

A Bayesian Perspective on EM Covariances

Ryan Giordano, March 2019

Outline

- **Prelude:** Linear response covariances.
- **Part 1:** The standard view of EM.
- **Part 2:** A Bayesian view of EM.
- **Part 3:** Covariance asymptotics.
- **Postscript:** Tools.

Prelude:

Linear response covariances.

The Bayesian Machinery

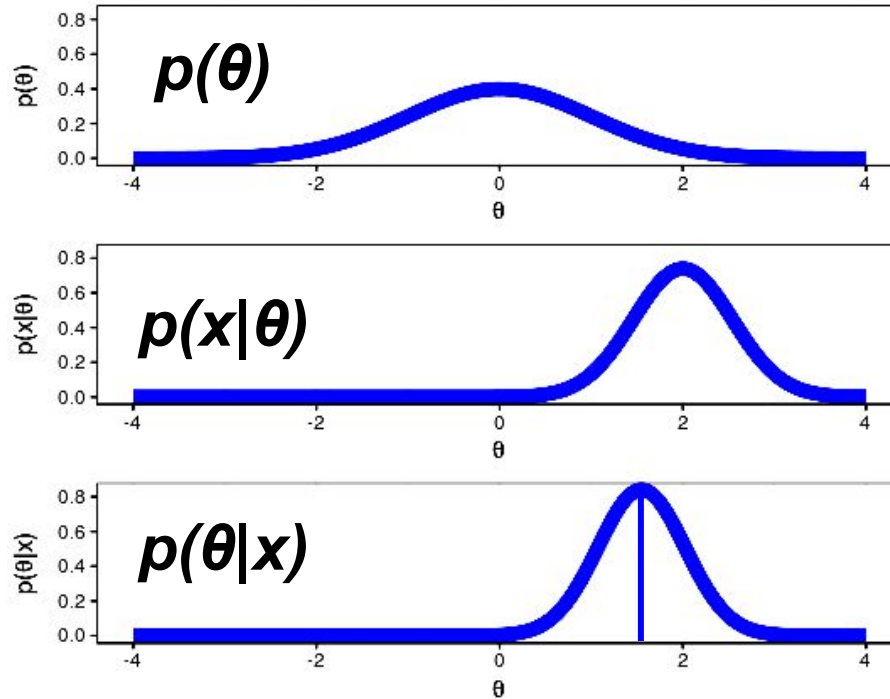
Prior

&

Likelihood

BAYES

Posterior



The Bayesian Machinery

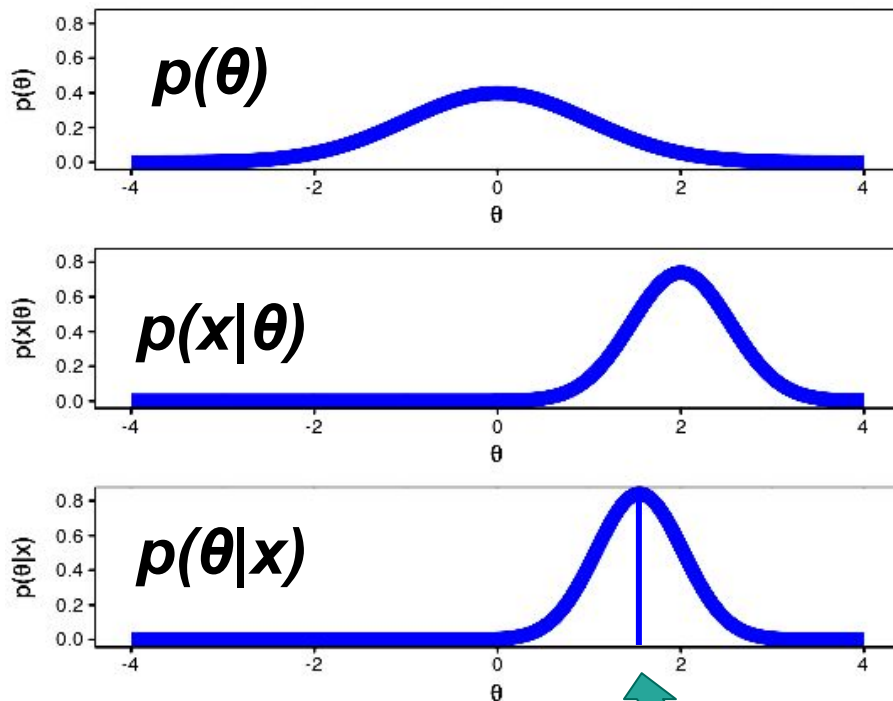
Prior

&

Likelihood

BAYES

Posterior



Summarize
posteriors
with a
“posterior
expectation”:

$$\mathbb{E}_{p(\theta|X)} [\theta]$$

Everything has “hyperparameters”

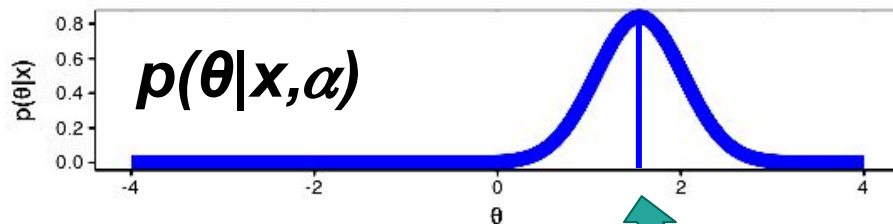
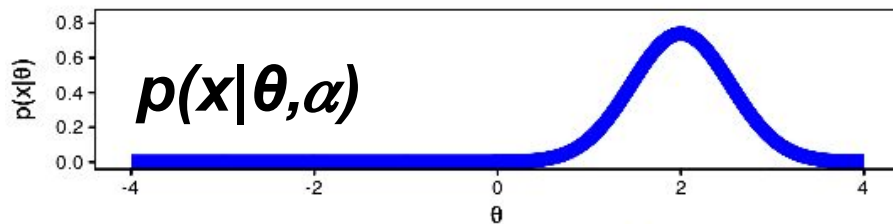
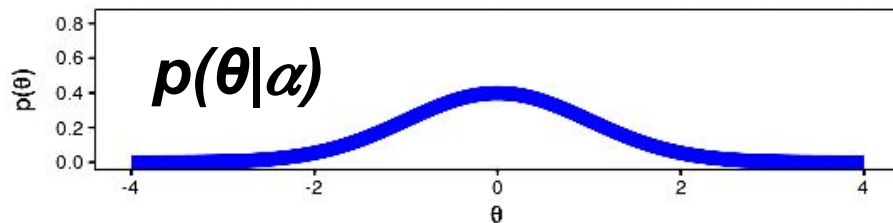
Prior

&

Likelihood

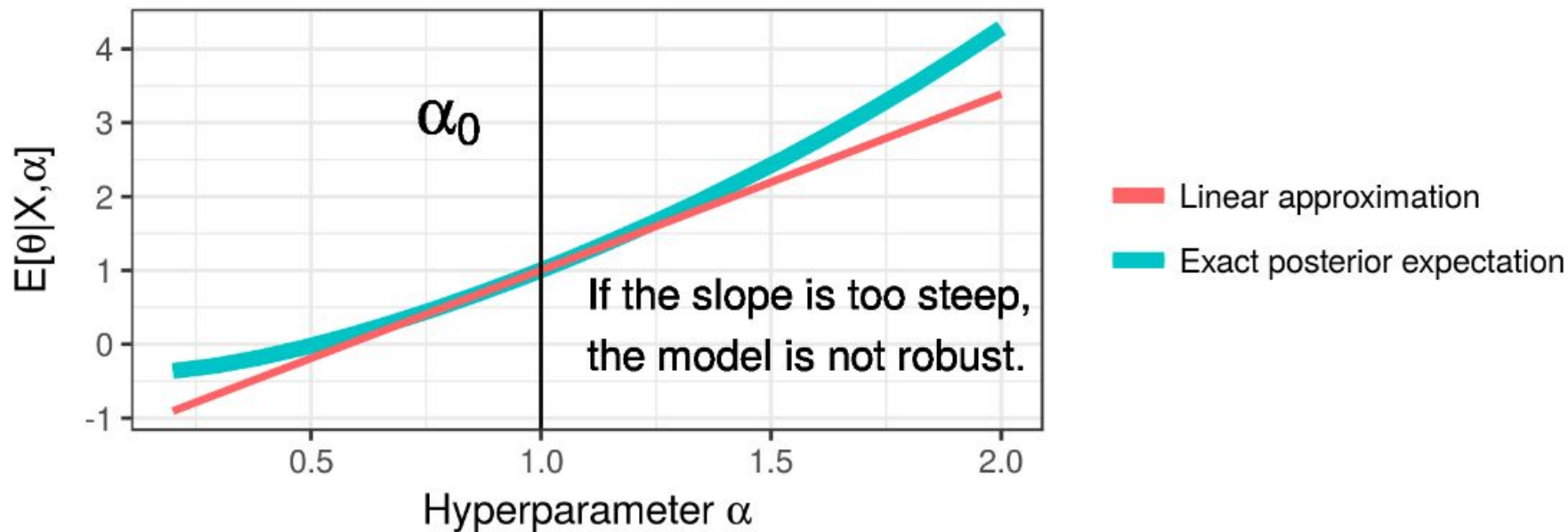
BAYES

Posterior



$$\mathbb{E}_{p(\theta|X, \alpha)} [\theta]$$

Sensitivity



Actual results calculated with
<https://github.com/rgiordan/StanSensitivity>

Sensitivity = Covariance

(exchange differentiation and integration)

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \underbrace{\frac{\partial}{\partial\alpha} \log p(\theta|\alpha, X)} \right)$$

Some nasty derivative

Classical Bayesian robustness:
Calculate the covariance to estimate the sensitivity

$$\frac{d\mathbb{E}[\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$



Linear response covariances:

Calculate the sensitivity to estimate the covariance.

