

Probabilistic Graphical Models

Lecture 12: Continuous-Variable Models

Volkan Cevher, Matthias Seeger
Ecole Polytechnique Fédérale de Lausanne

25/11/2011



- 1 Approximate Inference for Continuous Variables
- 2 The Sparse Linear Model
- 3 Sparsity Potentials
- 4 Super-Gaussian Bounding

The World Is Not Discrete

- Approximate Bayesian inference?
By far most activity for **discrete** variable models
 - Clean language of combinatorics on graphs. No numerical issues
 - Some relaxations are very fast (graph cuts)
 - Everything can be gridded, discretized, quantized in principle

The World Is Not Discrete

- Approximate Bayesian inference?
 - By far most activity for **discrete** variable models
 - Clean language of combinatorics on graphs. No numerical issues
 - Some relaxations are very fast (graph cuts)
 - Everything can be gridded, discretized, quantized in principle
- Viewed at useful scales, many problems are **continuous**.
 Quantization destroys structure useful for efficient computation.
 Trajectory of projectile? Planetary motion? Natural image?

Continuous

Newton mechanics

Differential equations

Integrals

Discrete

Quantum mechanics

Discretized finite differences

Ever larger sums

A Different World

Continuous inference needs more

Discrete inference

- Boils down to size of sums, hardness of graph (treewidth)
- True marginals easy to represent, “just” hard to compute

Continuous inference

- Distribution representation at least as important as graph
- Even local computations (often) not exact (\int for \sum)
- Numerical errors have to be controlled
- No ground truth even for smallish problems.
Local true marginals cannot be represented exactly

A Different World

Continuous inference: More flexibility, sometimes simpler

Discrete inference

- Most approaches today:
 - Recursive hyper-tree computations (smaller \sum)
 - Tractable combinatorial graph algorithms
- Smoothness? Non-local search directions? Global correlations?

Continuous inference

- Continuous optimization: Host of **different** approaches
- Global information from local computations (gradient, Hessian)
- Global correlation information over **all** variables (PCA)
- Continuous scientific computing well developed
 - Least squares estimation
 - Signal processing (Fourier transforms, ...)
 - PDEs

Sparse Linear Model

- Continuous variable inference:
 - Many different models for many different applications
 - Many (more or less) generic concepts
 - This lecture: Little time remaining . . .

Sparse Linear Model

- Continuous variable inference:
 - Many different models for many different applications
 - Many (more or less) generic concepts
 - This lecture: Little time remaining . . .
- Fortunately: **One** model
 - Surprisingly many applications (and growing)
 - Surprisingly many generic concepts can be demonstrated

Sparse Linear Model

- Continuous variable inference:
 - Many different models for many different applications
 - Many (more or less) generic concepts
 - This lecture: Little time remaining . . .
- Fortunately: **One** model
 - Surprisingly many applications (and growing)
 - Surprisingly many generic concepts can be demonstrated
- Workhorse for much of remaining lectures:
Sparse linear model
- Some important points we will skip
 - Multimodality of posteriors
 - Models with continuous and discrete variables

Sparsity

- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**

Sparsity

- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**

Uniform shrinkage: All weights are **smallish** \implies Gaussian F6

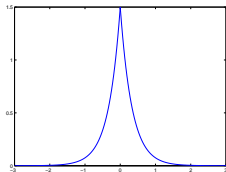
Selective shrinkage: Most weights are **tiny**, but some can be **tall**
 \implies Sparsity

Sparsity

- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**
 - Uniform shrinkage: All weights are **smallish** \implies Gaussian
 - Selective shrinkage: Most weights are **tiny**, but some can be **tall**
 \implies Sparsity
- How to implement sparsity?
 - Combinatorial search [:-(]

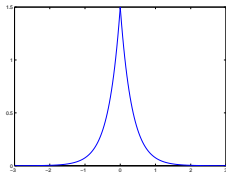
Sparsity

- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**
- Uniform shrinkage: All weights are **smallish** \implies Gaussian
- Selective shrinkage: Most weights are **tiny**, but some can be **tall**
 \implies Sparsity
- How to implement sparsity?
 - Combinatorial search [:-)]
 - **Super-Gaussian** distributions [:-)]



Sparsity

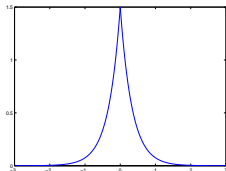
- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**
- Uniform shrinkage: All weights are **smallish** \implies Gaussian
- Selective shrinkage: Most weights are **tiny**, but some can be **tall**
 \implies Sparsity
- How to implement sparsity?
 - Combinatorial search [:-)]
 - **Super-Gaussian** distributions [:-)]
- We know sparsity for **efficiency** (SVMs, sparse matrices, ...)
 Here: Sparsity **captures signals better**
 (would be faster without)



Sparsity

- Statistics needs **regularization**: notions of **simplicity**
- Linear functions are simple if their weights are **small**
- Uniform shrinkage: All weights are **smallish** \implies Gaussian
- Selective shrinkage: Most weights are **tiny**, but some can be **tall**
 \implies Sparsity

- How to implement sparsity?
 - Combinatorial search [:-)]
 - Super-Gaussian** distributions [:-)]
- We know sparsity for **efficiency** (SVMs, sparse matrices, ...)



Here: Sparsity **captures signals better**
 (would be faster without)

- We know **sparse estimation** (Lasso, basis pursuit, ...)
 Here: **Bayesian inference** with sparsity distributions

Sparse Linear Model

Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

Sparse Linear Model

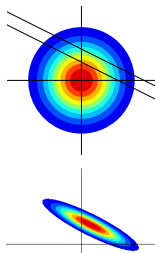
Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

Gaussian Prior $P(\mathbf{u})$

- Renders inference simple
Chosen often **only** for that
- Does not enforce $\mathbf{b}_j^T \mathbf{u} \approx 0$ strongly
- Does not allow **any** large $\mathbf{b}_j^T \mathbf{u}$



Sparse Linear Model

Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

- Whatever images are, they are not Gaussian!

