

2361-5

**School on Large Scale Problems in Machine Learning and Workshop on
Common Concepts in Machine Learning and Statistical Physics**

20 - 31 August 2012

**Large Scale Variational Bayesian Inference for Continuous Variable Models -
Lecture Notes**

Matthias SEEGER

*Laboratory for Probabilistic Machine Learning, EPFL, CH-1015 Lausanne
Switzerland*

Large Scale Variational Bayesian Inference for Continuous Variable Models

Matthias Seeger

Laboratory for Probabilistic Machine Learning
Ecole Polytechnique Fédérale de Lausanne

<http://lapmal.epfl.ch/>



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

- 1 Motivation
- 2 Variational Inference Relaxations
 - Super-Gaussian Bounding
 - Expectation Propagation
 - Gaussian KL Minimization
 - Conjugate Gradients Algorithm
- 3 Scalable Variational Inference
 - Scaling up Super-Gaussian Bounding
 - Penalized Least Squares
 - Gaussian Variances
- 4 Application Example

Outline

- 1 Motivation
- 2 Variational Inference Relaxations
 - Super-Gaussian Bounding
 - Expectation Propagation
 - Gaussian KL Minimization
 - Conjugate Gradients Algorithm
- 3 Scalable Variational Inference
 - Scaling up Super-Gaussian Bounding
 - Penalized Least Squares
 - Gaussian Variances
- 4 Application Example

Goals of Lecture

- Beyond point estimation:
Bayesian inference for non-Gaussian continuous variable models
- Beyond message passing:
Computational structure of variational inference relaxations
- The layer below:
Scalability through reductions to convex optimization and numerical mathematics

Image Reconstruction

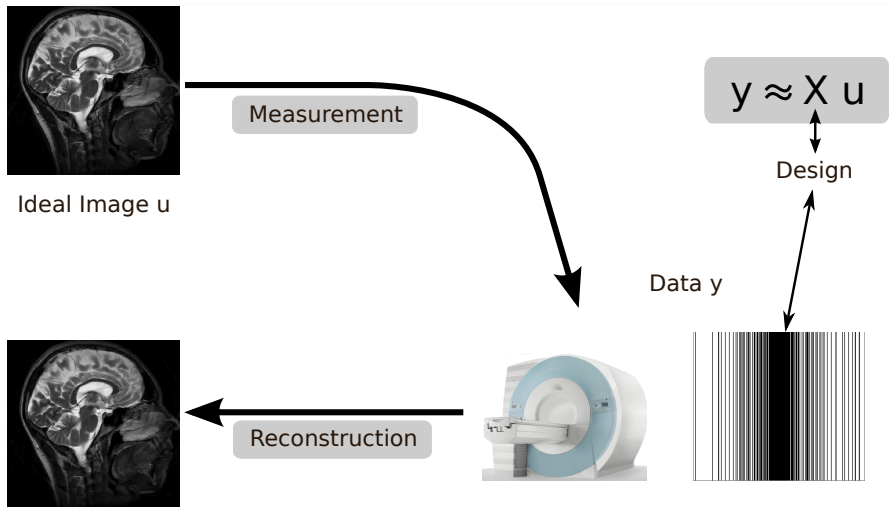
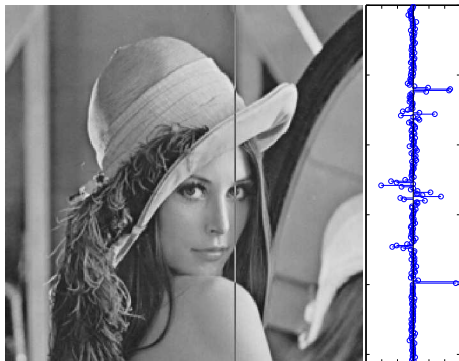
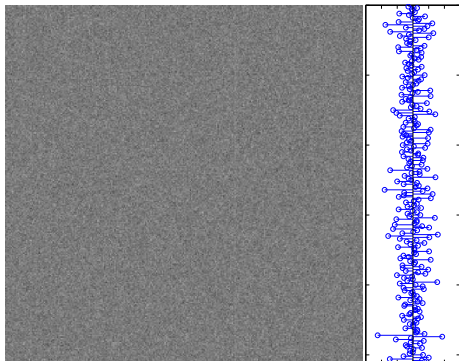


Image Statistics

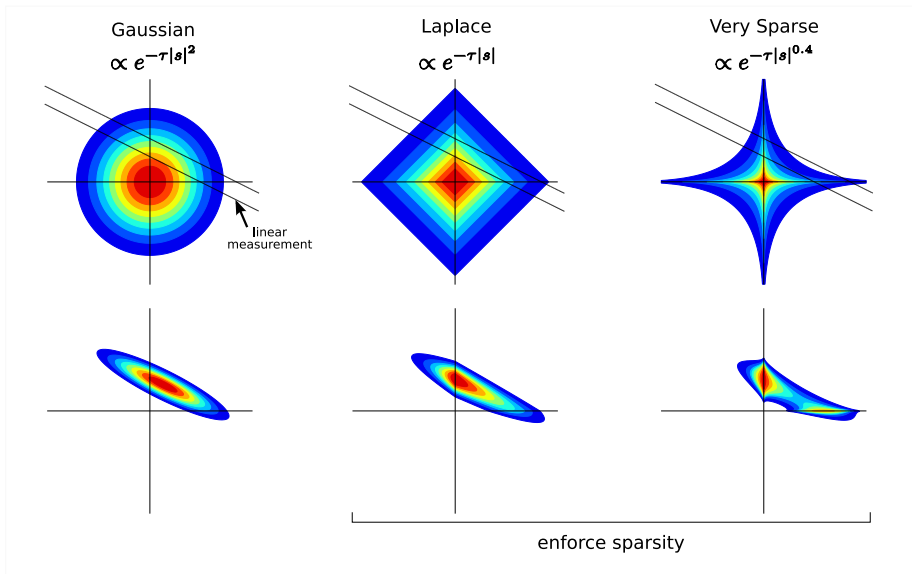
Whatever images are ...

they are not Gaussian!



Sparsity Priors

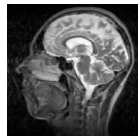
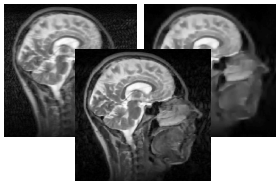
courtesy Florian Steinke



Posterior Distribution

- Likelihood $P(\mathbf{y}|\mathbf{u})$: Data fit

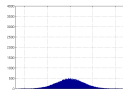
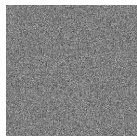
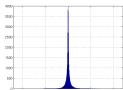
$$P(\mathbf{y}|\mathbf{u})$$



Posterior Distribution

- Likelihood $P(\mathbf{y}|\mathbf{u})$: Data fit
- Prior $P(\mathbf{u})$: Signal properties

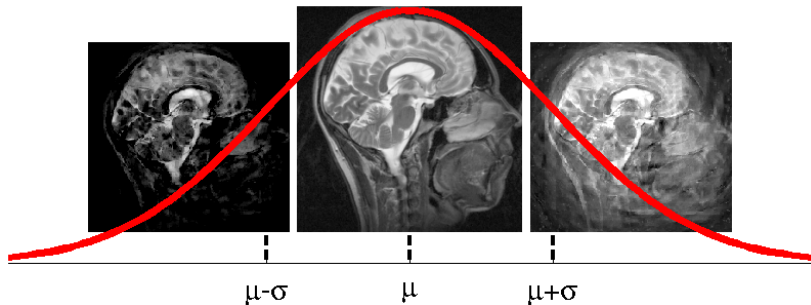
$$P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})$$



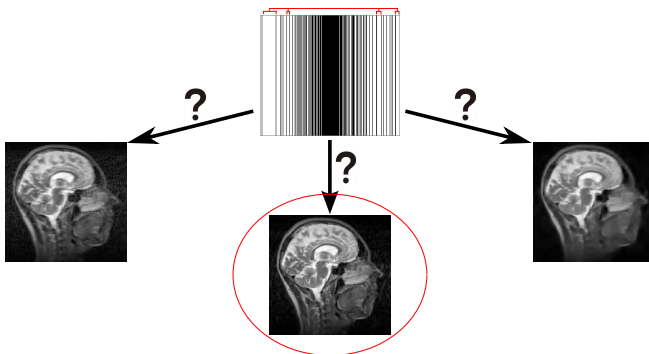
Posterior Distribution

- Likelihood $P(\mathbf{y}|\mathbf{u})$: Data fit
- Prior $P(\mathbf{u})$: Signal properties
- Posterior distribution $P(\mathbf{u}|\mathbf{y})$:
Consistent information summary

$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$



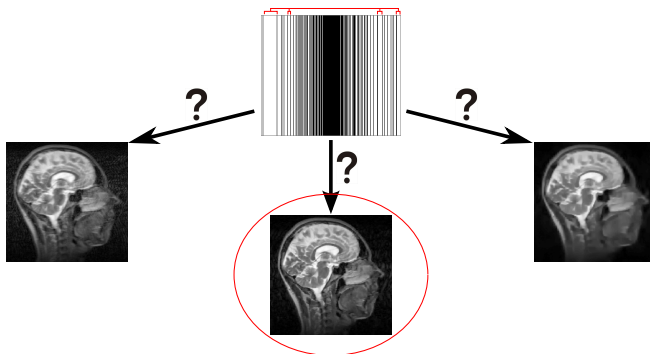
MAP Estimation



Maximum a Posteriori (MAP) Estimation

$$\mathbf{u}_* = \operatorname{argmax}_{\mathbf{u}} P(\mathbf{y}|\mathbf{u})P(\mathbf{u})$$

Why Move Beyond MAP?



Maximum a Posteriori (MAP) Estimation

$$\mathbf{u}_* = \operatorname{argmax}_{\mathbf{u}} P(\mathbf{y}|\mathbf{u})P(\mathbf{u})$$

Bayesian Calibration

y



k



u

www.wisdom.weizmann.ac.il/~levina

$$\mathbf{y} \approx \mathbf{k} \otimes \mathbf{u}$$

- Computer vision
 - Blind deconvolution
 - Calibrating camera parameters
- Magnetic resonance imaging
 - Autocalibrating parallel MRI

Bayesian Calibration

$$P(\mathbf{y}|\theta) = \int P(\mathbf{y}|\mathbf{u}, \theta)P(\mathbf{u}|\theta) d\mathbf{u}$$

Given raw data \mathbf{y} , no ground truth \mathbf{u} . Estimate model parameters θ .

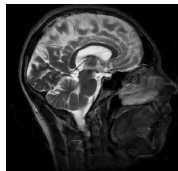
- Blind deconvolution (θ blur kernel)
- Multi-frame super-resolution (θ camera parameters, PSF)
- Image coding (θ codebook)
- Learning image priors ($P(\mathbf{u}) = P(\mathbf{u}|\theta)$)

Bayesian Experimental Design

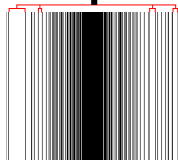
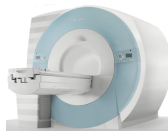
scan time \propto
 # phase encodes

$$y \approx X u$$

$$X \leftarrow ?$$

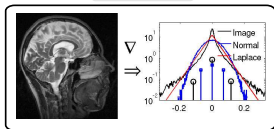


Reconstruction

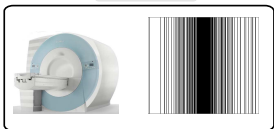


Bayesian Experimental Design

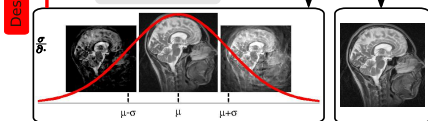
Prior $P(u)$



Data $P(y|u)$



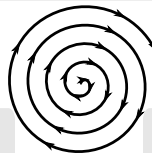
Posterior $P(u|y)$



Inference

Estimation

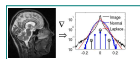
Measurement



Design
Decision

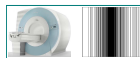
Posterior
Update

Sparse Linear Model

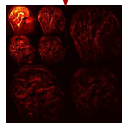


$$P(\mathbf{u}) \propto \prod_{i=1}^q t_i(s_i) =$$

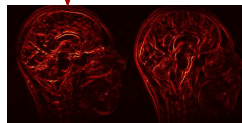
$$e^{-\tau_w \|\mathbf{B}_w \mathbf{u}\|_1} \times e^{-\tau_{tv} \|\mathbf{B}_{tv} \mathbf{u}\|_1}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$



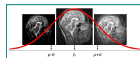
$$P(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I})$$



wavelet



gradient



$$P(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{u})P(\mathbf{y}|\mathbf{u})$$

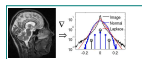
- \mathbf{X} , \mathbf{B} ? Fast operators of your choice (\mathbf{X} dictated by application)

Denoising: \mathbf{X} diagonal

Deconvolution: $\mathbf{X}\mathbf{u} = \mathbf{k} \otimes \mathbf{u}$

MRI reconstruction: $\mathbf{X} = \mathbf{I}_J \cdot \mathbf{F}$, \mathbf{F} DFT

Sparse Linear Model

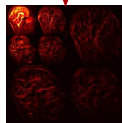


$$P(\mathbf{u}) \propto \prod_{i=1}^q t_i(s_i) =$$

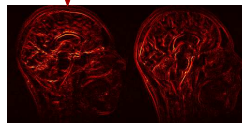
$$e^{-\tau_w \|\mathbf{B}_w \mathbf{u}\|_1} \times e^{-\tau_{tv} \|\mathbf{B}_{tv} \mathbf{u}\|_1}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$



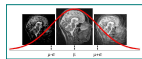
$$P(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I})$$



wavelet



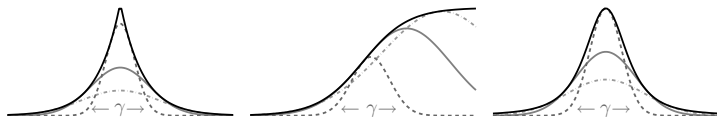
gradient



$$P(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{u})P(\mathbf{y}|\mathbf{u})$$

- \mathbf{X} , \mathbf{B} ? Fast operators of your choice (\mathbf{X} dictated by application)
- $t_i(s_i)$ Laplace here, but many other options

Nickisch, Seeger, ICML 2009



Outline

- 1 Motivation
- 2 Variational Inference Relaxations
 - Super-Gaussian Bounding
 - Expectation Propagation
 - Gaussian KL Minimization
 - Conjugate Gradients Algorithm
- 3 Scalable Variational Inference
 - Scaling up Super-Gaussian Bounding
 - Penalized Least Squares
 - Gaussian Variances
- 4 Application Example

Variational Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i), \quad Z = \int P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i) d\mathbf{u}$$

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable
- Integration tractable for **Gaussians** $Q(\mathbf{u}|\mathbf{y})$
⇒ Approximate $P(\mathbf{u}|\mathbf{y})$ by $Q(\mathbf{u}|\mathbf{y})!$

Variational approximation

Apply variational principle to fit master function $\log Z$

The Log Partition Function

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1}P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i), \quad Z = \int P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) d\mathbf{u}$$

Master function $\log Z$? Why this target?

