Idea:

- we have data $D$, params $\theta$ and missing data $\Delta$

- $P(D, \Delta | \theta)$ would be easy to work with

- $P(D | \theta) = \int P(D, \Delta | \theta) \, d\Delta$ is usually hard

  - eg $\log$ of sum

Algorithm: given $\theta^{(n)}$

E Step · form $Q(\theta ; \theta^{(n)}) = E_{\Delta | \theta^{(n)}} \left[ P(D, \Delta | \right.$

M Step · form $\theta^{(n+1)} = \arg\max_{\theta} Q(\theta ; \theta^{(n)}$

This takes a usefully simple form when $\log P(D, \Delta | \theta)$ is linear in

Example: dynamic model with only one clock interval

$D = Y_0^{(i)}$ ⟵ lots of different obs of $0$'th emission

$\Delta = \delta_{oj}^{(i)}$

$= \begin{cases} 1 & (x_0 = x_j) \\ 0 & (\text{otherwise}) \end{cases}$

$\log P(D, \Delta | \theta) = \log P(D | \Delta, \theta)$

$\qquad\qquad\qquad + \log P(\Delta | \theta)$

this term determined by emission mode

this is known prior

$$\log P(D_{1:\Delta}|\theta)\ P(Y_0|x_0,\theta)$$

$$= \sum_{i,j} S_j^i \left[ (Y_0 \text{ is } \mu(X_j^{(i)},\theta))\frac{\Sigma_{0i}^{-1}}{2}(Y_0, \Sigma^{(i)}(x),\theta) \right.$$

$$+ \quad \text{constant first term} \qquad \text{like kanosstritch}$$

here $S_j^i$ converges

so $\log P(Y_0|S_{oj},\theta)$

$$= (Y_0 - \mu(x_j,\theta))\frac{\Sigma^{-1}}{2}(Y_0 - \mu(x_j,\theta))$$

$$+ \quad \underbrace{\log K_n}$$

$\curvearrowright$ but this is not a
fn of $\theta, S$

$$\log P(D, \Delta | \theta)$$

$$= \sum_{i,j} \delta_j^i \left[ (Y_o^i - M(x_j, \theta))^T \frac{\Sigma^{-1}}{2} (Y_o^i - M(x_j, \theta)) \right.$$

$$+ \underline{\hspace{2cm}} \text{constant terms}$$

This acts like a switch

Now $$E[\delta_j^i]_{\delta_j^i | D, \theta}$$

$$= 1 \cdot P(\delta_j^i = 1 | D, \theta) + 0 \cdot$$

$$P(\delta_j^i = 1 | D, \theta) = \frac{P(Y_o^{(i)} | \delta_j^i, \theta) \, P(\delta_j^i | \theta)}{\sum_u P(Y_o^{(i)} | \delta_u^i, \theta) P(\delta_u^i | \theta)}$$

In this case, we get

$$P(\gamma_j^i \mid D, \theta) = \frac{\exp\left[(Y_0^{(i)} - \mu(x_j^i; \theta)) \sum_{2}^{l-1} \times (Y_0^{(i)} - \mu(x_j^i; \theta))\right] \times \pi_j}{\sum_u (\text{terms as above})}$$

So the E step is straightforward.

M-Step

- depends on $\mu(x_j^i, \theta)$
  (form of function)

eg. $\mu(x_j^i, \theta) = \theta \cdot X_j^i$

and this has a 1 in j'th location and zeros elsewhere

this case is one mean per state

- Now look at LLH as fn of j'th mean

$$\sum_i P(\delta_i^j | D, \theta^{(n)}) \cdot \left[ (Y_0^{(i)} - M_j)' \sum_2^{-1} (Y_0^{(i)} - M_j) \right]$$

+ other terms that don't depend on $M_j$

But this is just a weighted mean.

<u>case 2:</u>

$$P(Y_o^{\#} | X_o) \text{ is a table}$$

because $Y$ is discrete

Maximization is by weighted counts

<u>Example 2:</u>   sequences, multiple
  examples

$$P(Y_o^{(i)} \cdots Y_n^{(i)}, S_{oj}^i \cdots \delta_{nj}^i | \theta)$$

$$= P(Y_o^{(i)} \cdots Y_n^{(i)} | S_{oj}^i \cdots \delta_{nj}^i, \theta) \times$$

$$P(\delta_{oj}^i \cdots \delta_{nj}^i, \theta)$$

− we are assuming that dynamics
  are known, so second term
  is fixed

$$\log \; P(Y_0^{(i)} \cdots Y_n^{(i)} \mid S_{0j}^i \cdots \delta_{nj}^i, \Theta)$$

Switch $\triangle$

1 per clock tick.

