N \left(\log \sqrt{2\pi} + \log \sigma \right)}_{\text{does not depend on } a, b} - \frac{1}{2\sigma^2} \underbrace{\sum_{j=1}^N (y_j - ax_j - b)^2}_{:=S(a,b)}\end{aligned}$$

As before, $\arg \max_{a,b} \log \mathcal{L} = \arg \min_{a,b} S(a, b)$.

Linear regression

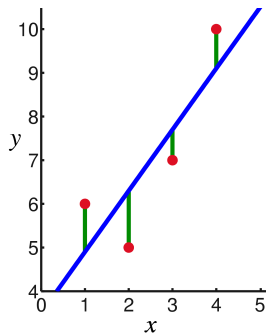
$$S(a, b) := \sum_{j=1}^N (y_j - ax_j - b)^2$$

Linear regression

$$S(a, b) := \sum_{j=1}^N (y_j - ax_j - b)^2$$

$$\frac{\partial S}{\partial a} = \sum_{j=1}^N 2x_j(ax_j + b - y_j) = 0$$

$$\frac{\partial S}{\partial b} = \sum_{j=1}^N 2(ax_j + b - y_j) = 0$$



Linear regression

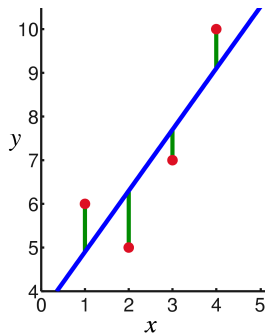
$$S(a, b) := \sum_{j=1}^N (y_j - ax_j - b)^2$$

$$\frac{\partial S}{\partial a} = \sum_{j=1}^N 2x_j(ax_j + b - y_j) = 0$$

$$\frac{\partial S}{\partial b} = \sum_{j=1}^N 2(ax_j + b - y_j) = 0$$

$$a = \frac{N \sum_{j=1}^N x_j y_j - \sum_{j=1}^N x_j \sum_{j=1}^N y_j}{N \sum_{j=1}^N x_j^2 - \left(\sum_{j=1}^N x_j \right)^2}$$

$$b = \frac{1}{N} \left(\sum_{j=1}^N y_j - a \sum_{j=1}^N x_j \right)$$



The takeaway message

The least squares method is a particular case of maximizing likelihood in cases where the probability density is Gaussian.

The more we parameters we have, the more cumbersome the analytical solutions are. Fortunately, we are not living in XVIII century anymore, we have computers!

Next we will try to build a geometric intuition on least squares, and see how can least squares problems be efficiently implemented.

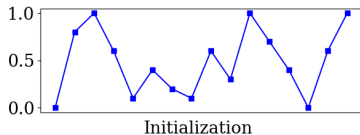
Table of Contents

- 1 Maximum likelihood through examples
- 2 Introduction to systems of linear equations**
- 3 Minimization of quadratic functions

Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

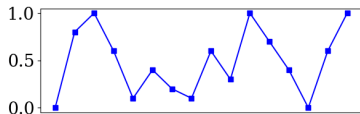
```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



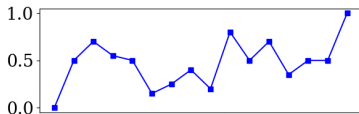
Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



Initialization

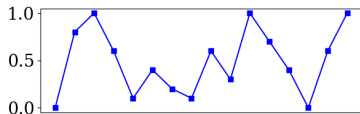


Iteration #1

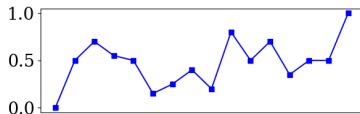
Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

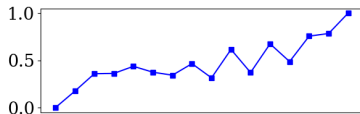
```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



Initialization



Iteration #1

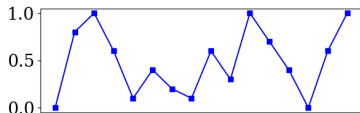


Iteration #5

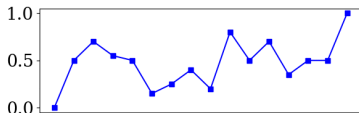
Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

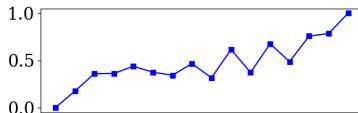
```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



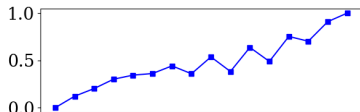
Initialization



Iteration #1



Iteration #5

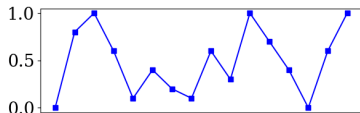


Iteration #10

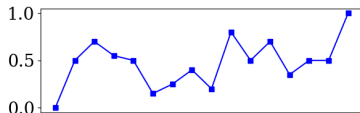
Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

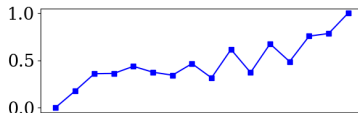
```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



