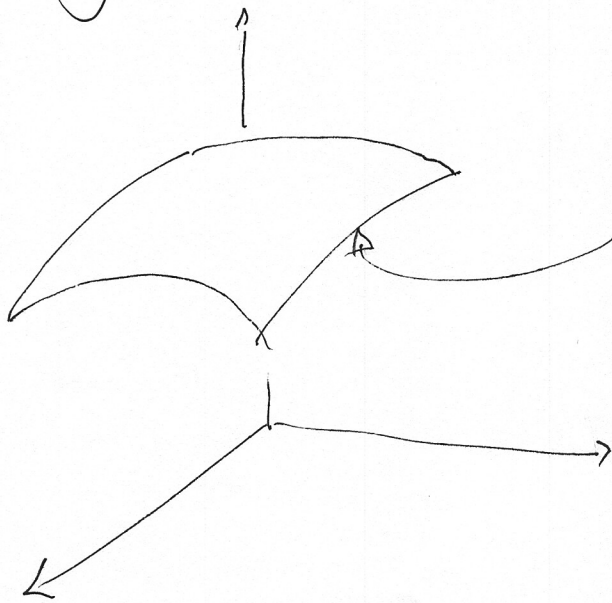# Constrained Optimization :

. Two kinds of problem — equality constraints

              — ineq constraints.

. These are quite different, in important ways.
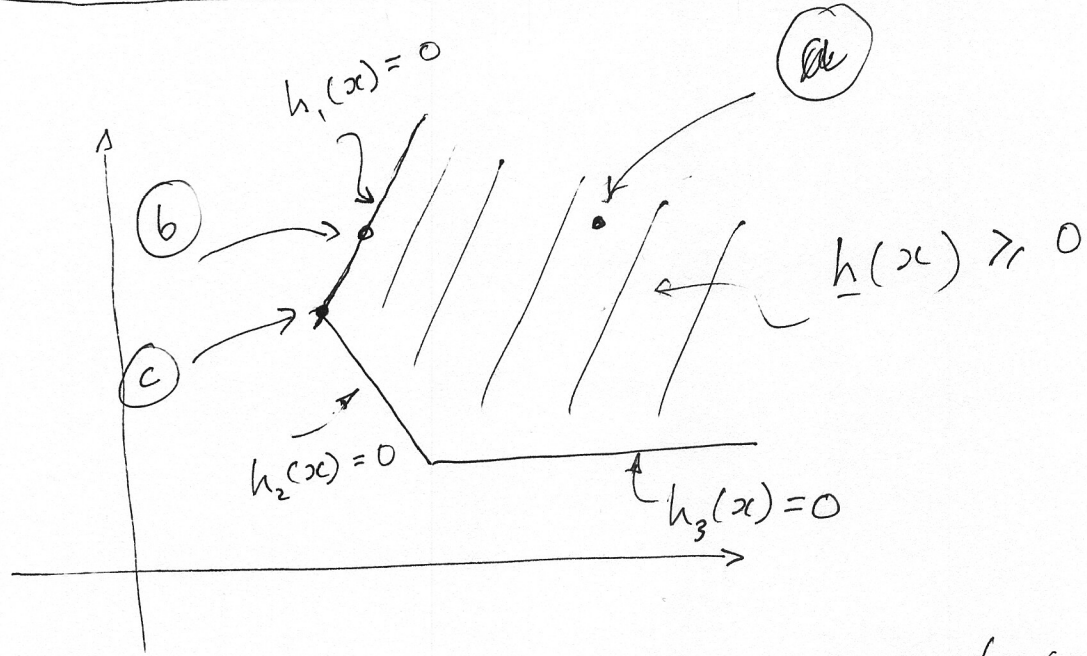
$g(x) = 0$