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$



Standard result from sensitivity analysis:

$$\hat{\theta}(t) = \operatorname{argmin}_{\theta} (f(\theta) + t\theta)$$

Standard result from sensitivity analysis:

$$\hat{\theta}(t) = \operatorname{argmin}_{\theta} (f(\theta) + t\theta)$$
$$\left. \frac{d\hat{\theta}}{dt} \right|_{t=0} = \left(\left. \frac{\partial^2 f(\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} \right)^{-1} \left. \vphantom{\frac{d\hat{\theta}}{dt}} \right\} \text{An inverse Hessian!}$$

Part 1:

The standard view of EM.

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$

Latent
variables: $Z = (Z_1, \dots, Z_N)$

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$ Observed, grows with N

Latent variables: $Z = (Z_1, \dots, Z_N)$ Unobserved, grows with N

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$ Unobserved, fixed dimension

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$

Latent variables: $Z = (Z_1, \dots, Z_N)$

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$

Generative process:

$p(\theta) = \Lambda$ given prior.

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$

Latent variables: $Z = (Z_1, \dots, Z_N)$

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$

Generative process:

$$p(Z|\theta) = \prod_{n=1}^N p(Z_n|\theta)$$

$p(\theta) = \Lambda$ given prior.

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$

Latent variables: $Z = (Z_1, \dots, Z_N)$

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$

Generative process:

$$p(Y|\theta, Z) = \prod_{n=1}^N p(Y_n|Z_n, \theta)$$

$$p(Z|\theta) = \prod_{n=1}^N p(Z_n|\theta)$$

$$p(\theta) = \text{A given prior.}$$

Parameters and data:

Data: $Y = (Y_1, \dots, Y_N)$

Latent
variables:

$$Z = (Z_1, \dots, Z_N)$$

Parameters: $\theta \in \Omega_\theta \subseteq \mathbb{R}^D$

Equivalent generative process:

$$p(Y|\theta) = \prod_{n=1}^N \int p(Y_n|Z_n, \theta) p(Z_n|\theta) dZ_n$$

$$p(Z|\theta) = \prod_{n=1}^N p(Z_n|\theta)$$

$p(\theta) = \Lambda$ given prior.

Why not estimate with $\hat{\theta}, \hat{Z} = \operatorname{argmax}_{\theta, Z} p(Y, Z|\theta)$ **?**

Why not estimate with $\hat{\theta}, \hat{Z} = \operatorname{argmax}_{\theta, Z} p(Y, Z|\theta)$ **?**

$$p(Y_n|\theta) = \int p(Y_n, Z_n|\theta) dZ_n \neq p(Y_n, \hat{Z}_n|\theta)$$

...unless $p(Z_n|Y_n, \theta)$ is concentrated.
(Roughly speaking.)

We do EM when:

Hard:

$$p(Y|\theta) = \prod_{n=1}^N \int p(Y_n|Z_n, \theta) p(Z_n|\theta) dZ_n$$

Easy:

$$p(Y|\theta, Z) = \prod_{n=1}^N p(Y_n|Z_n, \theta)$$

Dispersed:

$$p(Z_n|Y_n, \theta)$$

We do EM when:

Hard:

$$p(Y|\theta) = \prod_{n=1}^N \int p(Y_n|Z_n, \theta) p(Z_n|\theta) dZ_n$$

Easy:

$$p(Y|\theta, Z) = \prod_{n=1}^N p(Y_n|Z_n, \theta)$$

**Dispersed:
(and easy)**

$$p(Z_n|Y_n, \theta)$$

Notation for log probabilities.

$$\ell(Y|\theta) = \log p(Y|\theta)$$

(and in general)

Assume the MLE is nice.

Consistent:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \ell(Y|\theta)$$

$$\hat{\theta} \xrightarrow[N \rightarrow \infty]{} \theta_0$$

Assume the MLE is nice.

Consistent:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \ell(Y|\theta)$$

$$\hat{\theta} \xrightarrow[N \rightarrow \infty]{} \theta_0$$

Asymptotically normal:

$$\hat{\mathcal{I}}_{\theta\theta} := -\frac{1}{N} \left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}}$$

$$\hat{\Sigma} := \hat{\mathcal{I}}_{\theta\theta}^{-1}$$

$$\sqrt{N} \hat{\Sigma}^{-1/2} \left(\hat{\theta} - \theta_0 \right) \xrightarrow[N \rightarrow \infty]{} \mathcal{N}(0, I_D)$$

Hard to calculate
(by assumption)

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \overbrace{\ell(Y|\theta)}$$

$$\hat{\theta} \xrightarrow[N \rightarrow \infty]{} \theta_0$$

$$\hat{\mathcal{I}}_{\theta\theta} := -\frac{1}{N} \overbrace{\frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta}} \bigg|_{\hat{\theta}}$$

$$\hat{\Sigma} := \hat{\mathcal{I}}_{\theta\theta}^{-1}$$

$$\sqrt{N} \hat{\Sigma}^{-1/2} \left(\hat{\theta} - \theta_0 \right) \xrightarrow[N \rightarrow \infty]{} \mathcal{N}(0, I_D)$$

The “EM Identity”

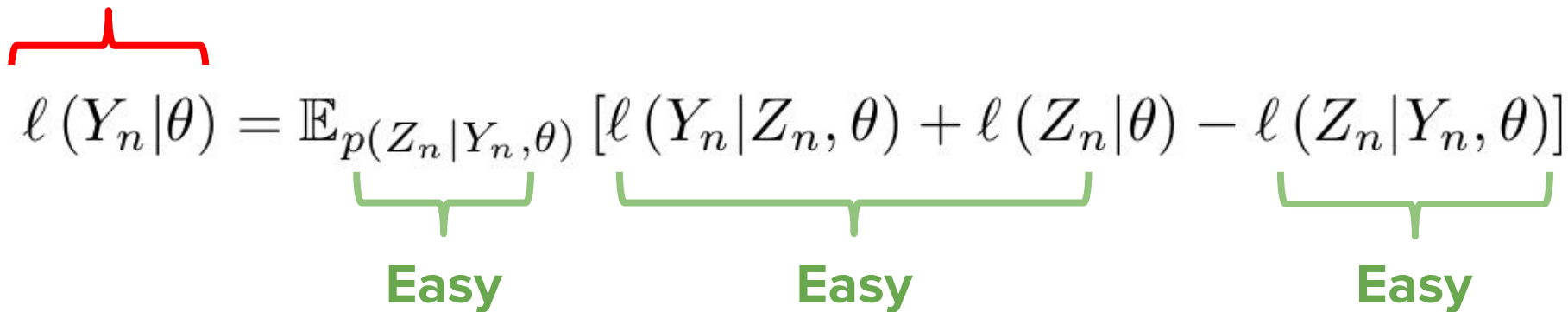
Hard



$$\ell(Y_n|\theta) = \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n,\theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)]$$

The “EM Identity”

Hard



The diagram illustrates the EM Identity equation: $\ell(Y_n|\theta) = \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n,\theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)]$. A red bracket above the left side of the equation is labeled "Hard". Three green brackets below the right side are each labeled "Easy". The first green bracket is under $\ell(Y_n|Z_n,\theta)$, the second is under $\ell(Z_n|\theta)$, and the third is under $\ell(Z_n|Y_n,\theta)$. The expectation operator $\mathbb{E}_{p(Z_n|Y_n,\theta)}$ is not bracketed.

$$\ell(Y_n|\theta) = \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n,\theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)]$$

Easy Easy Easy

Proof.

Apply Bayes' rule



$$\ell(Y_n|\theta) = \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n, \theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n, \theta)]$$

Proof.

Apply Bayes' rule



$$\begin{aligned}\ell(Y_n|\theta) &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n, \theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n, \theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [(\ell(Z_n|Y_n, \theta) + \ell(Y_n|\theta) - \ell(Z_n|\theta)) + \ell(Z_n|\theta) - \ell(Z_n|Y_n, \theta)]\end{aligned}$$

Proof.

Apply Bayes' rule



$$\begin{aligned}\ell(Y_n|\theta) &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n,\theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [(\ell(Z_n|Y_n,\theta) + \ell(Y_n|\theta) - \ell(Z_n|\theta)) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|\theta)]\end{aligned}$$

Proof.

Apply Bayes' rule



$$\begin{aligned}\ell(Y_n|\theta) &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n, \theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n, \theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [(\ell(Z_n|Y_n, \theta) + \ell(Y_n|\theta) - \ell(Z_n|\theta)) + \ell(Z_n|\theta) - \ell(Z_n|Y_n, \theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|\theta)] \\ &= \ell(Y_n|\theta)\end{aligned}$$

Proof.

$$\begin{aligned}\ell(Y_n|\theta) &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|Z_n,\theta) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [(\ell(Z_n|Y_n,\theta) + \ell(Y_n|\theta) - \ell(Z_n|\theta)) + \ell(Z_n|\theta) - \ell(Z_n|Y_n,\theta)] \\ &= \mathbb{E}_{p(Z_n|Y_n,\theta)} [\ell(Y_n|\theta)] \\ &= \ell(Y_n|\theta)\end{aligned}$$

This view of EM can simplify some EM proofs.