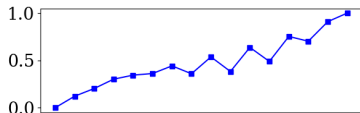
Initialization



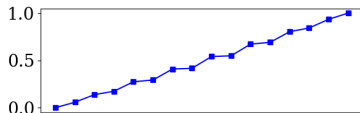
Iteration #1



Iteration #5



Iteration #10

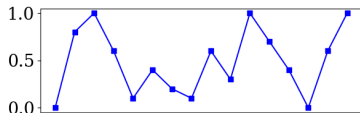


Iteration #50

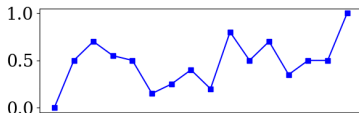
Smooth an array

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

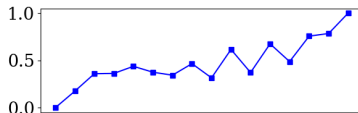
```
for _ in range(512):  
    x = [ x[0] ] + \  
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \  
        [ x[-1] ]
```



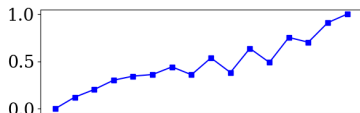
Initialization



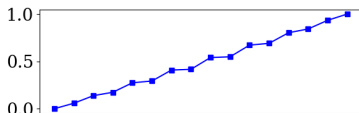
Iteration #1



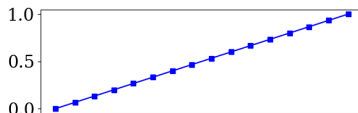
Iteration #5



Iteration #10



Iteration #50



Iteration #500

The Jacobi iterative method

Given an ordinary system of linear equations:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

The Jacobi iterative method

Given an ordinary system of linear equations:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

Let us rewrite it as follows:

$$\begin{aligned} x_1 &= \frac{1}{a_{11}}(b_1 - a_{12}x_2 - a_{13}x_3 - \cdots - a_{1n}x_n) \\ x_2 &= \frac{1}{a_{22}}(b_2 - a_{21}x_1 - a_{23}x_3 - \cdots - a_{2n}x_n) \\ &\vdots \\ x_n &= \frac{1}{a_{nn}}(b_n - a_{n1}x_1 - a_{n2}x_2 - \cdots - a_{n,n-1}x_{n-1}) \end{aligned}$$

The Jacobi iterative method

Let us start with an arbitrary vector $\vec{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})$,

The Jacobi iterative method

Let us start with an arbitrary vector $\vec{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})$,
we can define $\vec{x}^{(1)}$ as follows:

$$\begin{aligned}x_1^{(1)} &= \frac{1}{a_{11}}(b_1 - a_{12}x_2^{(0)} - a_{13}x_3^{(0)} - \dots - a_{1n}x_n^{(0)}) \\x_2^{(1)} &= \frac{1}{a_{22}}(b_2 - a_{21}x_1^{(0)} - a_{23}x_3^{(0)} - \dots - a_{2n}x_n^{(0)}) \\&\vdots \\x_n^{(1)} &= \frac{1}{a_{nn}}(b_n - a_{n1}x_1^{(0)} - a_{n2}x_2^{(0)} - \dots - a_{n,n-1}x_{n-1}^{(0)})\end{aligned}$$

The Jacobi iterative method

Let us start with an arbitrary vector $\vec{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})$,
we can define $\vec{x}^{(1)}$ as follows:

$$\begin{aligned}x_1^{(1)} &= \frac{1}{a_{11}}(b_1 - a_{12}x_2^{(0)} - a_{13}x_3^{(0)} - \dots - a_{1n}x_n^{(0)}) \\x_2^{(1)} &= \frac{1}{a_{22}}(b_2 - a_{21}x_1^{(0)} - a_{23}x_3^{(0)} - \dots - a_{2n}x_n^{(0)}) \\&\vdots \\x_n^{(1)} &= \frac{1}{a_{nn}}(b_n - a_{n1}x_1^{(0)} - a_{n2}x_2^{(0)} - \dots - a_{n,n-1}x_{n-1}^{(0)})\end{aligned}$$

Repeating the process k times, the solution can be approximated by the vector $\vec{x}^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})$.

Back to the array smoothing

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]

for _ in range(512):
    x = [ x[0] ] + \
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \
        [ x[-1] ]
```

$$\left\{ \begin{array}{lcl} x_0 & = & 0 \\ x_1 - x_0 & = & x_2 - x_1 \\ x_2 - x_1 & = & x_3 - x_1 \\ & \vdots & \\ x_{13} - x_{12} & = & x_{14} - x_{13} \\ x_{14} - x_{13} & = & x_{15} - x_{14} \\ x_{15} & = & 1 \end{array} \right.$$

The Gauß-Seidel iterative method

Jacobi:

$$x_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k-1)} \right), \quad \text{for } i = 1, 2, \dots, n$$

The Gauß-Seidel iterative method

Jacobi:

$$x_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k-1)} \right), \quad \text{for } i = 1, 2, \dots, n$$

Gauß-Seidel:

$$x_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} \right), \quad \text{for } i = 1, 2, \dots, n$$

The Gauß-Seidel iterative method

Jacobi:

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]

for _ in range(512):
    x = [ x[0] ] + \
        [ (x[i-1]+x[i+1])/2. for i in range(1, len(x)-1) ] + \
        [ x[-1] ]
```

Gauß-Seidel:

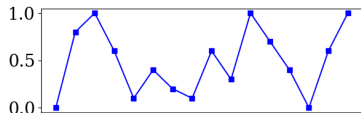
```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]

for _ in range(512):
    for i in range(1, len(x)-1):
        x[i] = ( x[i-1] + x[i+1] )/2.
```

