

# Probabilistic Graphical Models

## Inference for Continuous-Variable Models

Volkan Cevher, Matthias Seeger  
Ecole Polytechnique Fédérale de Lausanne

28/11/2011

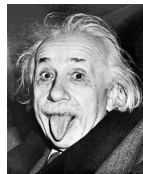


- 1 Introduction
- 2 Sparsity. Super-Gaussianity
- 3 Variational Bayesian Inference Relaxations

# Sparsity: A Fundamental Concept

... as simple as possible, but not simpler.

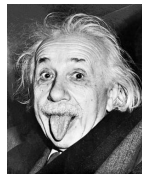
What do you mean with **simple**?



# Sparsity: A Fundamental Concept

... as simple as possible, but not simpler.

What do you mean with **simple**?

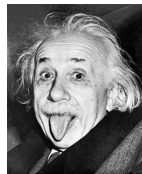


## Classical (Gaussian)

- **All** specified elements
- Use each of them **a little**

# Sparsity: A Fundamental Concept

... as simple as possible, but not simpler.



What do you mean with **simple**?

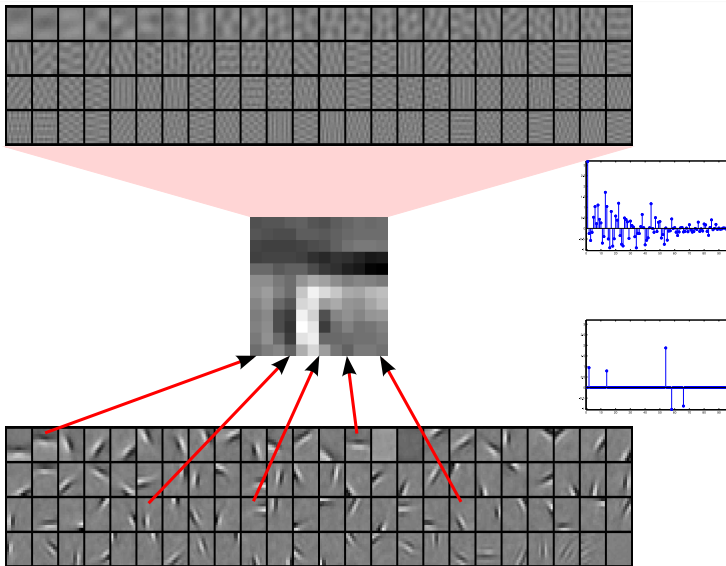
## Classical (Gaussian)

- **All** specified elements
- Use each of them **a little**

## Sparsity

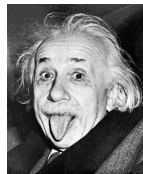
- As **few** elements as possible
- If at all, use them **big**

# Sparsity: A Fundamental Concept



# Sparsity: A Fundamental Concept

... as simple as possible, but not simpler.



What do you mean with **simple**?

## Classical (Gaussian)

- All specified elements
- Use each of them **very little**

## Sparsity

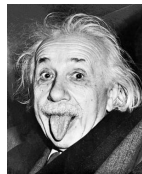
- As **few** elements as possible
- If at all, use them **big**

Classical linear framework: Shapes the way we think

- Nyquist/Shannon limit. Point spread function
- Aliasing. Ringing. Signal-to-noise ratio
- Linear measurements? Linear reconstruction!

# Sparsity: A Fundamental Concept

... as simple as possible, but not simpler.



What do you mean with **simple**?

