

## Least Squares

Suppose that we are given a linear system,  $Ax=b$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ . Assume that  $m > n$ . Moreover,  $\text{rank } A = n$ . (full column rank)

In this case, the system is usually inconsistent, that is, there's no  $x$  satisfying  $Ax=b$ .

One way of finding an approximate solution is to consider the residual vector  $r = Ax - b \in \mathbb{R}^m$  and to minimize the norm square of it. This leads to the least squares problem:

$$(LS) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|Ax - b\|_2^2.$$

$$\begin{aligned} \text{Note that } f(x) &= \|Ax - b\|_2^2 = (Ax - b)^T (Ax - b) \\ &= x^T A^T A x - x^T A b - b^T A x + b^T b \\ &= x^T A^T A x - 2(A^T b)^T x + b^T b. \end{aligned}$$

Hence,  $f$  is a quadratic function of  $n$  variables.

$$\nabla f(x) = 2A^T A x - 2A^T b \quad \text{and} \quad \nabla^2 f(x) = 2A^T A.$$

Let's check if  $A^T A$  is psd (even pd):

Note that for any vector  $v \in \mathbb{R}^n$ , we have

$$v^T A^T A v = (Av)^T (Av) = \|Av\|_2^2 \geq 0.$$

Hence  $A^T A$  is positive semidefinite. Then, any stationary point for (LS) is a global minimum point.

Now,  $\text{rank } A = n$  implies that the columns of  $A$ , say  $a_1, \dots, a_n \in \mathbb{R}^m$ , are linearly independent. Moreover for any vector  $v \in \mathbb{R}^n$  we have

$$Av = v_1 a_1 + \dots + v_n a_n \in \mathbb{R}^m.$$

By linear independence we know that  $Av = 0$  if and only if  $v = 0 \in \mathbb{R}^n$ .

Thus,  $v^T A^T A v = \|Av\|_2^2 = 0$  only for  $v = 0$ , and it's strictly positive for any other vector  $v \in \mathbb{R}^n$ .

This means that  $A^T A$  is positive definite. In particular, it has an inverse. This implies that  $\nabla f(x) = 0 \Leftrightarrow$

$$A^T A x = A^T b \Leftrightarrow (A^T A)^{-1} A^T A x = (A^T A)^{-1} A^T b$$

$$\Leftrightarrow \hat{x} = (A^T A)^{-1} A^T b.$$

Note that if we had  $m=n$  &  $A$  invertible, then this would simplify as

$$\hat{x} = A^{-1} \underbrace{(A^T)^{-1} A^T}_{I} b = A^{-1} b$$

as expected.

Ex: (Data fitting)

Suppose we are given a set of data points  $(s_i, t_i)$  for  $i=1, \dots, m$  where  $s_i \in \mathbb{R}^n$  and  $t_i \in \mathbb{R}$ . We want to fit this data to a linear form, that is, we look for a vector  $x \in \mathbb{R}^n$  s.t.  $s_i^T x = t_i, i=1, \dots, m$ .

In general,  $m$  is much larger than  $n$  and the system is overdetermined.

Then, one can use the least squares approach.

$$\text{let } S = \begin{bmatrix} s_1 & \dots & s_m \end{bmatrix} \in \mathbb{R}^{n \times m}, \quad t = (t_1, \dots, t_m)^T.$$

we want to minimize  $\|S^T x - t\|_2^2$ .

$\hat{x} = (SS^T)^{-1} S^T t$   
is the least square solution.

Ex: (Nonlinear data fitting)

let  $(x_i, y_i), i=1, \dots, m$  be data points such that  $s_i, t_i \in \mathbb{R}, \forall i$ .

Assume now, we want to fit data points in a polynomial form, that is,

$$(*) \quad y_i \simeq a_p x_i^p + a_{p-1} x_i^{p-1} + \dots + a_1 x_i + a_0, \quad i=1, \dots, m.$$

for some numbers  $a_0, \dots, a_p \in \mathbb{R}$ . The idea is to find the best  $a \in \mathbb{R}^{p+1}$  such that the residual is minimized.

$$\text{let } s_i = (x_i^p, x_i^{p-1}, \dots, x_i, 1)^T \in \mathbb{R}^{p+1}$$

we can write  $(*)$  as  $y_i \simeq s_i^T a$ .

we obtain a very similar case as in the previous example. let  $S = [s_1, \dots, s_m] \in \mathbb{R}^{(p+1) \times m}$ .

$$y = (y_1, \dots, y_m)^T \in \mathbb{R}^m$$

we want to minimize  $\|S^T \alpha - y\|_2^2$ , which is again a least square problem.

