

# Probabilistic Numerics

## IV – Linear Algebra

Philipp Hennig

Dobbiaco Summer School

21 June 2017



MAX-PLANCK-GESELLSCHAFT

Research Group for Probabilistic Numerics  
Max Planck Institute for Intelligent Systems  
Tübingen, Germany



Some of the presented work was supported by  
the Emmy Noether Programme of the DFG

## Yesterday:

- probabilities: a rigorous notion of uncertainty
- Gaussians provide a computationally efficient algebra of uncertainty
- classic methods for **Quadrature** and **ODEs** can be re-formulated as MAP/mean inference in Gaussian models. The associated covariances yield calibrated (worst-case) error bars
- these probabilistic formulations also provide additional functionality that would be tricky in the classic formulations

all slides from yesterday at ( $x = 1, 2, 3$ )  
<http://tinyurl.com/Dobbiaco-Hennig-x>

## Now:

- a probabilistic formulation of **iterative linear solvers**
- a connection to nonlinear optimization

# Problem Setting

The least-squares problem

$$A \begin{matrix} i \\ k \\ j \end{matrix} = \begin{matrix} x \\ k \\ j \end{matrix} = b \begin{matrix} i \\ j \end{matrix}$$

$$Ax = b \quad A \in \mathbb{R}^{N \times N} \quad b \in \mathbb{R}^N$$

- in machine learning and statistics:

$$A = \Phi\Phi^\top + \sigma^2 I \quad \text{symmetric positive definite}$$

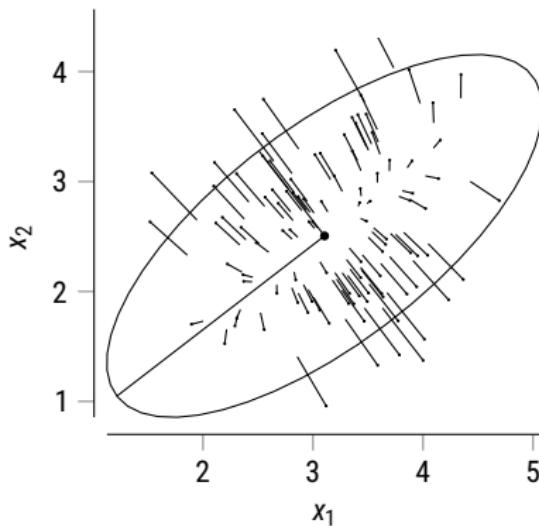
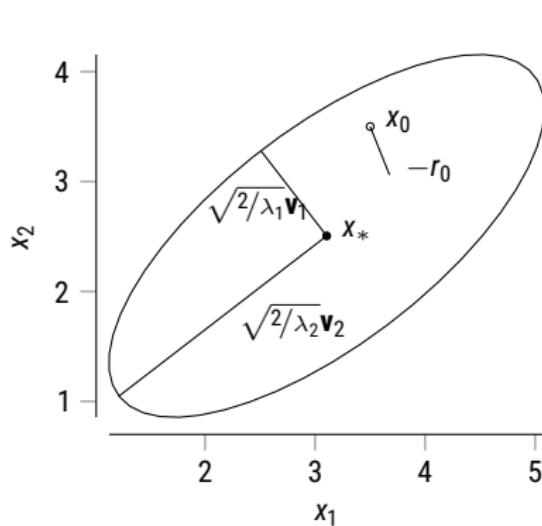
- aka. the **least-squares problem**, because then

$$x = \arg \min_{\tilde{x}} f(\tilde{x}) \quad \text{for} \quad f(\tilde{x}) = \frac{1}{2}\tilde{x}^\top A\tilde{x} - \tilde{x}^\top b$$

with  $\nabla_{\tilde{x}} f = A\tilde{x} - b$ . Define  $H := A^{-1}$ .

# A Pictorial View

SPD linear problems are quadratic optimization problems



$$Ax_* = b$$

$$A = [\mathbf{v}_1 \quad \mathbf{v}_2] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \end{bmatrix}$$

$$f(\tilde{x}) = 1/2 \tilde{x}^\top A \tilde{x} - \tilde{x}^\top b$$

$$r(x_i) := Ax_i - b = \nabla f(x_i)$$

# Method of Conjugate Gradients

our target for the day

□ Hestenes & Stiefel, 1952

```
1 procedure CG( $A(\cdot)$ ,  $b$ ,  $x_0$ )
2    $r_0 = Ax_0 - b$                                 // initial gradient
3    $d_0 = 0, \beta_0 = 0$                           // for notational clarity
4   for  $i = 1, \dots, N$  do
5      $d_i = -r_{i-1} + \beta_{i-1} p_{i-1}$           // compute direction
6      $Z_i = Ad_i$                                 // observe
7      $\alpha_i = -d_i^T r_{i-1} / d_i^T Z_i$         // optimal step-size
8      $s_i = \alpha_i d_i$                           // re-scale step
9      $y_i = \alpha_i Z_i$                           // re-scale observation
10     $x_i = x_{i-1} + s_i$                         // update estimate for  $x$ 
11     $r_i = r_{i-1} + y_i$                         // new gradient at  $x_i$ 
12     $\beta_i = r_i^T r_i / r_{i-1}^T r_{i-1}$       // compute conjugate correction
13  end for
14 end procedure
```

- produces **iterate sequence**  $x_0, x_1, \dots, x_i$ , so that  $x_N = x$
- cost dominated by line 6. But no required access to  $A$  per se, only  $A(\cdot)$

# A First Idea

a probabilistic linear solver

Idea: treat  $y_i = As_i$  as **observations** of  $A$ , **infer**  $H = A^{-1}$ .

```
1 procedure LINSOLVE_HOPEFUL( $A(\cdot)$ ,  $b$ ,  $p(A)$ )
2    $x_0 = H_0 b$                                      // initial guess
3    $r_0 = Ax_0 - b$                                // initial gradient
4   for  $i = 1, \dots, N$  do
5      $s_i = -H_{i-1} r_{i-1}$                       // compute direction
6      $y_i = As_i$                                 // observe
7      $x_i = x_{i-1} + s_i$                          // update estimate for  $x$ 
8      $r_i = r_{i-1} + y_i$                           // new gradient at  $x_i$ 
9      $H_i = \text{INFER}(H \mid Y_i, S_i, p(A))$     // estimate  $H$ . To be expanded below
10  end for
11 end procedure
```

- why not  $x_i = H_i b$ ? Because classic methods allow for generic  $x_0$ . Note

$$x_{i+1} = x_i - H_i(Ax_i - b) = x_i - H_i \left( \sum_{j \leq i} As_j + Ax_0 - b \right) = H_i b + (I - H_i A)x_0.$$

# A Refinement

analytic step-size give our generic probabilistic iterative solver

□ Hennig, SIOPT 2015

```
1 procedure LIN SOLVE_HOPEFUL( $A(\cdot)$ ,  $b$ ,  $p(A)$ )
2    $x_0 = H_0 b$                                      // initial guess
3    $r_0 = Ax_0 - b$                                // initial gradient
4   for  $i = 1, \dots, N$  do
5      $s_i = -H_{i-1} r_{i-1}$                       // compute direction
6      $y_i = As_i$                                 // observe
7      $x_i = x_{i-1} + s_i$                         // update estimate for  $x$ 
8      $r_i = r_{i-1} + y_i$                          // new gradient at  $x_i$ 
9      $H_i = \text{INFER}(H \mid Y_i, S_i, p(A))$     // estimate  $H$ . To be expanded below
10  end for
11 end procedure
```