The “EM Identity”

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \ell(Y|\theta) \quad \left. \vphantom{\operatorname{argmax}_{\theta \in \Omega_{\theta}}} \right\} \text{ Hard}$$

$$= \operatorname{argmax}_{\theta \in \Omega_{\theta}} \mathbb{E}_{p(Z|Y, \theta)} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \theta)]$$



Easy



Easy



Easy

The “EM algorithm”

Given iteration k , $\theta^{(k)}$:

Step 1k. “E step”:

1. Calculate $Q^{(k)}(\theta) = \mathbb{E}_{p(Z|Y, \theta^{(k)})} \left[\underbrace{\ell(Y|Z, \theta)}_{\text{Fixed}} + \underbrace{\ell(Z|\theta)}_{\text{A function of } \theta} - \underbrace{\ell(Z|Y, \theta^{(k)})}_{\text{Fixed (typically omitted)}} \right]$

The “EM algorithm”

Given iteration k , $\theta^{(k)}$:

Step 1k. “E step”:

A function of θ



1. Calculate $Q^{(k)}(\theta) = \mathbb{E}_{p(Z|Y, \theta^{(k)})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \theta^{(k)}) \right]$

Step 2k. “M step”:

Calculate $\theta^{(k+1)} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} Q^{(k)}(\theta)$

...repeat.

EM algorithm \neq EM identity

Under nice conditions, the EM algorithm solves the same optimization problem as the MLE.

$$\theta^{(k)} \xrightarrow[N \rightarrow \infty]{} \hat{\theta}$$

But what about covariances?

What about covariances?

We want the Hessian of the marginal log likelihood:

$$\frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \bigg|_{\hat{\theta}}$$

What about covariances?

We want the Hessian of the marginal log likelihood:

$$\left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}}$$

All we have is the Q function:

$$\hat{Q}(\theta) = \mathbb{E}_{p(Z|Y, \hat{\theta})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta}) \right]$$

What about covariances?

...but the Hessians are not the same.

$$\left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} \neq \left. \frac{\partial^2 \hat{Q}(\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}}$$


$$\hat{Q}(\theta) = \mathbb{E}_{p(Z|Y, \hat{\theta})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta}) \right]$$

...but the Hessians are not the same.

$$\left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} = \left. \frac{\partial^2}{\partial \theta \partial \theta} \mathbb{E}_{p(Z|Y, \theta)} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \theta)] \right|_{\hat{\theta}}$$

(by the EM identity)

...but the Hessians are not the same.

$$\begin{aligned} \left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} &= \left. \frac{\partial^2}{\partial \theta \partial \theta} \mathbb{E}_{p(Z|Y, \theta)} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \theta)] \right|_{\hat{\theta}} \\ &\neq \left. \frac{\partial^2}{\partial \theta \partial \theta} \mathbb{E}_{p(Z_n|Y_n, \hat{\theta})} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta})] \right|_{\hat{\theta}} \end{aligned}$$


(fix the “E step”)

...but the Hessians are not the same.

$$\begin{aligned}\frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \Big|_{\hat{\theta}} &= \frac{\partial^2}{\partial \theta \partial \theta} \mathbb{E}_{p(Z|Y, \theta)} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \theta)] \Big|_{\hat{\theta}} \\ &\neq \frac{\partial^2}{\partial \theta \partial \theta} \mathbb{E}_{p(Z_n|Y_n, \hat{\theta})} [\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta})] \Big|_{\hat{\theta}} \\ &= \frac{\partial^2 \hat{Q}(\theta)}{\partial \theta \partial \theta} \Big|_{\hat{\theta}}\end{aligned}$$

(by definition)

...but the Hessians are not the same.

$$\left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} \neq \left. \frac{\partial^2 \hat{Q}(\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}}$$

Standard work-arounds are kinda complicated*.

And what about uncertainty in Z ?

- Using EM to obtain asymptotic variance-covariance matrices: The SEM algorithm, Meng et al., 2001
- Direct calculation of the information matrix via the EM algorithm, Oakes, 1999
- The EM algorithm and extensions, McLachlan, G. and T. Krishnan, 2007

A Bayesian view of EM.

Let's calculate the posterior.

$$p(\theta, Z|Y) = \frac{p(Y|Z, \theta) p(Z|\theta) p(\theta)}{p(Y)}$$

Let's calculate the posterior.

Hard

$$p(\theta, Z|Y) = \frac{\overbrace{p(Y|Z, \theta) p(Z|\theta) p(\theta)}}{p(Y)}$$

Let's ~~calculate~~ approximate the posterior.

Hard

$$p(\theta, Z|Y) = \frac{\overbrace{p(Y|Z, \theta) p(Z|\theta) p(\theta)}}{p(Y)}$$

Let's ~~calculate~~ approximate the posterior.

$$p(\theta, Z|Y) = \frac{p(Y|Z, \theta) p(Z|\theta) p(\theta)}{p(Y)}$$

Variational Bayes (VB): find a $q(\theta, Z)$ that is

- (a) Easy to deal with and**
- (b) Close to $p(\theta, Z|Y)$ in some sense**

Variational Bayes.

Define a class of approximating distribution.

$$q(\theta, Z|\vartheta, \zeta) := \delta(\theta - \vartheta) q(Z|\zeta)$$

Variational Bayes.

Define a class of approximating distribution.

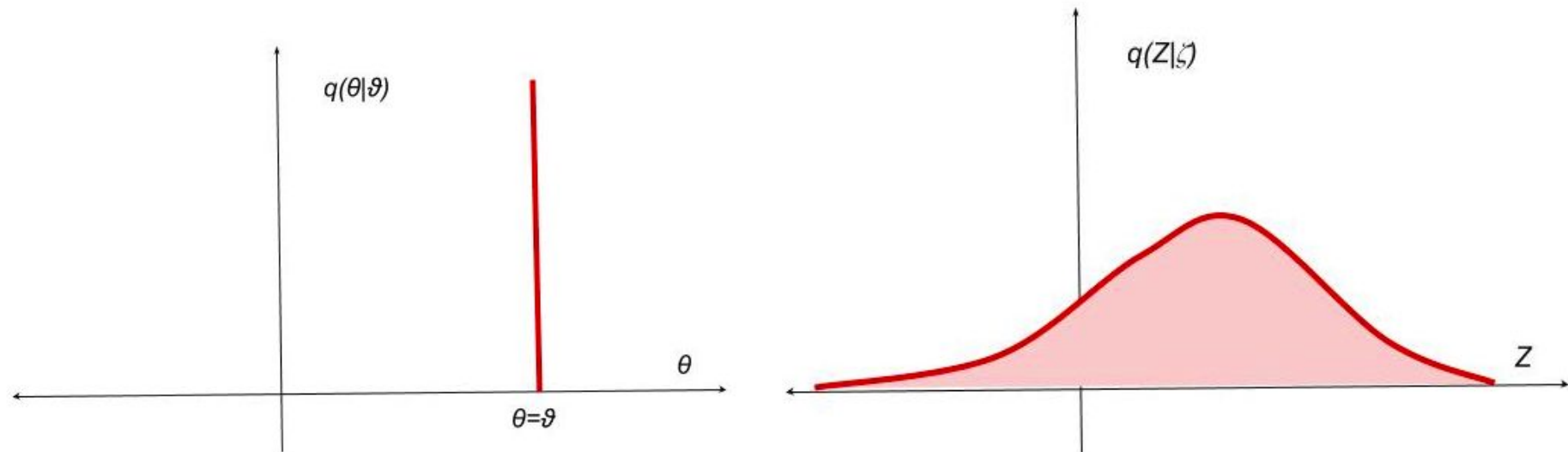
Degenerate at ϑ

$$q(\theta, Z | \vartheta, \zeta) := \overbrace{\delta(\theta - \vartheta)}^{\text{Degenerate at } \vartheta} \underbrace{q(Z | \zeta)}_{\text{Non-degenerate, in some parametric family:}}$$

Non-degenerate, in some
parametric family:

$$q(Z | \zeta) \in \mathcal{Q}$$

$$q(\theta, Z|\vartheta, \zeta) := \delta(\theta - \vartheta) q(Z|\zeta)$$



Variational Bayes.