minize $f$ subject to $g = 0$

$\equiv \nabla f$ is ~~flat~~ normal to the "surface" $g = 0$

( because otherwise, we could move along $g = 0$ in a way that reduces $f$ )
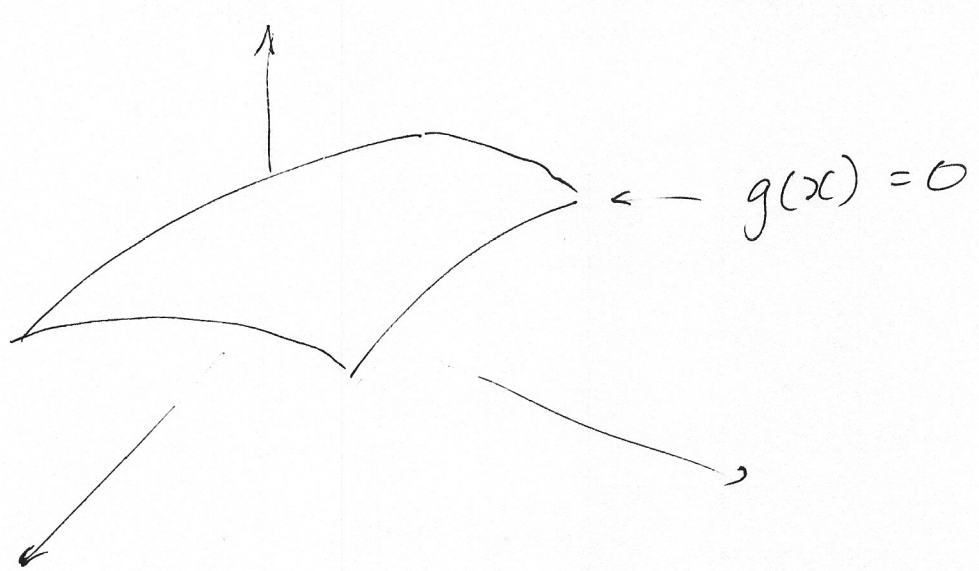
# Inequality Picture:

$h_1(x) = 0$

ⓐ

ⓑ

ⓒ

$\underline{h}(x) \geq 0$

$h_2(x) = 0$

$h_3(x) = 0$

- at ⓐ, constraints are irrelevant; locally, problem involves $\min f(x)$ w/o constraints

- at ⓑ, $h_1$ applies, but no others, so locally problem looks like
$$\min f(x) \quad \text{s.t.} \quad h_1(x) = 0$$

- at ⓒ, $h_1, h_2$ . . . .
$$\min f(x) \quad \text{st} \quad h_1(x) = 0$$
$$h_2(x) = 0$$

**BUT:** one step may mean picture <u>changes</u>!

# Equality constraints :

$\longleftarrow g(x) = 0$

- Simple picture :
  - 3D, one constraint $g(x) = 0$

- $\min f(x)$ st $g(x) = 0$

- answer occurs at points where

$$\nabla f \quad \text{is} \quad \text{normal to} \quad \{g(x) = 0\}$$

- Normal of implicit surface $g(x) = 0$

  is $\nabla g$

$$\therefore \quad \nabla f = \lambda \nabla g$$

$\lambda$ some unknown constant.

What if there are many constraints in N-D?

$$g_1(x) = 0, \quad g_2(x) = 0, \quad \text{etc} \ldots$$

$\nabla f$ is normal

$|||$

$$\nabla f \in \text{span}\{\nabla g_1, \nabla g_2, \nabla g_3, \ldots \}$$

$|||$

$$\nabla f = \lambda_1 \nabla g_1 + \lambda_2 \nabla g_2 + \ldots$$

equivalently, write $\underline{g} = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \end{bmatrix}$

$$J_g = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} & \cdots \\ \frac{\partial g_2}{\partial x_1} \\ \vdots \end{bmatrix}$$

then

$$\nabla f = \lambda^T J_g.$$

This justifies writing the _Lagrangian_

$$\mathcal{L} = f - \lambda^T g.$$

at minimum:

$$\nabla_x \mathcal{L} = \nabla f - \lambda^T J_g = 0$$

$$\nabla_\lambda \mathcal{L} = -g = 0$$

These conditions must be true at a minimum.

V. Important special cases for constrained optimization  ⑪

①    max    $\dfrac{x^{T}Ax}{2}$       s.t.   $x^{T}x = 1$

Lagrangian

$$\dfrac{x^{T}Ax}{2} - \lambda(x^{T}x - 1)$$

$\therefore$ $\boxed{Ax = \lambda x}$

$\underleftarrow{\quad}$ eigenvalue problem

$$\max \quad x^T \frac{A}{2} x \qquad st \quad x^T B x = 1$$

Lagrangian

$$x^T \frac{A}{2} x - \lambda(x^T B x - 1)$$

$$\therefore \boxed{Ax - \lambda Bx = 0}$$

$\uparrow$ generalized eigenvalue problem

Notice :     NOT    the same as

$$B^{-1} Ax - \lambda x = 0 \leftarrow \text{NAUGHTY!}$$

because $B^{-1}$ may not exist

• Any good Numerical linear Alg package can do these.