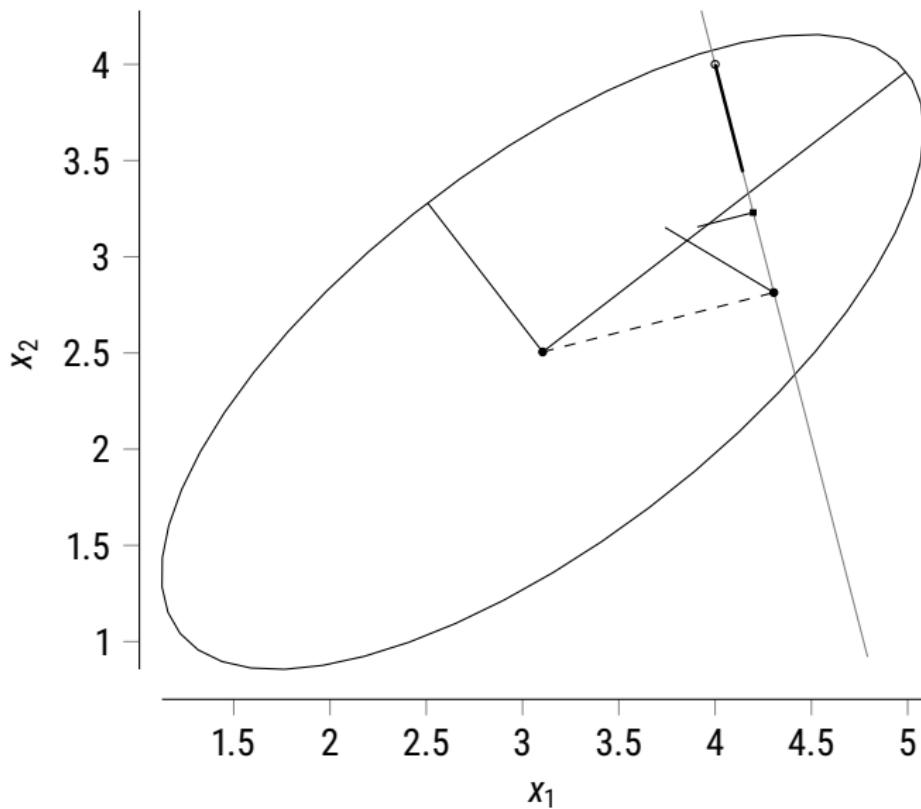
line 6 gives additional information before update (7), fixes step-size:

$$\begin{aligned}\frac{\partial f(x_{i-1} + \alpha_i d_i)}{\partial \alpha_i} &= \alpha_i d_i^T A d_i + d_i^T (A x_{i-1} - b) = \alpha_i d_i^T A d_i + d_i^T r_{i-1} \\ \Rightarrow \quad \alpha_i &= -\frac{d_i^T r_{i-1}}{d_i^T z_i}.\end{aligned}$$

# A Refinement

analytic step-size give our generic probabilistic iterative solver

Hennig, SIOPT 2015



# A Refinement

analytic step-size give our generic probabilistic iterative solver

□ Hennig, SIOPT 2015

```
1 procedure LINOLVE_PROB( $A(\cdot)$ ,  $b$ ,  $p(A)$ )
2    $x_0 = H_0 b$                                 // initial guess
3    $r_0 = Ax_0 - b$                           // initial gradient
4   for  $i = 1, \dots, N$  do
5      $d_i = -H_{i-1}r_{i-1}$                   // compute direction
6      $z_i = Ad_i$                             // observe
7      $\alpha_i = -d_i^T r_{i-1} / d_i^T z_i$     // optimal step-size
8      $s_i = \alpha_i d_i$                       // re-scale step
9      $y_i = \alpha_i z_i$                       // re-scale observation
10     $x_i = x_{i-1} + s_i$                     // update estimate for x
11     $r_i = r_{i-1} + y_i$                     // new gradient at  $x_i$ 
12     $H_i = \text{INFER}(H \mid Y_i, S_i, p(A))$  // estimate  $H$ 
13  end for
14 end procedure
```

# A Refinement

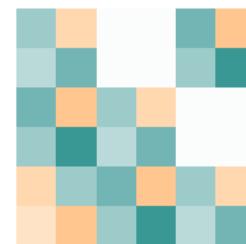
analytic step-size give our generic probabilistic iterative solver

□ Hennig, SIOPT 2015

```
1 procedure CG( $A(\cdot)$ ,  $b$ ,  $x_0$ )
2    $d_0 = 0, \beta_0 = 0$                                 // for notational clarity
3    $r_0 = Ax_0 - b$                                 // initial gradient
4   for  $i = 1, \dots, N$  do
5      $d_i = -r_{i-1} + \beta_{i-1}p_{i-1}$                 // compute direction
6      $z_i = Ad_i$                                     // observe
7      $\alpha_i = -d_i^T r_{i-1} / d_i^T z_i$             // optimal step-size
8      $s_i = \alpha_i d_i$                             // re-scale step
9      $y_i = \alpha_i z_i$                             // re-scale observation
10     $x_i = x_{i-1} + s_i$                           // update estimate for x
11     $r_i = r_{i-1} + y_i$                           // new gradient at  $x_i$ 
12     $\beta_i = r_i^T r_i / r_{i-1}^T r_{i-1}$         // compute conjugate correction
13  end for
14 end procedure
```

# Some Notation

unfortunately, computational linear algebra is tedious and often confusing



- vectorized matrices

$$A \in \mathbb{R}^{N \times N} \Rightarrow \vec{A} \in \mathbb{R}^{N^2 \times 1}. \quad [A]_{ij} = [\vec{A}]_{(ij)}$$

- the Kronecker product

$$A \in \mathbb{R}^{N_A \times M_A}, B \in \mathbb{R}^{N_B \times M_B} \Rightarrow A \otimes B \in \mathbb{R}^{N_A N_B \times M_A M_B}$$

$$[A \otimes B]_{(ij),(k\ell)} = [A]_{ik} \cdot [B]_{j\ell}$$

$$(A \otimes B) \vec{C} = \sum_{k\ell} A_{ik} C_{k\ell} B_{j\ell} = \overrightarrow{ACB^\top}$$

$$(A \otimes B) \cdot (C \otimes D) = (AC) \otimes (BD) \quad (A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

# Inference on Matrix Elements

general Gaussians on matrix objects work in principle, but are **way** too expensive

$$p(A) = \mathcal{N}(A; A_0, \Sigma_0) = \frac{1}{\sqrt{2\pi}|\Sigma|^{D/2}} \exp \left( (\overrightarrow{A - A_0})^\top \Sigma_0^{-1} (\overrightarrow{A - A_0}) \right)$$

$$p(S, Y | A) = \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) = \lim_{\beta \rightarrow 0} \mathcal{N}(Y; (I \otimes S^\top) \vec{A}, \beta \Lambda)$$

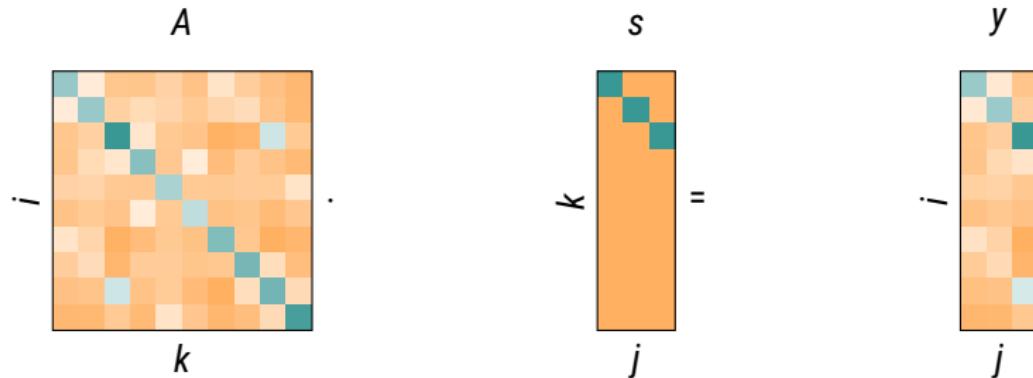
$$p_A(A | Y, S) = \mathcal{N}(A; A_M, \Sigma_M) \quad \text{with}$$

$$\vec{A}_M := \vec{A}_0 + \Sigma_0 (I \otimes S) ((I \otimes S^\top) \Sigma_0 (I \otimes S))^{-1} (\vec{Y} - (I \otimes S^\top) \vec{A}_0), \text{ and}$$

$$\Sigma_M := \Sigma_0 - \Sigma_0 (I \otimes S) \underbrace{((I \otimes S^\top) \Sigma_0 (I \otimes S))^{-1} (I \otimes S^\top) \Sigma_0}_{:= G_M \in \mathbb{R}^{NM \times NM}}$$

# Kronecker Products to the Rescue!

low-rank posteriors



$$p(A) = \mathcal{N}(A; A_0, V_0 \otimes W_0) \quad p(S, Y | A) = \delta(\vec{Y} - (I \otimes S^\top) \vec{A})$$