Estimate using a Kullback-Leibler (KL)-like divergence.

$$\hat{\vartheta}, \hat{\zeta} = \operatorname{argmin}_{\vartheta, \zeta} \widetilde{KL} (q(\theta, Z | \vartheta, \zeta) || p(\theta, Z | Y))$$

$$\widetilde{KL} (q (\theta, Z|\vartheta, \zeta) || p (\theta, Z|Y))$$

$$= \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\log q (Z|\zeta)] - \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\ell (\theta, Z|Y)]$$



**Entropy of θ is missing -- in
this sense it's not a real KL
divergence**

$$\widetilde{KL} (q (\theta, Z|\vartheta, \zeta) || p (\theta, Z|Y))$$

$$= \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\log q (Z|\zeta)] - \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\ell (\theta, Z|Y)]$$

$$= \mathbb{E}_{q(Z|\zeta)} [\log q (Z|\zeta) - \ell (\vartheta, Z|Y)]$$

$$\begin{aligned}
& \widetilde{KL} (q (\theta, Z|\vartheta, \zeta) || p (\theta, Z|Y)) \\
&= \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\log q (Z|\zeta)] - \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\ell (\theta, Z|Y)] \\
&= \mathbb{E}_{q(Z|\zeta)} [\log q (Z|\zeta) - \ell (\vartheta, Z|Y)] \\
&= \mathbb{E}_{q(Z|\zeta)} [\log q (Z|\zeta) - (\ell (Y|Z, \vartheta) + \ell (Z|\vartheta) - \ell (\vartheta) + \ell (Y))]
\end{aligned}$$

$$\begin{aligned}
& \widetilde{KL} (q (\theta, Z|\vartheta, \zeta) || p (\theta, Z|Y)) \\
&= \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\log q (Z|\zeta)] - \mathbb{E}_{q(\theta, Z|\vartheta, \zeta)} [\ell (\theta, Z|Y)] \\
&= \mathbb{E}_{q(Z|\zeta)} [\log q (Z|\zeta) - \ell (\vartheta, Z|Y)] \\
&= \mathbb{E}_{q(Z|\zeta)} [\log q (Z|\zeta) - (\ell (Y|Z, \vartheta) + \ell (Z|\vartheta) - \ell (\vartheta) + \ell (Y))] \\
&= -\mathbb{E}_{q(Z|\zeta)} [\ell (Y|Z, \vartheta) + \ell (Z|\vartheta) - \ell (\vartheta) - \log q (Z|\zeta)] + \ell (Y)
\end{aligned}$$

Contrast with the EM identity.

$$\ell (Y|\theta) = \mathbb{E}_{p(Z|Y, \theta)} [\ell (Y|Z, \theta) + \ell (Z|\theta) - \ell (Z|Y, \theta)]$$

Proposition: VB = EM

(Neal and Hinton, 1998)

Suppose that $p(Z|\theta, Y) \in \mathcal{Q}, \forall \theta \in \Omega_\theta$

Then the VB and EM optima are the same:

$$\hat{\vartheta} = \hat{\theta}$$

$$q(Z|\hat{\zeta}) = p(Z|Y, \hat{\theta})$$

Proposition: $\mathbf{VB} = \mathbf{EM}$

(Neal and Hinton, 1998)

**We can always find a parametric class \mathcal{Q} that satisfies this condition.
(Why?)**

$$p(Z|\theta, Y) \in \mathcal{Q}, \forall \theta \in \Omega_\theta$$

Proposition: VB = EM

(Neal and Hinton, 1998)

From now on q will be used for the conditional distribution of Z .

$$q \left(Z | \hat{\zeta} \right) = p \left(Z | Y, \hat{\theta} \right)$$

Proposition: VB = EM

(Neal and Hinton, 1998):

In fact, the EM algorithm is coordinate ascent in ϑ, ζ .

Proposition: VB = EM

(Neal and Hinton, 1998):

Given iteration k , $\vartheta^{(k)}, \zeta^{(k)}$

Step 1k. “E step”:

Calculate $\hat{\zeta}^{(k+1)} = \underset{\zeta}{\operatorname{argmin}} \widetilde{KL} \left(q \left(\theta, Z | \hat{\vartheta}^{(k)}, \zeta \right) || p \left(\theta, Z | Y \right) \right)$

Proposition: VB = EM

(Neal and Hinton, 1998):

Given iteration k , $\vartheta^{(k)}, \zeta^{(k)}$

Step 1k. “E step”:

Calculate $\hat{\zeta}^{(k+1)} = \underset{\zeta}{\operatorname{argmin}} \widetilde{KL} \left(q \left(\theta, Z | \hat{\vartheta}^{(k)}, \zeta \right) || p \left(\theta, Z | Y \right) \right)$

Step 2k. “M step”:

Calculate $\hat{\vartheta}^{(k+1)} = \underset{\vartheta}{\operatorname{argmin}} \widetilde{KL} \left(q \left(\theta, Z | \vartheta, \hat{\zeta}^{(k+1)} \right) || p \left(\theta, Z | Y \right) \right)$

...repeat.

Who uses coordinate ascent?



$$\hat{\vartheta}, \hat{\zeta} = \underset{\vartheta, \zeta}{\operatorname{argmin}} \widetilde{KL}(q(\theta, Z | \vartheta, \zeta) || p(\theta, Z | Y))$$

Part 3:

Covariance asymptotics.

Assumption : Bayesian CLT. Bernstein-von Mises (BVM) theorem

$$\mathcal{I}_{\theta\theta} = - \lim_{n \rightarrow \infty} \frac{1}{N} \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \Big|_{\hat{\theta}}$$

Covariance from the Laplace
approximation

Assumption : Bayesian CLT. Bernstein-von Mises (BVM) theorem


$$\mathcal{I}_{\theta\theta} = - \lim_{n \rightarrow \infty} \frac{1}{N} \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \Big|_{\hat{\theta}}$$

$$p \left(\sqrt{N} \mathcal{I}_{\theta\theta}^{1/2} \left(\theta - \hat{\theta} \right) | Y \right) \xrightarrow{N \rightarrow \infty} \mathcal{N} (0, I_D)$$

Assumption : Bayesian CLT.

Bernstein-von Mises (BVM) theorem

**The posterior on θ goes to
a degenerate distribution**


$$p \left(\sqrt{N} \mathcal{I}_{\theta\theta}^{1/2} \left(\theta - \hat{\theta} \right) | Y \right) \xrightarrow{N \rightarrow \infty} \mathcal{N} (0, I_D)$$

How good are the VB approximation's covariances?

$$\text{Cov}_{p(\theta|Y)}(\theta) \approx \frac{1}{N} \mathcal{I}_{\theta\theta}^{-1} = o_p\left(\sqrt{N}\right) \quad \text{Bayesian CLT}$$

How good are the VB approximation's covariances?

$$\text{Cov}_{p(\theta|Y)}(\theta) \approx \frac{1}{N} \mathcal{I}_{\theta\theta}^{-1} = o_p\left(\sqrt{N}\right) \quad \text{Bayesian CLT}$$

$$\text{Cov}_{q(\theta|\hat{v})}(\theta) \equiv 0$$

**Degenerate
approximation**

How good are the VB approximation's covariances?

$$\text{Cov}_{p(\theta|Y)}(\theta) \approx \frac{1}{N} \mathcal{I}_{\theta\theta}^{-1} = o_p\left(\sqrt{N}\right) \quad \text{Bayesian CLT}$$

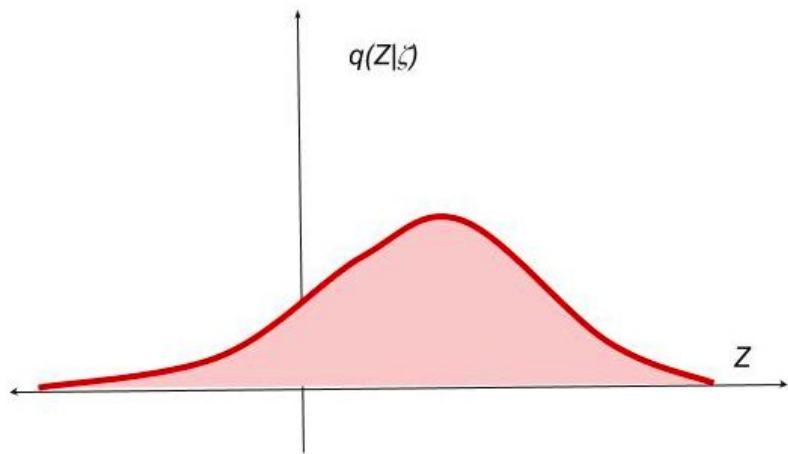
$$\text{Cov}_{q(\theta|\hat{v})}(\theta) \equiv 0$$

**Degenerate
approximation**

$$\text{Cov}_{p(\theta|Y)}(\theta) - \text{Cov}_{q(\theta|\hat{v})}(\theta) = o_p\left(\sqrt{N}\right) \quad \text{Consistent! (But trivial.)}$$

What about the variance of Z ?

$$\text{Cov}_{p(Z_n|Y)}(Z_n) \xrightarrow{N \rightarrow \infty} \text{Cov}_{p(Z_n|Y, \theta=\hat{\theta})}(Z_n) \neq 0$$



What about the variance of \mathbf{Z} ?

$$\text{Cov}_{p(Z_n|Y)}(Z_n) \xrightarrow{N \rightarrow \infty} \text{Cov}_{p(Z_n|Y, \theta = \hat{\theta})}(Z_n) \neq 0$$

\parallel

$$\text{Cov}_q(Z_n|\hat{\zeta})(Z_n)$$

What about the variance of \mathbf{Z} ?

$$\text{Cov}_{p(Z_n|Y)}(Z_n) \xrightarrow{N \rightarrow \infty} \text{Cov}_{p(Z_n|Y, \theta=\hat{\theta})}(Z_n) \neq 0$$
$$\parallel$$
$$\text{Cov}_{q(Z_n|\hat{\zeta})}(Z_n)$$

Therefore, trivially,

$$\text{Cov}_{p(Z_n|Y)}(Z_n) - \text{Cov}_{q(Z_n|\hat{\zeta})}(Z_n) = o_p(1)$$

What about the variance of \mathbf{Z} ?

$$\text{Cov}_{p(Z_n|Y)}(Z_n) \xrightarrow{N \rightarrow \infty} \text{Cov}_{p(Z_n|Y, \theta = \hat{\theta})}(Z_n) \neq 0$$
$$\parallel$$
$$\text{Cov}_{q(Z_n|\hat{\zeta})}(Z_n)$$

Only slightly less trivially:

$$\text{Cov}_{p(Z_n|Y)}(Z_n) - \text{Cov}_{q(Z_n|\hat{\zeta})}(Z_n) = o_p\left(\sqrt{N}\right)$$

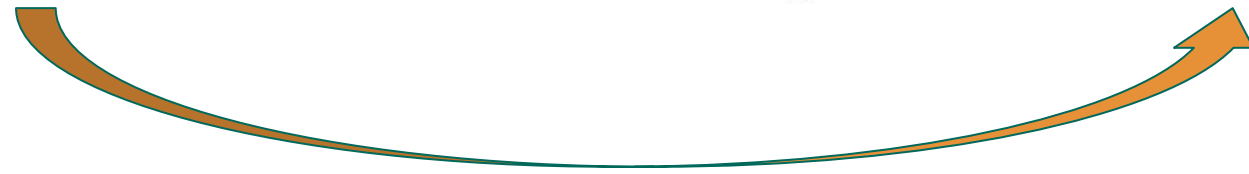
Can these naive approximations be improved?

