

CONJUGATE GRADIENT & NESTED SAMPLING

John Skilling

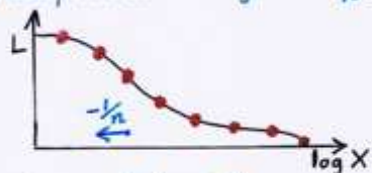
Cambridge, Sept 2008.

Nested sampling proceeds by compressing enclosed volume  $X$ .



New  $L^*$ , compression  $\Delta \log X = -\frac{1}{n}$ .

It plots  $L(x)$ , giving evidence  $Z = \int L dx$  and posterior, weight  $(z_i) = L_i \Delta X_i / Z$ .



It needs uniform samples within  $L(x) > L^*$ .

We have relied on traditional exploration (Gibbs etc.). This ignores gradients.

Computation of  $L(x)$  can easily get  $\nabla L$  also: a million times more info. for twice the cost.

But we seem to lose detailed balance.



Actually, we CAN do it!

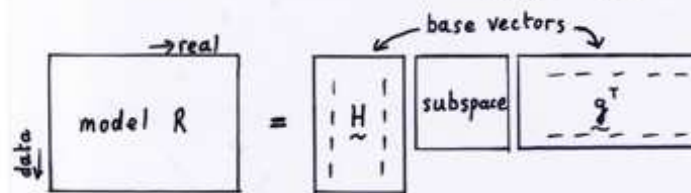
The archetype problem is  $\log L = -\frac{1}{2} X^2$ .  
 $X^2 = |Rz - D|^2$ ,  $D = \text{data}$ ,  $R = \text{"Opus"}$ .  
 $\nabla L = -R^T(Rz - D)$ ,  $\nabla \nabla L = -R^T R$ .

MATRIX MODELLING (secrets from 20 yrs ago)

Seek to simulate matrix operations (inversion ...) on

$$A = R^T R$$

Feedback  $\swarrow$   $\nwarrow$  "Opus" real-to-data  
 "Tropus" data-to-real



Wider subspace  $\rightarrow$  closer modelling

Generate vectors from a seed:

$$\begin{aligned} g_1 &= \text{seed} && \rightarrow H_1 \sim Rg_1 \\ g_2 &\sim \{g_1, R^T H_1\} && \rightarrow H_2 \sim \{H_1, Rg_2\} \\ g_3 &\sim \{g_1, g_2, R^T H_2\} && \rightarrow H_3 \sim \{H_1, H_2, Rg_3\} \\ & && \dots \end{aligned}$$

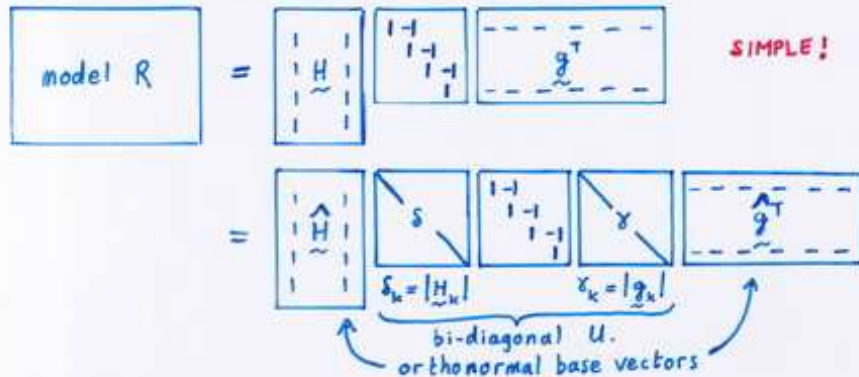
Orthogonal sets of vectors preferred.