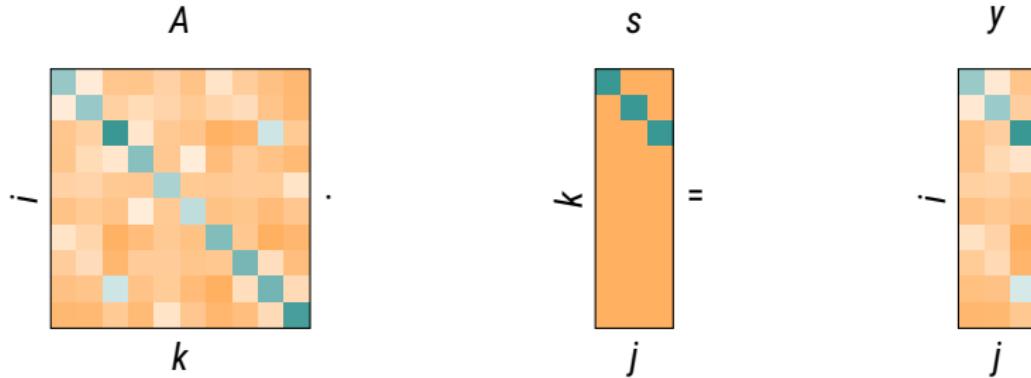
$$p_A(A | Y, S) = \mathcal{N}(A; A_M, V_0 \otimes W_M) \quad \text{with}$$

$$\vec{A}_M := \vec{A}_0 + (V_0 \otimes W_0)(I \otimes S)((I \otimes S^\top)(V_0 \otimes W_0)(I \otimes S))^{-1}(\vec{Y} - (I \otimes S^\top)\vec{A}_0)$$

$$\Sigma_M := (V_0 \otimes W_0) - (V_0 \otimes W_0)(I \otimes S)\underbrace{((I \otimes S^\top)(V_0 \otimes W_0)(I \otimes S))^{-1}}_{:= G_M \in \mathbb{R}^{NM \times NM}}(I \otimes S^\top)(V_0 \otimes W_0)(I \otimes S)$$

# Kronecker Products to the Rescue!

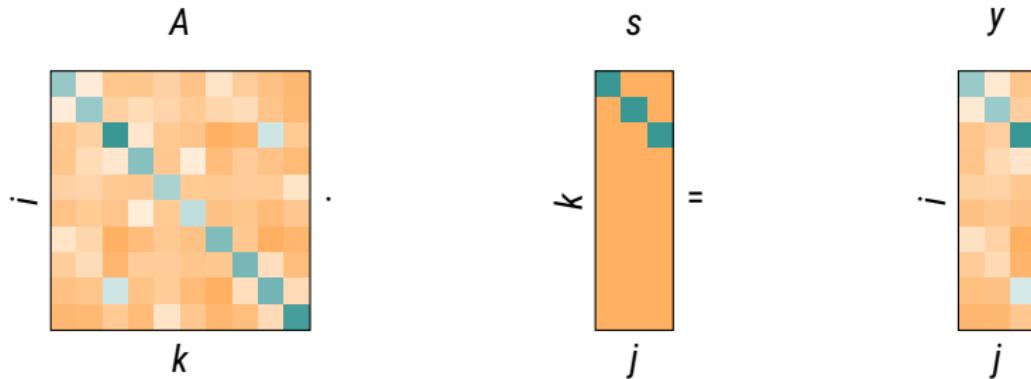
low-rank posteriors



$$\begin{aligned} p(A) &= \mathcal{N}(A; A_0, V_0 \otimes W_0) & p(S, Y | A) &= \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) \\ p_A(A | Y, S) &= \mathcal{N}(A; A_M, V_0 \otimes W_M) \quad \text{with} \\ \vec{A}_M &:= \vec{A}_0 + (V_0 \otimes W_0 S)(V_0 \otimes S^\top W_0 S)^{-1}(\vec{Y} - \vec{A}_0 \vec{S}), \text{ and} \\ \Sigma_M &:= (V_0 \otimes W_0) - (V_0 \otimes W_0 S)(V_0 \otimes S^\top W_0 S)^{-1}(V_0 \otimes S^\top W_0) \end{aligned}$$

# Kronecker Products to the Rescue!

low-rank posteriors



$$\begin{aligned} p(A) &= \mathcal{N}(A; A_0, V_0 \otimes W_0) & p(S, Y | A) &= \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) \\ p_A(A | Y, S) &= \mathcal{N}(A; A_M, V_M \otimes W_M) \quad \text{with} \\ \vec{A}_M &:= \vec{A}_0 + (I \otimes W_0 S (S^\top W_0 S)^{-1}) (\vec{Y} - \vec{A}_0 \vec{S}), \text{ and} \\ \Sigma_M &:= (V_0 \otimes W_0) - (V_0 \otimes W_0 S (S^\top W_0 S)^{-1} S^\top W_0) \end{aligned}$$

# Kronecker Products to the Rescue!

low-rank posteriors

$$\begin{aligned} p(A) &= \mathcal{N}(A; A_0, V_0 \otimes W_0) & p(S, Y | A) &= \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) \\ p_A(A | Y, S) &= \mathcal{N}(A; A_M, V_0 \otimes W_M) \quad \text{with} \\ A_M &:= A_0 + (Y - A_0 S)(S^\top W_0 S)^{-1} S^\top W_0, \text{ and} \\ \Sigma_M &:= V_0 \otimes \underbrace{(W_0 - W_0 S(S^\top W_0 S)^{-1} S^\top W_0)}_{=: W_M} \end{aligned}$$

# Kronecker Products to the Rescue!

low-rank posteriors

$$\begin{aligned} p(A) &= \mathcal{N}(A; A_0, V_0 \otimes W_0) & p(S, Y | A) &= \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) \\ p_A(A | Y, S) &= \mathcal{N}(A; A_M, V_0 \otimes W_M) \quad \text{with} \\ A_M &:= A_0 + (Y - A_0 S)(S^\top W_0 S)^{-1} S^\top W_0, \text{ and} \\ \Sigma_M &:= V_0 \otimes \underbrace{(W_0 - W_0 S(S^\top W_0 S)^{-1} S^\top W_0)}_{=: W_M} \end{aligned}$$

# Kronecker Products to the Rescue!

low-rank posteriors

$$\begin{aligned} p(A) &= \mathcal{N}(A; A_0, V_0 \otimes W_0) & p(S, Y | A) &= \delta(\vec{Y} - (I \otimes S^\top) \vec{A}) \\ p_A(A | Y, S) &= \mathcal{N}(A; A_M, V_0 \otimes W_M) \quad \text{with} \\ A_M &:= A_0 + (Y - A_0 S)(S^\top W_0 S)^{-1} S^\top W_0, \text{ and} \\ \Sigma_M &:= V_0 \otimes \underbrace{(W_0 - W_0 S (S^\top W_0 S)^{-1} S^\top W_0)}_{=: W_M} \end{aligned}$$

# What does this prior mean?

some intuition

- full support: if  $V, W$  are spd, then this is a **weighted Frobenius norm**

$$\begin{aligned}\vec{A}^\top (V \otimes W)^{-1} \vec{A} &= \vec{A}^\top \overline{(V^{-1} A W^{-1})^\top} \\ &= \text{tr}(A^\top V^{-1} A W^{-1}) \\ &= \|V^{-1/2} A W^{-1/2}\|_F^2 > 0 \quad \forall A \neq 0\end{aligned}$$

- this is **not** the same as drawing columns and rows independently. Instead:

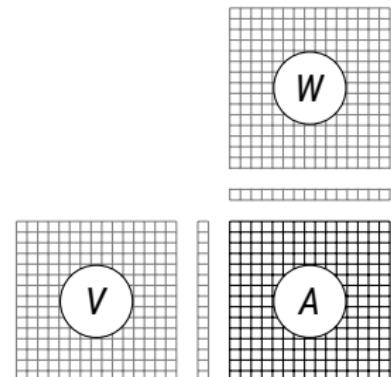
$$b_i \in \mathbb{R}^{N \times 1}, \text{ and } c_j \in \mathbb{R}^{1 \times N}, i = 1, \dots, N$$

$$b_i \sim \mathcal{N}(0, V_0) \quad c_j \sim \mathcal{N}(0, W_0)$$

$$A = B \odot C \quad \Rightarrow \quad \text{cov}(A_{ij}, A_{k\ell}) = [V_0 \otimes W_0]_{(ij),(k\ell)}$$

- $\text{var}(A_{ij}) = V_{A0,ii} \cdot W_{A0,jj}$ .

$$\begin{aligned}V_0, W_0 &\in \mathbb{R}^{N \times N} \\ V_0 \otimes W_0 &\in \mathbb{R}^{N^2 \times N^2}\end{aligned}$$