Let's try linear response covariances.

Linear response covariances reminder.

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$



$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$


$$\hat{\theta}(t) = \underset{\theta}{\operatorname{argmin}} (f(\theta) + t\theta)$$

$$\left. \frac{d\hat{\theta}}{dt} \right|_{t=0} = \left(\left. \frac{\partial^2 f(\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} \right)^{-1}$$

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$

Define a perturbation:

$$p(\theta, Z|Y, t) \propto p(\theta, Z|Y) \exp(t\theta)$$

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$

Define a perturbation:

$$p(\theta, Z|Y, t) \propto p(\theta, Z|Y) \exp(t\theta)$$

$$\left. \frac{d\mathbb{E}_{p(\theta, Z|Y, t)} [\theta]}{dt} \right|_{t=0} = \text{Cov}_{p(\theta, Z|Y)} (\theta)$$

$$\frac{d\mathbb{E} [\theta|\alpha, X]}{d\alpha} = \text{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$

Define a perturbation:

$$p(\theta, Z|Y, t) \propto p(\theta, Z|Y) \exp(t\theta)$$

$$\left. \frac{d\hat{\theta}(t)}{dt} \right|_{t=0} \approx \left. \frac{d\mathbb{E}_{p(\theta, Z|Y, t)} [\theta]}{dt} \right|_{t=0} = \text{Cov}_{p(\theta, Z|Y)}(\theta)$$

Sensitivity requires the Hessian of the optimum.

Fixed dimension

$$H := \left. \frac{d^2 \widehat{\text{KL}}(\eta)}{d\eta d\eta^\top} \right|_{\hat{\eta}} = \begin{pmatrix} H_{\theta\theta} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

**Grows
with N**

$$\eta := \begin{pmatrix} \vartheta \\ \zeta \end{pmatrix}$$

Linear response covariances for theta.

$$H := \left. \frac{d^2 \widehat{\text{KL}}(\eta)}{d\eta d\eta^\top} \right|_{\hat{\eta}} = \begin{pmatrix} H_{\theta\theta} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

Linear response covariances for theta.

$$H := \frac{d^2 \widehat{\text{KL}}(\eta)}{d\eta d\eta^\top} \bigg|_{\hat{\eta}} = \begin{pmatrix} H_{\theta\theta} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

$$\Sigma = H^{-1} = \begin{pmatrix} \Sigma_{\theta\theta} & \Sigma_{\theta\zeta} \\ \Sigma_{\zeta\theta} & \Sigma_{\zeta\zeta} \end{pmatrix}$$

Linear response covariances for theta.

$$H := \left. \frac{d^2 \widehat{\mathbf{KL}}(\eta)}{d\eta d\eta^\top} \right|_{\hat{\eta}} = \begin{pmatrix} H_{\theta\theta} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

$$\Sigma = H^{-1} = \begin{pmatrix} \Sigma_{\theta\theta} & \Sigma_{\theta\zeta} \\ \Sigma_{\zeta\theta} & \Sigma_{\zeta\zeta} \end{pmatrix}$$

☀ $\widehat{\text{Cov}}_{LR}(\theta) = \Sigma_{\theta\theta}$

Proposition 1:

With a flat prior, linear response posterior covariances for θ and classical frequentist covariance estimates are the same.

$$\widehat{\text{COV}}_{LR}(\theta) = \Sigma_{\theta\theta} = - \left(\frac{d^2 \ell(Y|\theta)}{d\theta d\theta^\top} \Big|_{\hat{\theta}} \right)^{-1}$$

Proposition 1 proof.

$$\widehat{\text{KL}}(\vartheta, \zeta) = -\mathbb{E}_{q(Z|\zeta)} [\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \log q(Z|\zeta)] + \ell(Y)$$

Definition, flat prior.

Proposition 1 proof.

$$\widehat{\text{KL}}(\vartheta, \zeta) = -\mathbb{E}_{q(Z|\zeta)} [\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \log q(Z|\zeta)] + \ell(Y)$$

Define

$$\hat{\zeta}(\vartheta) = \underset{\zeta}{\operatorname{argmin}} \widehat{\text{KL}}(\vartheta, \zeta)$$

$$q(Z|\hat{\zeta}(\theta)) = p(Z|Y, \theta)$$

Completeness of VB
approximation.

Proposition 1 proof.

$$\widehat{\text{KL}} \left(\vartheta, \hat{\zeta}(\vartheta) \right) = -\mathbb{E}_{q(Z|\hat{\zeta}(\vartheta))} \left[\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \log q \left(Z|\hat{\zeta}(\vartheta) \right) \right] + \ell(Y)$$

Definition.

Proposition 1 proof.

$$\begin{aligned}\widehat{\text{KL}}\left(\vartheta, \hat{\zeta}(\vartheta)\right) &= -\mathbb{E}_{q(Z|\hat{\zeta}(\vartheta))}\left[\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \log q\left(Z|\hat{\zeta}(\vartheta)\right)\right] + \ell(Y) \\ &= -\mathbb{E}_{p(Z|Y, \vartheta)}\left[\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \ell(Z|Y, \vartheta)\right] + \ell(Y)\end{aligned}$$

Completeness of VB
approximation.

Proposition 1 proof.

$$\begin{aligned}\widehat{\text{KL}}\left(\vartheta, \hat{\zeta}(\vartheta)\right) &= -\mathbb{E}_{q(Z|\hat{\zeta}(\vartheta))}\left[\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \log q\left(Z|\hat{\zeta}(\vartheta)\right)\right] + \ell(Y) \\ &= -\mathbb{E}_{p(Z|Y, \vartheta)}\left[\ell(Y|Z, \vartheta) + \ell(Z|\vartheta) - \ell(Z|Y, \vartheta)\right] + \ell(Y) \\ &= -\ell(Y|\vartheta) + \ell(Y).\end{aligned}$$

The EM identity.

Proposition 1 proof.

All these perturbed optimization problems are the same:

$$\hat{\vartheta}(t), \hat{\zeta}(t) = \operatorname{argmax}_{\vartheta, \zeta} - \widehat{\text{KL}}(\vartheta, \zeta) + t\vartheta$$

$$\hat{\hat{\vartheta}}(t) = \operatorname{argmax}_{\vartheta} - \widehat{\text{KL}}\left(\vartheta, \hat{\zeta}(\vartheta)\right) + t\vartheta$$

$$\hat{\theta}(t) = \operatorname{argmax}_{\theta} \ell(Y|\theta) + t\theta$$

$$\hat{\vartheta}(t) = \hat{\hat{\vartheta}}(t) = \hat{\theta}(t)$$

Recall our Q-function

$$\hat{Q}(\theta) = \mathbb{E}_{p(Z|Y, \hat{\theta})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta}) \right]$$

Recall our Q-function

$$\hat{Q}(\theta) = \mathbb{E}_{p(Z|Y, \hat{\theta})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta}) \right]$$

**This is the Hessian of
the Q-function**

$$H = \left(\begin{array}{cc} \overbrace{H_{\theta\theta}} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{array} \right)$$

Recall our Q-function

$$\hat{Q}(\theta) = \mathbb{E}_{p(Z|Y, \hat{\theta})} \left[\ell(Y|Z, \theta) + \ell(Z|\theta) - \ell(Z|Y, \hat{\theta}) \right]$$

**This is the Hessian of
the Q-function**

$$H = \begin{pmatrix} \overbrace{H_{\theta\theta}} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

**This is the linear
response covariance**

$$\Sigma = H^{-1} = \begin{pmatrix} \overbrace{\Sigma_{\theta\theta}} & \Sigma_{\theta\zeta} \\ \Sigma_{\zeta\theta} & \Sigma_{\zeta\zeta} \end{pmatrix}$$

$$\left. \frac{\partial^2 \ell(Y|\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}} \neq \left. \frac{\partial^2 \hat{Q}(\theta)}{\partial \theta \partial \theta} \right|_{\hat{\theta}}$$

$$\Sigma_{\theta\theta}^{-1} \neq H_{\theta\theta}^{-1}$$

This is the Hessian of
the Q-function

$$H = \begin{pmatrix} \overbrace{H_{\theta\theta}} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