Simplify the code to **kernel**

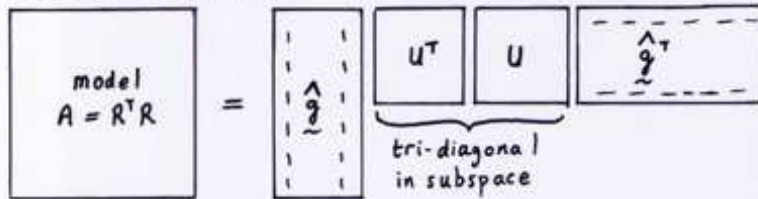
$$\begin{aligned} H_0 &= 0, & g_1 &= \text{seed} \\ \forall k: & H_k = H_{k-1} + \frac{Rg_k}{g_k^T g_k}, & g_{k+1} &= g_k - \frac{R^T H_k}{H_k^T H_k} \end{aligned}$$

Check:  $g_i \perp g_j$ ,  $H_i \perp H_j$

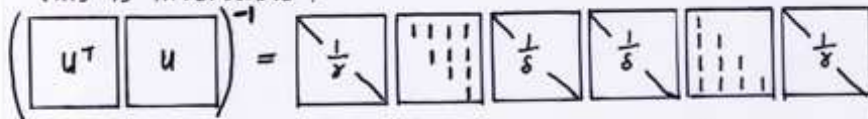
These "nicest" vectors have



Feedback matrix is



This is invertible:



Check: With seed  $\underline{b}$ , this models  $\underline{x} = A^{-1}\underline{b}$  by maximising  $q = 2\underline{x}^T \underline{b} - \underline{x}^T A \underline{x}$  in subspace.

Merge accumulation of  $\underline{x}$  and  $q$  into kernel. Get standard **conjugate gradient** for  $A^{-1}\underline{b}$ .

$\underline{g}_1 = \underline{b}, \underline{h}_0 = \underline{0}, \underline{x}_0 = \underline{0}, q_0 = 0$

$\forall k: \chi_k^2 = \underline{g}_k^T \underline{g}_k, \underline{h}_k = \underline{h}_{k-1} + \underline{g}_k / \delta_k^2, \underline{H} = R \underline{h}$

$\delta_k^2 = \underline{h}_k^T A \underline{h}_k, \underline{g}_{k+1} = \underline{g}_k - A \underline{h}_k / \delta_k^2$

$q_k = q_{k-1} + 1/\delta_k^2, \underline{x}_k = \underline{x}_{k-1} + \underline{h}_k / \delta_k^2$

- As subspace widens,  $0 = q_0 < q_1 < q_2 < \dots \leq Q = \underline{b}^T A^{-1} \underline{b}$
- $\underline{g}$  are gradients  $\nabla Q$  at successive maxima  $\underline{x}$ .
- Precursors  $\underline{h}$  of  $\underline{H} = R \underline{h}$  are conjugate ( $\underline{h}_i^T A \underline{h}_j = 0$ ) — hence name.

Modelling (J.S.)  $\iff$  Maximisation ("Establishment")



Regularisation (to stabilise small eigenvalues of A)

Instead of  $\max q = 2x^T b - x^T A x \leq Q = b^T A^{-1} b$

seek  $\max q' = 2x^T b - x^T B x \leq Q' = b^T B^{-1} b$ ,  $B = \alpha I + A$

Base vectors  $\hat{g}$  don't change, because subspace is same.

Scalars  $\gamma$  and  $\delta$  do change, but can be simulated "free".

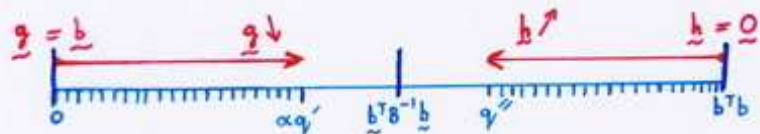
Get  $\gamma'$  and  $\delta'$ , hence  $x'$  and  $q'$ , for free.

Can also get for free

$$\max q'' = 2x^T A b - x^T A B x \leq Q'' = b^T A B^{-1} b$$

But  $\alpha Q' + Q'' = b^T b = \text{known}$ , so

$\alpha q''$  gives upper limit on  $b^T B^{-1} b$ .



Termination: Stop when 99% of  $b^T b$  accounted.

Extensions

Get controlled evaluation of  $(\alpha I + A)^{-1}$  on any  $b$

Also get  $(\beta I + A)^{-1}$ , even better controlled if  $\beta > \alpha$ .

Superpose: get  $(\beta_1 I + A)^{-1} + (\beta_2 I + A)^{-1}$ .