min

$$\frac{x^T x}{2} \qquad st \qquad Ax = b$$

(i.e. closest point on linear subspace to ) the origin

## Lagrangian :

$$\frac{x^T x}{2} - \lambda^T (Ax - b)$$

$$\therefore \qquad x - A^T \lambda = 0 \qquad \qquad \left( \nabla_x \mathcal{L} \right) \quad ⑬$$

So

$$A A^T \lambda = b \qquad \leftarrow \quad ⓐ$$

Alg    solve linear system ⓐ for $\lambda$,
then  subs  for  $x$.  in  ⑬

min

$$x^T \frac{A}{2} x + b^T x$$

s.t. $\quad Cx = d$

Lagrangian:

$$x^T \frac{A}{2} x + b^T x - \lambda^T (Cx - d)$$

$\nabla_x \mathcal{L} :$ 

$$Ax + b - C^T \lambda = 0$$

$\nabla_\lambda \mathcal{L} :$

$$Cx - d = 0$$

$$\begin{bmatrix} A & -C^T \\ C & 0 \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} -b \\ d \end{bmatrix}$$

↖ Solve this.

Notice how useful it has been to know the Lagrange multipliers.

## Algorithms for other cases

. ### Eliminating constraints

- Sometimes, we can parametrize the constraint set and move on that.

- not usually a good idea

eg. min $x^2 + (y-10)^2$

subject to

$y - \sin x = 0$

Notice the rich supply of local min

. we could rewrite as

$$\min x^2 + (\sin x - 10)^2 \qquad \text{w/ no constraints}$$

then try to min this.

. Notice when we do this, we are confining our steps to the constraint space:

    — <u>Problem</u> — don't see large scale structure of objective

. <u>Equivalent</u>:
    . take step <u>Tangent</u> to constraint space
    . project back.

( . i.e. make up local parametrization )

For example:

$$\min \quad x^T A x + b^T x$$

$$\text{s.t.} \quad \underline{\phi}(\underline{x}) = 0$$

$$\uparrow \text{ vector function.}$$

Now consider a step $\Delta x$

$$\underline{\phi}(\underline{x} + \Delta \underline{x}) \approx \underline{\phi}(\underline{x}) + J_\phi \cdot \Delta x.$$

So we could try:

Step :

1) ~~$\min \quad x^T A x + b^T x$~~

$$\text{st}$$

$$\min_{\Delta x} \quad (x_k + \Delta x)^T A (x_k + \Delta x) + b^T (x_k + \Delta x)$$

$$\text{st} \quad J_\phi \cdot \Delta x = 0$$

2) correct by finding $x_{k+1}$

$$\text{s.t.} \quad \phi(x_{k+1}) = 0, \quad \text{Start search at } x_k + \Delta x.$$

Again, not usually a great plan, because we may have a hard time taking big steps

## Quadratic penalty method

$$\min \quad f(\underline{x}) \qquad \text{st} \quad \underline{g}(\underline{x}) = 0$$

- approach by $\min \quad f(\underline{x}) + \frac{c}{2} \underline{g}^T \underline{g}$

- if $c$ is big, this forces $g^T g$ to be small

- advantage: — we could take steps off the constraint space
  — now its unconstrained.

· <u>Disadvantages</u>   (Big)

1)   big c   $\Rightarrow$   some big terms in Hessian

$$H = H_f + c\left[ \cancel{\nabla} J_g^T J_g + \cdots \right]$$

( so we should see terms that look like $\sum_{\substack{i \\ k}} \frac{\partial g_k}{\partial x_i}^2$ on Hess

<u>diag</u> _____ !

2)   at soln, g <u>isn't zero</u>

at   soln

$$\boxed{\nabla_x f + c\, g J_g = 0}$$

$\nabla_x f$ won't be zero, in general, so g can't be!

The method of multipliers
(also, Augmented Lagrangian)
method.

$$\text{min} \quad f(x) \qquad \text{st} \qquad g(x) = 0$$

form : Augmented Lagrangian

$$A(x;\lambda) = \underbrace{f(x) - \lambda^T g(x)}_{\text{Lagrangian}} + \underset{\uparrow}{\frac{c}{2}(g^T g)}$$

augmentation

- Now, assume we have an estimate $\lambda^K$ of the L.M.s

Minimize $A(x; \lambda^K)$ to get $x^{(k)}$

at $x^k$ we have

$$\nabla f(x^k) - \lambda^{k_T} g J_g + c g^T J_g = 0$$

Now, pattern match to conditions

$$\nabla_x \mathcal{L} = 0$$

$$\nabla_x \mathcal{L} = \nabla f - \lambda^T J_g$$

This suggests

$$\lambda^{k+1} = (\lambda^k - cg)$$

Notice: we could have a soln w/ g=0!