This is the linear
response covariance

$$\Sigma = H^{-1} = \begin{pmatrix} \overbrace{\Sigma_{\theta\theta}} & \Sigma_{\theta\zeta} \\ \Sigma_{\zeta\theta} & \Sigma_{\zeta\zeta} \end{pmatrix}$$

Theorem 1:

The linear response covariances add a root-N order of accuracy to the following covariances.

$$\text{Cov}_{p(\theta|Y)}(\theta) - \widehat{\text{Cov}}_{LR}(\theta) = o_p\left(\frac{1}{N}\right)$$

Theorem 1:

The linear response covariances add a root-N order of accuracy to the following covariances.

$$\text{Cov}_{p(\theta|Y)}(\theta) - \widehat{\text{Cov}}_{LR}(\theta) = o_p\left(\frac{1}{N}\right)$$

$$\text{Cov}_{p(\theta|Y)}(\theta, Z_n) - \widehat{\text{Cov}}_{LR}(\theta, Z_n) = o_p\left(\frac{1}{N}\right)$$

Theorem 1:

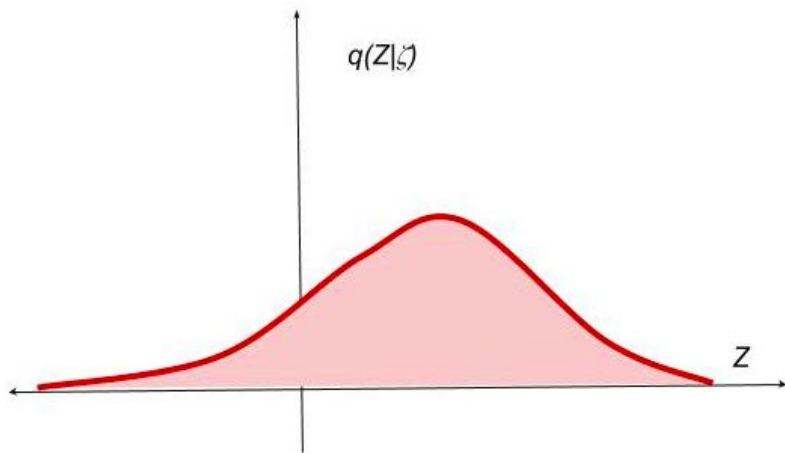
The linear response covariances add a root-N order of accuracy to the following covariances.

$$\text{Cov}_{p(\theta|Y)}(\theta) - \widehat{\text{Cov}}_{LR}(\theta) = o_p\left(\frac{1}{N}\right)$$

$$\text{Cov}_{p(\theta|Y)}(\theta, Z_n) - \widehat{\text{Cov}}_{LR}(\theta, Z_n) = o_p\left(\frac{1}{N}\right)$$

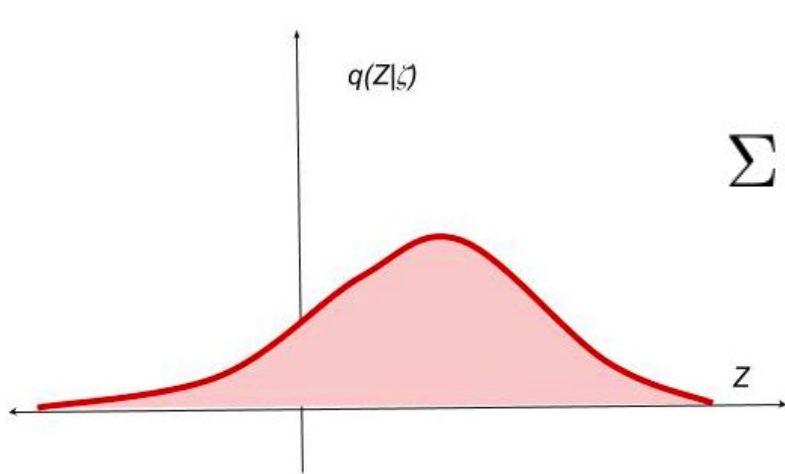
$$\text{Cov}_{p(\theta|Y)}(Z_n, Z_m) - \widehat{\text{Cov}}_{LR}(Z_n, Z_m) = o_p\left(\frac{1}{N}\right) \quad (\text{for } n \neq m)$$

What about Z?



What about Z?

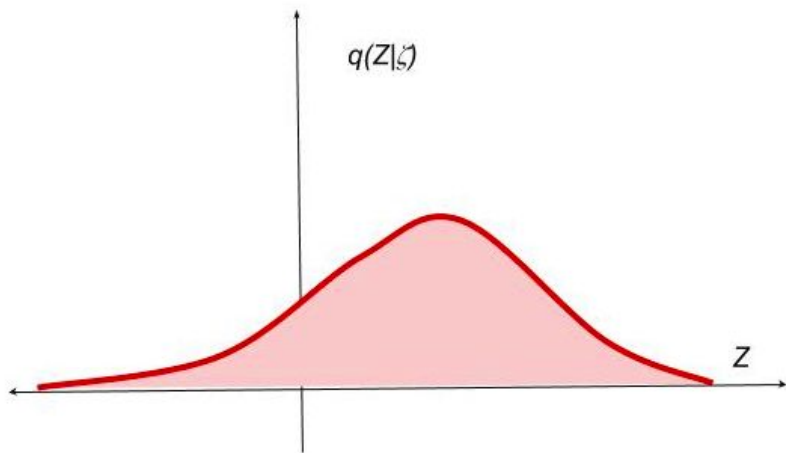
And what about these other Hessian terms?



$$\Sigma = H^{-1} = \left(\begin{array}{cc} \Sigma_{\theta\theta} & \overbrace{\Sigma_{\theta\zeta}} \\ \underbrace{\Sigma_{\zeta\theta}} & \underbrace{\Sigma_{\zeta\zeta}} \end{array} \right) \}$$

What about \mathbf{Z} ?

Results are best expressed in terms of the score function and its variance.

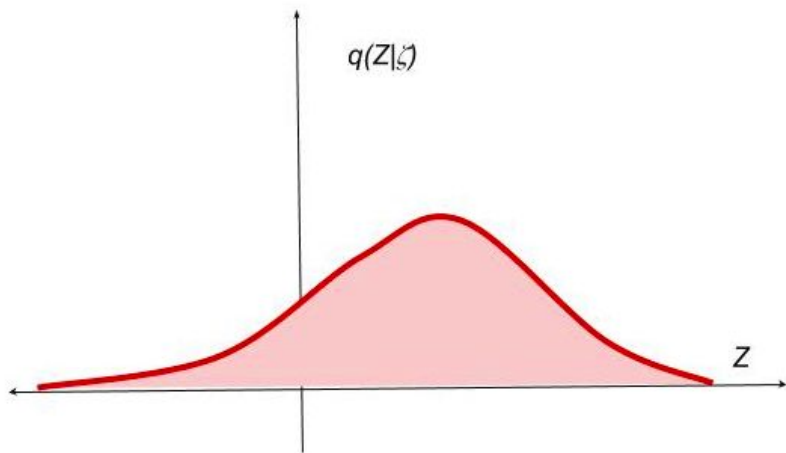


$$\hat{\gamma}(Z_n) := \left. \frac{\partial \log q(Z_n|\zeta)}{\partial \zeta} \right|_{\hat{\zeta}}$$

$$V_n := \text{Cov}_{q(Z_n|\hat{\zeta})}(\hat{\gamma}(Z_n))$$

What about \mathbf{Z} ?

If q is in the exponential family, $\gamma(Z)$ are sufficient statistics.



$$\hat{\gamma}(Z_n) := \left. \frac{\partial \log q(Z_n|\zeta)}{\partial \zeta} \right|_{\hat{\zeta}}$$

$$V_n := \text{Cov}_{q(Z_n|\hat{\zeta})}(\hat{\gamma}(Z_n))$$

Proposition 2:

The linear response covariances of the normalized score function are given by the inverse Hessian of the KL divergence.

$$\widehat{\text{Cov}}_{LR} \left(V_n^{-1} \hat{\gamma} (Z_n) \right) = \Sigma_{\zeta_n \zeta_n}$$

$$\hat{\gamma} (Z_n) := \left. \frac{\partial \log q (Z_n | \zeta)}{\partial \zeta} \right|_{\hat{\zeta}}$$
$$V_n := \text{Cov}_{q(Z_n | \hat{\zeta})} (\hat{\gamma} (Z_n))$$

Theorem 2:

The linear response covariances improves the constant, but not the rate, of the covariances of the score function.

Theorem 2:

The linear response covariances improves the constant, but not the rate, of the covariances of the score function.

Error when using the degenerate approximation for $p(\theta)$:

$$\text{Cov}_{p(\theta, Z|Y)}(\hat{\gamma}(Z)) - \text{Cov}_{q(Z|\hat{\zeta})}(\hat{\gamma}(Z)) = \left(\mathbb{E}_{q(Z_n|\hat{\zeta})} [\hat{\gamma}(Z)^3] + 1 \right) o_p\left(\frac{1}{\sqrt{N}}\right) + \text{Cov}_{p(\theta|Y)}\left(\mathbb{E}_{p(Z|Y, \theta)}[\hat{\gamma}(Z)]\right)$$

Theorem 2:

The linear response covariances improves the constant, but not the rate, of the covariances of the score function.