## Classical (Gaussian)

- **All** specified elements
- Use each of them **very little**

## Sparsity

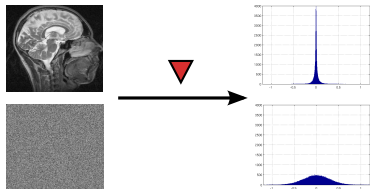
- As **few** elements as possible
- If at all, use them **big**

- Sparsity: A concept as basic as classical linear reconstruction
- Profound implications for how we (should) think about modelling, reconstruction, acquisition of real-world signals



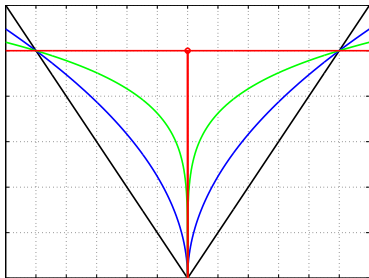
# Many Faces of Sparsity

- Image modelling
  - Processing
  - Reconstruction
  - Acquisition (sampling)
  - Computational neuroscience



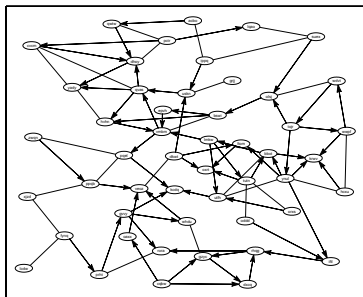
# Many Faces of Sparsity

- Image modelling
  - Processing
  - Reconstruction
  - Acquisition (sampling)
  - Computational neuroscience
- Relaxation of combinatorial optimization
  - Maximally sparse reconstruction



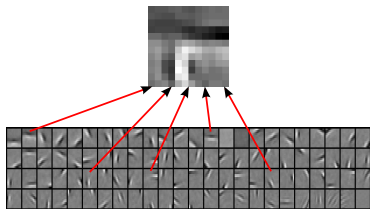
# Many Faces of Sparsity

- Image modelling
  - Processing
  - Reconstruction
  - Acquisition (sampling)
  - Computational neuroscience
- Relaxation of combinatorial optimization
  - Maximally sparse reconstruction
- Learning dependency structure
  - Meinshausen, Buehlmann
  - Graphical Lasso