- Physicist: Of course, it's the (negative) free energy!
- Probabilist: It generates posterior moments (cumulants)
- Variational definition of posterior distribution

$$E_{Q(\mathbf{u}|\mathbf{y})} \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i)}{Q(\mathbf{u}|\mathbf{y})} \right] \begin{cases} \operatorname{argmax}_{Q(\mathbf{u}|\mathbf{y})} & P(\mathbf{u}|\mathbf{y}) \\ \max_{Q(\mathbf{u}|\mathbf{y})} & \log Z \end{cases}$$

- Bayesian inference: **Optimization over distributions** Wainwright, Jordan, FTML 2008

Variational Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1}P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i), \quad Z = \int P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) d\mathbf{u}$$

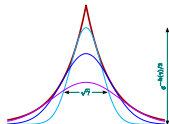
Variational approximation

Apply variational principle to fit master function $\log Z$

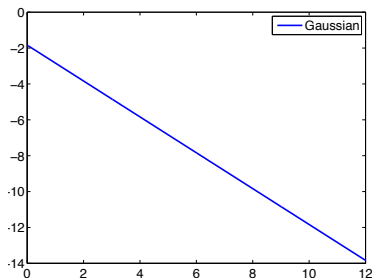
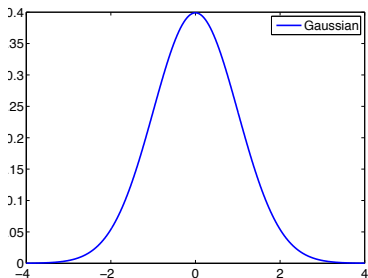
- **Super-Gaussian bounding**
- Expectation propagation
- Gaussian KL minimization

Super-Gaussian Potentials

$$t(s) = \max_{\gamma \geq 0} e^{-s^2/(2\gamma)} e^{-h(\gamma)/2}$$

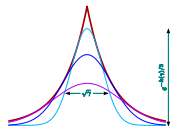


- $t(s)$ **even** and **positive**: Let's look at $s^2 \mapsto 2 \log t(s)$
- What's that for a Gaussian $t(s) = N(s|0, \sigma^2)$?
A **linear** (affine) function



Super-Gaussian Potentials

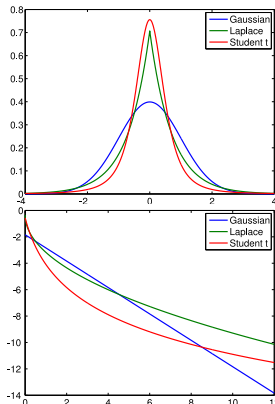
$$t(s) = \max_{\gamma \geq 0} e^{-s^2/(2\gamma)} e^{-h(\gamma)/2}$$



Sparsity potentials are **super-Gaussian**

$$s^2 \mapsto 2 \log t(s) \text{ is convex}$$

- Affine \rightarrow convex:
Shift mass to center and tails



Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $s^2 \mapsto 2 \log t(s)$ convex.

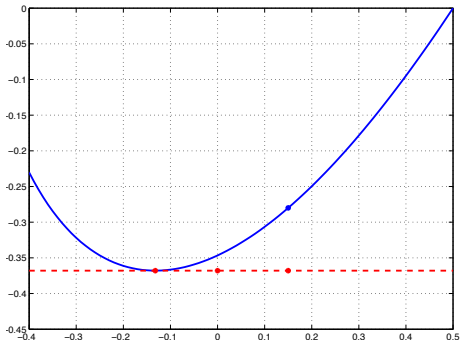
Convex function: Maximum of its affine lower bounds

Super-Gaussian function: Maximum of its Gaussian lower bounds

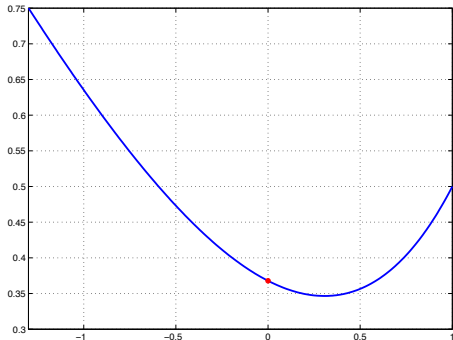
Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.



$$f(x) = \max_{\pi} \pi x - f^*(\pi)$$

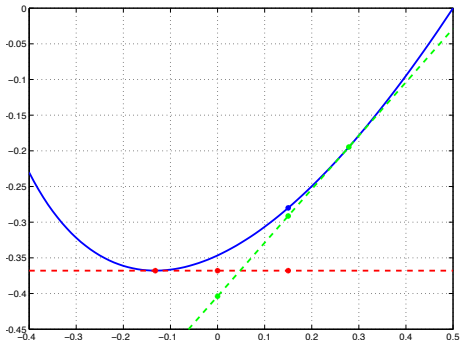


$$f^*(\pi) = \max_x \pi x - f(x)$$

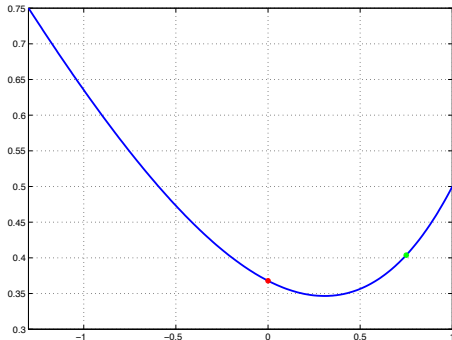
Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.



$$f(x) = \max_{\pi} x\pi - f^*(\pi)$$

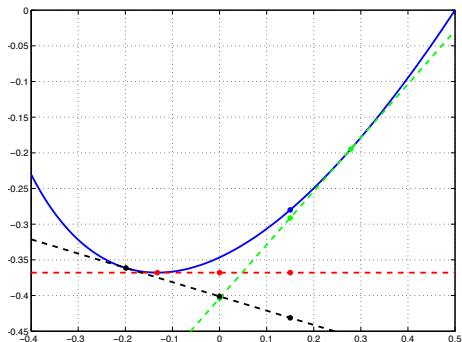


$$f^*(\pi) = \max_x x\pi - f(x)$$

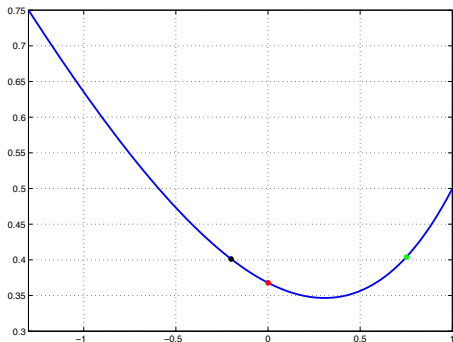
Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.



$$f(x) = \max_{\pi} \pi x - f^*(\pi)$$

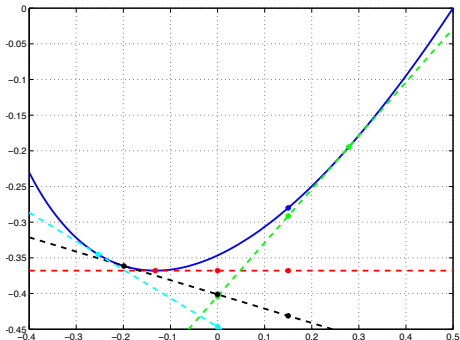


$$f^*(\pi) = \max_x \pi x - f(x)$$

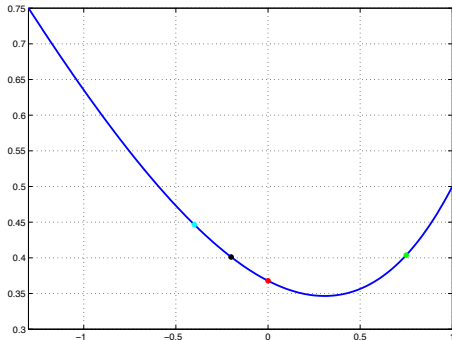
Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.



$$f(x) = \max_{\pi} \pi x - f^*(\pi)$$

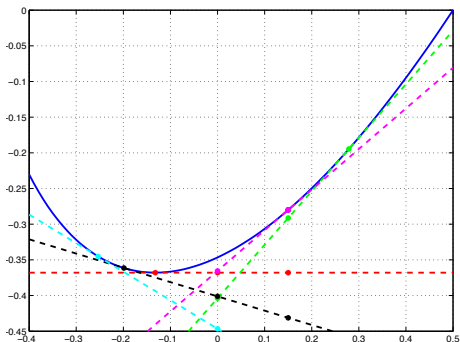


$$f^*(\pi) = \max_x \pi x - f(x)$$

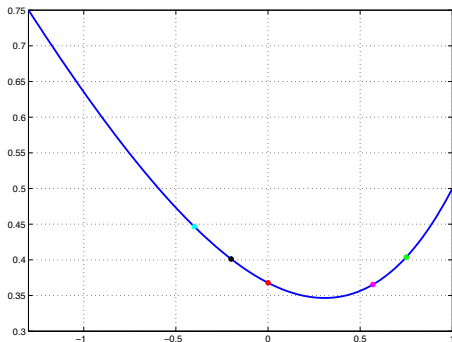
Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.



$$f(x) = \max_{\pi} x\pi - f^*(\pi)$$



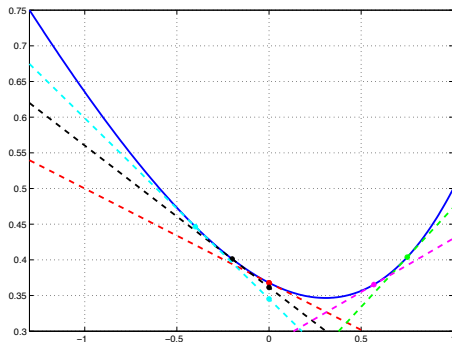
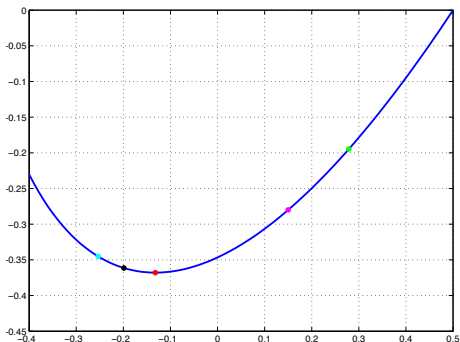
$$f^*(\pi) = \max_x \pi x - f(x)$$

Convex (Fenchel) Duality

Super-Gaussian:

$t(s)$ even, $\{x = s^2\} \mapsto \{f(x) = 2 \log t(s)\}$ convex.

F1



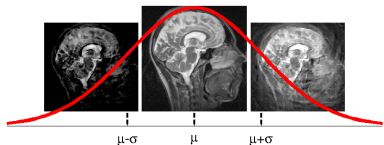
$$f(x) = \max_{\pi} x\pi - f^*(\pi)$$

$$t(s) = \max_{\gamma} e^{(-s^2/\gamma - h(\gamma))/2}$$

$$f^*(\pi) = \max_x \pi x - f(x)$$

$$h(\gamma) = \max_s -s^2/\gamma - 2 \log t(s)$$

Super-Gaussian Potentials



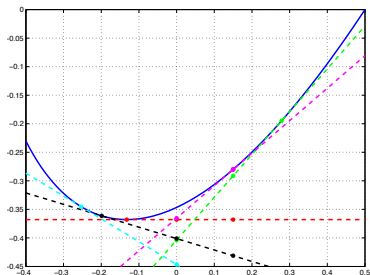
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Sparsity potentials are **super-Gaussian**

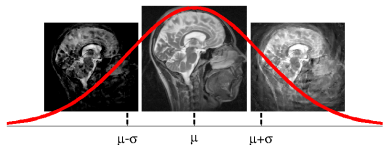
$$s_i^2 \mapsto 2 \log t_i(s_i) \text{ is convex}$$

Convex (Fenchel) duality

$$2 \log t_i(s_i) = \max_{\pi_i} s_i^2 \pi_i - f^*(\pi_i)$$



Super-Gaussian Bounding

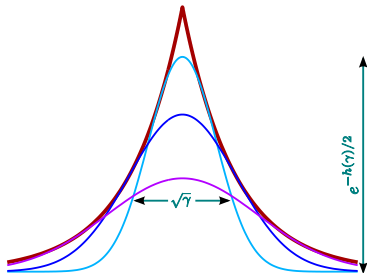


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

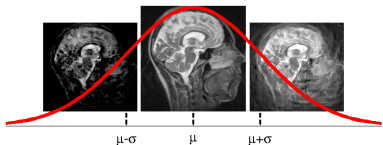
Sparsity potentials are **super-Gaussian**

$$t_i(s_i) = \max_{\gamma_i \geq 0} e^{-s_i^2 / (2\gamma_i) - h_i(\gamma_i)/2},$$

$$h(\gamma) := \sum_i h_i(\gamma_i), \quad \Gamma = \text{diag } \gamma$$



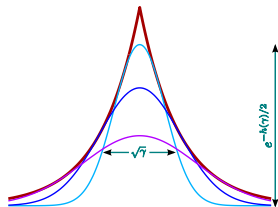
Super-Gaussian Bounding



$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

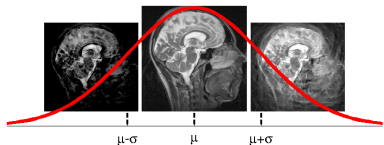
Exact representation

$$\begin{aligned} & \log Z \\ &= \log \int P(\mathbf{y}|\mathbf{u}) \max_{\gamma} e^{-(\mathbf{s}^T \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \end{aligned}$$



$$\begin{aligned} t_i(\mathbf{s}_i) &= \\ & \max_{\gamma_i \geq 0} e^{-s_i^2 / (2\gamma_i) - h_i(\gamma_i) / 2} \end{aligned}$$

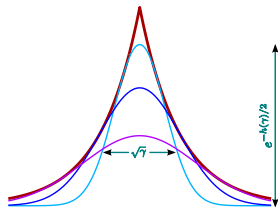
Super-Gaussian Bounding



$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

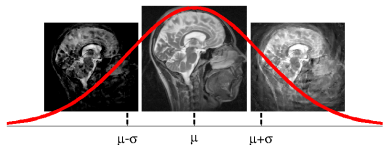
Lower bound

$$\begin{aligned} & \log Z \\ &= \log \int P(\mathbf{y}|\mathbf{u}) \max_{\gamma} e^{-(\mathbf{s}^T \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \\ &\geq \max_{\gamma} \log \int P(\mathbf{y}|\mathbf{u}) e^{-(\mathbf{s}^T \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \end{aligned}$$



$$\begin{aligned} t_i(\mathbf{s}_i) &= \\ & \max_{\gamma_i \geq 0} e^{-s_i^2 / (2\gamma_i) - h_i(\gamma_i) / 2} \end{aligned}$$

Super-Gaussian Bounding



$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Lower bound

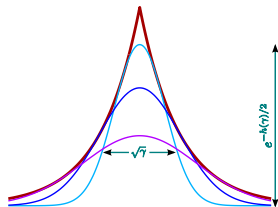
$\log Z$

$$\geq \max_{\gamma} \log \int P(\mathbf{y}|\mathbf{u}) e^{-(\mathbf{s}^T \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u}$$

$$= \max_{\gamma} \log Z_Q(\gamma) - h(\gamma)/2$$

Gaussian approximation

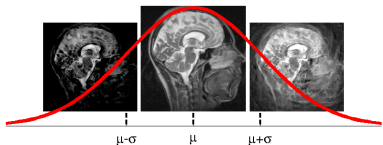
$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^T \Gamma^{-1} \mathbf{s}/2}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$



$$t_i(\mathbf{s}_i) =$$

$$\max_{\gamma_i \geq 0} e^{-\mathbf{s}_i^2 / (2\gamma_i) - h_i(\gamma_i)/2}$$

Super-Gaussian Bounding



$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

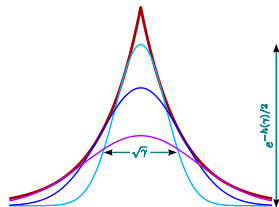
Variational problem: $Q(\mathbf{u}|\mathbf{y}) \approx P(\mathbf{u}|\mathbf{y})$

$$\min_{\gamma} \{ \phi(\gamma) = -2 \log Z_Q + h(\gamma) \}$$

Gaussian approximation

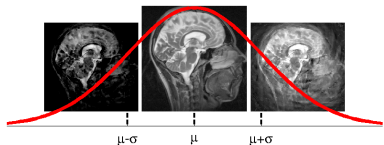
$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^T \Gamma^{-1} \mathbf{s} / 2}, \quad \mathbf{s} = \mathbf{B}\mathbf{u},$$

$$Z_Q = \int P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^T \Gamma^{-1} \mathbf{s} / 2} d\mathbf{u}$$



$$t_i(s_i) = \max_{\gamma_i \geq 0} e^{-s_i^2 / (2\gamma_i) - h_i(\gamma_i) / 2}$$

Super-Gaussian Bounding



$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

What did we do?