Sparse Linear Model

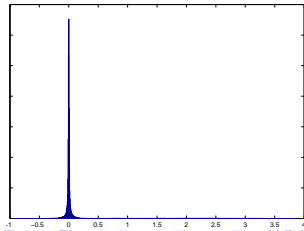
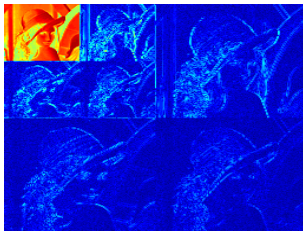
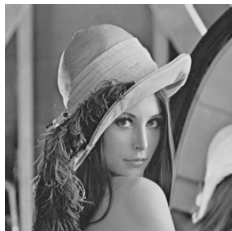
Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

- Whatever images are, they are not Gaussian!
 - Wavelet transform coefficients super-Gaussian

Simoncelli, SPIE 99



Sparse Linear Model

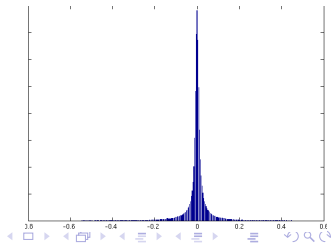
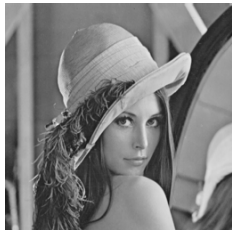
Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

- **Whatever images are, they are not Gaussian!**
 - Wavelet transform coefficients super-Gaussian
 - Spatial smoothness: Image gradient super-Gaussian

Simoncelli, SPIE 99



Sparse Linear Model

Linear Model

$$\mathbf{y} = \mathbf{X}\mathbf{u} + \varepsilon, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

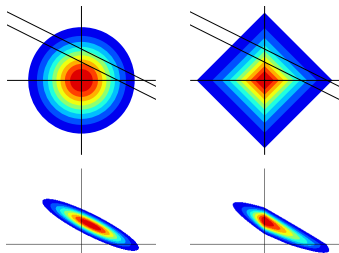
\mathbf{X} Design matrix
 $\mathbf{u} \in \mathbb{R}^n$ Latent variables
 $\mathbf{y} \in \mathbb{R}^m$ Responses

Laplace (Sparsity) Prior $P(\mathbf{u})$

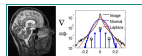
$\mathbf{s} = \mathbf{B}\mathbf{u}$ linear statistics

- Allows **few** s_j to be large
- Forces **most** $s_j \approx 0$

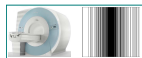
$$P(s_j) = \frac{\tau}{2} e^{-\tau|s_j|}, \quad \tau > 0$$



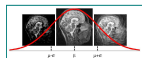
Sparse Linear Model



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

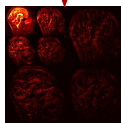


$$P(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{u}, \sigma^2\mathbf{I})$$



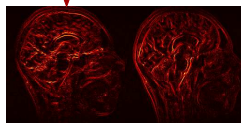
$$P(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{u})P(\mathbf{y}|\mathbf{u})$$

$$e^{-\tau_w \|\mathbf{B}_w \mathbf{u}\|_1} \times$$



wavelet

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



gradient

Gaussian Approximations

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

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?

Gaussian Approximations

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

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?
If all $t_i(s_i)$ were Gaussian ...

Gaussian Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?
If all $t_i(s_i)$ were Gaussian ...
- Approximate $P(\mathbf{u}|\mathbf{y})$ by Gaussian $Q(\mathbf{u}|\mathbf{y}; \mathbf{b}, \gamma)$

Gaussian Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?
If all $t_i(s_i)$ were Gaussian ...
- Approximate $P(\mathbf{u}|\mathbf{y})$ by Gaussian $Q(\mathbf{u}|\mathbf{y}; \mathbf{b}, \gamma)$
- Bad idea: $t_j(s_j)$ does not look like **any** Gaussian!

Gaussian Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?
If all $t_i(s_i)$ were Gaussian ...
- Approximate $P(\mathbf{u}|\mathbf{y})$ by Gaussian $Q(\mathbf{u}|\mathbf{y}; \mathbf{b}, \gamma)$
- Bad idea: $t_j(s_j)$ does not look like **any** Gaussian!
- Replace $t_j(s_j) \rightarrow e^{b_j s_j - s_j^2 / (2\gamma_j)}$,
then **adjust** \mathbf{b} , γ to fit **joint posterior**, not single $t_j(s_j)$!
 \Rightarrow Done by most algorithms: **Good idea**

Gaussian Approximations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Bayesian integration over $P(\mathbf{u}|\mathbf{y})$ intractable. Why?
If all $t_i(s_i)$ were Gaussian ...
- Approximate $P(\mathbf{u}|\mathbf{y})$ by Gaussian $Q(\mathbf{u}|\mathbf{y}; \mathbf{b}, \gamma)$
- Bad idea: $t_j(s_j)$ does not look like **any** Gaussian!
- Replace $t_j(s_j) \rightarrow e^{b_j s_j - s_j^2 / (2\gamma_j)}$,
then **adjust** \mathbf{b} , γ to fit **joint posterior**, not single $t_j(s_j)$!
 \Rightarrow Done by most algorithms: **Good idea**
- Criterion to minimize? Divergence $P(\mathbf{u}|\mathbf{y}) \leftrightarrow Q(\mathbf{u}|\mathbf{y})$?
 \Rightarrow Closer look at sparsity potentials $t_i(s_i)$

Selective Shrinkage and γ

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

Uniform shrinkage \Leftrightarrow Gaussian prior

Selective shrinkage \Leftrightarrow Sparsity prior (super-Gaussian)

$Q(\mathbf{u}|\mathbf{y})$ is **Gaussian**. Where is selective shrinkage?