Skewness (small when q is approximately symmetric).

$$\begin{aligned} \text{Cov}_{p(\theta, Z|Y)}(\hat{\gamma}(Z)) - \text{Cov}_{q(Z|\hat{\zeta})}(\hat{\gamma}(Z)) &= \overbrace{\left(\mathbb{E}_{q(Z_n|\hat{\zeta})} \left[\hat{\gamma}(Z)^3 \right] + 1 \right)} + o_p\left(\frac{1}{\sqrt{N}}\right) + \\ &\quad \underbrace{\text{Cov}_{p(\theta|Y)}\left(\mathbb{E}_{p(Z|Y, \theta)} \left[\hat{\gamma}(Z) \right] \right)}_{o_p\left(\frac{1}{\sqrt{N}}\right)} \end{aligned}$$

Theorem 2:

The linear response covariances improves the constant, but not the rate, of the covariances of the score function.

$$\text{Cov}_{p(\theta, Z|Y)}(\hat{\gamma}(Z)) - \text{Cov}_{q(Z|\hat{\zeta})}(\hat{\gamma}(Z)) = \left(\mathbb{E}_{q(Z_n|\hat{\zeta})} [\hat{\gamma}(Z)^3] + 1 \right) o_p\left(\frac{1}{\sqrt{N}}\right) + \text{Cov}_{p(\theta|Y)}\left(\mathbb{E}_{p(Z|Y, \theta)}[\hat{\gamma}(Z)]\right)$$

Error when using linear response covariances:

$$\text{Cov}_{p(\theta, Z|Y)}(\hat{\gamma}(Z)) - \widehat{\text{Cov}}_{LR}(\hat{\gamma}(Z)) = \left(\mathbb{E}_{q(Z_n|\hat{\zeta})} [\hat{\gamma}(Z)^3] + 1 \right) o_p\left(\frac{1}{\sqrt{N}}\right)$$

Theorem 3:

For covariances of functions of Z other than linear combinations of the score functions, linear response is inconsistent.

$$\text{Cov}_{p(\theta, Z|Y)}(h(Z_n)) - \widehat{\text{Cov}}_{LR}(h(Z_n)) = O_p(1)$$

Theorem 3:

For covariances of functions of Z other than linear combinations of the score functions, linear response is inconsistent.

$$\text{Cov}_{p(\theta, Z|Y)}(h(Z_n)) - \widehat{\text{Cov}}_{LR}(h(Z_n)) = O_p(1)$$

Practical workarounds:

- Increase the expressivity of q
- Use Monte Carlo instead of linear response

Tools.

“However, analytical evaluation of the second-order derivatives of the incomplete-data log likelihood may be difficult or at least tedious. Indeed, often it is for reasons of this nature that the EM algorithm is used to compute the MLE in the first place.”

- The EM Algorithm and Extensions, McLachlan (2008)

“However, analytical evaluation of the second-order derivatives of the incomplete-data log likelihood may be difficult or at least tedious. Indeed, often it is for reasons of this nature that the EM algorithm is used to compute the MLE in the first place.”

- The EM Algorithm and Extensions, McLachlan (2008)

Automatic differentiation makes everything easy.

Automatic e-steps (for conjugate Z distributions):

```
def log_likelihood(z, log_z, theta, y):  
    # ... return scalar log posterior value.  
  
    # Implemented with autodiff. Takes as an argument a function whose first  
    # two arguments are z and log_z and which evaluates the log probability..  
    get_e_step = get_gamma_e_step_funs(log_posterior)  
  
    y = load_data()  
    done = False  
    theta_k = # ... some init value  
    while not done:  
        # This is automatic.  
        e_z, log_z = get_e_step(theta_k, y)  
  
        # Your favorite M-step routine here.  
        theta_k = optimize(lambda theta: log_posterior(e_z, e_log_z, theta, y))  
  
        # Check convergence and set done
```

But why bother with an e-step?

```
def log_likelihood(z, log_z, theta, y):  
    # ... return scalar log posterior value.  
      
    # Implemented with autodiff. Takes as an argument a function whose first  
    # two arguments are z and log_z and which evaluates the log probability.  
    get_marginal_log_lik = get_gamma_marginal_log_lik(log_posterior)  
      
    y = load_data()  
    theta_hat = optimize(lambda theta: get_marginal_log_lik(theta, y))
```

This is what I actually do in practice.

And covariances are now easy.

```
def log_likelihood(z, log_z, theta, y):  
    # ... return scalar log posterior value.  
  
    # Implemented with autograd. Takes as an argument a function whose first  
    # two arguments are z and log_z and which evaluates the log probability.  
    get_marginal_log_lik = get_gamma_marginal_log_lik(log_posterior)  
  
    y = load_data()  
    theta_hat = optimize(lambda theta: get_marginal_log_lik(theta, y))  
  
    cov_theta_k = -1 * np.linalg.inv(  
        autograd.hessian(get_marginal_log_lik)(theta_hat, y))
```

(...and quite a lot of academic literature is obsolete.)

Paragami: “Parameter origami”

<https://github.com/rgiordan/paragami>

Converts parameter dictionaries and the functions that consume or return them between “folded” and “flat” representations.

All transformations are differentiable by autograd.

Paragami: “Parameter origami” <https://github.com/rgiordan/paragami>

Define “patterns” that describe your structured parameter sets.

```
[3] mvn_pattern = paragami.PatternDict(free_default=True)
     mvn_pattern['mean'] = paragami.NumericVectorPattern(length=dim)
     mvn_pattern['cov'] = paragami.PSDSymmetricMatrixPattern(size=dim)
```


Paragami: “Parameter origami” <https://github.com/rgiordan/paragami>

```
[36] true_mvn_par = dict()
      true_mvn_par['mean'] = mean_true
      true_mvn_par['cov'] = cov_true

      print('\nA dictionary of MVN parameters:\n{}'.format(
            true_mvn_par))

      mvn_par_free = mvn_pattern.flatten(true_mvn_par)
      print('\nA flat representation:\n{}'.format(
            mvn_pattern.flatten(true_mvn_par)))

      print('\nFolding recovers the original parameters:\n{}'.format(
            mvn_pattern.fold(mvn_par_free)))
```

A dictionary of MVN parameters:

```
{'cov': array([[1.1, 1. ],
               [1. , 1.1]]), 'mean': array([0.87367236, 0.21280422])}
```

A flat representation:

```
[ 0.87367236  0.21280422  0.04765509  0.95346259 -0.82797896]
```

Folding recovers the original parameters:

```
OrderedDict([('mean', array([0.87367236, 0.21280422])), ('cov', array([[1.1, 1. ],
                               [1. , 1.1]]))])
```

Vittles:

“Variational inference tools to leverage estimator sensitivity”

<https://github.com/rgiordan/vittles>

Calculates Taylor series approximations (to arbitrary order) of the dependence of optima on hyperparameters.

Linear response covariances are a special case.

Vittles:

“Variational inference tools to leverage estimator sensitivity”

<https://github.com/rgiordan/vittles>

```
def get_flat_kl_divergence(parameters):~
    ... # Return the KL divergence~
~

def get_posterior_moments_from_parameters(parameters):~
    ... # Return a vector of posterior moments~
~

optimal_parameters = optimize(get_flat_kl_divergence)~
~

get_lrvcov = vittles.LinearResponseCovariances(~
    ... get_flat_kl_divergence,~
    ... optimal_parameters)~
~

lrvcov = get_lrvcov.get_lr_covariance(~
    ... get_posterior_moments_from_parameters)~
```

Thank you for your attention!



Extra topics (in case anyone asks)

The Laplace approximation and linear response covariances.

Example: the Laplace approximation.

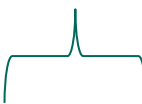
$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x)$$

Example: the Laplace approximation.

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x)$$

$$\mathbb{E}_{p(\theta|X)} [\theta] \approx \hat{\theta}$$

Cleverly chosen
perturbation


$$\hat{\theta}(\alpha) = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x) + \alpha\theta$$

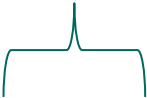
Cleverly chosen
perturbation

$$\hat{\theta}(\alpha) = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x) + \alpha \theta$$

Chosen so that this term
becomes θ

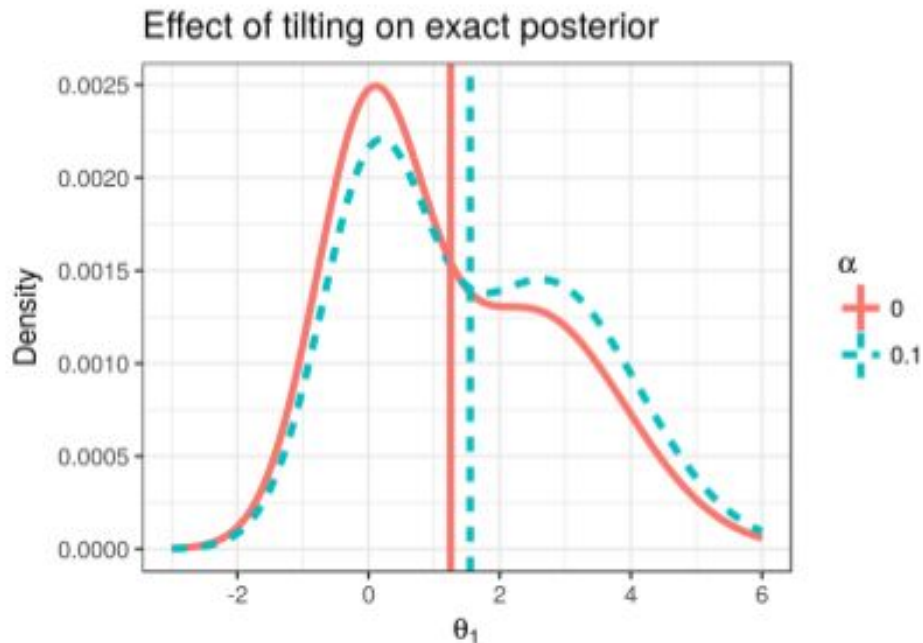
$$\frac{d\mathbb{E}[\theta|\alpha, X]}{d\alpha} = \operatorname{Cov}_{p(\theta|\alpha, X)} \left(\theta, \frac{\partial}{\partial \alpha} \log p(\theta|\alpha, X) \right)$$