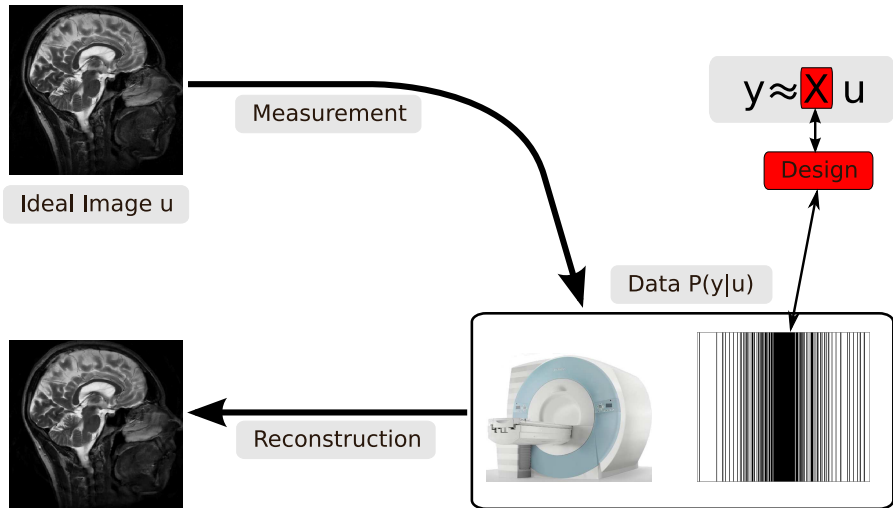


# Many Faces of Sparsity

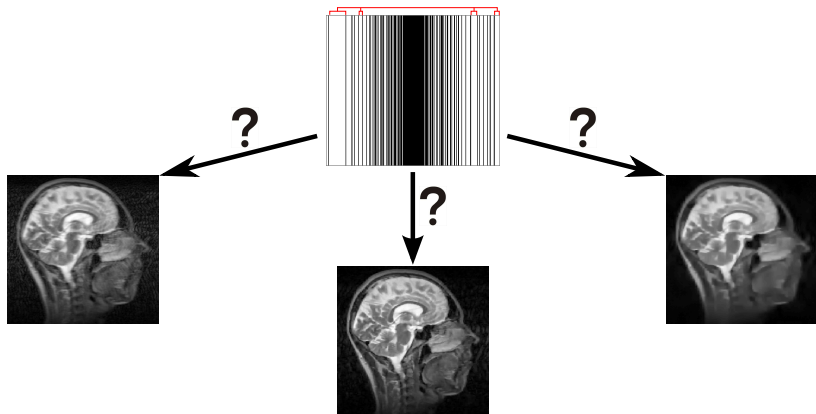
- Image modelling
  - Processing
  - Reconstruction
  - Acquisition (sampling)
  - Computational neuroscience
- Relaxation of combinatorial optimization
  - Maximally sparse reconstruction
- Learning dependency structure
  - Meinshausen, Buehlmann
  - Graphical Lasso
- Sparse coding
  - Olshausen, Field
  - Learning image priors



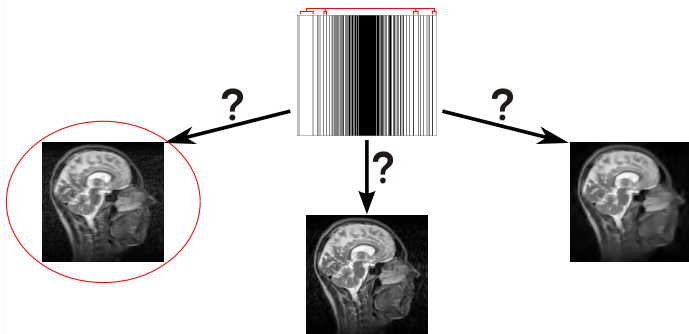
# Image Reconstruction



# Reconstruction is Ambiguous



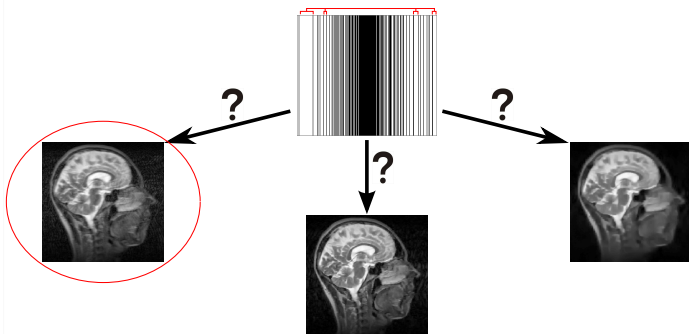
## Least Squares Estimation



## Least Squares Estimation (Linear Model)

$$\mathbf{u}_* = \operatorname{argmin}_{\mathbf{u}} \|\mathbf{y} - \mathbf{X}\mathbf{u}\|^2 \quad \text{s.t. } \|\mathbf{u}\|^2 \text{ small}$$

# Least Squares Estimation



## Least Squares Estimation (Linear Model)

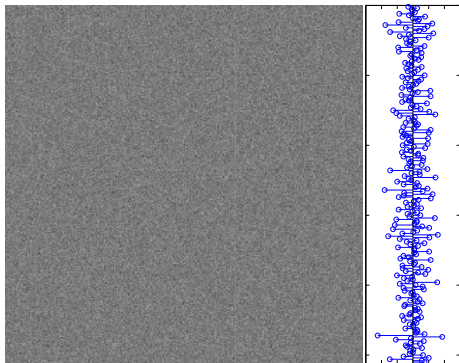
- Simple. Fast. Well understood
- Arbitrary decision (why **squares**?)



# Image Statistics

Whatever images are ...

they are not Gaussian!