Selective Shrinkage and γ

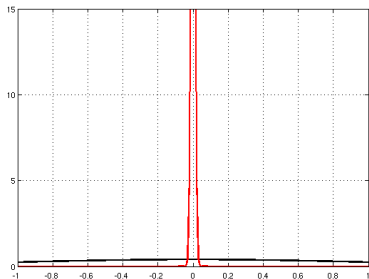
$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

Uniform shrinkage \Leftrightarrow Gaussian prior

Selective shrinkage \Leftrightarrow Sparsity prior (super-Gaussian)

$Q(\mathbf{u}|\mathbf{y})$ is **Gaussian**. Where is selective shrinkage?

- The γ_j allow for **selective shrinkage**
 - γ_j small: $|s_j|$ constrained to be small



Selective Shrinkage and γ

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_j s_j - s_j^2 / (2\gamma_j)}$$

Uniform shrinkage \Leftrightarrow Gaussian prior

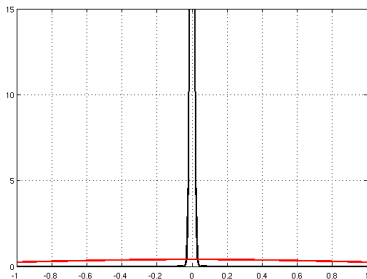
Selective shrinkage \Leftrightarrow Sparsity prior (super-Gaussian)

$Q(\mathbf{u}|\mathbf{y})$ is **Gaussian**. Where is selective shrinkage?

- The γ_j allow for **selective shrinkage**
 - γ_j small: $|s_j|$ constrained to be small
 - γ_j large: s_j rather unconstrained

Your exercise sheet:

$$\text{Var}_Q[s_j|\mathbf{y}] \leq \gamma_j$$



Selective Shrinkage and γ

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

Uniform shrinkage \Leftrightarrow Gaussian prior

Selective shrinkage \Leftrightarrow Sparsity prior (super-Gaussian)

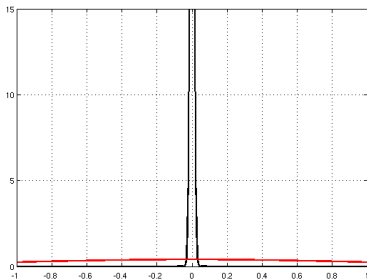
$Q(\mathbf{u}|\mathbf{y})$ is **Gaussian**. Where is selective shrinkage?

- The γ_j allow for **selective shrinkage**
 - γ_j small: $|s_j|$ constrained to be small
 - γ_j large: s_j rather unconstrained

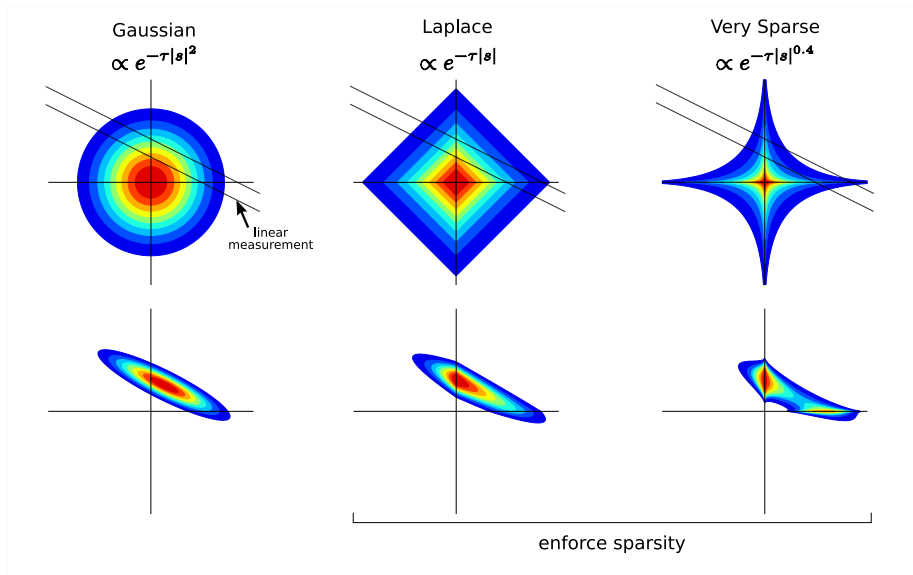
Your exercise sheet:

$$\text{Var}_Q[s_j|\mathbf{y}] \leq \gamma_j$$

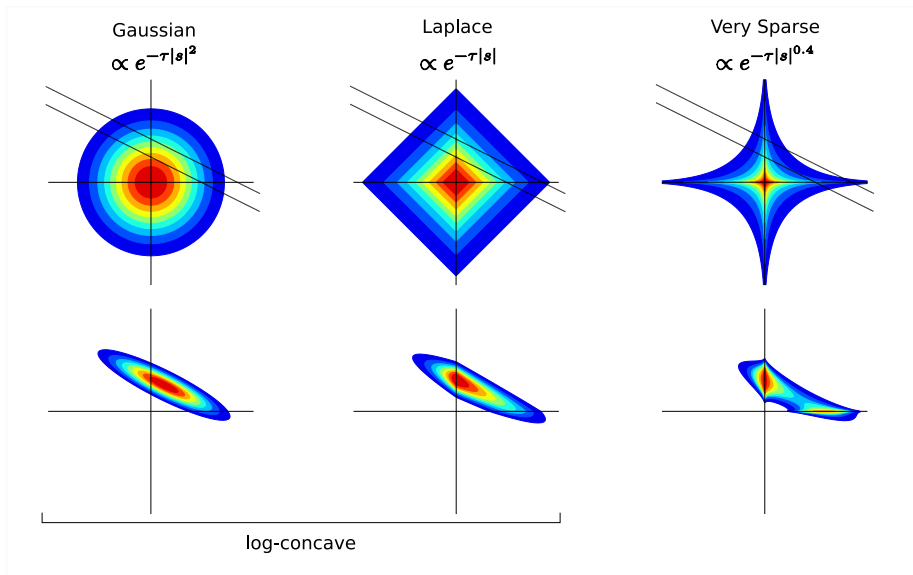
- Variational inference relaxation:
Update γ_j to implement selective shrinkage



Sparsity Priors



Sparsity Priors



Gaussian-Form Representations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Statistically, it's crucial that $t_j(s_j)$ are **not Gaussian**
- Computationally, we can only deal with **Gaussian** inference

What to do when you're stuck?