$$= \sum_j \left[ \log P(Y_0^{(i)} \mid X_0 = x_j, \Theta) \right] \cdot \delta_{0j}^{(i)}$$

$$+ \sum_j \left[ \log P(Y_1^{(i)} \mid X_\phi = x_j, \Theta) \right] S_{1j}^i$$

$$+ \quad \vdots$$

Now consider the E step

$$P(S_{\ell j}^i = 1 \mid Y_0^{(i)} \cdots Y_n^{(i)}, \Theta)$$

$$= \; P(X_\ell^i = x_j \mid Y_0^{(i)} \cdots Y_n^{(i)}, \Theta)$$

$$= \; \frac{P(X_\ell^i = x_j, Y_0^{(i)} \cdots Y_n^{(i)}, \Theta)}{P(Y_0^{(i)} \cdots Y_n^{(i)}, \Theta)}$$

we k

the

Constrained optimization:

$$\min \quad f(x) \qquad st \qquad c_i(x) = 0$$
$$g_i(x) \geqslant 0$$

## Lagrangian

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

(here $\lambda$ is a vector of constraints
whose elements csp to~~g~~ ineq $(\lambda^{(i)}$
or eq constraints $(\lambda^{(e)})$

## Necessary conditions       (KKT conds)

$$\nabla_x \mathcal{L} = 0 \qquad\qquad g_i(x) \geqslant 0$$

$$c_i(x) = 0 \qquad\qquad \lambda^{(i)} \geqslant 0$$

$$\lambda_i^{(e)} c_i = 0 \qquad\qquad \lambda_i^{(i)} g_i = 0$$

- Duality  Assume  inequality  constraints  only.

example:
$$\min f(x) \quad st \quad g_i(x) \geqslant 0$$

- Assume $\min -\frac{a^T x}{2}$ is st.convex $Ax = b$

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

define dual objective fn to be
from first condition:

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

i.e. $q(\lambda) = \max_x \mathcal{L}(x, \lambda)$.

substitute
on domain such that $\frac{\lambda^T AA \lambda}{2} - \lambda (AA x - b) \} = q(\lambda) \} = \alpha$

dual problem:
$$\xrightarrow{\quad \lambda \quad} \xrightarrow{\quad x \quad}$$

$$\max_x q(\lambda)$$

Knowledge of $\lambda$ wrt values is $\lambda \geqslant 0$
powerful!

Thm: $q$ is concave, domain is convex

(straightforward)

Thm: for feasible $x$, any $\lambda$

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

(straightforward)

Thm: suppose $x$ is soln of primal, $f$ and $-g_i$ are convex; then $\lambda$ such that $(x, \lambda)$ satisfies KKT is a soln of dual

~~Thm: with~~ other way round requires stronger technical conds

Thm: value of dual $\leq$ value of primal

Common application: in important cases, one may be able to write the dual directly.

## SVM

$$\min \quad \frac{w'w}{2}$$

$$\text{st} \quad y_i(w'x_i + b) \geqslant 1$$

Primal form, Separable

$$\mathcal{L}(w,\lambda) = \frac{w'w}{2} - \sum_i \lambda_i \left\{ \left[ y_i(w'x_i + b) \right] - 1 \right\}$$

$$\nabla_w \mathcal{L} = 0 = w - \sum_i \lambda_i \left\{ \left[ y_i x_i \right] \right\}$$

$$\nabla_b \mathcal{L} = 0 = - \sum_i \lambda_i y_i$$

Subst :

P

Subst

$$\mathcal{L}_D = \sum_i \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j \left[ y_i y_j (x_i^T x_j) \right]$$

Notice constraints

$$\sum_i \lambda_i y_i = 0$$

$$\lambda_i \geqslant 0$$

and we must <u>max</u> this in $\lambda$

<u>If</u> there is an fp for primal, the
max is soln to primal

i.e.     Value (Dual) = Value (Primal)

What if data is not separable?

$$\min \quad \frac{w'w}{2} + C \sum_i \xi_i$$

$$\text{st} \qquad y_i(w'x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

Primal prob

$\xi_i$ are __slack variables__

$$\mathcal{L}_p = \frac{w'w}{2} + C \sum_i \xi_i - \sum_i \lambda_i [y_i(w'x_i + b) - 1 + \xi_i]$$
$$- \sum_i \mu_i \xi_i$$

$$\nabla_w \mathcal{L}_p = w - \sum_i \lambda_i y_i x_i = 0$$

$$\nabla_b \mathcal{L}_p = 0 = -\sum_i \lambda_i y_i$$

$$\nabla_{\xi_i} \mathcal{L}_p = C - \lambda_i - \mu_i = 0 \quad \Big] \rightarrow \text{this gets rid of } \xi_i$$

Met/ we have

$$\mathcal{L}_D = \sum_i \lambda_i - \frac{1}{2} \sum_{ij} y_i y_j \lambda_i \lambda_j x_i' x_j$$

Subject to

$$\sum_i \lambda_i y_i = 0$$

$$0 \leq \lambda_i \leq C$$

Notice that $\xi_i$ can be interpreted as a <u>loss</u>

$$\text{hinge loss}(y, y_p) = \max\left(0, 1 - y_i y_p\right)$$

Methods:

## Quadratic penalty method

(assume equalities)

$$\min_{x} \quad f(x) + \frac{\mu}{2} \sum_i c_i^2(x) = Q_k(x)$$

and drive $\mu \to \infty$, resolve

Notice          at soln

$$\nabla_x Q_k \approx 0 = \nabla f + \sum_i (\mu_k c_i(x_k)) \nabla c_i(x)$$

By inspection, this would match
$\nabla_x \mathcal{L} = 0$,  if

$$-\mu_k c_i = \lambda_i^*$$

Which suggests that at conv $c_i = -\frac{\lambda_i^*}{\mu_k}$

This looks OK, because $\mu_K \to \infty$, but not exact. Also $M_K \to \infty$ creates major probs w/ Hessian

## Augmented Lagrangian method

Consider

$$\mathcal{L}_A (x, \lambda; \mu) = f - \sum_i \lambda_i c_i + \frac{\mu}{2} \sum_i c_i^2$$

- have an est of $\lambda^K$, $\mu_K$, get $x^*$

- at $x^*$ $\nabla_x \mathcal{L}_A = 0 = \nabla f - \sum_i (\lambda_i^K - \mu_K c_i) \nabla c_i$

- This suggests $\lambda_i^* \approx (\lambda_i^K - \mu_K c_i)$

and $c_i \approx -\frac{1}{\mu_K} . [\lambda_i^* - \lambda_i^K]$

which suggests moving $\lambda_i \to \lambda_i^*$

But we have a good est:

$$\lambda_i^* \approx \left( \lambda_i^k - \mu_k c_i \right)$$

so update ests, go again.

i) Method converges w/o increasing $\mu_k$ indefinitely

# Conjugate gradient

## We now have:

Start: $x_0$ , $r_0 = A x_0 - b$ , $P_0 = -r_0$

## Step:

$$x_{K+1} = x_K + \alpha_K P_K$$

$$\alpha_K = - \frac{r_K' A P_K}{P_K' A P_K}$$

$$P_{K+1} = -r_{K+1} + \beta_{K+1} P_K$$

$$\beta_{K+1} = \frac{P_K' A r_{K+1}}{P_K' A P_K}$$

We can make this more efficient

Conju
This gives

$$\frac{1}{2}\left[ (x_k + \alpha_k p_k)' A (x_k + \alpha_k p_k) \right]$$
$$- b_k (x_k + \alpha_k p_k)$$

Min is at:

$$- \frac{(A x_k - b)' p_k}{p_k' A p_k}$$

write

$$r_k = A x_k - b$$

so $\quad \alpha_k = - \frac{r_k' p_k}{p_k' A p_k}$

## gate gradient (simple form)

Start: $x_0$, $r_0 = Ax_0 - b$, $p_0 = -r_0$

Step:

$$x_{k+1} = x_k + \alpha_k p_k$$

$$\alpha_k = \frac{-r_k' p_k}{p_k' A p_k}$$

$$p_{k+1} = -r_{k+1} + \beta_{k+1} p_k$$

$$\beta_{k+1} = \frac{r_{k+1}' A p_k}{p_k^T A p_k}$$

$$r_{k+1} = r_k + \alpha_k A p_k$$

Conjugate gradient.

Cleaner form:

- By properties, we have

$$\alpha_{K+1} = \frac{r_K' r_K}{p_K' A p_K}$$

- Now $\alpha_K A p_K = r_{K+1} - r_K$

So $\beta_{K+1} = r_{K+1}' \left( \underbrace{r_{K+1} - r_K}_{\alpha_K} \right) \cdot \frac{1}{p_K' A p_K}$

$$= \frac{r_{K+1}' (r_{K+1} - r_K)}{r_K' r_K}$$

$$= \frac{r_{K+1}' r_{K+1}}{r_K' r_K} \qquad \text{(By properties)}$$

## Properties of conj. direction

$$r_k' \, p_i \;=\; 0, \qquad \forall i < k$$

(Show this by induction)

$$r_k' \, r_i \;=\; 0 \qquad \forall i < k$$

(thm 5.3 at end).

Conjugate direction in incremental form

Start with $x_0, P_0$

$$x_1 = x_0 + \alpha_0 P_0$$

now min wrt $\alpha_0$

to get

$$\frac{(Ax_0 - b)' P_0}{P_0' A P_0} = \alpha_0$$

write

$$r_k = (Ax_k - b)$$

and get

$$x_{k+1} = x_k + \alpha_k P_k$$

$$\alpha_k = \frac{r_k' P_k}{P_k' A P_k}$$

$$r_{k+1} = r_k + \alpha_k A P_k$$

# Conjugate Direction methods:

- a set of vectors, $P_0 \cdots P_n$ is conjugate for $A$ positive definite if

$$P_i' A P_j = 0 \qquad i \neq j$$

- Assume we wish to __Min__

$$\frac{x' A x}{2} - b' x$$

- useful because:

a) Solution to $Ax = b$
   for $A$ p.d.

b) $\underset{x}{\text{Min}} \quad \| Ux - b \|^2$ is like this

- now write

$$x = \alpha_0 P_0 + \alpha_1 P_1 + \cdots$$