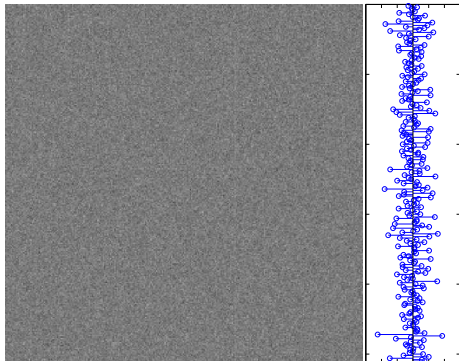


- Small noisy steps
- Gaussian random walker through pixel-land

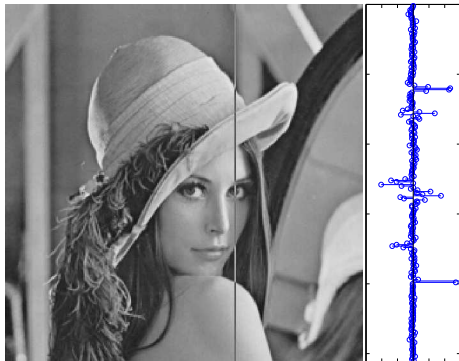
# Image Statistics

Whatever images are ...

they are not Gaussian!



- Small noisy steps
- Gaussian random walker through pixel-land



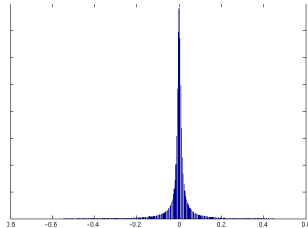
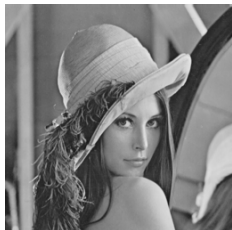
- Tiptoeing, edge jumping
- Gaussian won't do

# Image Statistics

Whatever images are ...

they are not Gaussian!

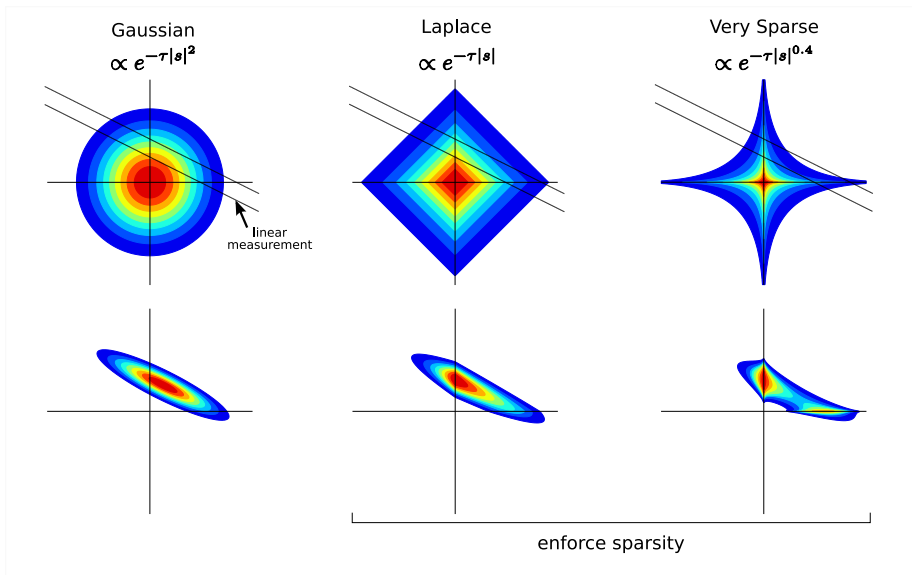
- Spatial smoothness: Image gradient super-Gaussian, **sparse**



Capture image properties in **prior distribution**  $P(\mathbf{u})$

## Sparsity Priors

courtesy Florian Steinke

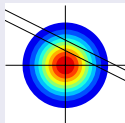


# Best of Both Worlds

$$P(\mathbf{u}) \propto \prod_{i=1}^q t_i(\mathbf{s}_i), \quad \mathbf{s} = \mathbf{B}\mathbf{u}, \quad t_i(\mathbf{s}_i) = e^{-\tau_i |\mathbf{s}_i|^2/2}$$