## ALM:

- start w $x^0$, $\lambda^0, c^0$

- min $A(x, \lambda^k) = f - \lambda^{k^T} g + \frac{c^0}{2} g^T g$

  to get $x^k$

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

  $c^{k+1} = r c^k$

  $\boxed{\text{often } 2}$

Q: How do we know its converged?

A: In ALM, usually nothing to do — we don't reject steps — but issue for future.

Q: do we have Hessian probs?

A: No, because $\lambda$ ests help. (formally, there is some bound on the $c$ required to get exact soln.)

# First glimple of Quality:

we have

$$\min \quad f(x) \qquad st \quad \underline{g(x)} = 0$$

$$\mathcal{L} = f(x) - \lambda^T g(x) = \mathcal{L}(x, \lambda)$$

we have solution when

$$\nabla_x \mathcal{L} = 0 \qquad \Bigg]$$
$$\nabla_\lambda \mathcal{L} = 0 \qquad \Bigg]$$

so solution is at a critical point of $\mathcal{L}$.

$\longrightarrow$ what kind of c.p. ?

(a) fix $\lambda = \hat{\lambda}$, then $\mathcal{L}(x, \hat{\lambda})$ is (locally) at a min

(b) but for fixed $x_a = \hat{x}$, $\mathcal{L}(\hat{x}, \lambda)$ is __linear__

(c) so at $\dfrac{x^c, \lambda^c}{\uparrow \text{solu}}$, $\dfrac{\partial^2 \mathcal{L}(x, \lambda)}{\partial \lambda_i \partial \lambda_j}$ is zero.

<u>i.e</u>   think   about

$$H_L = \text{Hessian of } \mathcal{L} \text{ in } x \text{ and } \lambda$$

at $x^c, \lambda^c$,   there are   some dirns
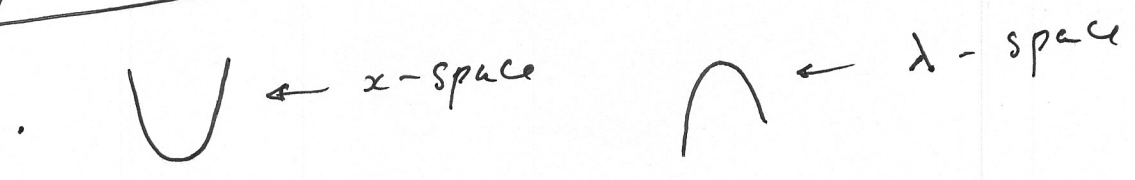(the $x$ dirns)   $\delta_x$   s.t.

$$\delta_x^T H \delta_x \geqslant 0$$

<u>AND</u>   some   dirns   ($\lambda$ dirns)

$\delta_\lambda$   st

$$\delta_\lambda^T H \delta_\lambda = 0$$

So   $x^c, \lambda^c$   <u>must</u>   be a

<u>Saddle point</u>

$\smile \leftarrow x\text{-space}$   $\frown \leftarrow \lambda\text{-space}$

This means we could think about

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

Notice:

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

(fairly easy: consider 1 constraint, then

$$\mathcal{L} = f(x) + \lambda g(x).$$

now for $q(\lambda)$ to be finite, we must have $\lambda g(x)$ bounded below; if $\lambda g(x)$ bound is greater than zero, no feasible point, so its less than zero; but then $\inf_x f(x) + \lambda g(x) \leq \inf_{st\ g(x)=0} f(x)$

and we're done
multiple dimensions follow

This is powerful because we could consider

$$\underbrace{\max_{\lambda} q(\lambda)}_{} \leq f(x^c).$$

if we have $q$, $\lambda^k$, $x^k$,

and $q(\lambda^k) - f(x^k)$ is small,

then $f(x^k) - f(x^c)$ is also small

this could help us track progress.