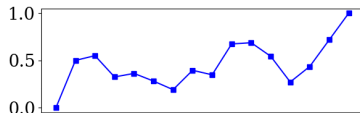
Smooth an array : Gauß-Seidel

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

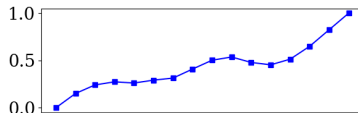
```
for _ in range(512):  
    for i in range(1, len(x)-1):  
        x[i] = ( x[i-1] + x[i+1] )/2.
```



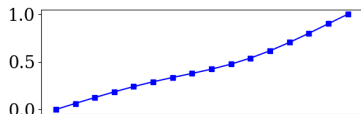
Initialization



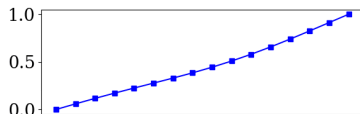
Iteration #1



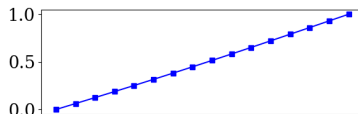
Iteration #3



Iteration #10



Iteration #15



Iteration #50

Equality of derivatives vs zero curvature

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

```
for _ in range(512):  
    for i in range(1, len(x)-1):  
        x[i] = ( x[i-1] + x[i+1] )/2.
```

$$\left\{ \begin{array}{lcl} x_0 & = & 0 \\ x_1 - x_0 & = & x_2 - x_1 \\ x_2 - x_1 & = & x_3 - x_1 \\ & \vdots & \\ x_{13} - x_{12} & = & x_{14} - x_{13} \\ x_{14} - x_{13} & = & x_{15} - x_{14} \\ x_{15} & = & 1 \end{array} \right.$$

Equality of derivatives vs zero curvature

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

```
for _ in range(512):  
    for i in range(1, len(x)-1):  
        x[i] = ( x[i-1] + x[i+1] )/2.
```

$$\left\{ \begin{array}{llllll} x_0 & & & & & & = 0 \\ -x_0 & +2x_1 & -x_2 & & & & = 0 \\ & & -x_2 & +2x_3 & -x_4 & & = 0 \\ & & & & \ddots & & \vdots \\ & & & & & -x_{12} & +2x_{13} & -x_{14} & = 0 \\ & & & & & & -x_{13} & +2x_{14} & -x_{15} & = 0 \\ & & & & & & & & x_{15} & = 1 \end{array} \right.$$

It also works for 3d surfaces

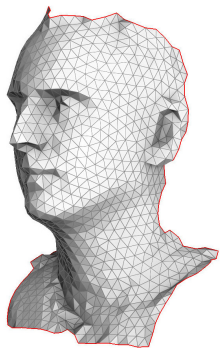
```
#include "model.h"

int main(void) {
    Model m("../input.obj"); // parse the input mesh

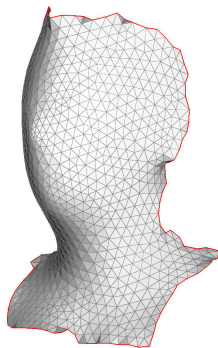
    // smooth the surface through Gauss-Seidel iterations
    for (int it=0; it<1000; it++)
        for (int v=0; v<m.nverts(); v++) // for all vertices
            if (!m.is_boundary_vert(v)) // fix the boundary
                m.point(v) = m.one_ring_barycenter(v);

    std::cout << m; // drop the result
    return 0;
}
```

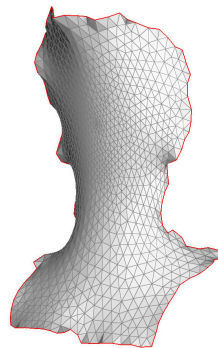
It also works for 3d surfaces



input



10 iterations



1000 iterations

Prescribe the right hand side

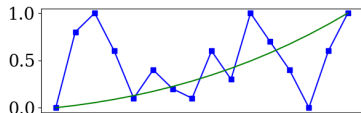
```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

```
for _ in range(512):  
    for i in range(1, len(x)-1):  
        x[i] = ( x[i-1] + x[i+1] - (i+15)/15**3 )/2.
```

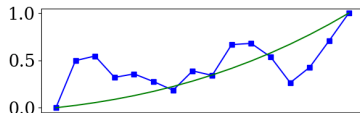
Prescribe the right hand side

```
x = [0, .8, 1, .6, .1, .4, .2, .1, .6, .3, 1, .7, .4, 0, .6, 1]
```

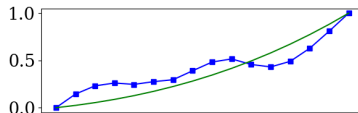
```
for _ in range(512):  
    for i in range(1, len(x)-1):  
        x[i] = ( x[i-1] + x[i+1] - (i+15)/15**3 )/2.
```



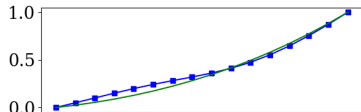
Initialization



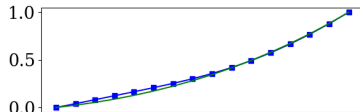
Iteration #1



Iteration #3



Iteration #10



Iteration #15



Iteration #50

Table of Contents

- 1 Maximum likelihood through examples
- 2 Introduction to systems of linear equations
- 3 Minimization of quadratic functions**

Matrices and numbers

What is a number a ?

```
float a;
```

Matrices and numbers

What is a number a ?

```
float a;
```

Is it a function $f(x) = ax : \mathbb{R} \rightarrow \mathbb{R}$?

Matrices and numbers

What is a number a ?