Generalise: get  $f(A) = \int_{\alpha}^{\infty} d\beta W(\beta) (\beta I + A)^{-1}$ ,  $W \geq 0$ .

for "simulatable" function  $f(t) = \int_{\alpha}^{\infty} \frac{W(\beta) d\beta}{\beta + t}$

Example:  $W(\beta) = (\beta - \alpha)^{-k}$ ,  $0 < k < 1$ , gives  $f(t) = (\alpha + t)^{-k}$

In particular,  $k = \frac{1}{2}$  gives  $\underline{s} = (\alpha I + A)^{-\frac{1}{2}} \underline{r}$

With  $r_i$  unit normal, covariance  $\langle \underline{s} \underline{s}^T \rangle = (\alpha I + A)^{-1}$ .

Example:  $k \rightarrow 0$  (with offset) gives  $f(t) = -\log(\alpha + t)$

In particular,  $\underline{r}^T \log(\alpha I + A) \underline{r}$  with  $r_i$  unit normal

$$= \sum r_i^2 \log(\alpha + \lambda_i) \quad \text{in eigencoords.}$$

with expectation  $\Sigma \log(\alpha + \lambda_i) = \log \det(\alpha I + A)$

To evaluate, diagonalise U to

$$U = \begin{bmatrix} & & & \\ & s & & \\ & & \begin{matrix} 1 & & \\ & \ddots & \\ & & 1 \end{matrix} & \\ & & & \gamma \end{bmatrix} = \begin{bmatrix} | & & | \\ \text{Lvec} & & \text{Rvec}^T \\ | & & | \end{bmatrix}$$

so that

$$f \left( \begin{bmatrix} U^T & U \end{bmatrix} \right) = \begin{bmatrix} | & & | \\ \text{Rvec} & & f(\lambda^2) \\ | & & | \end{bmatrix}$$

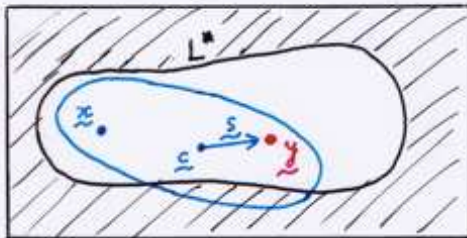
Upper and lower limits remain available.

NESTED SAMPLING

Space of  $\underline{x}$  has flat prior and convex

$$L(\underline{x}) = \log \text{Pr}(\text{data} | \underline{x})$$

Current position is  $\underline{x}$  with  $L(\underline{x}) = L_0 > L^*$ .



Seek new random point  $\underline{y}$  within  $L(\underline{y}) > L^*$ .

Proposal: Build local quadratic model, regularised.

$$\tilde{L}(\underline{y}) = L_0 + (\underline{y} - \underline{x})^T \underline{g} - \frac{1}{2} (\underline{y} - \underline{x})^T B (\underline{y} - \underline{x})$$

where  $\underline{g} = \nabla L$ ,  $B = \alpha I + \nabla \nabla L$  at  $\underline{x}$ .

Centre of model is

$$\underline{c} = \underline{x} + B^{-1} \underline{g} \quad \text{with} \quad L_c = L_0 + \frac{1}{2} \underline{g}^T B^{-1} \underline{g}$$

Propose  $\underline{y} = \underline{c} + \underline{s}$

within ellipsoid  $L(\underline{y}) = L_c + \frac{1}{2} \underline{s}^T B \underline{s} > L^*$ .

Offset  $\underline{s}$  constrained within ellipsoid

$$\underline{s}^T B \underline{s} < 2(L_c - L^*) \equiv a^2$$

Random point on surface is  $a \cdot B^{-\frac{1}{2}} \underline{\underline{r}}$

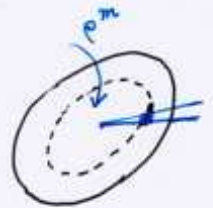
where  $\underline{\underline{r}}$  is random on sphere  $\underline{\underline{r}}^T \underline{\underline{r}} = 1$ .

Random point within volume is

$$\underline{s} = \rho \cdot a B^{-\frac{1}{2}} \underline{\underline{r}}$$

where  $\rho^m = \text{Uniform}(0, 1)$

( $m = \text{dimension}$ ).