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

- Simple. Fast
- Well understood

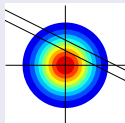


# Best of Both Worlds

$$P(\mathbf{u}) \propto \prod_{i=1}^q t_i(s_i), \quad \mathbf{s} = \mathbf{B}\mathbf{u}, \quad t_i(s_i) = e^{-\tau_i |s_i|}$$

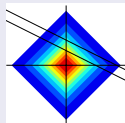
## Gaussian Prior $P(\mathbf{u})$

- Simple. Fast
- Well understood



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

- Better prior for real-world signals (images)

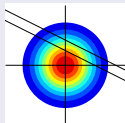


# Best of Both Worlds

$$P(\mathbf{u}) \propto \prod_{i=1}^q t_i(s_i), \quad \mathbf{s} = \mathbf{B}\mathbf{u}, \quad t_i(s_i) = e^{-\tau_i |s_i|}$$

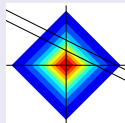
## Gaussian Prior $P(\mathbf{u})$

- Simple. Fast
- Well understood



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

- Better prior for real-world signals (images)



## Latent Gaussian Representations

- Gaussian scale mixtures
- Super-Gaussian potentials

$$t_i(s_i) = \int_{\gamma_i \geq 0} e^{-|s_i|^2 / (2\gamma_i)} f_i(\gamma_i) d\gamma_i$$

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

# Gaussian Scale Mixtures

- Mixture of Gaussians?  $K$ -means, EM, ...

$$P(X) = \sum_{j=1}^K \pi_j \mathcal{N}(X | \mu_j, \sigma^2)$$

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



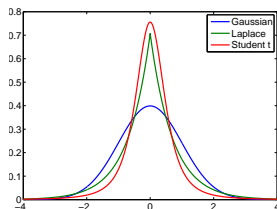
# Gaussian Scale Mixtures

- Mixture of Gaussians?  $K$ -means, EM, ...

$$P(X) = \sum_{j=1}^K \pi_j \mathcal{N}(X | \mu_j, \sigma^2)$$

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

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



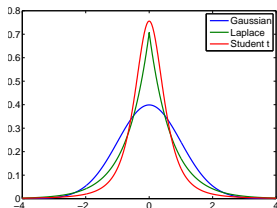
# Gaussian Scale Mixtures

- Mixture of Gaussians?  $K$ -means, EM, ...

$$P(X) = \sum_{j=1}^K \pi_j \mathcal{N}(X | \mu_j, \sigma^2)$$

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

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



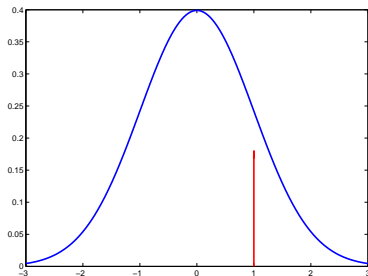
# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(\mathbf{0}, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(0, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

- Many distributions you know are **scale mixtures**
  - Gaussian [:-)].

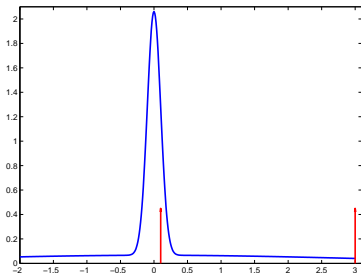


$$P(X) = N(X|0, \gamma)$$

# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(0, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

- Many distributions you know are **scale mixtures**
  - Gaussian [:-)]. Spike and slab

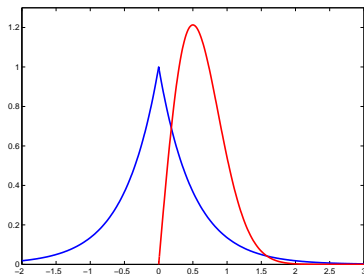


$$P(X) = \pi N(X|0, \gamma_1) + (1 - \pi)N(X|0, \gamma_2), \quad \gamma_1 \ll \gamma_2$$

# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(0, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