- Start with tight single potential bounds: $t_i(s_i) = \max_{\gamma_i \geq 0} \dots$
 \Rightarrow Auxiliary variables $\gamma \succeq \mathbf{0}$
- Plug into target function $\log Z$. Interchange $\int \dots d\mathbf{u} \leftrightarrow \max_{\gamma}$
 \Rightarrow Global **lower bound** on $\log Z$
- Lower bounds are log partition functions of **Gaussians** $Q(\mathbf{u}|\mathbf{y})$
 \Rightarrow Approximation family $\mathcal{Q} = \{Q(\mathbf{u}|\mathbf{y})\}$
- Divergence $Q(\mathbf{u}|\mathbf{y}) \leftrightarrow P(\mathbf{u}|\mathbf{y})$? Maximize lower bound!
 $\Rightarrow \phi(\gamma) = -2 \log Z_Q + h(\gamma)$

MAP Estimation and Variational Inference

MAP Estimation

$$\begin{aligned}
 & \max_{\mathbf{u}} \log P(\mathbf{u}|\mathbf{y})Z \\
 = & \max_{\mathbf{u}} \log N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) \max_{\gamma} e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} \\
 & \quad \quad \quad \parallel \\
 & \max_{\gamma} \max_{\mathbf{u}} \log N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2}
 \end{aligned}$$

Bayesian Inference

$$\begin{aligned}
 & \log Z \\
 = & \log \int N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) \max_{\gamma} e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} d\mathbf{u} \\
 & \quad \quad \quad \vee \\
 & \max_{\gamma} \log \int N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} d\mathbf{u}
 \end{aligned}$$

Coordinate Update Algorithm

- Simple algorithm: Update **single variables** γ_j

repeat

for $j \in \{1, \dots, q\}$ **do**

Update γ_j , based on marginal $Q(s_j|\mathbf{y})$

Gaussian propagation of pseudo-evidence change

end for

Refresh representation

until convergence

- Needs mean and variance of $Q(s_j|\mathbf{y})$ for each update
- Representation of $Q(\mathbf{u}|\mathbf{y})$: Backbone for Gaussian propagation.
Moderate size problems: Cholesky representation

Seeger, JMLR 2008

Variational Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1}P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i), \quad Z = \int P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i) d\mathbf{u}$$

Variational approximation

Apply variational principle to fit master function $\log Z$

- Super-Gaussian bounding
- **Expectation propagation**
- Gaussian KL minimization

Expectation Propagation

Opper, Winther, Phys. Rev. E 2001
Minka, UAI 2001

$$P(\mathbf{u}|\mathbf{y}) \approx Q(\mathbf{u}|\mathbf{y}; \gamma, \mathbf{b}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Best Gaussian approximation?

Expectation Propagation

Opper, Winther, Phys. Rev. E 2001
 Minka, UAI 2001

$$Q(\mathbf{u}|\mathbf{y}) \stackrel{\text{MM}}{\leftarrow} P(\mathbf{u}|\mathbf{y}) \Leftrightarrow Q(\mathbf{u}|\mathbf{y}) = N(\mathbb{E}[\mathbf{u}|\mathbf{y}], \text{Cov}[\mathbf{u}|\mathbf{y}])$$

- Best Gaussian approximation? **Moment matching** of $P(\mathbf{u}|\mathbf{y})$

Intractable conditions for $Q(\mathbf{u}|\mathbf{y})$

$$Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_j s_j - s_j^2 / (2\gamma_j)}$$

↑ MM

$$P(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)$$

Expectation Propagation

Opper, Winther, Phys. Rev. E 2001
 Minka, UAI 2001

$$Q(\mathbf{u}|\mathbf{y}) \stackrel{\text{MM}}{\leftarrow} P(\mathbf{u}|\mathbf{y}) \Leftrightarrow Q(\mathbf{u}|\mathbf{y}) = N(\mathbb{E}[\mathbf{u}|\mathbf{y}], \text{Cov}[\mathbf{u}|\mathbf{y}])$$

- Best Gaussian approximation? **Moment matching** of $P(\mathbf{u}|\mathbf{y})$

Intractable conditions for $Q(\mathbf{u}|\mathbf{y})$

$$Q(\mathbf{u}|\mathbf{y}) \propto e^{b_i s_i - s_i^2 / (2\gamma_i)} \times P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)}$$

↑ MM

$$P(\mathbf{u}|\mathbf{y}) \propto t_i(\mathbf{s}_i) \times P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} t_j(\mathbf{s}_j)$$

Expectation Propagation

Opper, Winther, Phys. Rev. E 2001
 Minka, UAI 2001

$$Q(\mathbf{u}|\mathbf{y}) \stackrel{\text{MM}}{\leftarrow} \hat{P}_i(\mathbf{u}) \quad \Leftrightarrow \quad Q(\mathbf{u}|\mathbf{y}) = N(\mathbb{E}_{\hat{P}_i}[\mathbf{u}], \text{Cov}_{\hat{P}_i}[\mathbf{u}])$$

- Best Gaussian approximation? **Moment matching** of $P(\mathbf{u}|\mathbf{y})$
- Tractable surrogate: Moment matching for single potentials

Self-consistency conditions for $Q(\mathbf{u}|\mathbf{y})$

$$\begin{aligned}
 Q(\mathbf{u}|\mathbf{y}) &\propto e^{b_i s_i - s_i^2 / (2\gamma_i)} \times P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)} \\
 &\quad \uparrow \text{MM} \\
 \hat{P}_i(\mathbf{u}) &\propto t_i(\mathbf{s}_i) \quad \times P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)}
 \end{aligned}$$

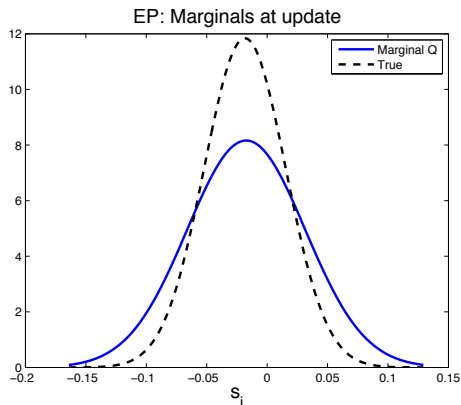
Expectation Propagation

$$Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_j s_j - s_j^2 / (2\gamma_j)}$$

Expectation Propagation

$$Q(s_i | \mathbf{y}) \propto \int P(\mathbf{y} | \mathbf{u}) \prod_j e^{b_j s_j - s_j^2 / (2\gamma_j)} d\{\mathbf{u} \setminus s_i\}$$

- 1 Marginal distribution:
 $Q(s_i | \mathbf{y})$



Expectation Propagation

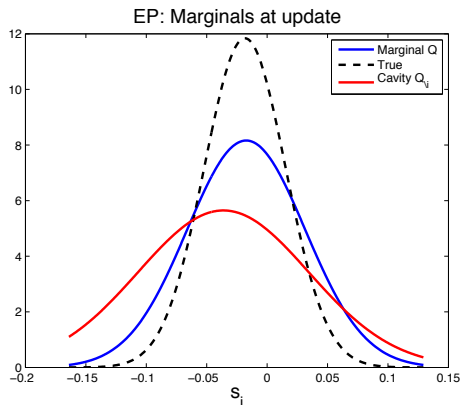
$$Q_{-i}(s_i) \propto \int P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)} d\{\mathbf{u} \setminus s_i\}$$

- 1 Marginal distribution:

$$Q(s_i|\mathbf{y})$$

- 2 Cavity distribution:

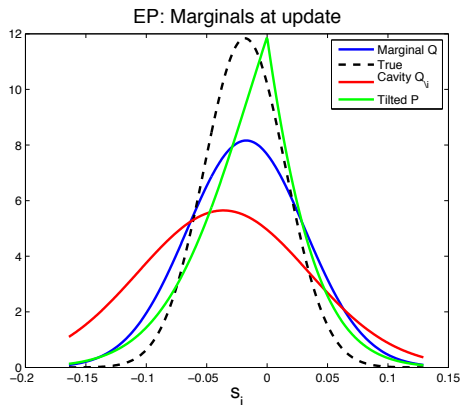
$$Q_{-i}(s_i) \propto Q(s_i|\mathbf{y}) / e^{b_i s_i - s_i^2 / (2\gamma_i)}$$



Expectation Propagation

$$\hat{P}_i(s_i) \propto t_i(s_i) \int P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)} d\{\mathbf{u} \setminus s_i\}$$

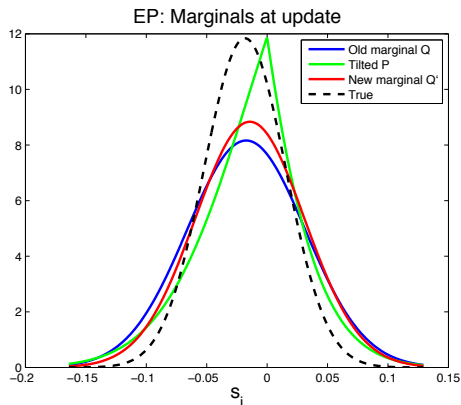
- 1 Marginal distribution:
 $Q(s_i|\mathbf{y})$
- 2 Cavity distribution:
 $Q_{-i}(s_i) \propto Q(s_i|\mathbf{y}) / e^{b_i s_i - s_i^2 / (2\gamma_i)}$
- 3 Tilted distribution:
 $\hat{P}_i(s_i) \propto Q_{-i}(s_i) t_i(s_i)$



Expectation Propagation

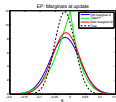
$$Q(s_i|\mathbf{y})' \propto e^{b_i' s_i - s_i^2 / (2\gamma_i')} \int P(\mathbf{y}|\mathbf{u}) \prod_{j \neq i} e^{b_j s_j - s_j^2 / (2\gamma_j)} d\{\mathbf{u} \setminus s_i\}$$

- 1 Marginal distribution:
 $Q(s_i|\mathbf{y})$
- 2 Cavity distribution:
 $Q_{-i}(s_i) \propto Q(s_i|\mathbf{y}) / e^{b_i s_i - s_i^2 / (2\gamma_i)}$
- 3 Tilted distribution:
 $\hat{P}_i(s_i) \propto Q_{-i}(s_i) t_i(s_i)$
- 4 Moment matching update:
 $Q(s_i|\mathbf{y})' = N(E_{\hat{P}_i}[s_i], \text{Var}_{\hat{P}_i}[s_i])$



Expectation Propagation

$$Q(s_i|\mathbf{y})' \xleftarrow{\text{MM}} Z_i^{-1}(Q(s_i|\mathbf{y})/e^{b_i s_i - s_i^2/(2\gamma_i)})t_i(s_i)$$



- Variational problem

$$\phi_{\text{EP}}(\gamma) = -2 \log Z_Q + \sum_i h_i^{\text{EP}}(\gamma_i, Q(s_i|\mathbf{y}))$$

$$h_i^{\text{EP}}(\gamma_i, Q(s_i|\mathbf{y})) = -2 \left(\log \mathbb{E}_{Q_{-i}}[t_i(s_i)] - \log \mathbb{E}_{Q_{-i}}[e^{b_i s_i - s_i^2/(2\gamma_i)}] \right)$$

- Arbitrary potentials
- Empirically very accurate
- Algorithmically difficult

Nickisch *et al.*, JMLR 2008

- Saddlepoint, not optimum of $\phi_{\text{EP}}(\gamma)$
- EP coordinate update algorithm lacks convergence proof
- Difficult to scale up

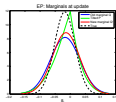
Seeger, Nickisch, AISTATS 2011

- Robust for log-concave potentials

Seeger, JMLR 2008

Expectation Propagation

$$Q(s_i|\mathbf{y})' \xleftarrow{\text{MM}} Z_i^{-1}(Q(s_i|\mathbf{y})/e^{b_i s_i - s_i^2/(2\gamma_i)})t_i(s_i)$$



Expectation propagation much more general

- Related to cavity methods (adaptive TAP) Opper, Winther, JMLR 2005
- Discrete graphical models (generalizes loopy belief propagation).
Tree expectation propagation Minka, Qi, NIPS 2004
- Dynamical systems: Natural generalization of moment matching
(assumed density) filtering Zoeter, Heskes, UAI 2002
- Inference in hybrid models (discrete and continuous)

research.microsoft.com/en-us/um/people/minka/papers/ep/roadmap.html

Variational Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1}P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i), \quad Z = \int P(\mathbf{y}|\mathbf{u}) \prod_i t_i(\mathbf{s}_i) d\mathbf{u}$$

Variational approximation

Apply variational principle to fit master function $\log Z$

- Super-Gaussian bounding
- Expectation propagation
- **Gaussian KL minimization**

Gaussian KL Minimization

Seeger, Dipl. 1999

Oppel, Archambeau, N. Comp. 2009

Variational inference

$$\log Z = \max_{Q(\mathbf{u}|\mathbf{y})} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right]$$

Gaussian KL Minimization

Seeger, Dipl. 1999

Opper, Archambeau, N. Comp. 2009

Why not use Gaussians directly?