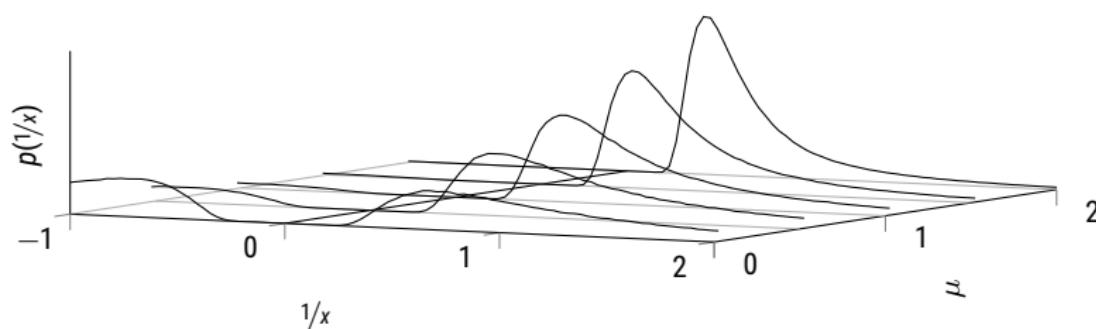


# A Prior Over what?

direct and inverse models

Hennig, SIOPT 2015

If  $A$  is spd, then  $Ax = b \Leftrightarrow x = Hb$ . Are we inferring  $A$  or  $H = A^{-1}$ ?



$$p(A) = \mathcal{N}(A; A_0, V_0 \otimes W_0)$$

$$p(Y, S | A) = \delta(Y - AS)$$

$$p(A | Y, S) = \mathcal{N}\left(A, \underbrace{A_0 + (Y - A_0 S)(S^\top W_0 S)^{-1} S^\top W_0}_{=: A_M}, V_0 \otimes (W_0 - W_0 S(S^\top W_0 S)^{-1} S^\top W_0)\right)$$

$$A_M^{-1} = A_0^{-1} + (S - A_0^{-1} Y)(S^\top W_0 A_0^{-1} Y)^{-1} S^\top W_0 A_0^{-1}$$

---

$$p(H) = \mathcal{N}(H; H_0, V_0 \otimes W_0)$$

$$p(Y, S | H) = \delta(S - HY)$$

$$p(H | Y, S) = \mathcal{N}\left(H, H_0 + (S - H_0 Y)(Y^\top W_0 Y)^{-1} Y^\top W_0, V_0 \otimes (W_0 - W_0 Y(Y^\top W_0 Y)^{-1} Y^\top W_0)\right)$$

# It only matters a little

direct and inverse models are not equivalent, but closely related

but wait for later!

$$A_M^{-1} = A_0^{-1} + (S - A_0^{-1}Y) \cdot (S^T W_0^A A_0^{-1} Y)^{-1} \cdot S^T W_0^A A_0^{-1}$$
$$H_M = H_0 + (S - H_0 Y) \cdot (Y^T W_0^H Y)^{-1} \cdot Y^T W_0^H$$

## Theorem:

Assume  $A_0$  and  $W_0$  are spd, and the search directions  $S$  are chosen linearly independent. Then, for spd  $A$ , the inverse  $A_M^{-1}$  exists.

Further, if  $H_0 = A_0^{-1}$  we have  $A_M^{-1} = H_M$  whenever  $S^T W_0^A A_0^{-1} = Y^T W_0^H$  (e.g. if  $W_0^A = A$ ,  $W_0^H = A_0^{-1}$ ).

**Proof:** If  $A_0$  is spd, its inverse exists.  $Y = AS$ , and products of spd matrices are spd. Thus,  $W_0 A_0^{-1} A$  is spd, hence  $S^T W_0 A_0^{-1} A S$  invertible. The rest of the theorem follows by inspection of the above equation.  $\square$

## Corollary:

If  $S$  are spd, for every **direct prior**  $p(A) = \mathcal{N}(A; A_0, V_0 \otimes W_0^A)$ , there exists an **equivalent inverse prior**  $p(H) = \mathcal{N}(H; H_0, \tilde{V}_0 \otimes W_0^H)$ , so that the probabilistic linear solver (using  $A_M^{-1}$  for line 12 in the first case) behaves identically.

# Projection Methods

every Gaussian iterative solver is a projection method

```
1 procedure LINSOLVE_PROP( $A(\cdot)$ ,  $b$ ,  $p(A)$ )
2    $x_0 = H_0 b$                                      // initial guess
3    $r_0 = Ax_0 - b$                                // initial gradient
4   for  $i = 1, \dots, N$  do
5      $d_i = -H_{i-1}r_{i-1}$                          // compute optimization direction
6      $z_i = Ad_i$                                  // observe
7      $\alpha_i = -d_i^T r_{i-1} / d_i^T z_i$           // optimal step-size
8      $s_i = \alpha_i d_i$                            // re-scale step
9      $y_i = \alpha_i z_i$                            // re-scale observation
10     $x_i = x_{i-1} + s_i$                          // update estimate for  $x$ 
11     $r_i = r_{i-1} + y_i$                          // new gradient at  $x_i$ 
12     $H_i = \text{INFER}(H \mid Y_i, S_i, p(A))$       // estimate  $H$ 
13  end for
14 end procedure
```

Definition:

By its general form, (regardless of line 12), LINSOLVE\_PROP selects  $x_i$  within a sequence  $\mathbb{K}_i = \text{span}\{d_1, \dots, d_i\}$ . If  $r_i \neq 0$ , then for all  $i > 0$ , there is a non-empty space  $\mathbb{L}_i$  such that  $r_i \perp \mathbb{L}_i$  (by construction,  $d_i \in \mathbb{L}_i$ ). Such an algorithm is called a **projection method**.

# Encoding Symmetry – The Symmetric Kronecker Product

We **know**  $A = A^T$ . But the prior does not reflect it!

$A - A^T = 0$  is an “observation” of a linear projection of  $A$ !

# Symmetric Kronecker Products

□ van Loan 2000. *The ubiquitous Kronecker product.* J of Computational and Applied Mathematics **123**, 85–100

- define  $\Pi_{\ominus}, \Pi_{\oplus} : \mathbb{R}^{N^2} \rightarrow \mathbb{R}^{N^2}$  by their elements

$$\Pi_{\ominus}(ij),(k\ell) := 1/2(\delta_{ik}\delta_{j\ell} + \delta_{i\ell}\delta_{jk}) \quad \Pi_{\oplus}(ij),(k\ell) := 1/2(\delta_{ik}\delta_{j\ell} - \delta_{i\ell}\delta_{jk})$$

- thus

$$\Pi_{\ominus} \vec{X} = \overrightarrow{1/2(X + X^T)} \quad \Pi_{\oplus} \vec{X} = \overrightarrow{1/2(X - X^T)}$$

- $\Pi_{\oplus}, \Pi_{\ominus}$  are **orthogonal projections**:  $\Pi_{\oplus}^T = \Pi_{\oplus}, \quad \Pi_{\ominus}^T = \Pi_{\ominus}$

$$\Pi_{\ominus}\Pi_{\ominus} = \Pi_{\ominus}, \quad \Pi_{\oplus}\Pi_{\oplus} = \Pi_{\oplus}, \quad \Pi_{\ominus}\Pi_{\oplus} = \Pi_{\oplus}\Pi_{\ominus} = \mathbf{0}, \quad \Pi_{\ominus} + \Pi_{\oplus} = I$$

- $\forall W : [\Pi_{\ominus}(W \otimes W)\Pi_{\oplus}]_{(ij),(k\ell)} = \frac{1}{4}(W_{ik}W_{j\ell} - W_{i\ell}W_{jk} - W_{ik}W_{j\ell} + W_{i\ell}W_{jk}) = 0$ .

- Thus define the **symmetric (antisymmetric) Kronecker product**  $\otimes$  ( $\otimes$ ):

$$W \otimes W = \underbrace{\Pi_{\ominus}(W \otimes W)\Pi_{\ominus}}_{=: W \otimes W} + \underbrace{\Pi_{\oplus}(W \otimes W)\Pi_{\oplus}}_{=: W \otimes W}$$

$$[W \otimes W]_{(ij),(k\ell)} = 1/2(W_{ik}W_{j\ell} + W_{i\ell}W_{jk})$$

$$(W \otimes W) \vec{X} = 1/2(WXW^T + WX^TW^T)$$

# Symmetric Gaussian Priors

encoding symmetry in the prior

□ Hennig, SIOPT 2015

- condition on symmetry

$$p(\Theta | A) = \delta(\Pi_{\Theta} A - 0) = \lim_{\beta \rightarrow 0} \mathcal{N}(0; \Pi_{\Theta} A, \beta I)$$