Cleverly chosen
perturbation


$$\hat{\theta}(\alpha) = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x) + \alpha\theta$$

$$\mathbb{E}_{p(\theta|X, \alpha)} [\theta] \approx \hat{\theta}(\alpha)$$

$$\hat{\theta}(\alpha) = \operatorname{argmax}_{\theta \in \Omega_{\theta}} \log p(\theta|x) + \alpha\theta$$

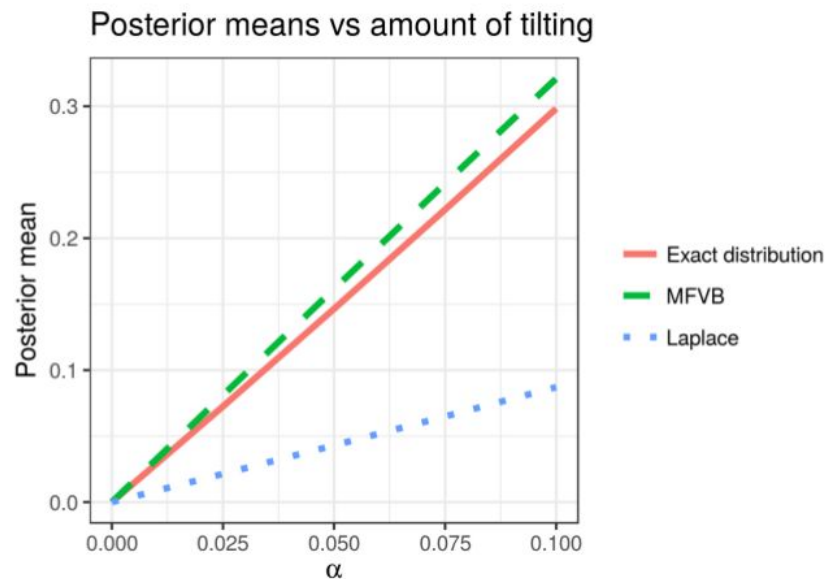
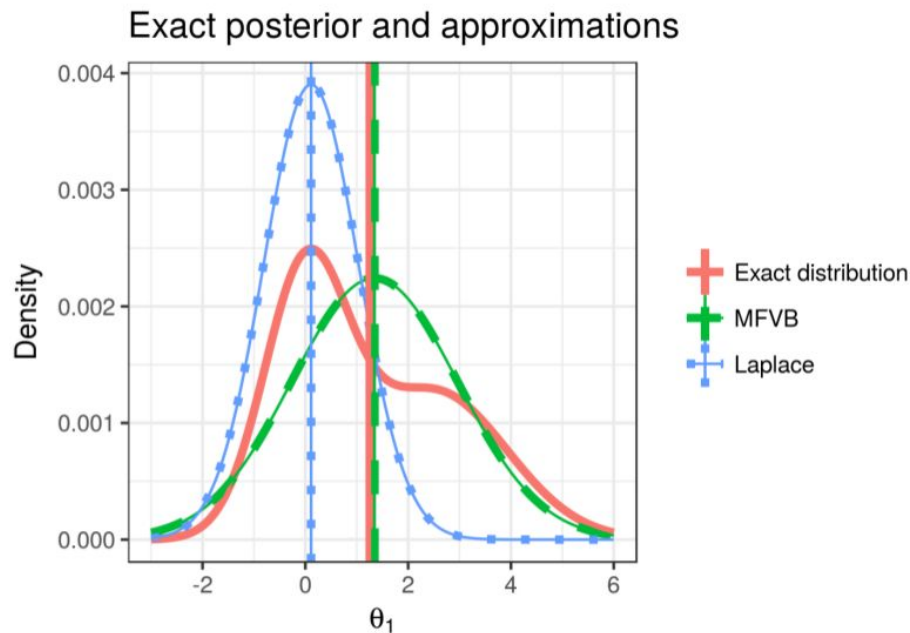


$$\text{Cov}_{p(\theta|X)}(\theta) = \left. \frac{d\mathbb{E}_{p(\theta|X,\alpha)}[\theta]}{d\alpha} \right|_{\alpha=0}$$

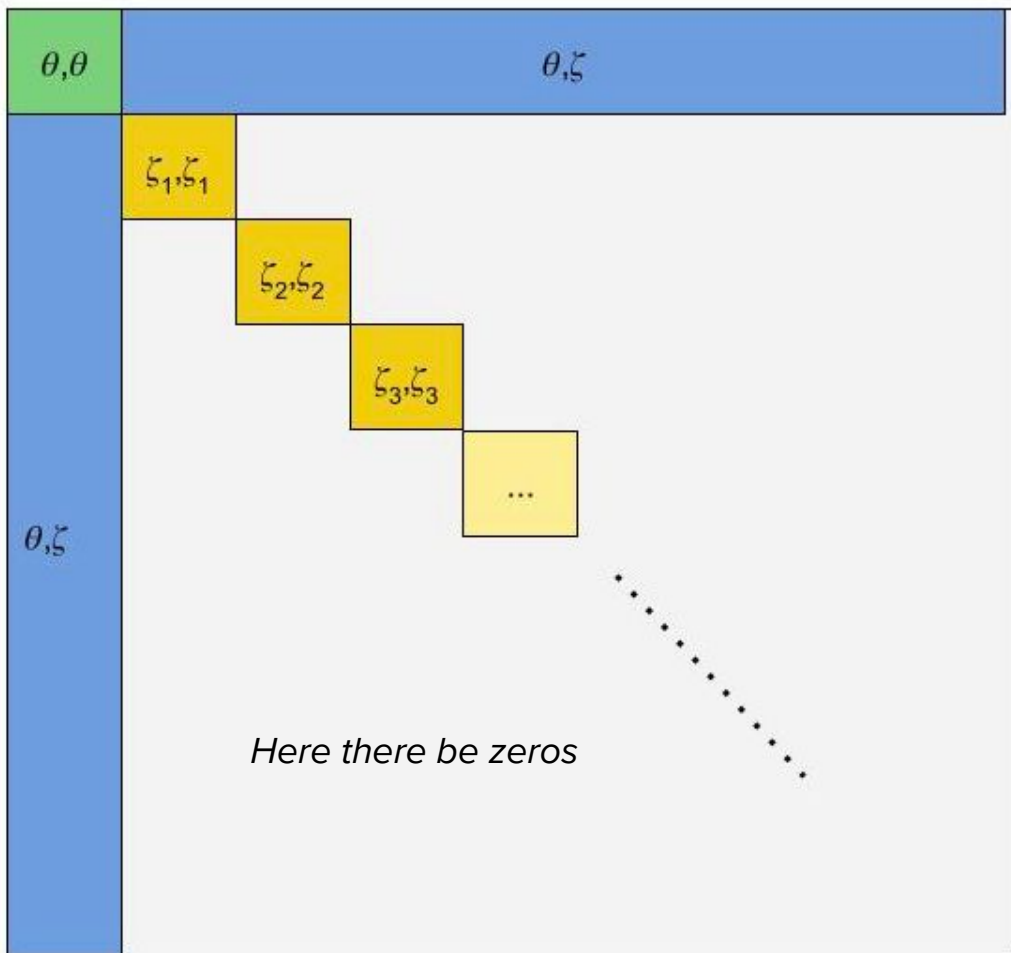
$$\begin{aligned} \text{Cov}_{p(\theta|X)}(\theta) &= \left. \frac{d\mathbb{E}_{p(\theta|X,\alpha)}[\theta]}{d\alpha} \right|_{\alpha=0} \\ &\approx \left. \frac{d\hat{\theta}(\alpha)}{d\alpha} \right|_{\alpha=0} \end{aligned}$$

$$\begin{aligned}
\text{Cov}_{p(\theta|X)}(\theta) &= \left. \frac{d\mathbb{E}_{p(\theta|X,\alpha)}[\theta]}{d\alpha} \right|_{\alpha=0} \\
&\approx \left. \frac{d\hat{\theta}(\alpha)}{d\alpha} \right|_{\alpha=0} \\
&= - \left(\left. \frac{\partial^2 \log p(\theta|x)}{\partial\theta\partial\theta} \right|_{\hat{\theta}} \right)^{-1}
\end{aligned}$$

In principle, linear response covariances can be calculated for *any optimization-based posterior approximation*.



Hessian picture



Fixed dimension

$$\begin{pmatrix} H_{\theta\theta} & H_{\theta\zeta} \\ H_{\zeta\theta} & H_{\zeta\zeta} \end{pmatrix}$$

Grows
with N