

This means we could think about

$$q(\lambda) = \inf_x \mathcal{L}(x, \lambda)$$

↑ recall: Lagrangian

Final  
 write  $x^c$  for  $\arg \min_x f(x)$  global min  
 st  $g(x) = 0$

Then we have

$$q(\lambda) \leq f(x^c)$$

(This is fairly easy. Think about a single constraint;  $\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$ )

Now: either  $q(\lambda) = -\infty$  or its finite;  
 if it's finite, then  $\lambda g(x)$  must be bounded below;  
 if the bound is 0, then  $g(x) = 0$  or  $\lambda = 0$   
 if  $g(x) = 0$ , then  $q$  can't be bigger than  $f(x^c)$   
 if  $\lambda = 0$ , same

now  
we  
so if bound is bigger than zero,  
can't have a soln.  
 $\lambda q(x) \leq 0$ , so ...

25a

Example Duals :

$$\min_x \quad \frac{x^T x}{2} \quad \text{st} \quad Ax + b = 0$$

$$L : \quad \frac{x^T x}{2} - \lambda^T (Ax + b)$$

now :  $\inf_x$  occurs when  $x - A^T \lambda = 0$  ( $\nabla_x L = 0$ )

so :

$$g(\lambda) = -\frac{\lambda^T A A^T \lambda}{2} - b^T \lambda$$

Notice :  $\max_{\lambda} g(\lambda)$  occurs when  $A A^T \lambda - b = 0$

(which we've seen before).

Idea:

- form dual and solve  $\max_{\lambda} q(\lambda)$

- This isn't usually a good idea, cause we can't get the dual.

Idea:

- go down in  $x$ , then up in  $\lambda$
- we've done something like that already! (AHM!)
- Sometimes called "Dual Ascent"

recall :

start:  $x^0, \lambda^0 = 0, c^0 > 0$

$$x^{k+1} = \underset{x}{\operatorname{arg\,min}} \quad A(x, \lambda^k, c^k) = f(x) - \lambda^{kT} g(x) + \frac{c^k}{2} g(x)^T g(x)$$

$$\lambda^{k+1} = \lambda^k - \frac{c^k}{2} g(x^{k+1})$$

$$c^{k+1} = \gamma c^k$$

Notice

$$A(x^{k+1}, \lambda^{k+1}, c^{k+1})$$

$$= A(x^{k+1}, \lambda^k, c^k) + \underbrace{\gamma \frac{c^k}{2} g(x^k)^T g(x^k)}_{\uparrow}$$

and this is +ve.



Another new idea: Splitting

• Example

want to  $\min_x \|Ax - b\|^2$   
tolerable.

for very large  $A$  big  $\downarrow A$

• don't want to do the linear algebra  
 $A^T A$  might be tough to work with

• instead, write

$$A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

and consider

$$\|A_1 x_1 - b_1\|^2 + \|A_2 x_2 - b_2\|^2$$

st  $x_1 = x_2$

① these two problems have solutions in the same place.

Now look at ALM

(30)

$$\text{Aug. Lag.} = \|A_1 x_1 - b_1\|^2 + \|A_2 x_2 - b_2\|^2 + \lambda^T (x_1 - x_2) + \frac{c}{2} (x_1 - x_2)^T (x_1 - x_2)$$

We could "split" the step obtaining a min in  $x$

so

$$x_1^{k+1} = \underset{x_1}{\text{argmin}} A(x_1, x_2^k, \lambda^k, c^k)$$

and this is a linear system which is smaller than the  $\|Ax - b\|^2$  system

$$x_2^{k+1} = \underset{x_2}{\text{argmin}} A(x_1^{k+1}, x_2, \lambda^k, c^k)$$

$\left. \begin{array}{l} \lambda^{k+1} \\ c^{k+1} \end{array} \right\} \text{ as usual.}$

Example An important working example

- Sparse regression.

recall linear regression.

- we have a dataset

$$(y_i, \underline{x}_i^T)$$

and we wish to predict best estimate of future  $y$  from new  $\underline{x}$

one approach:

linear model. d-dimensional

$$y_{\text{pred}}(\underline{x}) = \underline{x}_i^T \Theta$$

↑ parameters

How to choose  $\Theta$ ?

min least-squares error



$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \theta$$

and solve:

$$\min_{\theta} [Y - X\theta]^T [Y - X\theta]$$

Which is solved by

$$X^T X \theta = X^T Y$$

But there are problems with this story

(I) imagine  $d \gg N$   $\rightarrow$  many  $\theta$ !

II even if  $d \ll N$ , we could

have

$X^T X$  has small rank, and

this must lead to large prediction

errors

III

imagine some of the  $x$ -components are irrelevant — then we want the corresponding  $\theta$ -components to be 0 — otherwise they contribute error.

Practical Q:

can we force lots of  $\theta$ -components to be 0?

Obvious strategy doesn't work

(34)

• Obvious strategy — penalize  $\Theta^T \Theta$

so solve

$$\|Y - X\Theta\|_2^2 + \frac{\alpha}{2} \Theta^T \Theta$$

Solu :  $(X^T X + \alpha I) \Theta = X^T Y$

- notice if  $\alpha$  is big enough,  $\Pi$  ~~has~~ is solved
- Also  $I$ , but not convincingly
- We DON'T GET 0's
- $l_2$  regularization  
OR.

# Very simple example

(35)

$$y_i = c$$

$$\underline{x}_i^T = (1, \xi_i)$$

a sample from noise, mean = 0, var =  $\sigma^2$

then  $E[X^T X] = \begin{pmatrix} N & 0 \\ 0 & N\sigma^2 \end{pmatrix}$  and  $E[X^T Y] = \begin{pmatrix} Nc \\ 0 \end{pmatrix}$

↑ approx...

so if we get  $\begin{pmatrix} Nc \\ 0 \end{pmatrix}$  then  $\Theta = \begin{bmatrix} c \\ 0 \end{bmatrix}$

But imagine  $X^T Y = \begin{pmatrix} Nc \\ z \end{pmatrix}$

then  $\Theta = \begin{pmatrix} c \\ \frac{z}{N\sigma^2} \end{pmatrix}$

if we regularize w  $\lambda^2$  we get

$$\begin{pmatrix} \frac{c}{1+d} \\ \frac{z}{N\sigma^2+d} \end{pmatrix}$$

smaller, but not zero



Why:

- the penalty for  $\theta_i$  small, but not 0 is very small



- we would like a small  $\theta_i$  to be more expensive

→ rather than  $\frac{1}{2} \theta^T \theta$

use  $\sum_i |\theta_i| = \|\theta\|_1$

absolute value.

$l_1$  norm

so we could use

(37)

$$\min_{\theta} \|X\theta - Y\|^2 + \alpha \|\theta\|_1$$

(known as a lasso).

But how do we find a min?

Notice:  $\|\theta\|_1$  isn't differentiable ~~at~~  
(think  $\theta_i = 0$ )!

Notice:  
problems  $\rightarrow$  ignoring this will lead to  
we won't get zero's!