- generic Gaussian posterior from above, choosing  $\Sigma = W \otimes W$ ,

$$p(A | \Theta) = \mathcal{N}(A; \vec{A}_0 - \Sigma_0 \Pi_{\Theta}^T (\Pi_{\Theta} \Sigma_0 \Pi_{\Theta}^T)^{-1} (-\Pi_{\Theta} \vec{A}_0),$$

$$\Sigma_0 - \Sigma_0 \Pi_{\Theta}^T (\Pi_{\Theta} \Sigma_0 \Pi_{\Theta}^T)^{-1} \Pi_{\Theta} \Sigma_0)$$

$$\mathbb{E}(A) = \vec{A}_0 - (\Pi_{\Theta} + \Pi_{\Theta}) \Sigma_0 \Pi_{\Theta} (\Pi_{\Theta} \Sigma_0 \Pi_{\Theta})^{-1} \Pi_{\Theta} \vec{A}_0$$

$$= (\Pi_{\Theta} + \Pi_{\Theta}) \vec{A}_0 - \Pi_{\Theta} \vec{A}_0 = \Pi_{\Theta} \vec{A}_0$$

$$\text{cov}(A) = \Sigma_0 - \Sigma_0 \Pi_{\Theta}^T (\Pi_{\Theta} \Sigma_0 \Pi_{\Theta}^T)^{-1} \Pi_{\Theta} \Sigma_0 = \Pi_{\Theta} \Sigma_0 \Pi_{\Theta} = W \otimes W$$

$$\Rightarrow p(A | \Theta) = \mathcal{N}(A; \Pi_{\Theta} \vec{A}_0, W \otimes W)$$

- with likelihood  $p(Y, S | A) = \delta(Y - AS)$  get **rank-2M** posterior update

$$A_M = A_0 + (Y - A_0 S) (S^T WS)^{-1} S^T W + WS(S^T WS)^{-1} (Y - A_0 S)^T$$
$$- WS(S^T WS)^{-1} S^T (Y - A_0 S) (S^T WS)^{-1} S^T W$$

$$\Sigma_M = W_M \otimes W_M \quad \text{with (again)} \quad W_M := W_0 - W_0 S (S^T W_0 S) S^T W.$$

# Conjugate Direction Methods

every symmetric Gaussian iterative solver is a CD method

```
1 procedure LINSOLVE_PROB( $A(\cdot)$ ,  $b$ ,  $p(A)$ )
2    $x_0 = H_0 b$                                      // initial guess
3    $r_0 = Ax_0 - b$                                // initial gradient
4   for  $i = 1, \dots, N$  do
5      $d_i = -H_{i-1}r_{i-1}$                          // compute optimization direction
6      $z_i = Ad_i$                                  // observe
7      $\alpha_i = -d_i^T r_{i-1} / d_i^T z_i$           // optimal step-size
8      $s_i = \alpha_i d_i$                            // re-scale step
9      $y_i = \alpha_i z_i$                            // re-scale observation
10     $x_i = x_{i-1} + s_i$                          // update estimate for x
11     $r_i = r_{i-1} + y_i$                          // new gradient at  $x_i$ 
12     $H_i = \text{INFER}(H \mid Y_i, S_i, p(A))$       // estimate  $H$ 
13  end for
14 end procedure
```

Theorem: (proof on next slide)

If  $A$  is symmetric, and line 12 produces a **symmetric** estimator  $H_i$ , then  
LINSOLVE\_PROB is a **conjugate directions method**. That is,  $s_i^T A s_j = 0$  if  $i \neq j$ .

Conjugate directions methods converge to the correct  $x$  in at most  $N$  steps [Nocedal & Wright, Thm. 5.1] (they are *linearly consistent*). They also ensure  $r_k \perp s_{j < k}$  [op.cit., Thm. 5.2]. (they are *orthogonal projection methods*).

# Proof of Theorem

by induction

Hennig, Osborne, Girolami, upcoming

By induction: For the base case<sup>1</sup>  $i = 2$ , i.e. after the first iteration of the loop, we have (recall that  $\alpha_1 = -d_1^T r_0 / d_1^T A d_1$ ). The symmetry of the estimator  $H_i$  is used at \*

$$\begin{aligned} d_1^T A d_2 &= -d_1^T A(H_1 r_1) &= -d_1^T A(H_1(y_1 + r_0)) &= -d_1^T A(s_1 + H_1 r_0) \\ &= -\alpha_1 d_1^T A d_1 - d_1^T A H_1 r_0 &= d_1^T r_0 - \alpha_1^{-1} s_1^T A H_1 r_0 &= d_1^T r_0 - \alpha_1^{-1} y_1^T H_1 r_0 \\ &\stackrel{*}{=} d_1^T r_0 - \alpha_1^{-1} s_1^T r_0 &= d_1^T r_0 - d_1^T r_0 = 0. \end{aligned}$$

For the inductive step, assume  $\{d_0, \dots, d_{i-1}\}$  are pairwise  $A$ -conjugate. Consider any  $k < i$ , and use this assumption twice to find

$$\begin{aligned} d_k^T A d_i &= -d_k^T A(H_i r_i) &&= -d_k^T A H_i \left( \sum_{j \leq i} y_j + r_0 \right) \\ &= -d_k^T A \left( \sum_{j \leq i} s_j + H_i r_0 \right) &&= -d_k^T A \left( \sum_{j \leq i} \alpha_j d_j + H_i r_0 \right) \\ &= -\alpha_k d_k^T A d_k - d_k^T A(H_i r_0) &&= d_k^T r_{k-1} - d_k^T r_0 \\ &= d_k^T \left( \sum_{j < k} y_j + r_0 \right) - d_k^T r_0 = 0 &&= \sum_{j < k} \alpha_j d_k^T A d_j = 0. \end{aligned}$$

---

<sup>1</sup>For this proof, it does not actually matter how the first direction  $d_1$  is chosen.

# Connection to Conjugate Gradients

main result of this lecture

Hennig, Osborne, Girolami, upcoming

Theorem: (proof only on request)

Assume  $A$  is spd,  $H_i$  is symmetric for all  $i \geq 0$ , and LINSOLVE\_PROB does not terminate before step  $k < N$ . Then, if

$$d_i = -H_{i-1}r_{i-1} \in \text{span}\{s_1, \dots, s_{i-1}, y_1, \dots, y_{i-1}, r_{i-1}\}.$$

it holds that

$$r_i \perp r_j \quad \forall 0 \leq i \neq j \leq k,$$

and there exist  $\gamma_i \in \mathbb{R}_{\setminus 0}$  for all  $i < k$  so that line 12 of LINSOLVE\_PROB is

$$d_i = -H_{i-1}r_{i-1} = \gamma_i \left( -r_{i-1} + \frac{\beta_i}{\gamma_{i-1}} d_{i-1} \right), \quad \text{with } \beta_i := \frac{r_{i-1}^T r_{i-1}}{r_{i-2}^T r_{i-2}}.$$

**Then, LINSOLVE\_PROB is equivalent to the method of conjugate gradients** (with  $x_0 = H_0 b$ ).

# Probabilistic Conjugate Gradients

main result of this lecture

■ Hennig, Osborne, Girolami, upcoming

## Theorem:

Consider a prior  $p(H) = \mathcal{N}(H; H_0, W_H \otimes W_H)$ . For all choices of  $(H_0, W_H)$  with  $H_0 = \alpha I$  and  $W_H = \beta I + \gamma H$ , for  $\alpha \in \mathbb{R}$ ,  $\beta, \gamma \in \mathbb{R}_+$ , LINSOLVE\_PROB is equivalent to the method of conjugate gradients, in the sense that it produces the exact same sequence  $\{x_i\}$ . The same is true for MAP inference using  $p(A) = \mathcal{N}(A; A_0, W_A \otimes W_A)$  with scalar  $A_0 = \alpha I$  and  $W_A = \beta I + \gamma A$ .