```
float a;
```

Is it a function $f(x) = ax : \mathbb{R} \rightarrow \mathbb{R}$?

```
float f(float x) {  
    return a*x;  
}
```

Matrices and numbers

What is a number a ?

```
float a;
```

Is it a function $f(x) = ax : \mathbb{R} \rightarrow \mathbb{R}$?

```
float f(float x) {  
    return a*x;  
}
```

Or is it $f(x) = ax^2 : \mathbb{R} \rightarrow \mathbb{R}$?

Matrices and numbers

What is a number a ?

```
float a;
```

Is it a function $f(x) = ax : \mathbb{R} \rightarrow \mathbb{R}$?

```
float f(float x) {  
    return a*x;  
}
```

Or is it $f(x) = ax^2 : \mathbb{R} \rightarrow \mathbb{R}$?

```
float f(float x) {  
    return x*a*x;  
}
```

Matrices and numbers

The same goes for matrices, what is a matrix $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$?

Matrices and numbers

The same goes for matrices, what is a matrix $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$?

```
float A[m][n];
```

Matrices and numbers

The same goes for matrices, what is a matrix $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$?

```
float A[m][n];
```

Is it $f(x) = Ax : \mathbb{R}^2 \rightarrow \mathbb{R}^2$?

```
vector<float> f(vector<float> x) {  
    return vector<float>{a11*x[0] + a12*x[1], a21*x[0] + a22*x[1]};  
}
```

Matrices and numbers

The same goes for matrices, what is a matrix $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$?

```
float A[m][n];
```

Is it $f(x) = Ax : \mathbb{R}^2 \rightarrow \mathbb{R}^2$?

```
vector<float> f(vector<float> x) {  
    return vector<float>{a11*x[0] + a12*x[1], a21*x[0] + a22*x[1]};  
}
```

Or $f(x) = x^T Ax = \sum_i \sum_j a_{ij} x_i x_j : \mathbb{R}^2 \rightarrow \mathbb{R}$?