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Equivalent to

$$\min_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} D[Q(\mathbf{u}|\mathbf{y}) \| P(\mathbf{u}|\mathbf{y})]$$

Gaussian KL Minimization

Seeger, Dipl. 1999

Opper, Archambeau, N. Comp. 2009

Why not use Gaussians directly?

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Working out the variational problem:

$$-2 \log Z \leq \min_{\gamma, \mathbf{b}} 2 \mathbb{E}_Q \left[\log \frac{Q(\mathbf{u}|\mathbf{y})}{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)} \right]$$

Gaussian KL Minimization

Seeger, Dipl. 1999

Oppel, Archambeau, N. Comp. 2009

Why not use Gaussians directly?

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Working out the variational problem:

$$-2 \log Z \leq \min_{\gamma, \mathbf{b}} 2 \mathbb{E}_Q \left[\log Z_Q^{-1} \prod_j \frac{e^{b_j s_j - s_j^2 / (2\gamma_j)}}{t_j(s_j)} \right]$$

Gaussian KL Minimization

Seeger, Dipl. 1999

Opper, Archambeau, N. Comp. 2009

Why not use Gaussians directly?

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Working out the variational problem:

$$-2 \log Z \leq \min_{\gamma, \mathbf{b}} -2 \log Z_Q + \sum_j 2 \mathbb{E}_Q [-\log t_j(s_j) - s_j^2 / (2\gamma_j) + \mathbf{b}_j s_j]$$

Gaussian KL Minimization

Seeger, Dipl. 1999

Opper, Archambeau, N. Comp. 2009

Why not use Gaussians directly?

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Working out the variational problem:

$$-2 \log Z \leq \min_{\gamma, \mathbf{b}} -2 \log Z_Q + \sum_j h_j^{\text{KL}}(\gamma_j, \mathbf{b}_j; Q(s_j|\mathbf{y}))$$

Gaussian KL Minimization

$$\min_{\gamma, \mathbf{b}} -2 \log Z_Q + \sum_j h_j^{\text{KL}}(\gamma_j, \mathbf{b}_j; Q(s_j | \mathbf{y}))$$

Comparison to super-Gaussian bounding:

- More general (t_j need not be super-Gaussian; \mathbf{b} parameters)
- More difficult to solve
 - h_j^{KL} depends on $Q(s_j | \mathbf{y})$, so on all of γ, \mathbf{b}
 - Non-convex in general
 - No large scale algorithm so far
- Tighter bound on log partition function $\log Z$

Exercise

Gaussian KL Minimization

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y}) \in \mathcal{Q}_{\text{tract}}} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right],$$

$$\mathcal{Q}_{\text{tract}} = \left\{ Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) e^{\mathbf{b}^T \mathbf{s} - \frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} \right\}$$

- Why this form? Why not **any** Gaussian $Q(\mathbf{u}|\mathbf{y}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$?
- Any Gaussian maximizer lies in $\mathcal{Q}_{\text{tract}}$

Seeger, Dipl. 1999; Exercise

$$Q^*(\mathbf{u}|\mathbf{y}) \in \underset{\text{Gaussian}}{\text{argmin}} D[Q(\mathbf{u}|\mathbf{y}) \| P(\mathbf{u}|\mathbf{y})]$$

$$\Rightarrow \text{Cov}_{Q^*}[\mathbf{u}|\mathbf{y}]^{-1} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T (\text{diag } \gamma_*)^{-1} \mathbf{B}$$

Gaussian KL Minimization

$$\log Z \geq \max_{Q(\mathbf{u}|\mathbf{y})=N(\boldsymbol{\mu},\boldsymbol{\Sigma})} \mathbb{E}_Q \left[\log \frac{P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j)}{Q(\mathbf{u}|\mathbf{y})} \right]$$

- Log-concave potentials $t_j(s_j)$:

Problem jointly convex in $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}^{\frac{1}{2}}$

Challis, Barber, AISTATS 2011

- Reduced Cholesky parameterizations
- Factorization assumptions

- Coordinate update algorithm

Seeger, Dipl. 1999

- Open problem:

Scalable algorithm for Q_{tract} parameterization ($\boldsymbol{\gamma}$, \mathbf{b})

What About Large Models?

repeat

for $j \in \{1, \dots, q\}$ **do**

Update γ_j , based on marginal $Q(s_j|\mathbf{y})$

Gaussian propagation of pseudo-evidence change

end for

Refresh representation

until convergence

- Needs mean and variance of $Q(s_j|\mathbf{y})$ for each update
- Moderate size problems: Cholesky representation
- Moderate-sized high-resolution image: $n = 65536$ pixels.

Storage: 32G (single matrix)

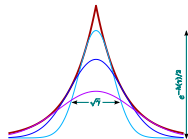
Time for Cholesky decomposition: $\approx 3\text{h}$ (if enough memory)

Out of the question

Gaussian Computations

$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = ?$$



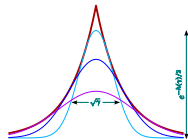
$$Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}}$$

F2

Gaussian Computations

$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = ?$$

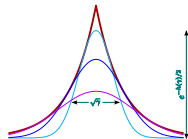


$$Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{u}^T \mathbf{B}^T \Gamma^{-1} \mathbf{B}\mathbf{u}}$$

Gaussian Computations

$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = ?$$

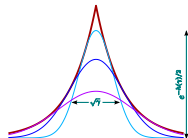


$$Q(\mathbf{u}|\mathbf{y}) \propto e^{-\frac{1}{2} (\sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \mathbf{u}^T \mathbf{B}^T \Gamma^{-1} \mathbf{B}\mathbf{u})}$$

Gaussian Computations

$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = ?$$



$$Q(\mathbf{u}|\mathbf{y}) \propto e^{-\frac{1}{2} \mathbf{u}^T (\sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \Gamma^{-1} \mathbf{B}) \mathbf{u} + \dots}$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = \mathbf{A}^{-1}, \quad \mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \Gamma^{-1} \mathbf{B}$$

If $\mathbf{v} = \mathbf{A}^{-1} (\mathbf{B}^T \delta_j)$:

$$E_Q[s_j|\mathbf{y}] = \mathbf{v}^T (\sigma^{-1} \mathbf{X}^T \mathbf{y}), \quad \text{Var}_Q[s_j|\mathbf{y}] = \mathbf{v}^T (\mathbf{B}^T \delta_j)$$

Iterative Solvers

$$\mathbb{E}_Q[\mathbf{u}|\mathbf{y}] = \mathbf{A}^{-1}(\sigma^{-2}\mathbf{X}^T\mathbf{y}), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

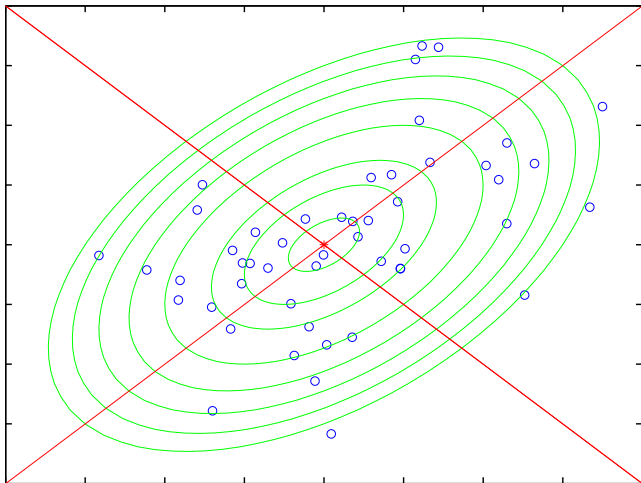
- Can multiply with \mathbf{A} rapidly:
 \mathbf{X}, \mathbf{X}^T : FFT. \mathbf{B}, \mathbf{B}^T : Simple filters
- Solve systems by iterating over matrix-vector multiplications
- Equivalent to linear least squares estimation

$$\mathbb{E}_Q[\mathbf{u}|\mathbf{y}] = \underset{\mathbf{u}}{\operatorname{argmin}} \sigma^{-2}\|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \mathbf{s}^T\mathbf{\Gamma}^{-1}\mathbf{s}$$

Minimizing Quadratic Functions

Positive definite \mathbf{A} :

$$\mathbf{x}_* = \operatorname{argmin}\{q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}\} \Leftrightarrow \mathbf{A} \mathbf{x}_* = \mathbf{b}$$



Minimizing Quadratic Functions

$$q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad \mathbf{g}(\mathbf{x}) = \nabla q(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}$$

Require: Operator \mathbf{A} . Initial \mathbf{x}_0

for $k = 1, 2, \dots$ **do**

Pick **search direction** \mathbf{d}_k , based on $\mathbf{g}_{k-1} = \mathbf{g}(\mathbf{x}_{k-1})$, $\{\mathbf{d}_l : l < k\}$

Line minimization:

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k, \quad \alpha_k = \operatorname{argmin}_{\alpha} q(\mathbf{x}_{k-1} + \alpha \mathbf{d}_k)$$

end for

Conjugate Directions

- Why, of course down **as steep as possible**

$$q(\mathbf{x}_{k-1} + d\mathbf{x}) = q(\mathbf{x}_{k-1}) + \underbrace{\mathbf{g}_{k-1}^T(d\mathbf{x})}_{\text{Smallest: } d\mathbf{x} \propto -\mathbf{g}_{k-1}} + O(\|d\mathbf{x}\|^2)$$

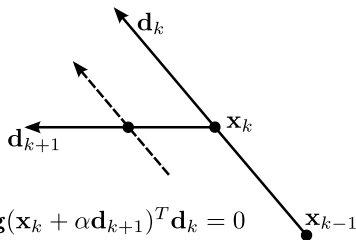
Steepest descent: $\mathbf{d}_k = -\mathbf{g}_{k-1}$

- Wrong:** For steepest descent: $\mathbf{d}_{k+1}^T \mathbf{d}_k = 0$
 \Rightarrow Improvements from previous iterations rapidly tempered with

New gradients \perp **old** directions?

\Rightarrow Retains previous efforts:

$$\mathbf{g}_k^T \mathbf{d}_k = 0 \rightarrow \mathbf{g}_{k+1}^T \mathbf{d}_k = 0 \dots$$



Conjugate Directions

$$\mathbf{d}_k^T \mathbf{A} \mathbf{d}_j = 0 \text{ for all } j < k$$

$$\mathbf{g}(\mathbf{x}_k + \alpha \mathbf{d}_{k+1})^T \mathbf{d}_k = 0$$

Towards Conjugate Gradients

Details: **Handout**

- 1 Directions conjugate: Gradient $\mathbf{g}_k \perp$ **all** previous directions:
 $\mathbf{g}_k^T \mathbf{d}_j = 0$ for all $j \leq k$
- 2 After n steps we are done: $\mathbf{g}_n = \mathbf{0}$
- 3 Construct conjugate directions by recurrence:
 $\mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$
- 4 All gradients are orthogonal: $\mathbf{g}_k^T \mathbf{g}_j = 0, j < k$
 [Bit of misnomer: **Directions** are conjugate]
- 5 What is α_k ? From line minimization:

$$\alpha_k = \frac{\|\mathbf{g}_{k-1}\|^2}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}$$

- 6 What is β_k ? The great synthesis!

$$\beta_k = \frac{\|\mathbf{g}_k\|^2}{\|\mathbf{g}_{k-1}\|^2}$$

Conjugate Gradients Algorithm

Require: Operator \mathbf{A} . Initial \mathbf{x}_0 . $\mathbf{g}_0 = \mathbf{Ax}_0 - \mathbf{b}$

for $k = 1, 2, \dots$ (no more than n) **do**

$$\rho_{k-1} = \|\mathbf{g}_{k-1}\|^2$$

if $k = 1$ **then**

$$\mathbf{d}_1 = -\mathbf{g}_0$$

else

$$\beta_{k-1} = \rho_{k-1} / \rho_{k-2}; \mathbf{d}_k = -\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1}$$

end if

$$\mathbf{q}_k = \mathbf{Ad}_k; \alpha_k = \rho_{k-1} / (\mathbf{d}_k^T \mathbf{q}_k)$$

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k; \mathbf{g}_k = \mathbf{g}_{k-1} + \alpha_k \mathbf{q}_k$$

Check for convergence (say $\|\mathbf{g}_k\| < \varepsilon \|\mathbf{b}\|$)

end for

Conjugate Gradients Algorithm

Let $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{d}_1, \dots, \mathbf{d}_k\}$. Then:

$$\mathbf{g}_k^T \mathbf{d}_j = 0, j \leq k \quad \Rightarrow \quad \mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$$

But $\mathcal{K}_k = \mathbf{x}_0 + \text{span}\{\mathbf{A}^j \mathbf{g}_0 \mid j < k\}$

\Rightarrow Optimal with k ($\mathbf{A} \cdot$) multiplications!

- $\mathcal{K}_k \subset \mathcal{K}_{k+1} \subset \dots, \mathbf{x}_* \in \mathcal{K}_n$ (Cayley/Hamilton)
- What about $k \ll n$ for huge n ? Depends on eigenspectrum of \mathbf{A} .
 $\mathbf{x}_k \approx \mathbf{x}_*$ in surprisingly many cases in practice
 \Rightarrow Krylov subspace view key to convergence analysis Exercise
- **Preconditioning:** $\mathbf{M} = \mathbf{C}\mathbf{C}^T \approx \mathbf{A}$, but **easy** to solve systems with
 - Work on $(\mathbf{C}^{-T} \mathbf{A} \mathbf{C}^{-1}) \mathbf{C} \mathbf{x} = \mathbf{C}^{-T} \mathbf{b}$
 \Rightarrow Better spectral properties \rightarrow Faster convergence
 - CG as before, with one $(\mathbf{M}^{-1} \cdot)$ per iteration
 - Art of iterative linear solvers

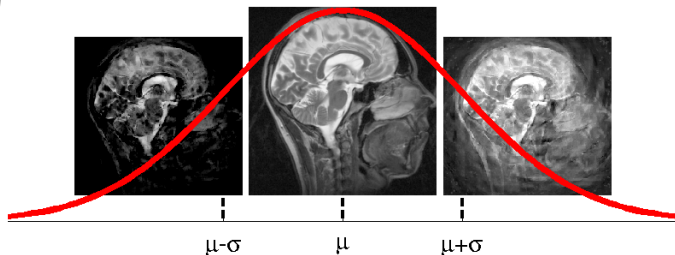
Outline

- 1 Motivation
- 2 Variational Inference Relaxations
 - Super-Gaussian Bounding
 - Expectation Propagation
 - Gaussian KL Minimization
 - Conjugate Gradients Algorithm
- 3 Scalable Variational Inference
 - Scaling up Super-Gaussian Bounding
 - Penalized Least Squares
 - Gaussian Variances
- 4 Application Example

Recap

- Sparse linear models:
Inverse problems, Bayesian calibration and sampling optimization
- Super-Gaussian bounding:
From local max-of-Gaussian representations to global bound
- Expectation propagation:
Tractable self-consistency by local moment matching
- Gaussian KL minimization:
Tighter, but more difficult than super-Gaussian bounding
- Conjugate gradients:
Large scale linear solvers by iterated matrix-vector multiplications

Need for Scalability



Bayesian inference over full images (256×256)?