Gaussian-Form Representations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Statistically, it's crucial that $t_j(s_j)$ are **not Gaussian**
- Computationally, we can only deal with **Gaussian** inference

What to do when you're stuck? **Add new variables!**

Represent $t_j(s_j)$ as **latent** Gaussian

Gaussian-Form Representations

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_j t_j(s_j) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_j e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- Statistically, it's crucial that $t_j(s_j)$ are **not Gaussian**
- Computationally, we can only deal with **Gaussian** inference

What to do when you're stuck? **Add new variables!**

Represent $t_j(s_j)$ as **latent** Gaussian

- Gaussian scale mixtures
- Super-Gaussian potentials

$$t_j(s_j) = \int_{\gamma_i > 0} e^{-s_j^2 / (2\gamma_i)} f_i(\gamma_i) d\gamma_i$$

$$t_j(s_j) = \max_{\gamma_i > 0} e^{-s_j^2 / (2\gamma_i)} g_i(\gamma_i)$$

F8

Gaussian Scale Mixtures

- Mixture of Gaussians: Typically over means

$$P(X) = \sum_{j=1}^k \pi_k \mathcal{N}(X | \mu_k, \sigma^2)$$

$t_i(s_i)$ unimodal: Means are not the issue

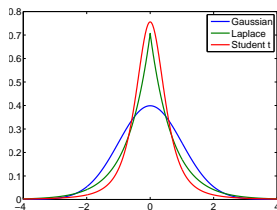
Gaussian Scale Mixtures

- Mixture of Gaussians: Typically over means

$$P(X) = \sum_{j=1}^k \pi_k \mathcal{N}(X | \mu_k, \sigma^2)$$

$t_i(s_i)$ unimodal: Means are not the issue

- What makes $t_i(s_i)$ non-Gaussian: **Shape**
 - More mass close to origin
 - More mass in tails (far from origin)
 - Less mass at moderate distances
- ⇒ Mass at different **scales**



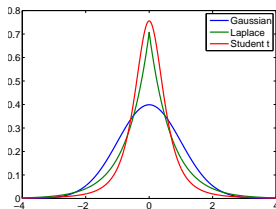
Gaussian Scale Mixtures

- Mixture of Gaussians: Typically over means

$$P(X) = \sum_{j=1}^k \pi_k \mathcal{N}(X | \mu_k, \sigma^2)$$

$t_i(s_j)$ unimodal: Means are not the issue

- What makes $t_i(s_j)$ non-Gaussian: **Shape**
 - More mass close to origin
 - More mass in tails (far from origin)
 - Less mass at moderate distances
- ⇒ Mass at different **scales**
- Why not mix over the scales?

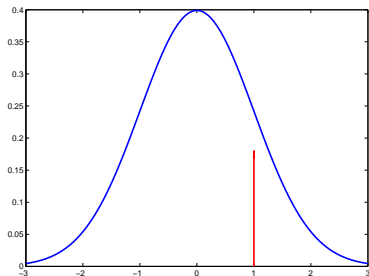


Gaussian Scale Mixtures

- $X = \rho Y, Y \sim N(0, 1), \rho \sim P(\rho)I_{\{\rho>0\}}$

Gaussian Scale Mixtures

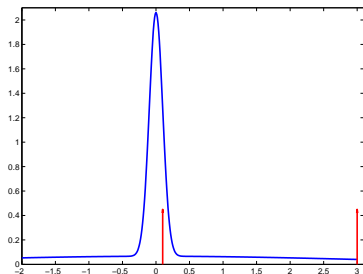
- $X = \rho Y$, $Y \sim N(0, 1)$, $\rho \sim P(\rho)I_{\{\rho>0\}}$
- Many distributions you know are **scale mixtures**
 - Gaussian [:-)].



$$P(X) = N(X|0, \rho^2)$$

Gaussian Scale Mixtures

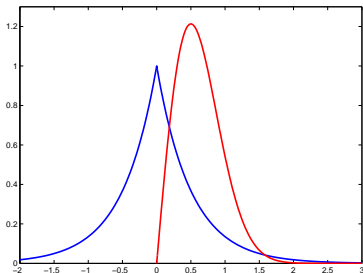
- $X = \rho Y$, $Y \sim N(0, 1)$, $\rho \sim P(\rho)I_{\{\rho>0\}}$
- Many distributions you know are **scale mixtures**
 - Gaussian [:-)]. Spike and slab



$$P(X) = \pi N(X|0, \rho_1^2) + (1 - \pi)N(X|0, \rho_2^2), \quad \rho_1 \ll \rho_2$$

Gaussian Scale Mixtures

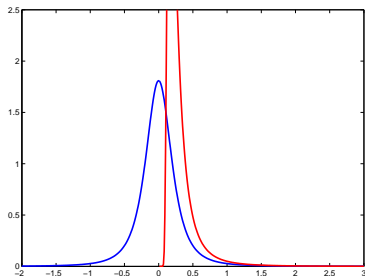
- $X = \rho Y$, $Y \sim N(0, 1)$, $\rho \sim P(\rho)I_{\{\rho>0\}}$
- Many distributions you know are **scale mixtures**
 - Gaussian [:-)]. Spike and slab
 - Exponential power ($\alpha \leq 2$)



$$P(X) \propto e^{-\tau|X|^\alpha}, \quad \alpha \in (0, 2], \tau > 0$$

Gaussian Scale Mixtures

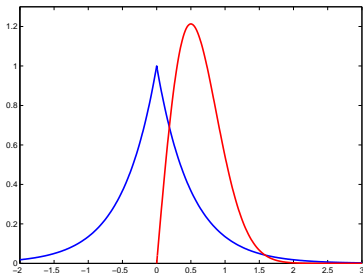
- $X = \rho Y$, $Y \sim N(0, 1)$, $\rho \sim P(\rho)I_{\{\rho>0\}}$
- Many distributions you know are **scale mixtures**
 - Gaussian [:-)]. Spike and slab
 - Exponential power ($\alpha \leq 2$)
 - Student's t



$$P(X) \propto (1 + (\tau/\nu)s^2)^{-(\nu+1)/2}, \quad \tau, \nu > 0$$

Gaussian Scale Mixtures

- $X = \rho Y$, $Y \sim N(0, 1)$, $\rho \sim P(\rho)I_{\{\rho>0\}}$
- Many distributions you know are **scale mixtures**
 - Gaussian [-:-]. Spike and slab
 - Exponential power ($\alpha \leq 2$)
 - Student's t