## Proof:

Use theorems on previous slides. First: conjugate directions by symmetry. Then to show:  $H_i$  resulting from the above models satisfies assumption of theorem on previous slide., i.e.  $r_{i-1}$  is mapped to  $s_i = -H_{i-1}r_{i-1}$  in the span of  $\{S, Y, r_{i-1}\}$ . Consider the image of  $H_i$  as defined above. For the model on  $H$ , if  $S_{:i-1}$  denotes  $(s_1, \dots, s_{i-1})$  and analogously for  $Y_{:i-1}$ , then  $H_{i-1}$  maps  $v \in \mathbb{R}^N$  to  $\text{span}\{H_0 v, S_{:i-1}, H_0 Y_{:i-1}, W_H Y_{:i-1}\}$ . Hence, if  $H_0$  is scalar and  $W_H = \beta I + \gamma H$ , then (since  $HY = S$ ),  $r_{i-1}$  is mapped to  $\text{span}\{r_{i-1}, S_{:i-1}, Y_{:i-1}\}$ . For models on  $A$ ,  $H_{i-1}$  maps  $r_{i-1}$  to  $\text{span}\{A_0^{-1}r_{i-1}, A_0^{-1}W_A S_{:i-1}, A_0^{-1}Y_{:i-1}, S_{:i-1}\}$ . For a scalar  $A_0$  and  $W_A = \beta I + \gamma A$ , using  $AS = Y$ , this is the span of  $\{r_{i-1}, Y_{:i-1}, S_{:i-1}\}$ . □

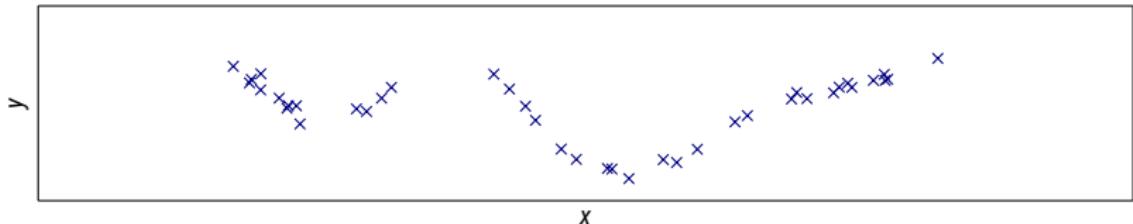
The big picture:

- Gaussian measures on matrices can be used to construct **projection methods** by MAP estimation
- **Symmetry** can be encoded analytically, and gives **conjugate direction methods**
- There is a nontrivial subset of symmetric priors that give **conjugate gradients**

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

□ Bartels & Hennig, AISTATS 2016



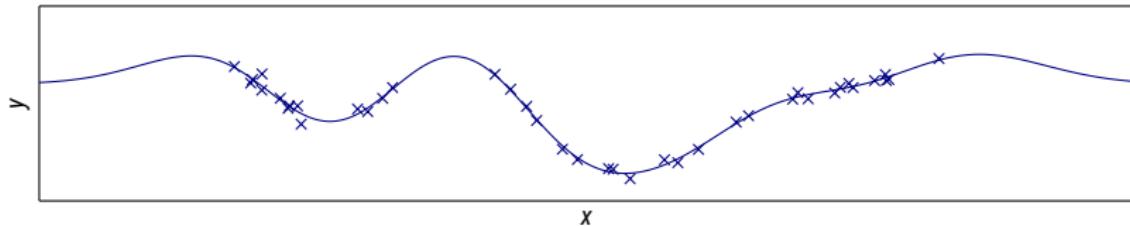
- $\Rightarrow \mu(x) = k_{xx}k_{XX}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_x\|^2 + \|f\|_k^2]$  where

$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

□ Bartels & Hennig, AISTATS 2016



- $\Rightarrow \mu(x) = k_{xx}k_{xx}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_x\|^2 + \|f\|_k^2]$  where

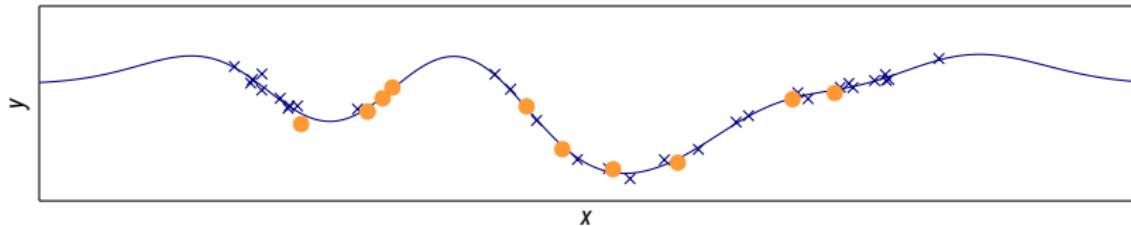
$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

- $X \sim p(X) = \prod_i p(x_i)$ .  $f(x) \sim \mathcal{GP}(0, k)$

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

□ Bartels & Hennig, AISTATS 2016



- $\Rightarrow \mu(x) = k_{xX}k_{XX}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_X\|^2 + \|f\|_k^2]$  where

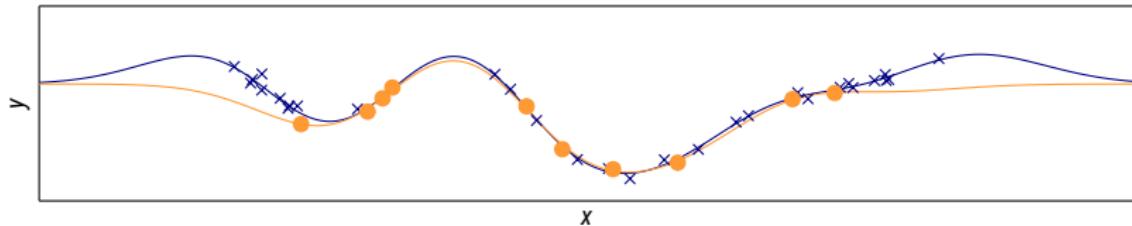
$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

- $X \sim p(X) = \prod_i p(x_i)$ .  $f(x) \sim \mathcal{GP}(0, k)$
- subset of size  $M$ . Equivalent to Gauss-Jordan for  $M$  steps, get approximation to full LSq-estimate  $\tilde{\mu}$

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

□ Bartels & Hennig, AISTATS 2016



- $\Rightarrow \mu(x) = k_{xX}k_{XX}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_X\|^2 + \|f\|_k^2]$  where

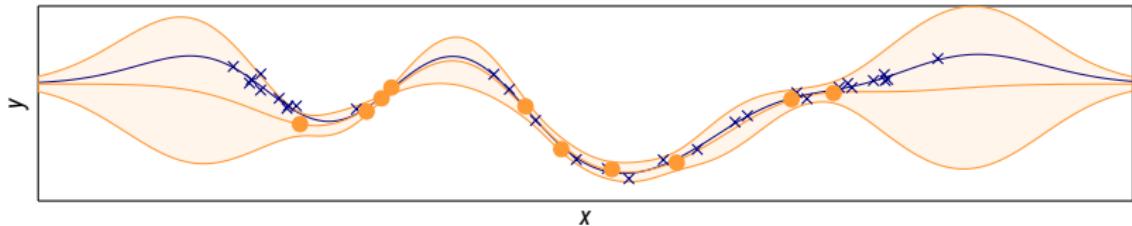
$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

- $X \sim p(X) = \prod_i p(x_i)$ .  $f(x) \sim \mathcal{GP}(0, k)$
- subset of size  $M$ . Equivalent to Gauss-Jordan for  $M$  steps, get approximation to full LSq-estimate  $\tilde{\mu}$

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

□ Bartels & Hennig, AISTATS 2016



- $\Rightarrow \mu(x) = k_{xX}k_{XX}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_X\|^2 + \|f\|_k^2]$  where

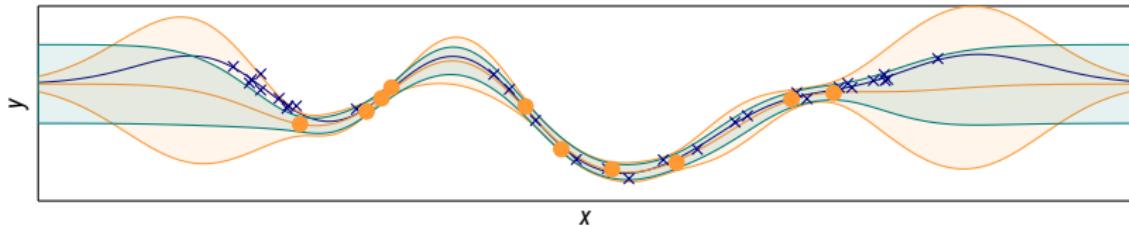
$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

- $X \sim p(X) = \prod_i p(x_i)$ .  $f(x) \sim \mathcal{GP}(0, k)$
- subset of size  $M$ . Equivalent to Gauss-Jordan for  $M$  steps, get approximation to full LSq-estimate  $\tilde{\mu}$
- choose  $p(k_{XX}^{-1}) = \mathcal{N}(\alpha_0 I, \beta^2 I \otimes I)$ . Use statistical regularity assumptions to get calibrate posterior to get upper bound on approximation error  $|\mu - \tilde{\mu}|$

# Uncertainty Calibration for Iterative Solvers

an example of statistically calibrated error estimates

Bartels & Hennig, AISTATS 2016



- $\Rightarrow \mu(x) = k_{xX}k_{XX}^{-1}Y = \arg \min_{f \in \mathcal{H}_k} [\|Y - f_X\|^2 + \|f\|_k^2]$  where