$\Rightarrow \mathbf{u} \in \mathbb{C}^{65536}, \gamma \in \mathbb{R}^{196096}$

Need for Scalability

repeat

for $j \in \{1, \dots, q\}$ **do**

Update γ_j , based on marginal $Q(s_j|\mathbf{y})$

Gaussian propagation of pseudo-evidence change

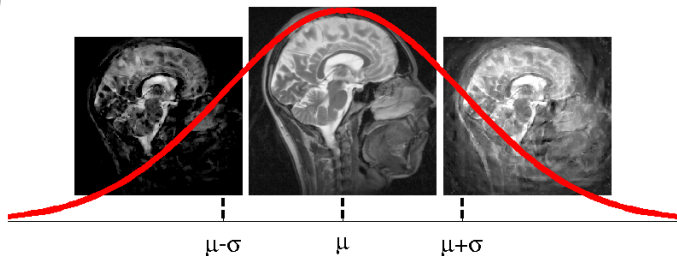
end for

Refresh representation

until convergence

- Needs mean and variance of $Q(s_j|\mathbf{y})$ for **each update**

Need for Scalability



Bayesian inference over full images (256×256)?

$\Rightarrow \mathbf{u} \in \mathbb{C}^{65536}, \gamma \in \mathbb{R}^{196096}$

- Coordinate update algorithm

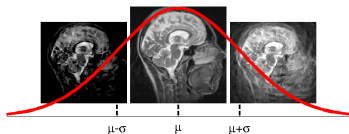
- Linear system $\mathbf{A}^{-1} \mathbf{r}$ for each update
- At least 196096 systems (visit each γ_i once)

(most previous methods)

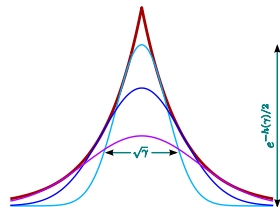
(needs conjugate gradients)

Out of the question

Properties of Super-Gaussian Bounding



$$\min_{\gamma} -2 \log \int P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u} + h(\gamma)$$

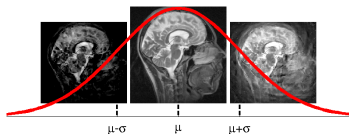


Super-Gaussian bounding stands out

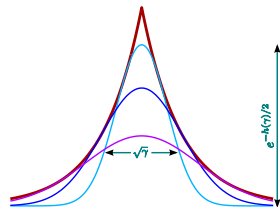
Seeger, Nickisch, SIAM IS 2011

- **Convex problem iff MAP estimation is convex**
- Can be solved at much larger scales than others

Properties of Super-Gaussian Bounding



$$\min_{\gamma} -2 \log \int P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u} + h(\gamma)$$



Super-Gaussian bounding stands out

Seeger, Nickisch, SIAM IS 2011

- Convex problem iff MAP estimation is convex
- **Can be solved at much larger scales than others**

Why is that?

MAP estimation will help solving it!

Towards Scalable Variational Inference

MAP Estimation

$$\begin{aligned}
 & \max_{\mathbf{u}} \log P(\mathbf{u}|\mathbf{y})Z \\
 = & \max_{\mathbf{u}} \log N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) \max_{\gamma} e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} \\
 & \quad \quad \quad \parallel \\
 & \max_{\gamma} \max_{\mathbf{u}} \log N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2}
 \end{aligned}$$

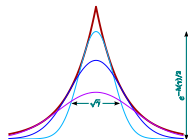
Bayesian Inference

$$\begin{aligned}
 & \log Z \\
 = & \log \int N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) \max_{\gamma} e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} d\mathbf{u} \\
 & \quad \quad \quad \vee \\
 & \max_{\gamma} \log \int N(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I}) e^{-(\mathbf{s}^T\Gamma^{-1}\mathbf{s}+h(\gamma))/2} d\mathbf{u}
 \end{aligned}$$

Towards Scalable Variational Inference

$$\min_{\gamma} -2 \log \int P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u} + h(\gamma)$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = \mathbf{A}^{-1}, \quad \mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \Gamma^{-1} \mathbf{B}$$



- Harder than MAP estimation. But why?
- Convert integration to optimization

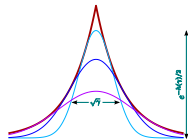
F4

$$\int P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u} \stackrel{!}{=} |2\pi \mathbf{A}^{-1}|^{1/2} \max_{\mathbf{u}_*} P(\mathbf{y}|\mathbf{u}_*) e^{-\frac{1}{2} \mathbf{s}_*^T \Gamma^{-1} \mathbf{s}_*}$$

Towards Scalable Variational Inference

$$\min_{\gamma} -2 \log \int P(\mathbf{y}|\mathbf{u}) e^{-\frac{1}{2} \mathbf{s}^T \Gamma^{-1} \mathbf{s}} d\mathbf{u} + h(\gamma)$$

$$\text{Cov}_Q[\mathbf{u}|\mathbf{y}] = \mathbf{A}^{-1}, \quad \mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \Gamma^{-1} \mathbf{B}$$



- Harder than MAP estimation. **Because of $\log |\mathbf{A}|$.**

Super-Gaussian bounding

$$\min_{\gamma, \mathbf{u}_*} \left\{ \underbrace{\phi(\mathbf{u}_*, \gamma) = \sigma^{-2} \|\mathbf{y} - \mathbf{X} \mathbf{u}_*\|^2 + \mathbf{s}_*^T \Gamma^{-1} \mathbf{s}_* + h(\gamma)}_{\text{MAP criterion } \phi_U(\mathbf{u}_*, \gamma)} + \log |\mathbf{A}| \right\}$$

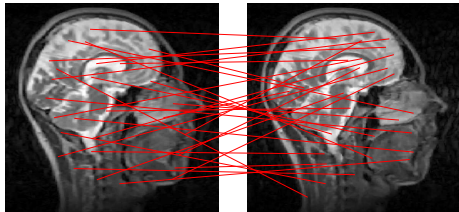
Decoupling by Convex Duality

$$-2 \log Z \leq \min_{\gamma, \mathbf{u}_*} \log |\mathbf{A}(\gamma)| + \phi_{\mathbf{U}}(\mathbf{u}_*, \gamma)$$

- Dependencies in posterior $P(\mathbf{u}|\mathbf{y})$
 \Rightarrow Difficult **coupling term** $\log |\mathbf{A}|$ in criterion ϕ

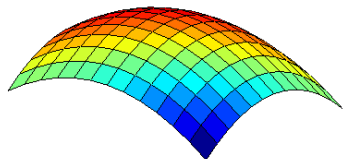
$$\mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \boldsymbol{\Gamma}^{-1} \mathbf{B}$$

- $\mathbf{A} \mapsto \log |\mathbf{A}|$ concave
- $\boldsymbol{\Gamma}^{-1} \mapsto \log |\mathbf{A}|$ concave



Decoupling by Convex Duality

$$\min_{\gamma, \mathbf{u}_*} \phi(\mathbf{u}_*, \gamma) = \min_{\gamma, \mathbf{u}_*} \underbrace{\log |\mathbf{A}(\gamma^{-1})|}_{\text{concave}} + \underbrace{\phi_U(\mathbf{u}_*, \gamma)}_{\text{convex}}$$

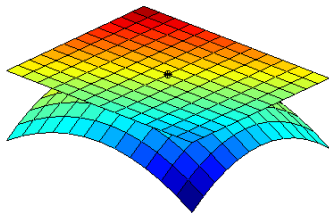


Decoupling by Convex Duality

$$\min_{\gamma, \mathbf{u}_*} \phi(\mathbf{u}_*, \gamma) = \min_{\gamma, \mathbf{u}_*} \underbrace{\log |\mathbf{A}(\gamma^{-1})|}_{\text{concave}} + \underbrace{\phi_U(\mathbf{u}_*, \gamma)}_{\text{convex}}$$

Convex (Fenchel) duality

$$\log |\mathbf{A}(\gamma^{-1})| = \min_{\mathbf{z}} \mathbf{z}^T (\gamma^{-1}) - g^*(\mathbf{z})$$

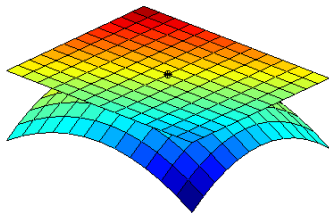


Decoupling by Convex Duality

$$\log |\mathbf{A}(\gamma^{-1})| + \phi_U(\mathbf{u}_*, \gamma) = \min_{\mathbf{z}} \underbrace{\mathbf{z}^T(\gamma^{-1}) + \phi_U(\mathbf{u}_*, \gamma) - g^*(\mathbf{z})}_{\phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) \text{ (convex, decoupled)}}$$

Convex (Fenchel) duality

$$\log |\mathbf{A}(\gamma^{-1})| = \min_{\mathbf{z}} \mathbf{z}^T(\gamma^{-1}) - g^*(\mathbf{z})$$



Scalable Double Loop Algorithm

Double loop algorithm

Seeger *et.al.*, NIPS 2009; insp. by Wipf *et.al.*, NIPS 2008

- Inner loop optimization: $\min_{\gamma} \min_{\mathbf{u}_*} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) + g^*(\mathbf{z})$ [fixed \mathbf{z}]

F5

$$\min_{\mathbf{u}_*} \min_{\gamma} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 + \mathbf{z}^T (\gamma^{-1}) + \mathbf{s}_*^T \mathbf{\Gamma}^{-1} \mathbf{s}_* + h(\gamma)$$

Scalable Double Loop Algorithm

Double loop algorithm

Seeger *et al.*, NIPS 2009; insp. by Wipf *et al.*, NIPS 2008

- Inner loop optimization: $\min_{\gamma} \min_{\mathbf{u}_*} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) + g^*(\mathbf{z})$ [fixed \mathbf{z}]
Smoothed MAP Reconstruction

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 - 2 \sum_{i=1}^q \log t_i \left(\sqrt{z_i + \mathbf{s}_{*i}^2} \right), \quad z_i > 0$$

Scalable Double Loop Algorithm

Double loop algorithm

Seeger *et al.*, NIPS 2009; insp. by Wipf *et al.*, NIPS 2008

- Inner loop optimization: $\min_{\gamma} \min_{\mathbf{u}_*} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) + g^*(\mathbf{z})$ [fixed \mathbf{z}]
Smoothed MAP Reconstruction
- Outer loop update: $\min_{\mathbf{z}} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma)$ [fixed (\mathbf{u}_*, γ)]

$$\text{Tangent : } \mathbf{z} \leftarrow \nabla_{\gamma^{-1}} \log |\mathbf{A}|, \quad \mathbf{A} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T \mathbf{\Gamma}^{-1} \mathbf{B}$$

Scalable Double Loop Algorithm

Double loop algorithm

Seeger *et.al.*, NIPS 2009; insp. by Wipf *et.al.*, NIPS 2008

- Inner loop optimization: $\min_{\gamma} \min_{\mathbf{u}_*} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) + g^*(\mathbf{z})$ [fixed \mathbf{z}]
Smoothed MAP Reconstruction
- Outer loop update: $\min_{\mathbf{z}} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma)$ [fixed (\mathbf{u}_*, γ)]
Gaussian (Co)Variances

$$\mathbf{z} \leftarrow \nabla_{\gamma^{-1}} \log |\mathbf{A}| = \text{diag}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T) = (\text{Var}_Q[\mathbf{s}_i | \mathbf{y}])$$

Reductions

Computational primitives driving large scale inference

1 Penalized least squares (\approx MAP estimation)

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 - 2 \sum_{i=1}^q \log t_i \left(\sqrt{z_i + \mathbf{s}_{*i}^2} \right)$$

- MAP special case: $z_i = 0$
- Scalable algorithms en masse (thanks to MAP “gold rush”)

2 Gaussian variances

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

- More difficult
- Methods from numerical maths, spatial statistics, solid state physics

Summary

Variational problem: $\min_{\gamma} -2 \log Z_Q + h(\gamma)$

- 1 Variational representation of **Gaussian** log partition function

$$-2 \log Z_Q \doteq \min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X} \mathbf{u}_*\|^2 + (\mathbf{s}_*^2)^T (\gamma^{-1}) + \log |\mathbf{A}|$$

Summary

Variational problem: $\min_{\gamma} -2 \log Z_Q + h(\gamma)$

- 1 Variational representation of **Gaussian** log partition function

$$-2 \log Z_Q \doteq \min_{\mathbf{z}, \mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X} \mathbf{u}_*\|^2 + (\mathbf{z} + \mathbf{s}_*^2)^T (\gamma^{-1}) - g^*(\mathbf{z})$$

- 2 Choose computationally favourable ordering of updates

$$\min_{\mathbf{z}} \left(\min_{\mathbf{u}_*, \gamma} \sigma^{-2} \|\mathbf{y} - \mathbf{X} \mathbf{u}_*\|^2 + (\mathbf{z} + \mathbf{s}_*^2)^T (\gamma^{-1}) + h(\gamma) \right) - g^*(\mathbf{z})$$

Variances \mathbf{z} expensive: Fix them as long as sensible

- 3 Convergence guarantee: Tangential bound to $\log |\mathbf{A}|$ Wipf et al., NIPS 2008

Factorization Assumptions

$$\min_{\mathbf{z}} \left(\min_{\mathbf{u}_*} \min_{\gamma} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) \right), \quad \mathbf{z} \leftarrow \text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T)$$

- Real time? As fast as MAP estimation?
Gaussian variances can be a real problem

Factorization Assumptions

$$\min_{\mathbf{z}} \left(\min_{\mathbf{u}_*} \min_{\gamma} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) \right)$$

- Real time? As fast as MAP estimation?
Gaussian variances can be a real problem
- Factorization assumptions:

$$Q(\mathbf{u}|\mathbf{y}) = \prod_{i=1}^n Q(u_i|\mathbf{y}) \quad \Rightarrow \quad Z_Q = \prod_{i=1}^n Z_{Q(u_i|\mathbf{y})}$$

\Rightarrow Update \mathbf{z} in $O(q) = O(n)$

Factorization Assumptions

$$\min_{\mathbf{u}_*} \left(\min_{\mathbf{z}} \min_{\gamma} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma) \right)$$

- Real time? As fast as MAP estimation?
Gaussian variances can be a real problem
- Factorization assumptions:

$$Q(\mathbf{u}|\mathbf{y}) = \prod_{i=1}^n Q(u_i|\mathbf{y}) \quad \Rightarrow \quad Z_Q = \prod_{i=1}^n Z_{Q(u_i|\mathbf{y})}$$

\Rightarrow Update \mathbf{z} in $O(q) = O(n)$

- One penalized least squares problem (coupled regularizer)

Factorization Assumptions

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 + \left(\min_{\mathbf{z}} \min_{\gamma} \mathcal{R}_{\mathbf{z}}(\mathbf{u}_*, \gamma) \right)$$

- Real time? As fast as MAP estimation?
Gaussian variances can be a real problem
- Factorization assumptions:

$$Q(\mathbf{u}|\mathbf{y}) = \prod_{i=1}^n Q(u_i|\mathbf{y}) \quad \Rightarrow \quad Z_Q = \prod_{i=1}^n Z_{Q(u_i|\mathbf{y})}$$

\Rightarrow Update \mathbf{z} in $O(q) = O(n)$

- One penalized least squares problem (coupled regularizer)
- Convexity properties are retained

Factorization Assumptions

$$Q(\mathbf{u}|\mathbf{y}) = \prod_{i=1}^n Q(u_i|\mathbf{y})$$

• Advantages

- Essentially as fast as MAP estimation
- Gaussian KL minimization convex
- It might just work . . .