- Duality between $P(X)$ and $P(\rho)$
- For the Laplace:

West, Biom. 87

$$\begin{aligned} \frac{\tau}{2} e^{-\tau|s|} &= \mathbb{E}[N(s|0, \gamma)], \quad \gamma \sim (\tau^2/2) e^{-(\tau^2/2)\gamma} \\ &= \int N(s|0, \gamma) P(\gamma) d\gamma \quad [\text{scale_mix_plot}] \end{aligned}$$

Super-Gaussian Potentials

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

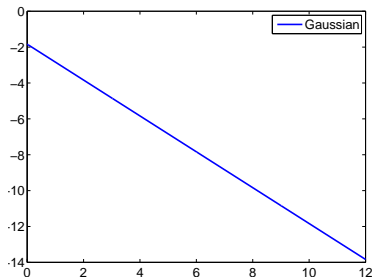
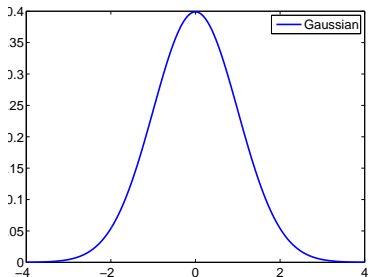
- $t_i(s_i)$ is **even**: Let's look at $s_i^2 \mapsto t_i(s_i)$
- $t_i(s_i)$ is **positive**: Let's look at $s_i^2 \mapsto 2 \log t_i(s_i)$
- What's that for a Gaussian $t_i(s_i) = N(s_i|0, \sigma_i^2)$?

F11

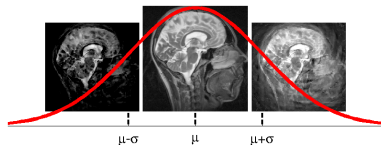
Super-Gaussian Potentials

$$P(\mathbf{u}|\mathbf{y}) = Z^{-1} P(\mathbf{y}|\mathbf{u}) \prod_i t_i(s_i) \approx Q(\mathbf{u}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{u}) \prod_i e^{b_i s_i - s_i^2 / (2\gamma_i)}$$

- $t_i(s_i)$ is **even**: Let's look at $s_i^2 \mapsto t_i(s_i)$
 $t_i(s_i)$ is **positive**: Let's look at $s_i^2 \mapsto 2 \log t_i(s_i)$
- What's that for a Gaussian $t_i(s_i) = N(s_i|0, \sigma_i^2)$?
 A **linear** (affine) function



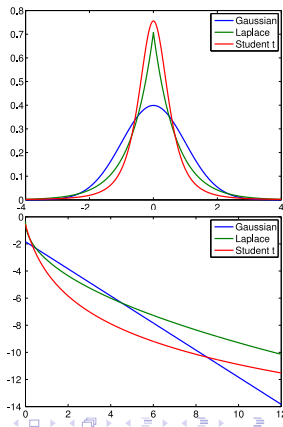
Super-Gaussian Potentials



Sparsity potentials are **super-Gaussian** F_{12}

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

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



Convex (Legendre) Duality

Super-Gaussian:

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

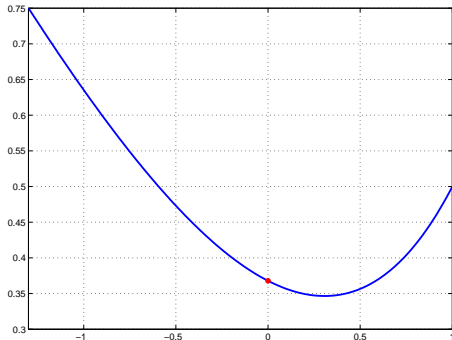
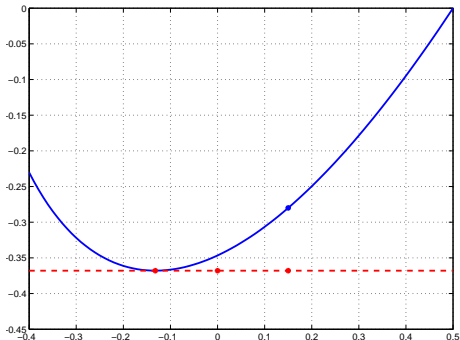
Remember Jensen's inequality?

Convex (Legendre) Duality

Super-Gaussian:

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

Remember Jensen's inequality? F13



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

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

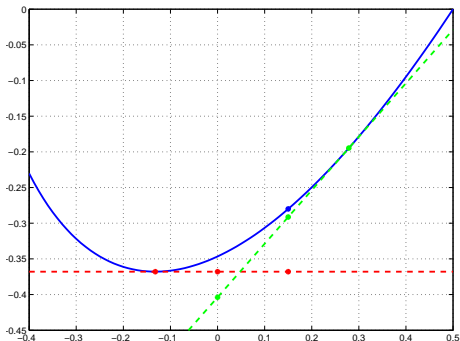
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

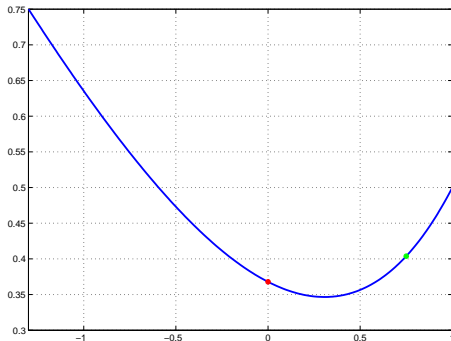
Convex (Legendre) Duality

Super-Gaussian:

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



Remember Jensen's inequality?