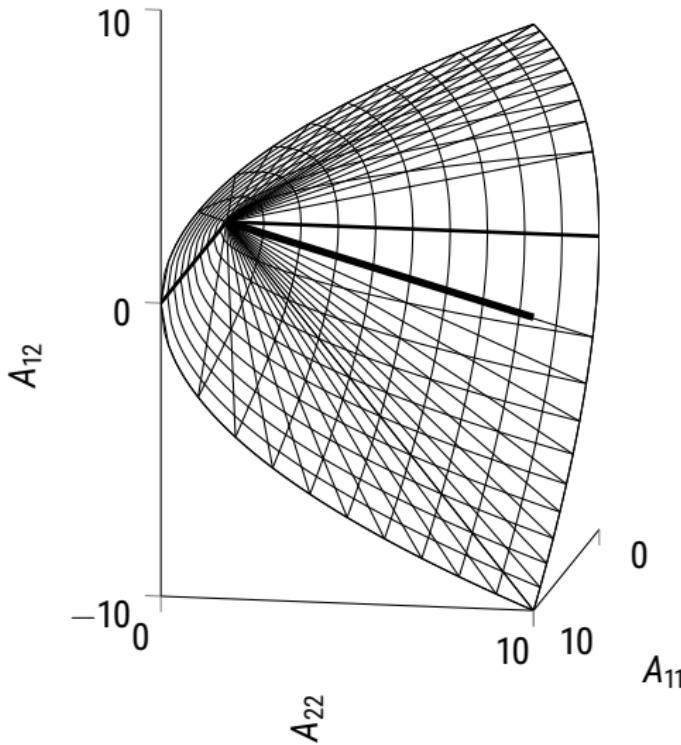
$$[k_{XX}]_{ij} = k(x_i, x_j) \quad \text{e.g.} \quad k(a, b) = \exp\left(-\frac{(a - b)^2}{2\lambda^2}\right)$$

- $X \sim p(X) = \prod_i p(x_i)$ .  $f(x) \sim \mathcal{GP}(0, k)$
- subset of size  $M$ . Equivalent to Gauss-Jordan for  $M$  steps, get approximation to full LSq-estimate  $\tilde{\mu}$
- choose  $p(k_{XX}^{-1}) = \mathcal{N}(\alpha_0 I, \beta^2 I \otimes I)$ . Use statistical regularity assumptions to get calibrate posterior to get upper bound on approximation error  $|\mu - \tilde{\mu}|$
- not the same as GP posterior variance!

# What about Positive Definiteness

positive definites can only be modelled approximately

Hennig, Osborne, Girolami, CUP, upcoming

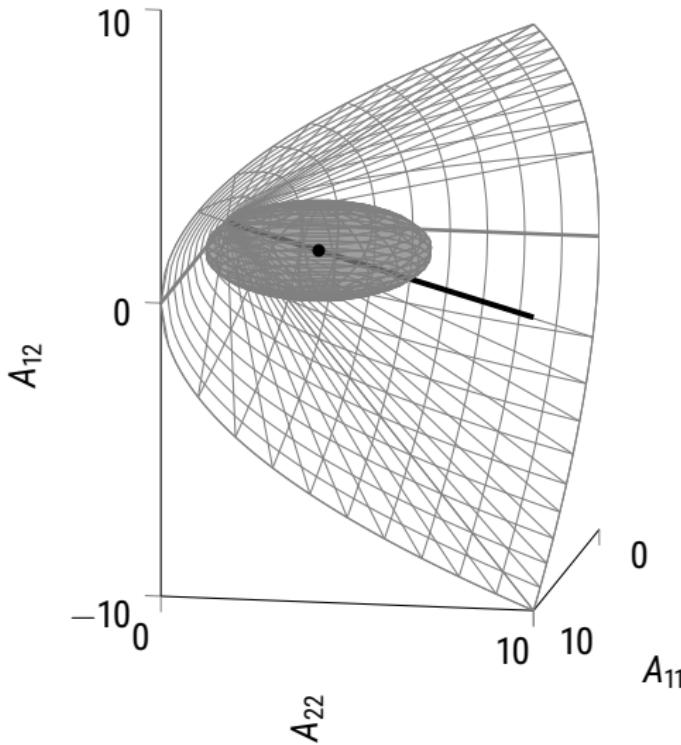


- the positive definite cone is not a linear sub-space of  $\mathbb{R}^{N^2}$

# What about Positive Definiteness

positive definites can only be modelled approximately

Hennig, Osborne, Girolami, CUP, upcoming

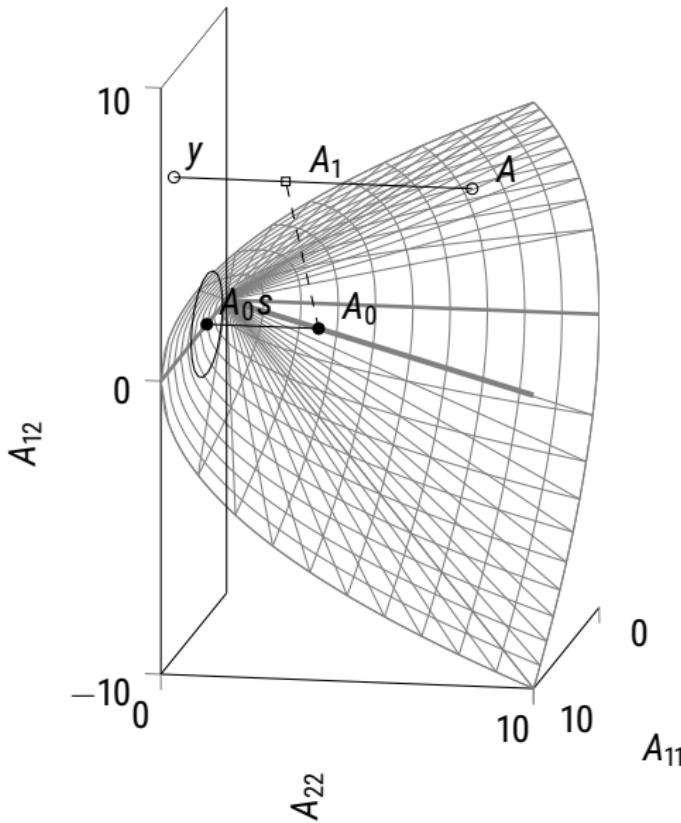


- the positive definite cone is not a linear sub-space of  $\mathbb{R}^{N^2}$
- Gaussian priors can not be exactly conditioned on pd-ness

# What about Positive Definiteness

positive definites can only be modelled approximately

Hennig, Osborne, Girolami, CUP, upcoming

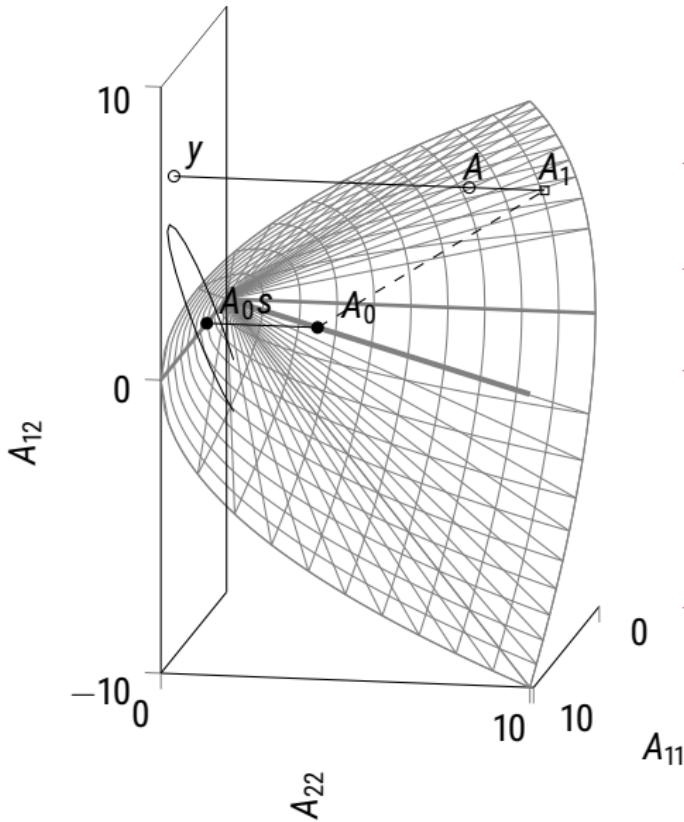


- the positive definite cone is not a linear sub-space of  $\mathbb{R}^{N^2}$
- Gaussian priors can not be exactly conditioned on pd-ness
- Theorem** (proof on request): For  $p(A) = \mathcal{N}(\beta I, \alpha^2 I \otimes I)$ , the posterior mean  $A_M$  may lie outside the cone. But for any  $y$ , there is an  $\alpha^2$  so  $A_M$  is spd.