Challis, Barber, AISTATS 2011

• Drawbacks

- True posterior tightly and strongly coupled:
Expect better results without factorizations
- Bayesian experimental design (active sampling) relies on covariances

Penalized Least Squares

Computational primitives driving large scale inference

- 1 Penalized least squares (\approx MAP estimation)

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 - 2 \sum_{i=1}^q \log t_i \left(\sqrt{z_i + \mathbf{s}_{*i}^2} \right)$$

Penalized Least Squares

Computational primitives driving large scale inference

- 1 Penalized least squares (\approx MAP estimation)

$$\min_{\mathbf{u}_*} \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 + \sum_{i=1}^q \psi_i(\mathbf{s}_{*i}), \quad \mathbf{s}_* = \mathbf{B}\mathbf{u}_*$$

Iteratively Reweighted Least Squares

$$\min_{\mathbf{u}_*} \left\{ \phi_{\mathbf{z}}(\mathbf{u}_*) = \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 + \sum_{i=1}^q \psi_i(\mathbf{s}_{*i}) \right\}, \quad \mathbf{s}_* = \mathbf{B}\mathbf{u}_*$$

- ψ_i twice differentiable: Newton-Raphson optimization
- Taylor approximation at $\mathbf{u}_* = \mathbf{u}_k$:

$$\phi_{\mathbf{z}}(\mathbf{u}_*) \approx \sigma^{-2} \|\mathbf{y} - \mathbf{X}\mathbf{u}_*\|^2 + \mathbf{s}_*^T (\text{diag } \mathbf{h}_k) \mathbf{s}_* - 2\mathbf{g}_k^T \mathbf{s}_* + C_k$$

- Newton search direction by conjugate gradients:

$$\mathbf{d}_k = \mathbf{A}_k^{-1} \left(\sigma^{-2} \mathbf{X}^T \mathbf{y} + \mathbf{B}^T \mathbf{g}_k \right), \quad \mathbf{A}_k = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \mathbf{B}^T (\text{diag } \mathbf{h}_k) \mathbf{B}$$

- Line search in $O(q)$:

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \alpha_* \mathbf{d}_k, \quad \alpha_* = \underset{\alpha > 0}{\text{argmin}} \phi_{\mathbf{z}}(\mathbf{u}_k + \alpha \mathbf{d}_k)$$

Iteratively Reweighted Least Squares

- Advantages

- Rapid (quadratic) convergence
- Reuse code for conjugate gradients algorithm

- Drawbacks

- Two nested loops: Difficult to fine-tune
- MAP estimation: ψ_i may not be twice differentiable
- For our setup: Systems

$$\mathbf{u} = \left(\mathbf{X}^T \mathbf{X} + \rho \mathbf{B}^T \mathbf{B} \right)^{-1} \mathbf{r}, \quad \rho > 0$$

can be solved analytically.

Augmented Lagrangian Solvers

$$\min_{\mathbf{u}} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 - \sum_j \log e^{-\tau|s_j|}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$

- Consider MAP estimation problem: Not differentiable

Augmented Lagrangian Solvers

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{B}\mathbf{u}\|_1, \quad \kappa = \tau\sigma^2$$

- Consider MAP estimation problem: Not differentiable

Augmented Lagrangian Solvers

$$\min_{\mathbf{u}, \mathbf{s}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{s}\|_1 \quad \text{s.t. } \mathbf{s} = \mathbf{B}\mathbf{u}$$

- Consider MAP estimation problem: Not differentiable
- Rewrite: Operator splitting.
⇒ Would be simple without constraint

Augmented Lagrangian Solvers

$$\max_{\mathbf{b}} \min_{\mathbf{u}, \mathbf{s}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{s}\|_1 + \lambda \mathbf{b}^T (\mathbf{B}\mathbf{u} - \mathbf{s})$$

- Consider MAP estimation problem: Not differentiable
- Rewrite: Operator splitting.
⇒ Would be simple without constraint
- Dualize constraint (Lagrange multipliers \mathbf{b})

Augmented Lagrangian Solvers

$$\max_{\mathbf{b}} \min_{\mathbf{u}, \mathbf{s}} \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{s}\|_1 + \lambda \mathbf{b}^T (\mathbf{B}\mathbf{u} - \mathbf{s}) + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - \mathbf{s}\|^2}_{\text{saddle function}}$$

- Consider MAP estimation problem: Not differentiable
- Rewrite: Operator splitting.
 ⇒ Would be simple without constraint
- Dualize constraint (Lagrange multipliers \mathbf{b})
- **Augmented** Lagrangian technique (additional smoothing)

F6

Alternating Direction Method of Multipliers

$$\max_{\mathbf{b}} \min_{\mathbf{u}, \mathbf{s}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{s}\|_1 + \lambda \mathbf{b}^T (\mathbf{B}\mathbf{u} - \mathbf{s}) + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - \mathbf{s}\|^2$$

Alternating Direction Method of Multipliers

Iterate:

- Linear least squares (fixed \mathbf{s} , \mathbf{b})

$$\mathbf{u} \leftarrow \operatorname{argmin} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - \mathbf{s} + \mathbf{b}\|^2$$

- Proximal map (fixed \mathbf{u} , \mathbf{b})

$$\mathbf{s} \leftarrow \operatorname{argmin} \kappa \|\mathbf{s}\|_1 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - \mathbf{s} + \mathbf{b}\|^2$$

- Lagrange multiplier update (fixed \mathbf{u} , \mathbf{s})

$$\mathbf{b} \leftarrow \mathbf{b} + \mathbf{B}\mathbf{u} - \mathbf{s}$$

The Proximal Map

Moreau

$$\text{prox}_f(\mathbf{r}) := \underset{\mathbf{s}}{\text{argmin}} f(\mathbf{s}) + \frac{1}{2} \|\mathbf{s} - \mathbf{r}\|^2$$

- Recall: $\mathbf{s} \leftarrow \text{prox}_{(\kappa/\lambda)\|\cdot\|_1}(\mathbf{B}\mathbf{u} + \mathbf{b})$

F7

The Proximal Map

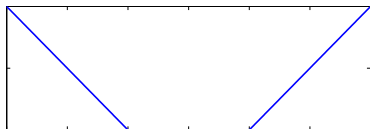
Moreau

$$\text{prox}_f(\mathbf{r}) := \underset{\mathbf{s}}{\text{argmin}} f(\mathbf{s}) + \frac{1}{2} \|\mathbf{s} - \mathbf{r}\|^2$$

- Recall: $\mathbf{s} \leftarrow \text{prox}_{(\kappa/\lambda)\|\cdot\|_1}(\mathbf{B}\mathbf{u} + \mathbf{b})$
- Simple for decoupling $f(\mathbf{s}) = \sum_i f_i(s_i)$:
 $\text{prox}_f(\mathbf{r}) = [\text{prox}_{f_i}(r_i)]$
- For Laplace (ℓ_1): **Soft thresholding**:

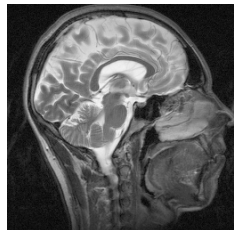
$$\text{prox}_{\alpha|\cdot|}(r) = \frac{\max\{|r| - \alpha, 0\}}{|r|} r$$

\Rightarrow Sparsity in \mathbf{s}



MRI Reconstruction

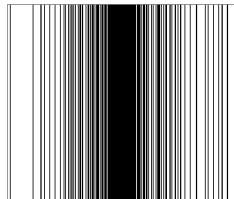
$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{B}\mathbf{u}\|_1$$



MRI Reconstruction

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \kappa \|\mathbf{B}\mathbf{u}\|_1$$

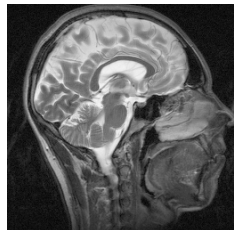
- $\mathbf{X} = \mathbf{I}_J \cdot \mathbf{F}$, \mathbf{F} DFT of size n , $J \subset \{1, \dots, n\}$
- Blocks of \mathbf{B} :
Orthonormal (wavelets), FIR filters (Δ_x, Δ_y)



MRI Reconstruction

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - (\mathbf{s} - \mathbf{b})\|^2$$

- $\mathbf{X} = \mathbf{I}_J \cdot \mathbf{F}$, \mathbf{F} DFT of size n , $J \subset \{1, \dots, n\}$
- Blocks of \mathbf{B} :
Orthonormal (wavelets), FIR filters (Δ_x, Δ_y)
- Linear least squares:



F8

$$\left(\mathbf{X}^H \mathbf{X} + \lambda \mathbf{B}^T \mathbf{B} \right) \mathbf{u} = \mathbf{r} := \mathbf{X}^H \mathbf{y} + \lambda \mathbf{B}^T (\mathbf{s} - \mathbf{b})$$

MRI Reconstruction

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - (\mathbf{s} - \mathbf{b})\|^2$$

- $\mathbf{X} = \mathbf{I}_J \cdot \mathbf{F}$, \mathbf{F} DFT of size n , $J \subset \{1, \dots, n\}$
- Blocks of \mathbf{B} :
Orthonormal (wavelets), FIR filters (Δ_x, Δ_y)
- Linear least squares:

$$\left(\mathbf{F}^H \mathbf{I}_{\cdot, J} \mathbf{I}_{J, \cdot} \mathbf{F} + \lambda \mathbf{B}^T \mathbf{B} \right) \mathbf{u} = \mathbf{r}$$



MRI Reconstruction

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - (\mathbf{s} - \mathbf{b})\|^2$$

- $\mathbf{X} = \mathbf{I}_{J,\cdot} \mathbf{F}$, \mathbf{F} DFT of size n , $J \subset \{1, \dots, n\}$
- Blocks of \mathbf{B} :
 Orthonormal (wavelets), FIR filters (Δ_x, Δ_y)
- Linear least squares:



$$\mathbf{F}^H \left(\mathbf{I}_{\cdot, J} \mathbf{I}_{J, \cdot} + \underbrace{\lambda \mathbf{F} \mathbf{B}^T \mathbf{B} \mathbf{F}^H}_{\text{diagonal}} \right) \mathbf{F} \mathbf{u} = \mathbf{r}$$

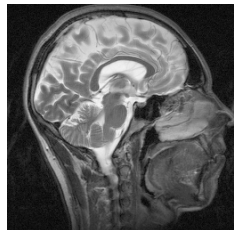
MRI Reconstruction

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{B}\mathbf{u} - (\mathbf{s} - \mathbf{b})\|^2$$

- $\mathbf{X} = \mathbf{I}_{J,\cdot} \mathbf{F}$, \mathbf{F} DFT of size n , $J \subset \{1, \dots, n\}$
- Blocks of \mathbf{B} :
Orthonormal (wavelets), FIR filters (Δ_x, Δ_y)
- Linear least squares:

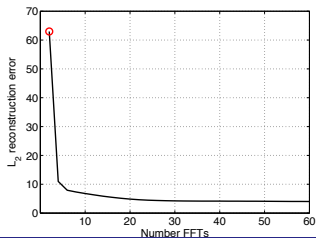
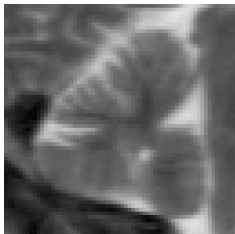
$$\underbrace{(\mathbf{I}_{\cdot,J} \mathbf{I}_{J,\cdot} + \mathbf{D})}_{\text{diagonal}} \mathbf{F}\mathbf{u} = \mathbf{F}\mathbf{r}$$

\Rightarrow Two fast Fourier transforms only!



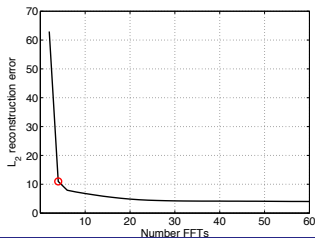
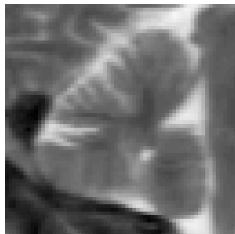
MRI Reconstruction

courtesy Mateusz Malinowski



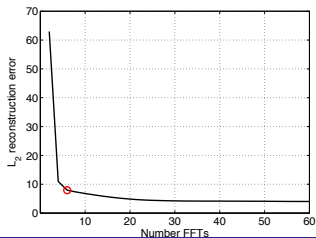
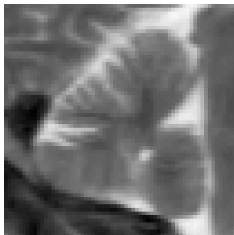
MRI Reconstruction

courtesy Mateusz Malinowski



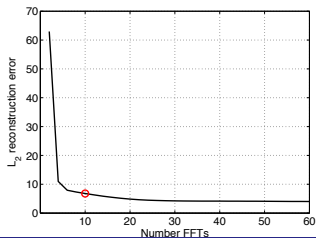
MRI Reconstruction

courtesy Mateusz Malinowski



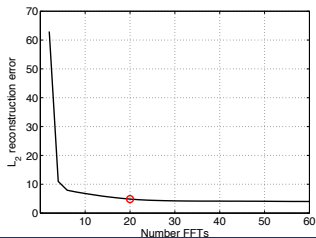
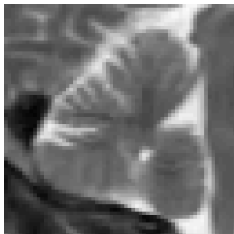
MRI Reconstruction

courtesy Mateusz Malinowski



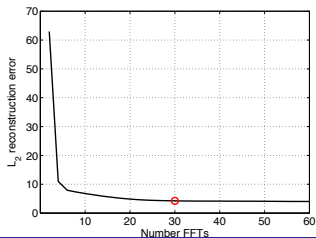
MRI Reconstruction

courtesy Mateusz Malinowski



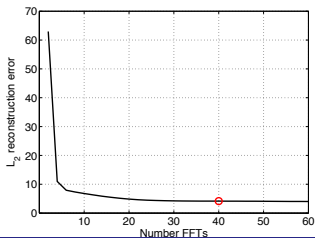
MRI Reconstruction

courtesy Mateusz Malinowski



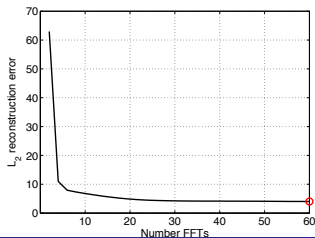
MRI Reconstruction

courtesy Mateusz Malinowski



MRI Reconstruction

courtesy Mateusz Malinowski



ADMM versus IRLS

Inner loop problems are smooth (twice differentiable)