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

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

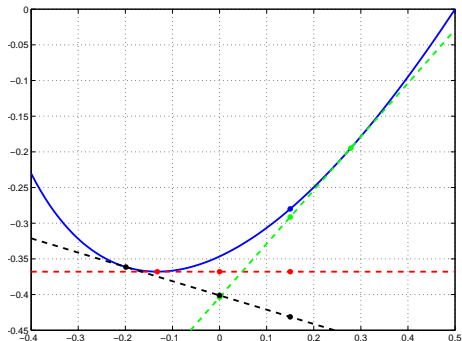
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

Convex (Legendre) Duality

Super-Gaussian:

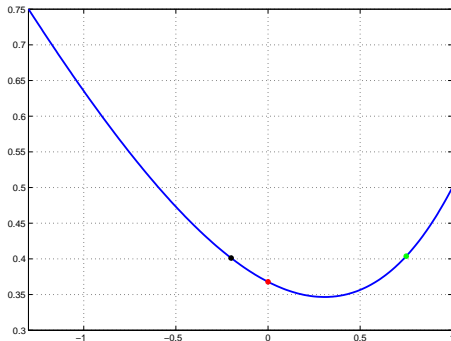
$t(s)$ even, $s^2 \mapsto \log t(s)$ convex.



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

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

Remember Jensen's inequality?

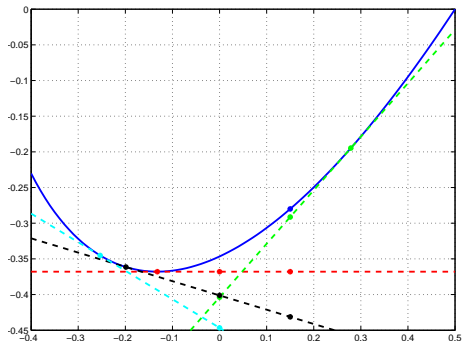


$$f^*(\pi) = \max_x \pi x - f(x)$$

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

Convex (Legendre) Duality

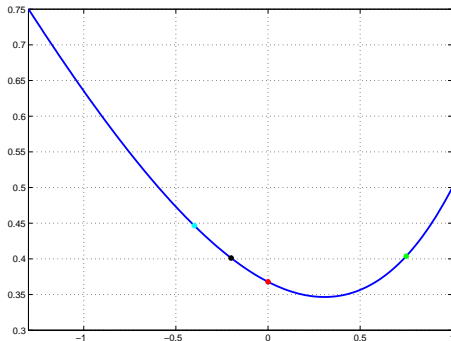
Super-Gaussian:

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

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

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

Remember Jensen's inequality?



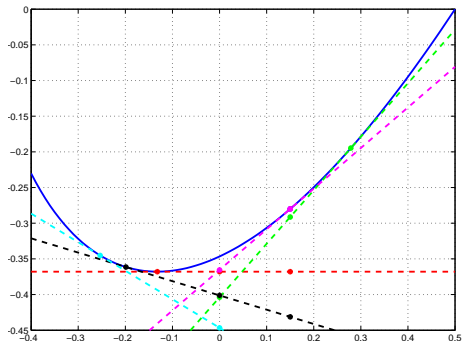
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

Convex (Legendre) Duality

Super-Gaussian:

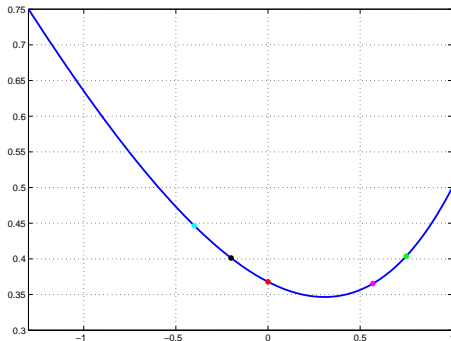
$t(s)$ even, $s^2 \mapsto \log t(s)$ convex.



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

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

Remember Jensen's inequality?



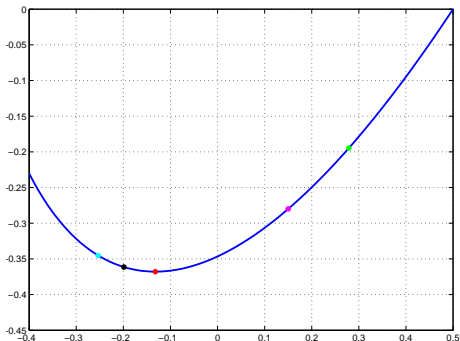
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

Convex (Legendre) Duality

Super-Gaussian:

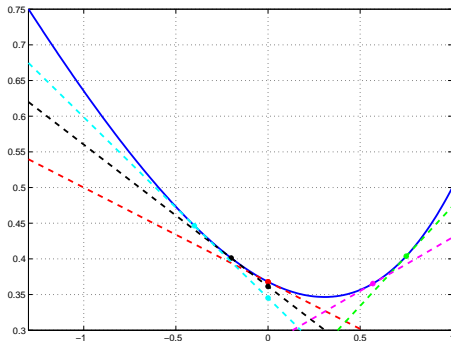
$t(s)$ even, $s^2 \mapsto \log t(s)$ convex.



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

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

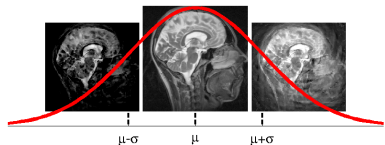
Remember Jensen's inequality?