Propose

$$\underline{y} = \underline{x} + B^{-1} \underline{g} + \rho a B^{-\frac{1}{2}} \underline{\underline{r}}$$

$$B = \alpha I + A$$

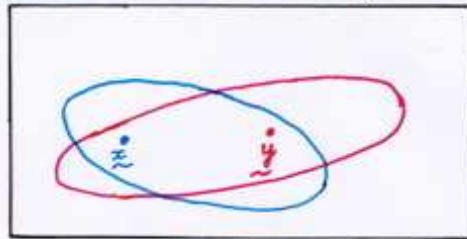
$$a = \sqrt{2(L_0 + \frac{1}{2} \underline{g}^T B^{-1} \underline{g} - L^*)}$$

Volume of destination ellipsoid

$$V = \frac{\pi^{m/2}}{(m/2)!} \frac{a^m}{\sqrt{\det B}}$$

Acceptance ?

If  $L(\underline{y})$  outside  $L^*$ , reject.  
 If OK, build new model around  $\underline{y}$ .



If original  $\underline{x}$  outside new ellipsoid, reject.  
 If OK, accept in ratio of volumes (detailed balance).

**But**  $V = \exp(m \times \dots)$ , so volume ratio not  $O(1)$ .

Fixup: Adjust  $\alpha$  to maintain

$$\frac{\log V}{m} = \log a - \frac{1}{2m} \log \det B = \log a - \frac{1}{2} \underline{r}^T \log B \underline{r}$$

= typical length in model.

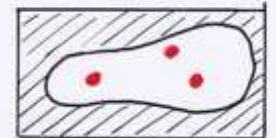
Exploration loop:

Propose  $\underline{y} = \underline{x} + B^{-1} \underline{g} + \rho \alpha B^{-\frac{1}{2}} \underline{r}$  at current  $\alpha$ .  
 If  $L(\underline{y})$  outside  $L^*$ , reject.  
 If OK, build new ellipsoid of same  $V$  (adjust  $\alpha$ ).  
 If new ellipsoid does not contain  $\underline{x}$ , reject.  
 If OK, accept  $\underline{y}$ .

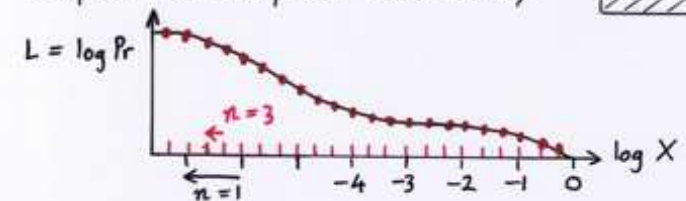
Oddity: Need to impose dimension  $m$  in  $\rho = \text{Uniform}^{1/m}$ .

Multiple samples

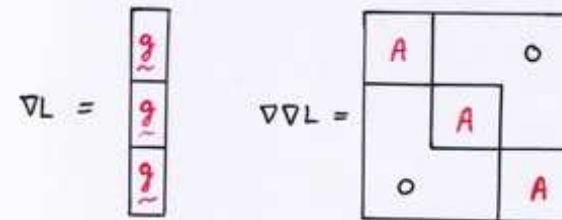
Nested sampling can use several samples to compress more slowly



e.g.  $n=3$ .



Samples can evolve separately  
 — or together in an ensemble with



Random detail apart, the  $\alpha$ 's stay as before, and algorithm proceeds as before, in bigger space.

Each individual sample evolves as before, except

$$\rho = \text{Uniform}^{1/3m}$$



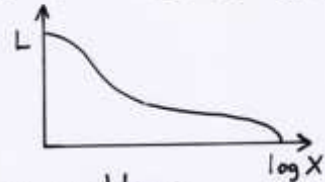
Sampling is biased outward, giving slower compression.



## SUMMARY

Nested sampling was designed for evidence value.

Actually, it plots  $L(x)$ .



Evidence and posterior  
are by-products.

Can cope with multi-phase problems.

Exploration need not rely on guessed directions,  
mostly orthogonal to useful directions.

Can use gradients and curvatures.

Control by volume, not point-wise values.

Compression can be arbitrarily slow (accuracy)  
or arbitrarily fast (speed).

High-dimensional problems need not be slow.

$\text{CPU} \propto \text{size}$

---

J.S. Sept 2008