```
float f(vector<float> x) {  
    return x[0]*a11*x[0] + x[0]*a12*x[1] +  
           x[1]*a21*x[0] + x[1]*a22*x[1];  
}
```

What is a positive number?

We have a great tool called the predicate “greater than” $>$.

What is a positive number?

We have a great tool called the predicate “greater than” $>$.

Definition

The real number a is positive if and only if for all non-zero real $x \in \mathbb{R}$, $x \neq 0$ the condition $ax^2 > 0$ is satisfied.

What is a positive number?

We have a great tool called the predicate “greater than” $>$.

Definition

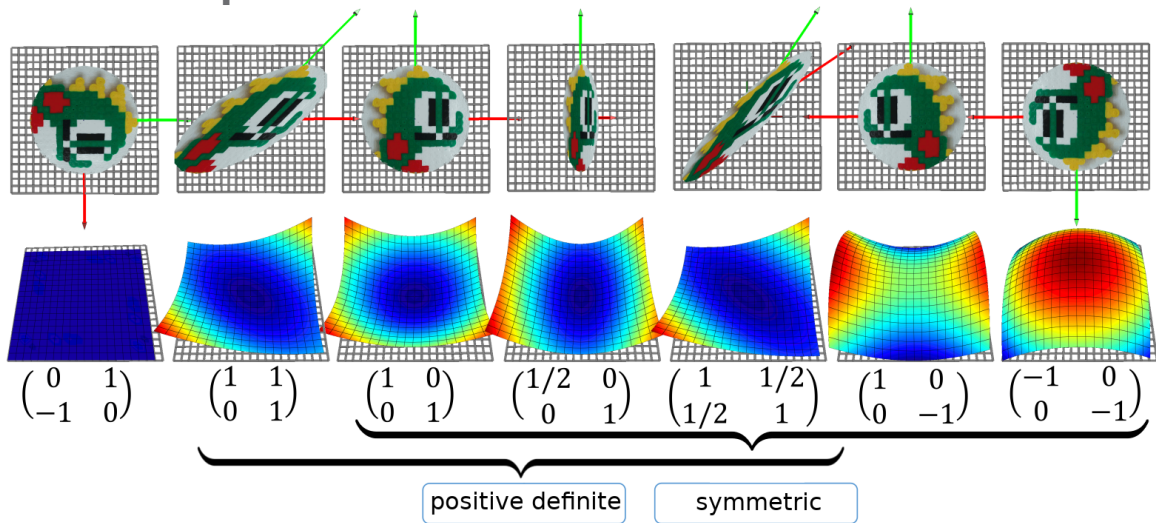
The real number a is positive if and only if for all non-zero real $x \in \mathbb{R}$, $x \neq 0$ the condition $ax^2 > 0$ is satisfied.

This definition looks pretty awkward, but it applies perfectly to matrices:

Definition

The square matrix A is called positive definite if for any non-zero x the condition $x^T A x > 0$ is met, i.e. the corresponding quadratic form is strictly positive everywhere except at the origin.

What is a positive number?



Minimizing a 1d quadratic function

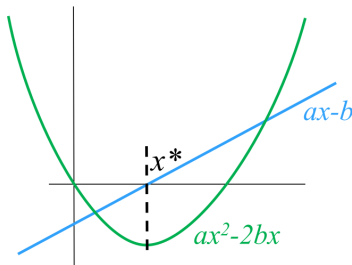
Let us find the minimum of the function $f(x) = ax^2 - 2bx$ (with a positive).

$$\frac{d}{dx}f(x) = 2ax - 2b = 0$$

Minimizing a 1d quadratic function

Let us find the minimum of the function $f(x) = ax^2 - 2bx$ (with a positive).

$$\frac{d}{dx}f(x) = 2ax - 2b = 0$$



In 1d, the solution x^* of the equation $ax - b = 0$ solves the minimization problem $\arg \min_x (ax^2 - 2bx)$ as well.

Differentiating matrix expressions

The first theorem states that 1×1 matrices are invariant w.r.t the transposition:

Theorem

$$x \in \mathbb{R} \Rightarrow x^T = x$$

The proof is left as an exercise.

Differentiating matrix expressions

For a 1d function bx we know that $\frac{d}{dx}(bx) = b$, but what happens in the case of a real function of n variables?

Theorem

$$\nabla b^\top x = \nabla x^\top b = b$$

Differentiating matrix expressions

For a 1d function bx we know that $\frac{d}{dx}(bx) = b$, but what happens in the case of a real function of n variables?

Theorem

$$\nabla b^\top x = \nabla x^\top b = b$$

$$\nabla(b^\top x) = \begin{bmatrix} \frac{\partial(b^\top x)}{\partial x_1} \\ \vdots \\ \frac{\partial(b^\top x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial(b_1 x_1 + \dots + b_n x_n)}{\partial x_1} \\ \vdots \\ \frac{\partial(b_1 x_1 + \dots + b_n x_n)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} = b$$

Differentiating matrix expressions

For a 1d function ax^2 we know that $\frac{d}{dx}(ax^2) = 2ax$, but what about quadratic forms?

Theorem

$$\nabla(x^\top Ax) = (A + A^\top)x$$

Note that if A is symmetric, then $\nabla(x^\top Ax) = 2Ax$.

The proof is straightforward, let us express the quadratic form as a double sum:

$$x^\top Ax = \sum_i \sum_j a_{ij} x_i x_j$$

Differentiating matrix expressions

$$\frac{\partial(x^\top Ax)}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\sum_{k_1} \sum_{k_2} a_{k_1 k_2} x_{k_1} x_{k_2} \right) =$$

Differentiating matrix expressions

$$\begin{aligned}\frac{\partial(x^\top Ax)}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\sum_{k_1} \sum_{k_2} a_{k_1 k_2} x_{k_1} x_{k_2} \right) = \\ &= \frac{\partial}{\partial x_i} \left(\underbrace{\sum_{k_1 \neq i} \sum_{k_2 \neq i} a_{ik_2} x_{k_1} x_{k_2}}_{k_1 \neq i, k_2 \neq i} + \underbrace{\sum_{k_2 \neq i} a_{ik_2} x_i x_{k_2}}_{k_1 = i, k_2 \neq i} + \underbrace{\sum_{k_1 \neq i} a_{k_1 i} x_{k_1} x_i}_{k_1 \neq i, k_2 = i} + \underbrace{a_{ii} x_i^2}_{k_1 = i, k_2 = i} \right) =\end{aligned}$$

Differentiating matrix expressions

$$\begin{aligned}\frac{\partial(x^\top Ax)}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\sum_{k_1} \sum_{k_2} a_{k_1 k_2} x_{k_1} x_{k_2} \right) = \\ &= \frac{\partial}{\partial x_i} \left(\underbrace{\sum_{k_1 \neq i} \sum_{k_2 \neq i} a_{ik_2} x_{k_1} x_{k_2}}_{k_1 \neq i, k_2 \neq i} + \underbrace{\sum_{k_2 \neq i} a_{ik_2} x_i x_{k_2}}_{k_1 = i, k_2 \neq i} + \underbrace{\sum_{k_1 \neq i} a_{k_1 i} x_{k_1} x_i}_{k_1 \neq i, k_2 = i} + \underbrace{a_{ii} x_i^2}_{k_1 = i, k_2 = i} \right) = \\ &= \sum_{k_2 \neq i} a_{ik_2} x_{k_2} + \sum_{k_1 \neq i} a_{k_1 i} x_{k_1} + 2a_{ii} x_i =\end{aligned}$$