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

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

Super-Gaussian Potentials



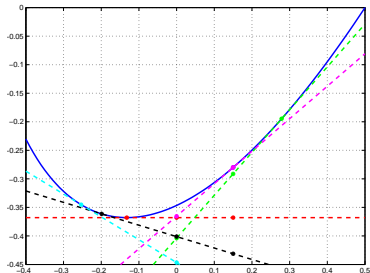
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Sparsity potentials are **super-Gaussian**

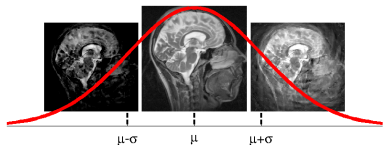
$$s_i^2 \mapsto 2 \log t_i(s_i) \quad \text{convex}$$

Convex (Legendre) duality

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



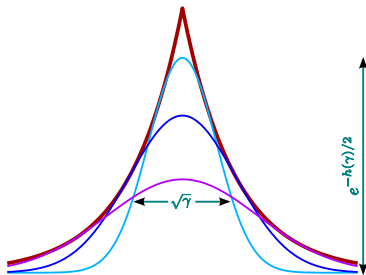
Super-Gaussian Potentials



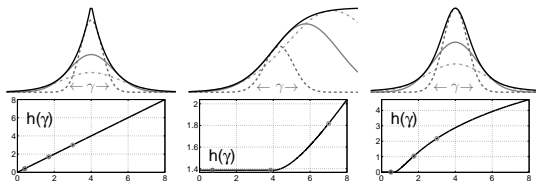
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Sparsity potentials are **super-Gaussian**

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

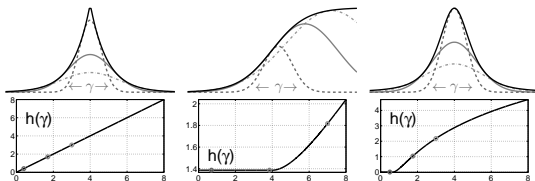


Super-Gaussian Potentials



- $t(s) = \hat{t}(s)e^{\kappa s}$ **super-Gaussian** iff
 - $\hat{t}(s)$ even function (for some κ ; $\kappa = 0$ if $t(s)$ itself even)
 - $s^2 \mapsto \log \hat{t}(s)$ convex, decreasing

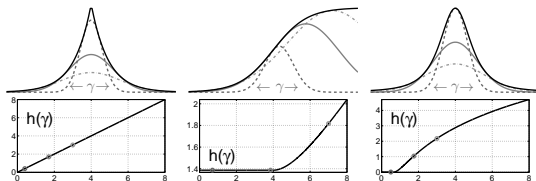
Super-Gaussian Potentials



- $t(s) = \hat{t}(s)e^{\kappa s}$ **super-Gaussian** iff
 - $\hat{t}(s)$ even function (for some κ ; $\kappa = 0$ if $t(s)$ itself even)
 - $s^2 \mapsto \log \hat{t}(s)$ convex, decreasing
- Bernoulli (logistic) $t(s) = (1 + e^{-y\tau s})^{-1}$, $y \in \{\pm 1\}$?

F15

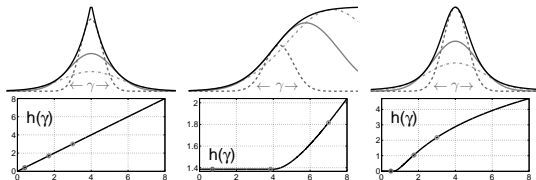
Super-Gaussian Potentials



- $t(s) = \hat{t}(s)e^{\kappa s}$ **super-Gaussian** iff
 - $\hat{t}(s)$ even function (for some κ ; $\kappa = 0$ if $t(s)$ itself even)
 - $s^2 \mapsto \log \hat{t}(s)$ convex, decreasing
- Bernoulli (logistic) $t(s) = (1 + e^{-y^\tau s})^{-1}$, $y \in \{\pm 1\}$?
 $t(s) = \hat{t}(s)e^{(y^\tau/2)s}$, $g(x) \doteq -\log \cosh((y^\tau/2)x^{1/2})$
- All scale mixtures are super-Gaussian

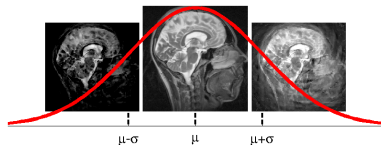
Palmer *et.al.*, NIPS 2005

Super-Gaussian Potentials



- $t(s) = \hat{t}(s)e^{\kappa s}$ **super-Gaussian** iff
 - $\hat{t}(s)$ even function (for some κ ; $\kappa = 0$ if $t(s)$ itself even)
 - $s^2 \mapsto \log \hat{t}(s)$ convex, decreasing
 - Bernoulli (logistic) $t(s) = (1 + e^{-y^\tau s})^{-1}$, $y \in \{\pm 1\}$?
 $t(s) = \hat{t}(s)e^{(y^\tau/2)s}$, $g(x) \doteq -\log \cosh((y^\tau/2)x^{1/2})$
 - All scale mixtures are super-Gaussian
- Palmer *et.al.*, NIPS 2005
- F4b
- Some closure properties: $\{t_i(s_i)\}$ super-Gaussian, $\alpha_i > 0$
 - $\prod_i t_i(s_i)^{\alpha_i}$ super-Gaussian
 - $\sum_i \alpha_i t_i(s_i)$ super-Gaussian
- F15b

Super-Gaussian Bounding

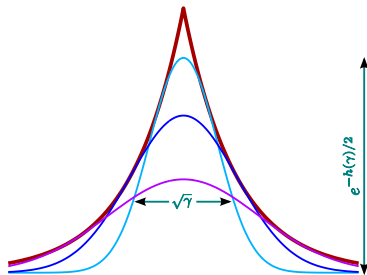


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

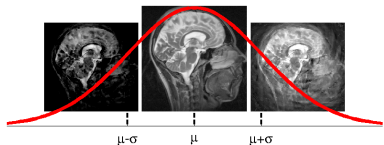
Sparsity potentials are **super-Gaussian**