- Many distributions you know are **scale mixtures**
  - Gaussian [:-)]. Spike and slab
  - Exponential power ( $\alpha \leq 2$ )

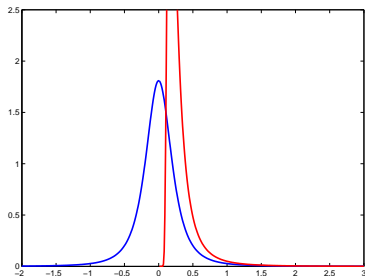


$$P(X) \propto e^{-\tau|X|^\alpha}, \quad \alpha \in (0, 2], \tau > 0$$

# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(0, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

- Many distributions you know are **scale mixtures**
  - Gaussian [:-)]. Spike and slab
  - Exponential power ( $\alpha \leq 2$ )
  - Student's t

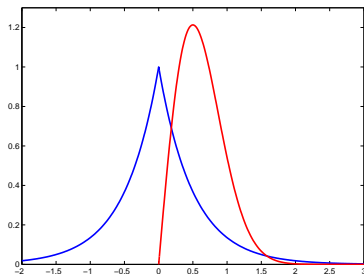


$$P(X) \propto (1 + (\tau/\nu)|X|^2)^{-(\nu+1)/2}, \quad \tau, \nu > 0$$

# Gaussian Scale Mixtures

$$X = \sqrt{\gamma}Y: Y \sim N(0, 1), \gamma \sim f(\gamma)\mathbf{I}_{\{\gamma \geq 0\}}$$

- Many distributions you know are **scale mixtures**
  - Gaussian [:-)]. Spike and slab
  - Exponential power ( $\alpha \leq 2$ )
  - Student's t



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

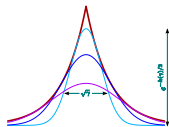
West, Biometrika 87

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



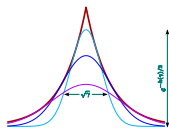
# Super-Gaussian Potentials

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



# Super-Gaussian Potentials

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

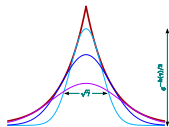


- $t(s)$  **even** and **positive**: Let's look at  $|s|^2 \mapsto 2 \log t(s)$
- What's that for a Gaussian  $t(s) = N(|s||0, \sigma^2)$ ?

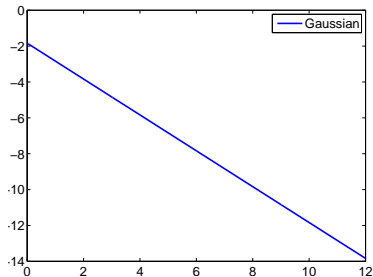
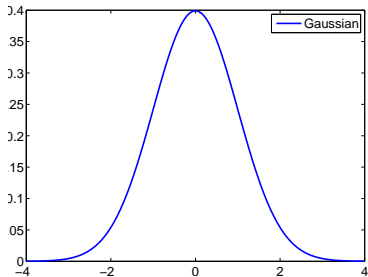
F7

# Super-Gaussian Potentials

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

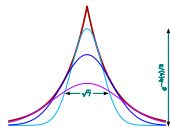


- $t(s)$  **even** and **positive**: Let's look at  $|s|^2 \mapsto 2 \log t(s)$
- What's that for a Gaussian  $t(s) = N(|s||0, \sigma^2)$ ?  
A **linear** (affine) function



# Super-Gaussian Potentials

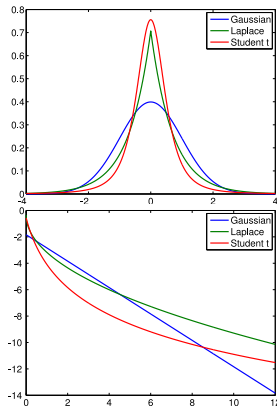
$$t(s) = \max_{\gamma \geq 0} e^{-|s|^2/(2\gamma)} g(\gamma)$$