Differentiating matrix expressions

$$\begin{aligned}\frac{\partial(x^\top Ax)}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\sum_{k_1} \sum_{k_2} a_{k_1 k_2} x_{k_1} x_{k_2} \right) = \\ &= \frac{\partial}{\partial x_i} \left(\underbrace{\sum_{k_1 \neq i} \sum_{k_2 \neq i} a_{ik_2} x_{k_1} x_{k_2}}_{k_1 \neq i, k_2 \neq i} + \underbrace{\sum_{k_2 \neq i} a_{ik_2} x_i x_{k_2}}_{k_1 = i, k_2 \neq i} + \underbrace{\sum_{k_1 \neq i} a_{k_1 i} x_{k_1} x_i}_{k_1 \neq i, k_2 = i} + \underbrace{a_{ii} x_i^2}_{k_1 = i, k_2 = i} \right) = \\ &= \sum_{k_2 \neq i} a_{ik_2} x_{k_2} + \sum_{k_1 \neq i} a_{k_1 i} x_{k_1} + 2a_{ii} x_i = \\ &= \sum_{k_2} a_{ik_2} x_{k_2} + \sum_{k_1} a_{k_1 i} x_{k_1} =\end{aligned}$$

Differentiating matrix expressions

$$\begin{aligned}
 \frac{\partial(x^{\top}Ax)}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\sum_{k_1} \sum_{k_2} a_{k_1 k_2} x_{k_1} x_{k_2} \right) = \\
 &= \frac{\partial}{\partial x_i} \left(\underbrace{\sum_{k_1 \neq i} \sum_{k_2 \neq i} a_{ik_2} x_{k_1} x_{k_2}}_{k_1 \neq i, k_2 \neq i} + \underbrace{\sum_{k_2 \neq i} a_{ik_2} x_i x_{k_2}}_{k_1 = i, k_2 \neq i} + \underbrace{\sum_{k_1 \neq i} a_{k_1 i} x_{k_1} x_i}_{k_1 \neq i, k_2 = i} + \underbrace{a_{ii} x_i^2}_{k_1 = i, k_2 = i} \right) = \\
 &= \sum_{k_2 \neq i} a_{ik_2} x_{k_2} + \sum_{k_1 \neq i} a_{k_1 i} x_{k_1} + 2a_{ii} x_i = \\
 &= \sum_{k_2} a_{ik_2} x_{k_2} + \sum_{k_1} a_{k_1 i} x_{k_1} = \\
 &= \sum_j (a_{ij} + a_{ji}) x_j \qquad \Rightarrow \nabla(x^{\top}Ax) = (A + A^{\top})x
 \end{aligned}$$

Minimum of a quadratic form and the linear system

Recall that for $a > 0$ solving the equation $ax = b$ is equivalent to the quadratic function $\arg \min_x (ax^2 - 2bx)$ minimization.

Minimum of a quadratic form and the linear system

Recall that for $a > 0$ solving the equation $ax = b$ is equivalent to the quadratic function $\arg \min_x (ax^2 - 2bx)$ minimization.

To minimize a quadratic form $\arg \min_{x \in \mathbb{R}^n} (x^\top Ax - 2b^\top x)$ with a symmetric positive definite matrix A , equate the derivative to zero:

$$\nabla(x^\top Ax - 2b^\top x) = [0 \quad \dots \quad 0]^\top.$$

Minimum of a quadratic form and the linear system

Recall that for $a > 0$ solving the equation $ax = b$ is equivalent to the quadratic function $\arg \min_x (ax^2 - 2bx)$ minimization.

To minimize a quadratic form $\arg \min_{x \in \mathbb{R}^n} (x^\top Ax - 2b^\top x)$ with a symmetric positive definite matrix A , equate the derivative to zero:

$$\nabla(x^\top Ax - 2b^\top x) = [0 \quad \dots \quad 0]^\top.$$

The Hamilton operator is linear: $\nabla(x^\top Ax) - 2\nabla(b^\top x) = [0 \quad \dots \quad 0]^\top$.

Minimum of a quadratic form and the linear system

Recall that for $a > 0$ solving the equation $ax = b$ is equivalent to the quadratic function $\arg \min_x (ax^2 - 2bx)$ minimization.

To minimize a quadratic form $\arg \min_{x \in \mathbb{R}^n} (x^\top Ax - 2b^\top x)$ with a symmetric positive definite matrix A , equate the derivative to zero:

$$\nabla(x^\top Ax - 2b^\top x) = [0 \quad \dots \quad 0]^\top.$$

The Hamilton operator is linear: $\nabla(x^\top Ax) - 2\nabla(b^\top x) = [0 \quad \dots \quad 0]^\top$.
Apply the differentiation theorems:

$$(A + A^\top)x - 2b = [0 \quad \dots \quad 0]^\top.$$

Minimum of a quadratic form and the linear system

Recall that for $a > 0$ solving the equation $ax = b$ is equivalent to the quadratic function $\arg \min_x (ax^2 - 2bx)$ minimization.

To minimize a quadratic form $\arg \min_{x \in \mathbb{R}^n} (x^\top Ax - 2b^\top x)$ with a symmetric positive definite matrix A , equate the derivative to zero:

$$\nabla(x^\top Ax - 2b^\top x) = [0 \quad \dots \quad 0]^\top.$$

The Hamilton operator is linear: $\nabla(x^\top Ax) - 2\nabla(b^\top x) = [0 \quad \dots \quad 0]^\top$.
Apply the differentiation theorems:

$$(A + A^\top)x - 2b = [0 \quad \dots \quad 0]^\top.$$

Recall that A is symmetric: $Ax = b$.

Back to the linear regression

Given two points (x_1, y_1) and (x_2, y_2) , find the line that passes through: $y = \alpha x + \beta$.

Back to the linear regression