## The Gradient Method

Consider the unconstrained minimization problems:  $\left( \begin{array}{l} \text{minimize } f(x) \\ x \in \mathbb{R}^n \end{array} \right)$

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ , continuously differentiable

We may be able to solve it by finding the stationary points (after showing existence) and selecting the minimum point among them. However, this is not always possible:

- (1)  $\nabla f(x) = 0$  may be difficult to solve
- (2)  $\nabla f(x) = 0$  may have inf.<sup>y</sup> many solutions...

There are also iterative methods:

- Pick  $x_0 \in \mathbb{R}^n$  (arbitrarily, or a guess...)  $\rightarrow$  ?
- Pick a "descent direction"  $d_k$
- Find a stepsize  $t_k$  s.t.  $f(x_k + t_k d_k) < f(x_k)$   $\leftarrow$  ?
- $x_{k+1} \leftarrow x_k + t_k d_k$
- Stop if a stopping criteria is satisfied.  $\leftarrow$  ?

Defn:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , continuously differentiable. A vector  $0 \neq d \in \mathbb{R}^n$  is a descent direction of  $f$  at  $x$  if the directional derivative of  $f$  in direction  $d$  at  $x$  is negative, that is,

$$f'(x; d) = \lim_{t \downarrow 0} \frac{f(x+td) - f(x)}{t} < 0.$$

\* Note that  $f'(x; d) = \nabla f(x)^T d$ .

For gradient method, the negative of the gradient is selected as the descent direction, that is,  $d_k = -\nabla f(x_k)$ .

Note that  $f'(x_k; d_k) = \nabla f(x_k)^T (-\nabla f(x_k)) = -\|\nabla f(x_k)\|^2 < 0$ .

Indeed, this is the "steepest" descent direction, that is, over all  $d \in \mathbb{R}^n$  with  $\|d\|=1$ ,  $-\nabla f(x_k)$  yields the minimum directional derivative

$$f'(x_k, -\nabla f(x_k)) = \min_{d \in \mathbb{R}^n} \{ f'(x_k, d) \mid \|d\|=1 \}.$$

• How to select the stepsize?

- constant stepsize:  $t_k = \bar{t} \quad \forall k$

easy to apply but how to choose? If too small, then may converge slow.  
If not small, may not satisfy  $f(x_k + t_k d_k) < f(x_k)$ .

- exact line search: minimize  $f$  along the ray  $x_k + t \cdot d_k$

i.e. solve  $\underset{t \geq 0}{\text{minimize}} f(x_k + t d_k)$

This may be difficult (may not be possible) to solve...

- backtracking: start with an initial guess  $t_k = s$ , iterate  $t_k \leftarrow \beta t_k$  for some  $0 < \beta < 1$ .

until we have  $f(x_k) - f(x_k + t_k d_k) \geq -\alpha t_k \nabla f(x_k)^T d_k$  for some  $\alpha \in (0, 1)$ .

• Stopping condition For a predetermined error  $\varepsilon > 0$ , if  $\|\nabla f(x_{k+1})\| \leq \varepsilon$ , then STOP.

Return  $x_{k+1}$ .



Ex: Exact line search for quadratic functions. Let  $A \in \mathbb{R}^{n \times n}$  be p.d.

$b \in \mathbb{R}^n$ ,  $c \in \mathbb{R}$ . Consider  $f(x) = x^T A x + 2b^T x + c$

Let  $d \in \mathbb{R}^n$  be a descent direction (may not be the steepest descent) at a point  $x$ . We want to find the stepsize by exact line search.

$\left( \begin{array}{l} \text{minimize } f(x+td) \\ \text{s.t. } t \geq 0 \end{array} \right)$  The decision variable is "t". Let's write the objective function, as a function of t:

$$g(t) := f(x+td) = (x+td)^T A (x+td) + 2b^T (x+td) + c$$

$$= x^T A x + t x^T A d + t d^T A x + t^2 d^T A d + 2b^T x + 2t b^T d + c$$

$$= t^2 (d^T A d) + 2t(x^T A d) + 2t b^T d + x^T A x + 2b^T x + c$$

$$= (d^T A d) t^2 + 2(x^T A d + b^T d) t + x^T A x + 2b^T x + c \quad \left( \begin{array}{l} \text{a quadratic} \\ \text{function of } t \end{array} \right)$$

$$g'(t) = 2(d^T A d)t + 2(x^T A d + b^T d)$$

$$= 2(d^T A d)t + 2d^T (Ax + b)$$

(as  $A^T = A$ )

$$\Rightarrow t^* = \frac{-d^T (Ax + b)}{d^T A d} \quad \text{Note that } \nabla f(x) = 2Ax + 2b = 2(Ax + b).$$

$$\text{Hence, } t^* = \frac{-d^T \nabla f(x)}{2d^T A d}.$$