## Simple duals:

① ~~$\dfrac{x^T A x}{2}$~~

① $\min \dfrac{x^T x}{2}$     st $Ax + b = 0$

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

Now $\inf_x L(x, \lambda)$ occurs when

$$x - A^T \lambda = 0$$

~~so $q(\lambda) = \lambda^T A A^T \lambda$~~

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

(subs. into $L$)

(It's not always this easy)

Notice $\max_\lambda q(\lambda)$

occurs when

$$A A^T \lambda - b = 0$$

(i.e. at soln).

# Interesting example

## Variational calc.

Problem: find a P.D.F. that has a fixed set of Expectations

$$\left\{ \; E_p(f_i) = m_i \; \longleftarrow \text{known number})\right.$$

(i.e.

while maximizing entropy.

→ useful modelling idea. We observe good estimates of some expectations in data, and want model to respect these. But we know nothing else, so max entropy.

so

$$\max -\int p \log p \; dx$$

$$\text{s.t.} \quad \int p \; dx = 1$$

$$\int p \cdot f_i \; dx = m_i$$

Variational problem with constraint.

# Lagrangian

$$\mathscr{L}(p) = -\int p \log p \, dx$$

$$- \lambda_0 \left[ \int p \, dx - 1 \right]$$

$$- \sum_i \lambda_i \left[ \int p \, f_i \, dx - m_i \right]$$

we want to form 2 gradients

$\nabla_x \mathscr{L}$ is easy

$\nabla_p \mathscr{L}$ follows the case we saw earlier.

$\left( i.e \quad at \quad p^*, \left[ \dfrac{d}{d\varepsilon} \mathscr{L}(p^* + \varepsilon \varphi) \right] \bigg|_{\varepsilon = 0} = 0 \quad for \; \underline{any} \; \varphi . \right)$

$$\frac{d}{d\varepsilon} \mathscr{L}(p^* + \varepsilon \varphi) \bigg|_{\varepsilon = 0} = \int \phi \left[ -\log p^* - 1 - \lambda_0 - \sum_i \lambda_i f_i \right] dx .$$

this must be zero for $\underline{any}$ $\varphi$,

so

$$p^* \; \alpha \; e^{-\lambda_0} \cdot e^{-\sum_i \lambda_i f_i(x)}$$

This class of model <u>used</u> to be
called a <u>maximum entropy model</u>

## Fitting:

- (Old way)

  - adjust $\lambda_i$ so that

    $$\int p^* f_i \, dx = M_i$$

    $\uparrow$

    $\frac{1}{N} \sum_j f_i(x_j)$ — an estimate from <u>data</u> of this expectation

  - and $\lambda_0$ so that

    $$\int p^* \, dx = 1.$$

But,

Imagine we have a model of the form

$$p^*(x) = e^{-\lambda_0} \, e^{-\sum_i \lambda_i f_i(x)}$$

and we fit with Max likelihood

We must solve

$$\max_{} \sum_j \log p^*(x_j) \qquad \text{s.t.} \int p^*(x) dx = 1$$

(problem in $\lambda_0, \lambda_i$)

Now $\int p^*(x) dx = 1 = \int e^{-\lambda_0} e^{-\sum_i \lambda_i f_i(x)} dx$

$$= e^{-\lambda_0} \int e^{-\sum_i \lambda_i f_i(x)} dx$$

So $\lambda_0 = \log\left[ \int e^{-\lambda_i f_i(x)} dx \right] = \log Z(\lambda_i)$

So we must solve:

$$\max_{\lambda_i} \sum_j \left[ -\log Z - \sum_i \lambda_i f_i(x_j) \right] = Q(\lambda)$$

$$\frac{\partial Q}{\partial \lambda_k} = \sum_j \left[ -\frac{1}{Z} \cdot \frac{\partial Z}{\partial \lambda_k} - f_k(x_j) \right]$$

$$Z = e^{\lambda_0} = \left[ \int e^{-\sum_i \lambda_i f_i(x)} dx \right]$$

$$\frac{\partial Z}{\partial \lambda_k} = -\int e^{-\sum_i \lambda_i f_i(x)} \cdot f_k(x) \, dx$$

$$\text{so} \quad -\frac{1}{Z} \frac{\partial Z}{\partial \lambda_k} = \int e^{-\lambda_0} \cdot e^{-\sum_i \lambda_i f_i(x)} \cdot f_k(x) \, dx$$

So we must solve

$$\boxed{\int p^* f_k(x) dx = \frac{1}{N} \sum_j f_k(x_j)}$$

Exact expectations          empirical expectations