- Use ADMM only if linear least squares step has direct solution (e.g., two FFTs)
- ADMM simpler to code and run (no CG inside)
- With complex data terms (many parts): ADMM easier to parallelize
- IRLS has much better convergence rate

Gaussian Variances

Computational primitives driving large scale inference

2 Gaussian variances

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

A Difficult Problem

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

- Gaussian variances much more difficult than Gaussian means
 - All means? One linear system
 - All variances? n linear systems!
 - Situation for Gaussian loopy belief propagation (LBP)
 - If LBP converges: Means are exact Weiss *et.al.*, JCOMP 2001
 - Variances are wrong in general: Malioutov *et.al.*, JMLR 2006
 Major part of computation (typically) not done by LBP
 - Some tractable cases
 - Tree-structured graphical model: Both means and variances in $O(n)$
 - \mathbf{A} admits sparse Cholesky factorization: van Gerven *et.al.*, Neuroimage 2010
 Variances by Takahashi equation
- Our \mathbf{A} is none of these (densely coupled)

A Difficult Problem

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

- Other fields need them as well
 - Electronic structure calculations
 - Uncertainty quantifications for PDEs
 - Gaussian MRFs for remote sensing

Low Rank Approximations

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T), \quad \mathbf{A} = \sigma^{-2}\mathbf{X}^T\mathbf{X} + \mathbf{B}^T\mathbf{\Gamma}^{-1}\mathbf{B}$$

- Pick $\mathbf{V} \in \mathbb{R}^{n \times L}$, $L \ll n$:

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T) \approx \text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{V}\mathbf{V}^T\mathbf{B}^T) = \sum_{l=1}^L (\mathbf{B}\mathbf{A}^{-1}\mathbf{v}_l) \circ (\mathbf{B}\mathbf{v}_l)$$

- Solve L linear systems instead of n
- How to choose \mathbf{V} ?

Hadamard Vectors

Bekas *et al.*, App. Num. Math. 2007

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T) \approx \text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{V}\mathbf{V}^T\mathbf{B}^T)$$

- Hadamard matrix:

$$\mathbf{H}_n \in \{-1, +1\}^{n \times n}, \quad \mathbf{H}_n^T \mathbf{H}_n = n\mathbf{I}$$

- Pick L columns:

$$\mathbf{V} = \frac{1}{\sqrt{L}}(\mathbf{H}_n)_{:,L}$$

- Deterministic estimator.

Intuition: Maximize smallest angle between any pair of rows of \mathbf{V}

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T) \approx \text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{V}\mathbf{V}^T\mathbf{B}^T)$$

- Draw L independent samples $\mathbf{q}_l \sim N(\mathbf{0}, \mathbf{A}^{-1})$:

$$\text{diag}^{-1}(\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T) = \mathbb{E}[(\mathbf{B}\mathbf{q}_1)^2] \approx \frac{1}{L} \sum_{l=1}^L (\mathbf{B}\mathbf{q}_l)^2$$

Optimal Monte Carlo estimator (no knowledge about \mathbf{A})

- One linear system per sample:

$$\mathbf{w}_l \sim N(\mathbf{0}, \mathbf{A}), \quad \mathbf{q}_l = \mathbf{A}^{-1} \mathbf{w}_l$$

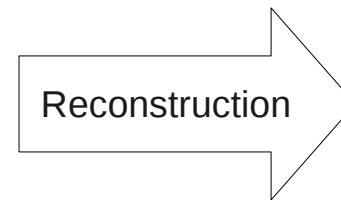
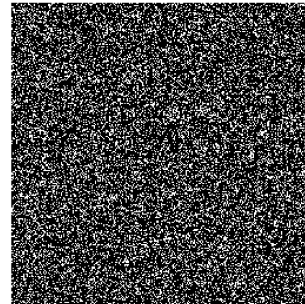
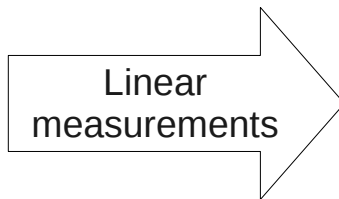
F9

Outline

- 1 Motivation
- 2 Variational Inference Relaxations
 - Super-Gaussian Bounding
 - Expectation Propagation
 - Gaussian KL Minimization
 - Conjugate Gradients Algorithm
- 3 Scalable Variational Inference
 - Scaling up Super-Gaussian Bounding
 - Penalized Least Squares
 - Gaussian Variances
- 4 Application Example

Ill-Posed Inverse Problems

- Undersampled image reconstruction



u

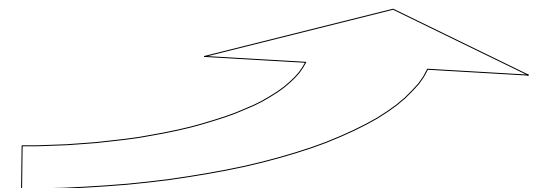
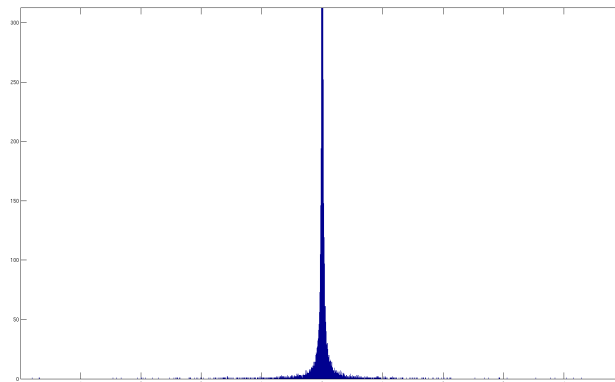
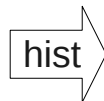
$$y = Xu + \varepsilon$$

\hat{u}

- Resolve ambiguities from prior knowledge (transform sparsity, ...)



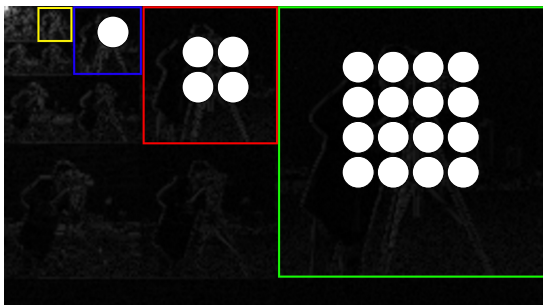
$$s = Wu$$



$P(u)$

Model: Factorial vs. Structured

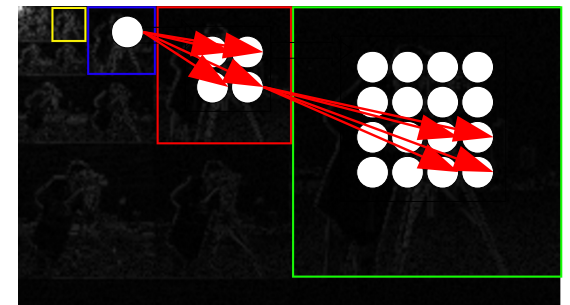
- Independent sparsity potentials
- No structure beyond component-wise sparsity
- Discrete tree-structured backbone
- **Mixtures** of sparsity potentials



$$P(\mathbf{u}) = \prod_j t_j(s_j)$$



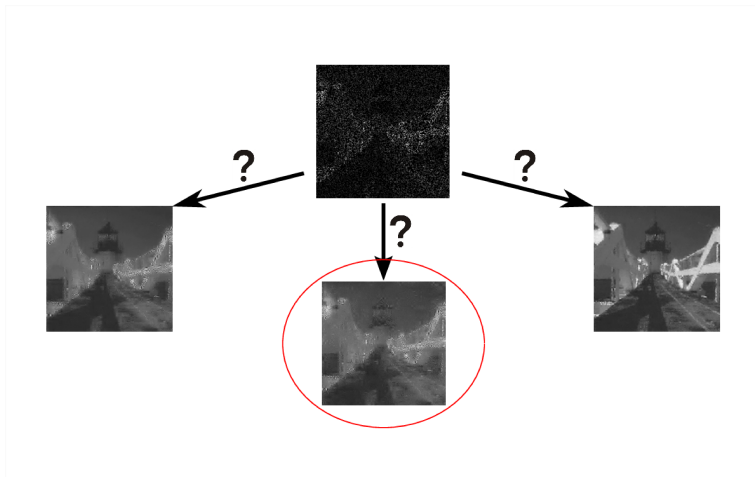
$$\mathbf{s} = \mathbf{W}\mathbf{u}$$



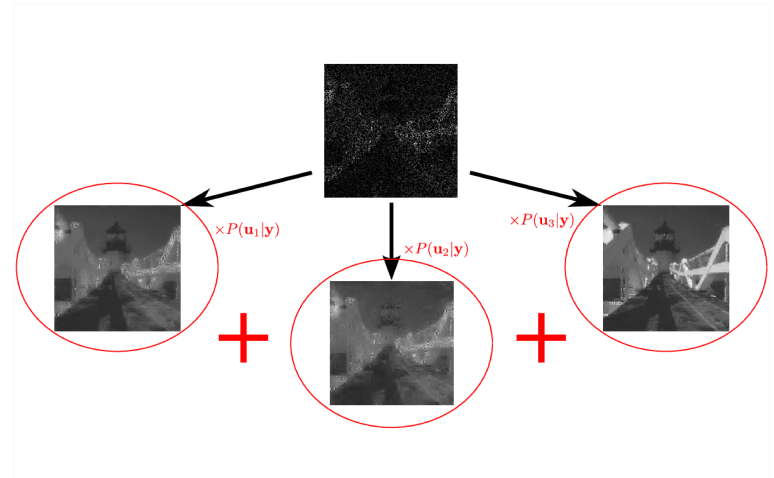
$$P(\mathbf{u}) = \sum_{\boldsymbol{\delta}} \prod_j t_j(s_j; \delta_j) P(\boldsymbol{\delta})$$

Method: MAP vs. Inference

- Estimate single maximum point
- Many fast algorithms
- Bayesian inference over posterior distribution
- Integrate, don't just maximize



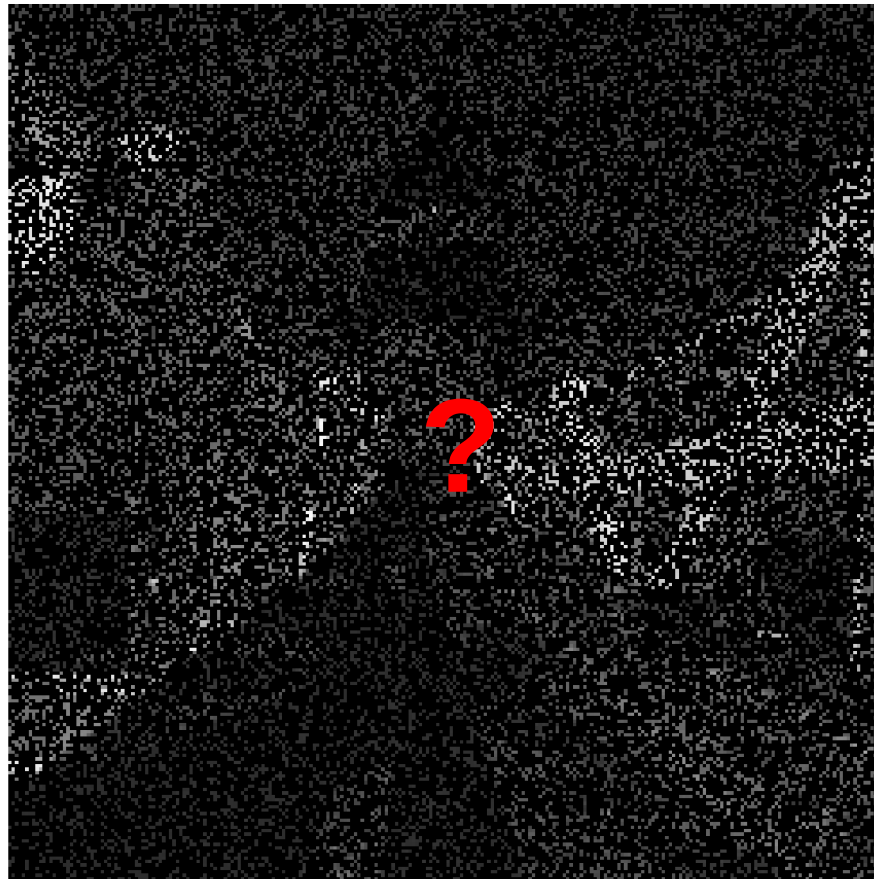
$$\hat{\mathbf{u}} = \operatorname{argmax}_{\mathbf{u}} P(\mathbf{y}|\mathbf{u})P(\mathbf{u})$$







$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u})P(\mathbf{u})}{\int P(\mathbf{y}|\mathbf{u})P(\mathbf{u})d\mathbf{u}}$$

Example: Image Inpainting

- Problem: 75% of pixels randomly removed



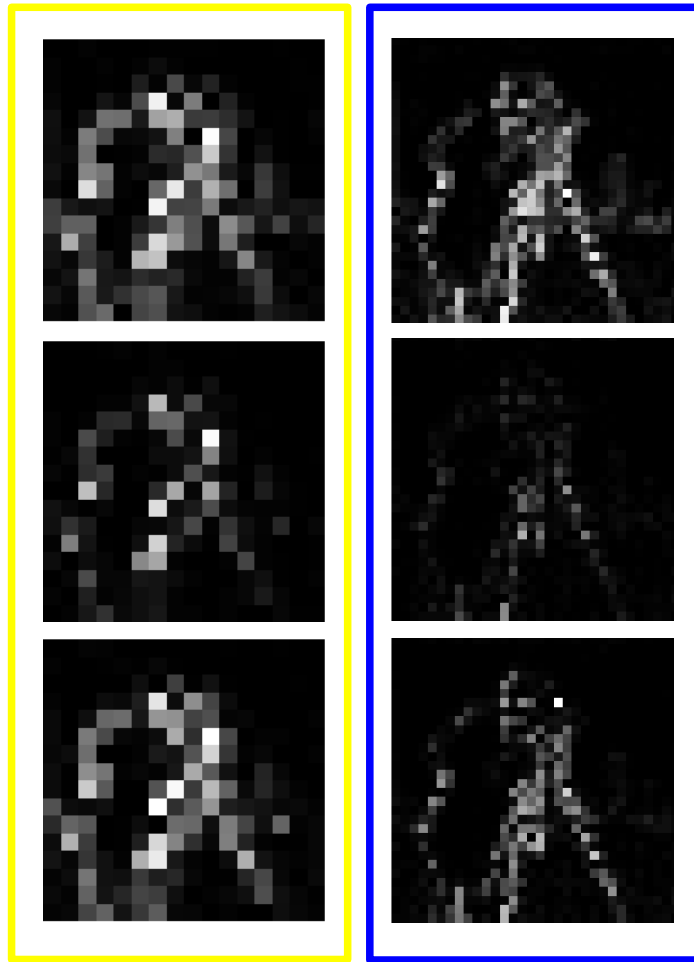
MAP	Inference	
 A grayscale image of a lighthouse at night, heavily obscured by a dense, noisy pattern of white pixels, representing the MAP (Maximum A Posteriori) estimate for the Factorial method.	 A grayscale image of a lighthouse at night, showing a smoother and more coherent reconstruction compared to the MAP image, representing the Inference result for the Factorial method.	Factorial
 A grayscale image of a lighthouse at night, heavily obscured by a dense, noisy pattern of white pixels, representing the MAP (Maximum A Posteriori) estimate for the Structured method.	 A grayscale image of a lighthouse at night, showing a reconstruction that is significantly clearer and more detailed than the MAP image, representing the Inference result for the Structured method.	Structured