# What about Positive Definiteness

positive definites can only be modelled approximately

Hennig, Osborne, Girolami, CUP, upcoming



- the positive definite cone is not a linear sub-space of  $\mathbb{R}^{N^2}$
- Gaussian priors can not be exactly conditioned on pd-ness
- Theorem** (proof on request): For  $p(A) = \mathcal{N}(\beta I, \alpha^2 I \otimes I)$ , the posterior mean  $A_M$  may lie outside the cone. But for any  $y$ , there is an  $\alpha^2$  so  $A_M$  is spd.
- Theorem** (proof on request): For  $p(A) = \mathcal{N}(I, \beta I, A \otimes A)$ , the posterior mean  $A_M$  always lies in the cone.

# What else?

In linear algebra, PN is tightly constrained by complexity

more in Hennig, Osborne, Girolami, upcoming

Some things one may want to try

- noise on matrix elements  $p(Y | A, S) = \mathcal{N}(Y; AS, \Lambda)$  ?
- breaks Kronecker structure of covariance, requires nontrivial extensions
- hierarchical inference on scale/structure of prior?
- possible using Gauss-Gamma prior.  
But CG does not produce iid. observations. More on request

PN affords a (one) unifying view of linear algebra methods, but significant additional functionality typically requires serious work.

There can be a trade-off between different desired probabilistic functionality: calibrated uncertainty is easier if the point-estimate is not efficient. Noisy observations require approximate inference. Etc.

The big picture:

- Gaussian measures on matrices can be used to construct **projection methods** by MAP estimation
- **Symmetry** can be encoded analytically, and gives **conjugate direction methods**
- There is a nontrivial subset of symmetric priors that give **conjugate gradients**
- There are many **cumbersome restrictions** imposed by the extremely high-dimensional space of matrix elements and the tight restrictions on computational **complexity**.

in **nonlinear Optimization**, many symmetries are broken.

# Nonlinear Optimization

## Problem Setting

$$f : \mathbb{R}^N \rightarrow \mathbb{R} \quad \nabla f(x_*) = 0$$

- local approximation to second order, assuming Hessian  $[B(x)]_{ij} := \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ .

$$f(x - x_i) = f(x_i) + (x - x_i)^\top \nabla f(x_i) + \frac{1}{2}(x - x_i)^\top B(x_i)(x - x_i) + \mathcal{O}(\|x - x_i\|^3)$$

# Nonlinear Optimization

## Problem Setting

$$f : \mathbb{R}^N \rightarrow \mathbb{R} \quad \nabla f(x_*) = 0$$

- local approximation to second order, assuming Hessian  $[B(x)]_{ij} := \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ .

$$f(x - x_i) = f(x_i) + (x - x_i)^\top \nabla f(x_i) + \frac{1}{2}(x - x_i)^\top B(x_i)(x - x_i) + \mathcal{O}(\|x - x_i\|^3)$$

- if  $B$  is spd (i.e.  $f$  is convex), efficient local optimization by **Newton iteration**

$$x_{i+1} \approx \arg \min f(x - x_i) = x_i - B^{-1}(x_i) \nabla f(x_i)$$

# Nonlinear Optimization

## Problem Setting

$$f : \mathbb{R}^N \rightarrow \mathbb{R} \quad \nabla f(x_*) = 0$$

- local approximation to second order, assuming Hessian  $[B(x)]_{ij} := \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ .

$$f(x - x_i) \approx f(x_i) + (x - x_i)^\top \nabla f(x_i) + \frac{1}{2}(x - x_i)^\top B(x_i)(x - x_i) + \mathcal{O}(\|x - x_i\|^3)$$

- if  $B$  is spd (i.e.  $f$  is convex), efficient local optimization by **Newton iteration**

$$x_{i+1} \approx \arg \min f(x - x_i) = x_i - B^{-1}(x_i) \nabla f(x_i)$$

- but computing  $B^{-1}(x_i)$  is  $\mathcal{O}(N^3)$ . So we have to **infer** it.

# Nonlinear Optimization

## Problem Setting

$$f : \mathbb{R}^N \rightarrow \mathbb{R} \quad \nabla f(x_*) = 0$$

- local approximation to second order, assuming Hessian  $[B(x)]_{ij} := \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$ .

$$f(x - x_i) \approx f(x_i) + (x - x_i)^\top \nabla f(x_i) + \frac{1}{2}(x - x_i)^\top B(x_i)(x - x_i) + \mathcal{O}(\|x - x_i\|^3)$$

- if  $B$  is spd (i.e.  $f$  is convex), efficient local optimization by **Newton iteration**

$$x_{i+1} \approx \arg \min f(x - x_i) = x_i - B^{-1}(x_i) \nabla f(x_i)$$

- but computing  $B^{-1}(x_i)$  is  $\mathcal{O}(N^3)$ . So we have to **infer** it.
- accessible:  $\{\nabla f(x_i)\}_{i=0, \dots}$

# Quasi Newton Methods

cf. □ Dennis & Moré, *Quasi-Newton Methods, motivation and theory*, SIREV 1977

- ensure **secant equation**

$$y_i := \nabla f(x_i) - \nabla f(x_{i-1}) = \hat{B}_i(x_i - x_{i-1}) =: \hat{B}_i s_i \quad \text{or} \quad \hat{H}_i y_i = s_i$$

- many options! E.g. The **Dennis family** (1971) of direct symmetric updates:

$$B_{i+1} = B_i + \frac{(y_i - B_i s_i)c_i^\top + c_i(y_i - B_i s_i)^\top}{c_i^\top s_i} - \frac{c_i s_i^\top (y_i - B_i s_i) c_i^\top}{(c_i^\top s_i)^2}$$

with members

[cf. □ Hennig, SIOPT 2015]

Symmetric Rank-1 (SR1)<sup>Davidon, 1959</sup>

$$c = y - B_{i-1}s$$

Powell<sup>1970</sup> Symmetric Broyden<sup>1965</sup>

$$c = s$$

Greenstadt<sup>1970</sup>

$$c = B_{i-1}s$$

Davidon<sup>1959</sup> Fletcher Powell<sup>1963</sup>

$$c = y$$

Broyden<sup>1969</sup> Fletcher<sup>1970</sup> Goldfarb<sup>1970</sup> Shanno<sup>1970</sup>

$$c = y + \sqrt{\frac{y^\top s}{s^\top B_{i-1}s}} B_{i-1}s$$

# A Probabilistic Interpretation of the Dennis Family

Symmetric Gaussian Priors

□ Hennig, SIOPT 2015

- direct models:  $p(B \mid B_{i-1}) = \mathcal{N}(B, B_i, W_i \otimes W_i)$  and  $p(y_i \mid s_i, B_i) = \delta(y_i - B_i s_i)$ :

Symmetric Rank-1 (SR1) <sup>Davidon, 1959</sup>

$$W_i = \alpha(\bar{B} - B_i)$$

Powell <sup>1970</sup> Symmetric Broyden <sup>1965</sup>

$$W_i = \alpha I$$

Greenstadt <sup>1970</sup>

$$c = B_{i-1}s$$

Davidon <sup>1959</sup> Fletcher Powell <sup>1963</sup>

$$W_i = \alpha\bar{B}$$

where  $\bar{B} = \int_0^1 B(x_{i-1} + ts_i) dt$ ,  $\alpha > 0$ .

- inverse models:  $p(H \mid H_{i-1}) = \mathcal{N}(H, H_i, V_i \otimes V_i)$  and  $p(y_i \mid s_i, H_i) = \delta(s_i - H_i y_i)$ :

Broyden <sup>1969</sup> Fletcher <sup>1970</sup> Goldfarb <sup>1970</sup> Shanno <sup>1970</sup>  $V_i = \alpha\bar{H}$

# Equivalence in the Linear Case

[cf. Nazareth, SINUM 1979]

Theorem (simplified, Nazareth, 1979)

If  $f$  is a convex (positive definite) quadratic function, **all** members of the Dennis family are equivalent to the method of conjugate gradients.

- Gaussian models can be used to infer matrix elements
- the potential model spaces are severely limited by computational considerations
- nevertheless, class of priors consistent with existing methods is surprisingly large
- algebraic constraints also make desirable functionality hard to achieve
- **nonlinear** optimization methods are closely related, but many **symmetries are broken**

These slides can be found at  
<http://tinyurl.com/Dobbiaco-Hennig-4>

Beyond classics, let's build new functionality for contemporary challenges!