Sparsity potentials are **super-Gaussian**

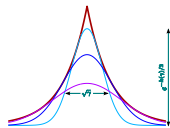
$$|s|^2 \mapsto 2 \log t(s) \text{ is convex}$$

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



# Super-Gaussian Potentials

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

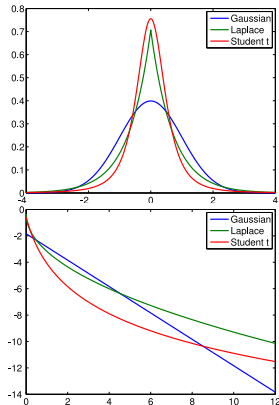


Sparsity potentials are **super-Gaussian**

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

- Affine  $\rightarrow$  convex:  
Shift mass to center and tails
- Scale mixtures are super-Gaussian

Palmer *et al.*,  
NIPS 2005



# Scale Mixtures are Super-Gaussian

$$\text{Gaussian scale mixture : } t(\mathbf{s}) = \int_{\geq 0} e^{-|\mathbf{s}|^2/(2\gamma)} f(\gamma) d\gamma$$

- $t(\mathbf{s})$  **even** and **positive**:  
 $x := |\mathbf{s}|^2 \Rightarrow t(\mathbf{s}) = e^{g(x)}$

F9

# Scale Mixtures are Super-Gaussian

Gaussian scale mixture : 
$$t(s) = \int_{\geq 0} e^{-|s|^2/(2\gamma)} f(\gamma) d\gamma$$

- $t(s)$  **even** and **positive**:  
 $x := |s|^2 \Rightarrow t(s) = e^{g(x)}$
- Super-Gaussian?  $|s|^2 \mapsto 2 \log t(s)$  is convex.  
 Show that  $g(x)$  is convex

# Scale Mixtures are Super-Gaussian

$$\text{Gaussian scale mixture : } t(s) = \int_{\geq 0} e^{-|s|^2/(2\gamma)} f(\gamma) d\gamma$$

- $t(s)$  **even** and **positive**:  
 $x := |s|^2 \Rightarrow t(s) = e^{g(x)}$
- Super-Gaussian?  $|s|^2 \mapsto 2 \log t(s)$  is convex.  
 Show that  $g(x)$  is convex
- Log-convexity: Closed under summation

Boyd, Vandenberghe, 2002

$$\psi(x, \gamma) \text{ convex } \forall \gamma \in \mathcal{C} \Rightarrow \log \int_{\mathcal{C}} e^{\psi(x, \gamma)} d\gamma \text{ convex}$$



# Scale Mixtures are Super-Gaussian

$$\text{Gaussian scale mixture : } t(s) = \int_{\geq 0} e^{-|s|^2/(2\gamma)} f(\gamma) d\gamma$$

- $t(s)$  **even** and **positive**:  
 $x := |s|^2 \Rightarrow t(s) = e^{g(x)}$
- Super-Gaussian?  $|s|^2 \mapsto 2 \log t(s)$  is convex.  
 Show that  $g(x)$  is convex
- Log-convexity: Closed under summation

Boyd, Vandenberghe, 2002

$$\psi(x, \gamma) \text{ convex } \forall \gamma \in \mathcal{C} \Rightarrow \log \int_{\mathcal{C}} e^{\psi(x, \gamma)} d\gamma \text{ convex}$$

- Apply to  $g(x)$ :

$$g(x) = \log \int_{\geq 0} e^{-x/(2\gamma)} f(\gamma) d\gamma = \log \int_{\geq 0} e^{-x/(2\gamma) + \log f(\gamma)} d\gamma$$

# Group Sparsity

$$t_i(\mathbf{s}_i) = \max_{\gamma_i \geq 0} e^{-|\mathbf{s}_i|^2 / (2\gamma_i)} g_i(\gamma_i)$$

- $t_i(\mathbf{s}_i)$  depends on absolute value  $|\mathbf{s}_i|$  only
- Can just as well plug in vector norm  $\|\mathbf{s}_i\|$