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

$$h(\gamma) := \sum_i h_i(\gamma_i)$$



Super-Gaussian Bounding

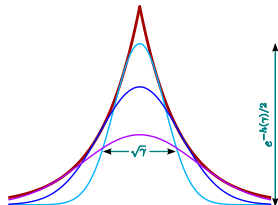


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Exact representation

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

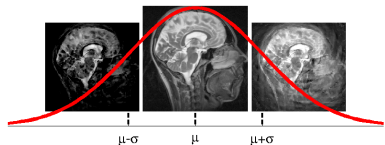
F16



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

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

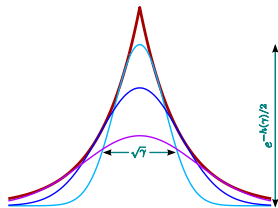
Super-Gaussian Bounding



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

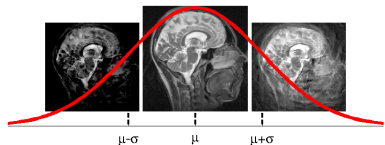
Lower bound

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



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

Super-Gaussian Bounding



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

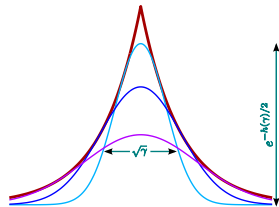
Lower bound

$\log Z$

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

Gaussian approximation

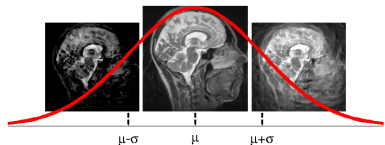
$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^T \Gamma^{-1} \mathbf{s}/2}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$



$$t_i(s_i) =$$

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

Super-Gaussian Bounding



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

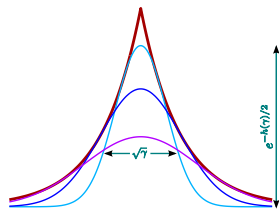
Variational problem: $Q(\mathbf{u}|\mathbf{y}) \approx P(\mathbf{u}|\mathbf{y})$

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

Gaussian approximation

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

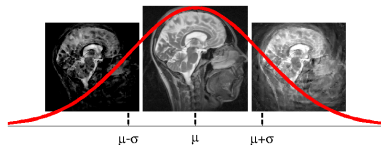
$$Z_Q = \int P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^T \Gamma^{-1} \mathbf{s} / 2} d\mathbf{u}$$



$$t_i(s_i) =$$

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

Super-Gaussian Bounding

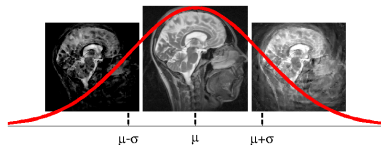


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

What did we do?

- Start with tight single potential bounds: $t_i(s_i) = \max_{\gamma_i > 0} \dots$
 \Rightarrow Auxiliary variables $\gamma \succ \mathbf{0}$

Super-Gaussian Bounding

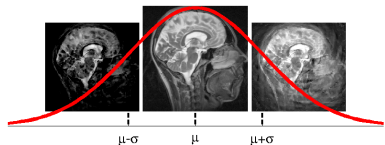


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

What did we do?

- Start with tight single potential bounds: $t_i(s_i) = \max_{\gamma_i > 0} \dots$
 \Rightarrow Auxiliary variables $\gamma \succ \mathbf{0}$
- Plug into target function $\log Z$. Interchange $\int \dots d\mathbf{u} \leftrightarrow \max_{\gamma}$
 \Rightarrow Global lower bound on $\log Z$ (not tight)

Super-Gaussian Bounding

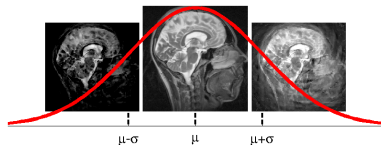


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

What did we do?

- Start with tight single potential bounds: $t_i(s_i) = \max_{\gamma_i > 0} \dots$
 \Rightarrow Auxiliary variables $\gamma \succ \mathbf{0}$
- Plug into target function $\log Z$. Interchange $\int \dots d\mathbf{u} \leftrightarrow \max_{\gamma}$
 \Rightarrow Global lower bound on $\log Z$ (not tight)
- Lower bounds **are** log partition functions of Gaussians $Q(\mathbf{u}|\mathbf{y})$
 \Rightarrow Approximation family $\mathcal{Q} = \{Q(\mathbf{u}|\mathbf{y})\}$

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: $t_i(s_i) = \max_{\gamma_i > 0} \dots$
 \Rightarrow Auxiliary variables $\gamma \succ \mathbf{0}$
- Plug into target function $\log Z$. Interchange $\int \dots d\mathbf{u} \leftrightarrow \max_{\gamma}$
 \Rightarrow Global lower bound on $\log Z$ (not tight)
- 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 Divergence $\phi(\gamma) = -2 \log Z_Q + h(\gamma)$

Coordinate Descent 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

Exercise sheet

Coordinate Descent 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

- **Representation** of $Q(\mathbf{u}|\mathbf{y})$: Backbone for Gaussian propagation
Moderate size problems: Cholesky representation

Exercise sheet

Coordinate Descent 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

- **Representation** of $Q(\mathbf{u}|\mathbf{y})$: Backbone for Gaussian propagation
Moderate size problems: Cholesky representation
- Large scale problems?
This algorithm is too slow (not scalable)

Exercise sheet

MAP Estimation and Variational Inference

MAP Estimation

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

Bayesian Inference

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

Wrap-Up

- Continuous-variable approximate inference: A different game
- Sparse linear model:
Combinatorial properties with continuous variables
- Gaussian distributions (possibly graph-structured):
Major backbone for continuous-variable inference
- Gaussian-form representations:
 - Scale mixtures
 - Super-Gaussian potentials
- Super-Gaussian bounding:
From local potential bounds to global log partition function bound