Multi-Scale Wavelet Analysis

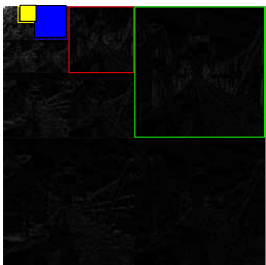
$$s = W u$$

$$\hat{s}_{MAP} = W \hat{u}_{MAP}$$

$$\hat{s}_{INF} = W \hat{u}_{INF}$$

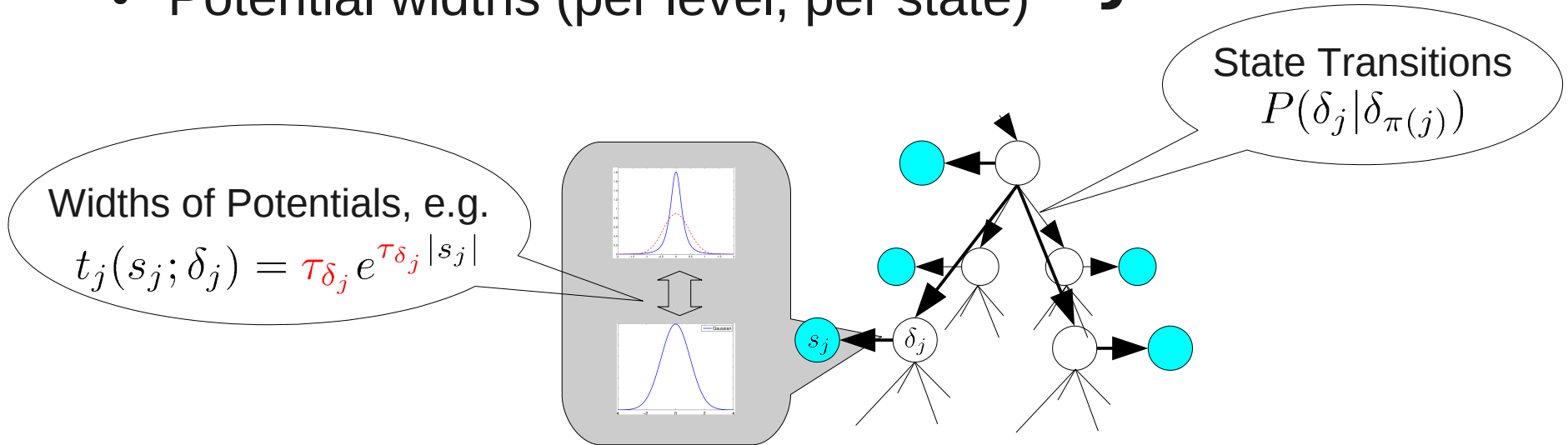


- Quad-tree: Correspondence between scales
- Energy percolates down the tree
- Better recovery where it really matters



Learning Structured Models

- Rich parametrization
 - Transition probabilities (per level)
 - Potential widths (per level, per state)
- } θ



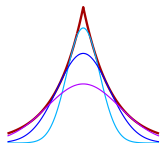
- Parameters learned automatically from raw data
 - Simple closed form updates
 - No expensive cross validation

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

1 Super-Gaussian bounding

F10

$$\log P(\mathbf{y}) \geq \max_{\gamma} \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{Q(\mathbf{u}|\delta; \gamma)}_{\text{Gaussian}} e^{-\frac{1}{2}h(\gamma; \delta)} P(\delta) d\mathbf{u}$$



Variational Bayesian Inference

Ko, Seeger, ICML 2012

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

- 1 Super-Gaussian bounding
- 2 Factorize: $Q(\mathbf{u}, \delta|\mathbf{y}) = Q(\mathbf{u}|\mathbf{y})Q(\delta|\mathbf{y})$

F11

$$\log P(\mathbf{y}) \geq \max_{\gamma, Q(\delta|\mathbf{y})} \log \int \underbrace{P(\mathbf{y}|\mathbf{u})Q(\mathbf{u}|\langle\delta\rangle_Q; \gamma)}_{\text{Gaussian}} d\mathbf{u} \\ - D[Q(\delta|\mathbf{y}) \| P(\delta)] - \frac{1}{2}h(\gamma; \langle\delta\rangle_Q)$$



Variational Bayesian Inference

Ko, Seeger, ICML 2012

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

- 1 Super-Gaussian bounding
- 2 Factorize: $Q(\mathbf{u}, \delta|\mathbf{y}) = Q(\mathbf{u}|\mathbf{y})Q(\delta|\mathbf{y})$

$$\begin{aligned} \log P(\mathbf{y}) &\geq \log Z_Q(\langle \delta \rangle_Q, \gamma) \\ &\quad - D[Q(\delta|\mathbf{y}) \| P(\delta)] - \frac{1}{2}h(\gamma; \langle \delta \rangle_Q) \end{aligned}$$



Variational Bayesian Inference

Ko, Seeger, ICML 2012

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

- 1 Super-Gaussian bounding
- 2 Factorize: $Q(\mathbf{u}, \delta|\mathbf{y}) = Q(\mathbf{u}|\mathbf{y})Q(\delta|\mathbf{y})$
- 3 Variational representation of $\log Z_Q$

$$-2 \log P(\mathbf{y}) \leq \min_{\gamma, Q(\delta|\mathbf{y})} \min_{\mathbf{z}, \mathbf{u}_*} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma, \langle \delta \rangle_Q) + 2D[Q(\delta|\mathbf{y}) \| P(\delta)]$$

Variational Bayesian Inference

Ko, Seeger, ICML 2012

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

- 1 Super-Gaussian bounding
- 2 Factorize: $Q(\mathbf{u}, \delta|\mathbf{y}) = Q(\mathbf{u}|\mathbf{y})Q(\delta|\mathbf{y})$
- 3 Variational representation of $\log Z_Q$

$$-2 \log P(\mathbf{y}) \leq \min_{\mathbf{z}} \left(\min_{Q(\delta|\mathbf{y}), \mathbf{u}_*, \gamma} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma, \langle \delta \rangle_Q) + 2D[Q(\delta|\mathbf{y}) \| P(\delta)] \right)$$

- Inner loop problem: Alternating minimization
 - Penalized least squares over \mathbf{u}_* (eliminate γ)

F12a

Variational Bayesian Inference

Ko, Seeger, ICML 2012

$$\log P(\mathbf{y}) = \log \sum_{\delta} \int P(\mathbf{y}|\mathbf{u}) \underbrace{P(\mathbf{u}|\delta)}_{\text{mixture}} \underbrace{P(\delta)}_{\text{tree}} d\mathbf{u}$$

- 1 Super-Gaussian bounding
- 2 Factorize: $Q(\mathbf{u}, \delta|\mathbf{y}) = Q(\mathbf{u}|\mathbf{y})Q(\delta|\mathbf{y})$
- 3 Variational representation of $\log Z_Q$

$$-2 \log P(\mathbf{y}) \leq \min_{\mathbf{z}} \left(\min_{Q(\delta|\mathbf{y}), \mathbf{u}_*, \gamma} \phi_{\mathbf{z}}(\mathbf{u}_*, \gamma, \langle \delta \rangle_Q) + 2D[Q(\delta|\mathbf{y}) \| P(\delta)] \right)$$

- Inner loop problem: Alternating minimization
 - Penalized least squares over \mathbf{u}_* (eliminate γ)
 - Belief propagation on tree for $Q(\delta|\mathbf{y})$ (eliminate γ)

F12b

No (Variational) Inference Without ...



**Gaussian
Variances**

$$z \leftarrow \text{diag}^{-1}(BA^{-1}B^T)$$



**Linear
Systems**

$$u_* \leftarrow \text{argmin} u_*^T A u_* - 2r^T u_*$$

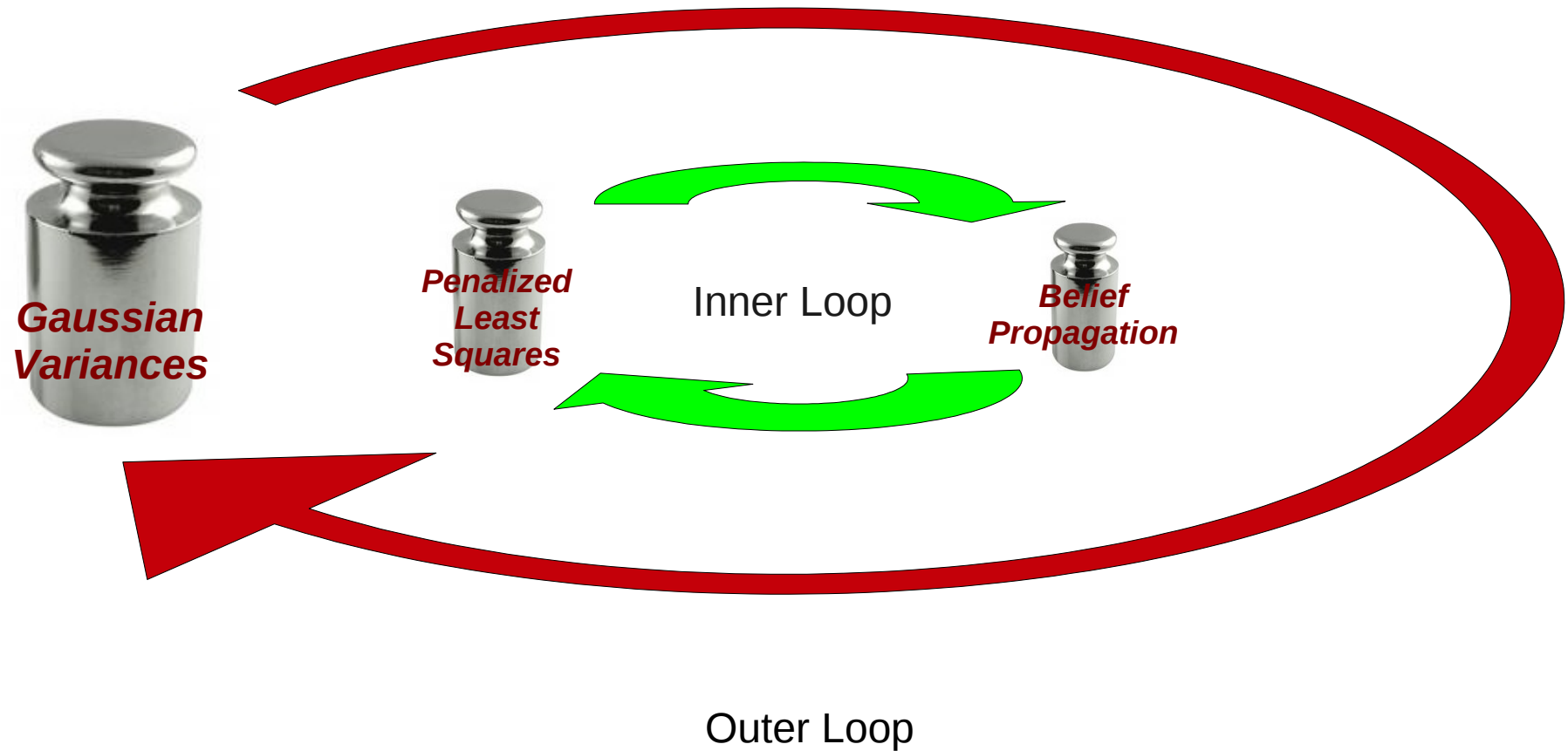






Peanuts
 $O(q)$




Why Inference Algorithms Can Be Slow



Attaining Scalability



MAP	Inference	
 A grayscale image showing a close-up of a person's face with a striped pattern. The image is heavily degraded with significant salt-and-pepper noise, making the details difficult to discern.	 A grayscale image showing a close-up of a person's face with a striped pattern. The image is heavily degraded with significant salt-and-pepper noise, making the details difficult to discern.	Factorial
 A grayscale image showing a close-up of a person's face with a striped pattern. The image is heavily degraded with significant salt-and-pepper noise, making the details difficult to discern.	 A grayscale image showing a close-up of a person's face with a striped pattern. The image is heavily degraded with significant salt-and-pepper noise, making the details difficult to discern.	Structured

MAP	Inference	
 A grayscale image of a woman wearing a hat, showing significant salt-and-pepper noise. The noise is most prominent in the darker areas of the hat and hair.	 A grayscale image of the same woman and hat, showing a smoother result with less noise than the MAP method. The features are more clearly defined.	Factorial
 A grayscale image of the same woman and hat, showing significant salt-and-pepper noise, similar to the top-left image.	 A grayscale image of the same woman and hat, showing a result that is very similar to the Inference result for the Factorial method, with reduced noise and clear features.	Structured

Large Scale Variational Inference

- Inference beyond MAP estimation
 - Robust solutions for ill-posed problems
 - Model calibration from raw data
 - Bilinear models (blind deconvolution, dictionary learning)
 - Decision making (experimental design)
- Algorithms beyond belief propagation
 - Exploit workhorses from computational mathematics
 - Reductions to convex optimization
 - Randomized techniques may be part of the solution
- What about expectation propagation?
 - More difficult to scale up
 - Talk at workshop

Seeger, Nickisch, AISTATS 2011

Try This At Home

glm-ie: Toolbox by Hannes Nickisch

`mloss.org/software/view/269/`

- Generalized sparse linear models
- MAP reconstruction and variational Bayesian inference (double loop algorithm for super-Gaussian bounding)
- Matlab 7.x, GNU Octave 3.2.x

Physics and Bayesian Machine Learning

- You love physics?
 - Predictive models for real-world phenomena
 - Intuitive analysis of large complex systems
 - Statistical evaluation by clever experiments

You'll love Bayesian machine learning!

- Most Bayesian concepts come from (statistical) physics
- Challenges of our time
 - Huge, extremely complex datasets
 - Connectivity at all scales
 - Human-level recognition and decision-making

Needs physics thinking more than ever

Bayesian Machine Learning @ EPFL

Care for new vistas? Postdoc, PhD, Internship . . .