Given two points (x_1, y_1) and (x_2, y_2) , find the line that passes through: $y = \alpha x + \beta$.

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \end{cases}$$

Back to the linear regression

Given two points (x_1, y_1) and (x_2, y_2) , find the line that passes through: $y = \alpha x + \beta$.

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \end{cases} \quad \underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \end{bmatrix}}_{:=A} \underbrace{\begin{bmatrix} \alpha \\ \beta \end{bmatrix}}_{:=x} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}}_{:=b} \quad \Rightarrow x^* = A^{-1}b$$

Back to the linear regression

Given two points (x_1, y_1) and (x_2, y_2) , find the line that passes through: $y = \alpha x + \beta$.

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \end{cases}$$

$$\underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \end{bmatrix}}_{:=A} \underbrace{\begin{bmatrix} \alpha \\ \beta \end{bmatrix}}_{:=x} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}}_{:=b} \quad \Rightarrow x^* = A^{-1}b$$

Now add a **third** point:

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \\ \alpha x_3 + \beta = y_3 \end{cases}$$

Back to the linear regression

Given two points (x_1, y_1) and (x_2, y_2) , find the line that passes through: $y = \alpha x + \beta$.

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \end{cases} \quad \underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \end{bmatrix}}_{:=A} \underbrace{\begin{bmatrix} \alpha \\ \beta \end{bmatrix}}_{:=x} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}}_{:=b} \quad \Rightarrow x^* = A^{-1}b$$

Now add a **third** point:

$$\begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \\ \alpha x_3 + \beta = y_3 \end{cases} \quad \underbrace{\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix}}_{:=A(3 \times 2)} \underbrace{\begin{bmatrix} \alpha \\ \beta \end{bmatrix}}_{:=x(2 \times 1)} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}}_{:=b(3 \times 1)}$$

A is rectangular, and thus it is not invertible. Oops!

Back to the linear regression

No biggie, let us rewrite the system:

$$\alpha \underbrace{[x_1 \ x_2 \ x_3]^T}_{:=\vec{i}} + \beta \underbrace{[1 \ 1 \ 1]^T}_{:=\vec{j}} = [y_1 \ y_2 \ y_3]^T$$

Back to the linear regression

No biggie, let us rewrite the system:

$$\alpha \underbrace{[x_1 \ x_2 \ x_3]^T}_{:=\vec{i}} + \beta \underbrace{[1 \ 1 \ 1]^T}_{:=\vec{j}} = [y_1 \ y_2 \ y_3]^T$$

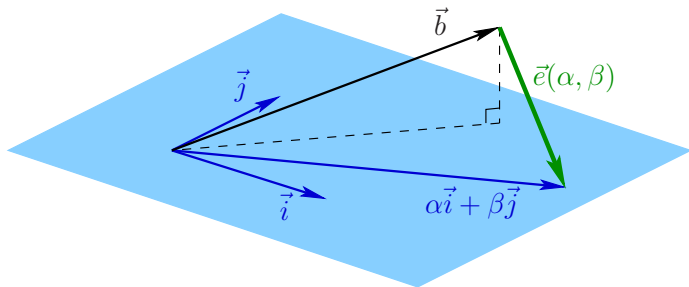
$$\alpha \vec{i} + \beta \vec{j} = \vec{b}.$$

Back to the linear regression

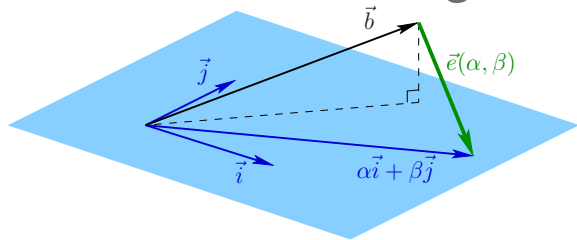
No biggie, let us rewrite the system:

$$\alpha \underbrace{\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}^\top}_{:=\vec{i}} + \beta \underbrace{\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^\top}_{:=\vec{j}} = \begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix}^\top \quad \alpha\vec{i} + \beta\vec{j} = \vec{b}.$$

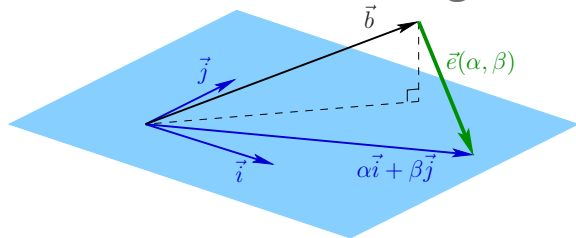
Solve for $\arg \min_{\alpha, \beta} \|\vec{e}(\alpha, \beta)\|$, where $\vec{e}(\alpha, \beta) := \alpha\vec{i} + \beta\vec{j} - \vec{b}$:



Back to the linear regression



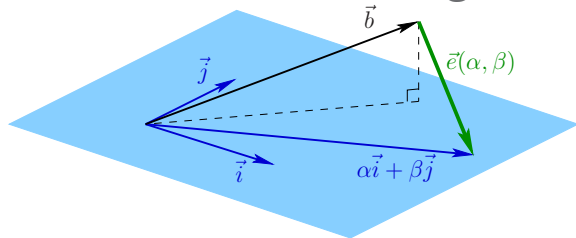
Back to the linear regression



The $\|\vec{e}(\alpha, \beta)\|$ is minimized when $\vec{e}(\alpha, \beta) \perp \text{span}\{\vec{i}, \vec{j}\}$:

$$\begin{cases} \vec{i}^\top \vec{e}(\alpha, \beta) = 0 \\ \vec{j}^\top \vec{e}(\alpha, \beta) = 0 \end{cases}$$

Back to the linear regression

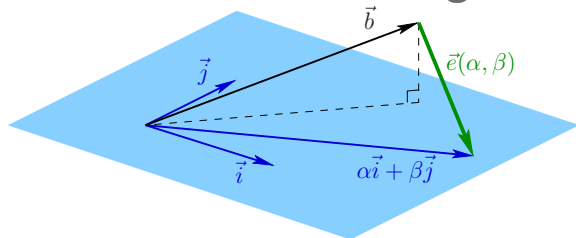


The $\|\vec{e}(\alpha, \beta)\|$ is minimized when $\vec{e}(\alpha, \beta) \perp \text{span}\{\vec{i}, \vec{j}\}$:

$$\begin{cases} \vec{i}^\top \vec{e}(\alpha, \beta) = 0 \\ \vec{j}^\top \vec{e}(\alpha, \beta) = 0 \end{cases}$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 1 \end{bmatrix} \left(\alpha \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \beta \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Back to the linear regression



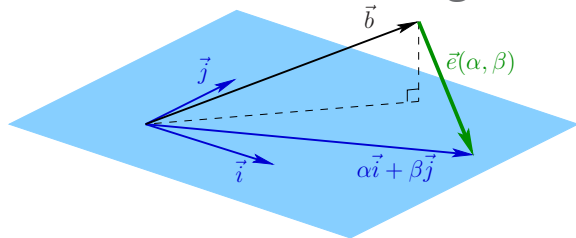
The $\|\vec{e}(\alpha, \beta)\|$ is minimized when $\vec{e}(\alpha, \beta) \perp \text{span}\{\vec{i}, \vec{j}\}$:

$$\begin{cases} \vec{i}^\top \vec{e}(\alpha, \beta) = 0 \\ \vec{j}^\top \vec{e}(\alpha, \beta) = 0 \end{cases}$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 1 \end{bmatrix} \left(\alpha \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \beta \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$A^\top (Ax - b) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Back to the linear regression



The $\|\vec{e}(\alpha, \beta)\|$ is minimized when $\vec{e}(\alpha, \beta) \perp \text{span}\{\vec{i}, \vec{j}\}$:

$$\begin{cases} \vec{i}^\top \vec{e}(\alpha, \beta) = 0 \\ \vec{j}^\top \vec{e}(\alpha, \beta) = 0 \end{cases}$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \\ 1 & 1 & 1 \end{bmatrix} \left(\alpha \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \beta \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$A^\top (Ax - b) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

In a general case the matrix $A^\top A$ can be invertible!

$$A^\top Ax = A^\top b.$$

Some nice properties of $A^\top A$

Theorem

$A^\top A$ is symmetric.

Some nice properties of $A^\top A$

Theorem

$A^\top A$ is symmetric.

It is very easy to show:

$$(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A.$$

Some nice properties of $A^\top A$

Theorem

$A^\top A$ is symmetric.

It is very easy to show:

$$(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A.$$

Theorem

$A^\top A$ positive semidefinite: $\forall x \in \mathbb{R}^n \quad x^\top A^\top A x \geq 0$.

Some nice properties of $A^\top A$

Theorem

$A^\top A$ is symmetric.

It is very easy to show:

$$(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A.$$

Theorem

$A^\top A$ positive semidefinite: $\forall x \in \mathbb{R}^n \quad x^\top A^\top A x \geq 0$.

It follows from the fact that $x^\top A^\top A x = (Ax)^\top Ax \geq 0$. Moreover, $A^\top A$ is positive definite in the case where A has linearly independent columns (rank A is equal to the number of the variables in the system).

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\arg \min \|Ax - b\|^2$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\arg \min \|Ax - b\|^2 = \arg \min (Ax - b)^\top (Ax - b) =$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\arg \min \|Ax - b\|^2 = \arg \min (Ax - b)^\top (Ax - b) = \arg \min (x^\top A^\top - b^\top)(Ax - b) =$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\begin{aligned}\arg \min \|Ax - b\|^2 &= \arg \min (Ax - b)^\top (Ax - b) = \arg \min (x^\top A^\top - b^\top)(Ax - b) = \\ &= \arg \min (x^\top A^\top Ax - b^\top Ax - x^\top A^\top b + \underbrace{b^\top b}_{\text{const}}) =\end{aligned}$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\begin{aligned}\arg \min \|Ax - b\|^2 &= \arg \min (Ax - b)^\top (Ax - b) = \arg \min (x^\top A^\top - b^\top)(Ax - b) = \\ &= \arg \min (x^\top A^\top Ax - b^\top Ax - x^\top A^\top b + \underbrace{b^\top b}_{\text{const}}) = \\ &= \arg \min (x^\top A^\top Ax - 2b^\top Ax) =\end{aligned}$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\begin{aligned}\arg \min \|Ax - b\|^2 &= \arg \min (Ax - b)^\top (Ax - b) = \arg \min (x^\top A^\top - b^\top)(Ax - b) = \\ &= \arg \min (x^\top A^\top Ax - b^\top Ax - x^\top A^\top b + \underbrace{b^\top b}_{\text{const}}) = \\ &= \arg \min (x^\top A^\top Ax - 2b^\top Ax) = \arg \min (x^\top \underbrace{(A^\top A)}_{:=A'} x - 2 \underbrace{(A^\top b)}_{:=b'}^\top x)\end{aligned}$$

Least squares in more than two dimensions

The same reasoning applies, here is an algebraic way to show it:

$$\begin{aligned}\arg \min \|Ax - b\|^2 &= \arg \min (Ax - b)^\top (Ax - b) = \arg \min (x^\top A^\top - b^\top)(Ax - b) = \\ &= \arg \min (x^\top A^\top Ax - b^\top Ax - x^\top A^\top b + \underbrace{b^\top b}_{\text{const}}) = \\ &= \arg \min (x^\top A^\top Ax - 2b^\top Ax) = \arg \min (x^\top \underbrace{(A^\top A)}_{:=A'} x - 2 \underbrace{(A^\top b)}_{:=b'}^\top x)\end{aligned}$$

The takeaway message

The least squares problem $\arg \min \|Ax - b\|^2$ is equivalent to minimizing the quadratic function $\arg \min (x^\top A' x - 2b'^\top x)$ with (in general) a symmetric positive definite matrix A' . This can be done by solving a linear system $A'x = b'$.



Dmitry Sokolov

Least squares for programmers
— with color plates —