# Group Sparsity

$$t(\mathbf{s}_i) = \max_{\gamma_i \geq 0} e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} g_i(\gamma_i)$$

- $t_i(\mathbf{s}_i)$  depends on absolute value  $|\mathbf{s}_i|$  only
- Can just as well plug in vector norm  $\|\mathbf{s}_i\|$
- Useful for complex values:  $|\mathbf{s}_i| = \|(\Re \mathbf{s}_i, \Im \mathbf{s}_i)^T\|$

# Group Sparsity

$$t(\mathbf{s}_i) = \max_{\gamma_i \geq 0} e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} g_i(\gamma_i)$$

- $t_i(\mathbf{s}_i)$  depends on absolute value  $|\mathbf{s}_i|$  only
- Can just as well plug in vector norm  $\|\mathbf{s}_i\|$
- Useful for complex values:  $|\mathbf{s}_i| = \|(\Re \mathbf{s}_i, \Im \mathbf{s}_i)^T\|$
- Useful to **structure** sparsity: Joint penalization of **groups**  
 $\Rightarrow \ell_1 - \ell_2$  norms, group Lasso, and all that ...

# Group Sparsity

$$t(\mathbf{s}_i) = \max_{\gamma_i \geq 0} e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} g_i(\gamma_i)$$

- $t_i(\mathbf{s}_i)$  depends on absolute value  $|\mathbf{s}_i|$  only
- Can just as well plug in vector norm  $\|\mathbf{s}_i\|$
- Useful for complex values:  $|\mathbf{s}_i| = \|(\Re \mathbf{s}_i, \Im \mathbf{s}_i)^T\|$
- Useful to **structure** sparsity: Joint penalization of **groups**  
 $\Rightarrow \ell_1 - \ell_2$  norms, group Lasso, and all that ...
- Latent Gaussian representations: Just parameter tying  
 $e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} \propto N(\mathbf{s}_i | \mathbf{0}, \gamma_i \mathbf{1})$

# Group Sparsity

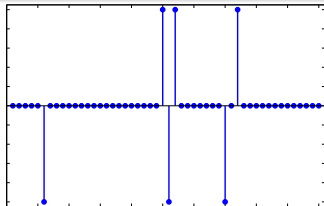
$$t(\mathbf{s}_i) = \max_{\gamma_i \geq 0} e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} g_i(\gamma_i)$$

- $t_i(\mathbf{s}_i)$  depends on absolute value  $|\mathbf{s}_i|$  only
- Can just as well plug in vector norm  $\|\mathbf{s}_i\|$
- Useful for complex values:  $|\mathbf{s}_i| = \|(\Re \mathbf{s}_i, \Im \mathbf{s}_i)^T\|$
- Useful to **structure** sparsity: Joint penalization of **groups**  
 $\Rightarrow \ell_1 - \ell_2$  norms, group Lasso, and all that ...
- Latent Gaussian representations: Just parameter tying  
 $e^{-\|\mathbf{s}_i\|^2 / (2\gamma_i)} \propto N(\mathbf{s}_i | \mathbf{0}, \gamma_i \mathbf{1})$

# Sparsity vs. Super-Gaussianity

## Sparse $\mathbf{s}$

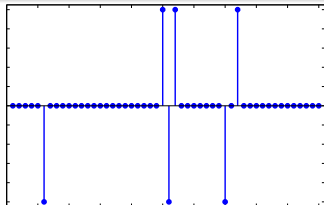
- Many/most  $s_j = 0$



# Sparsity vs. Super-Gaussianity

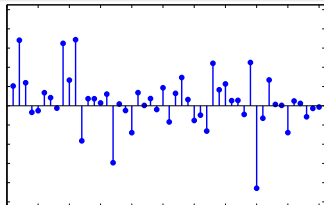
## Sparse $\mathbf{s}$

- Many/most  $s_i = 0$



## Super-Gaussian $\mathbf{s}$

- Super-Gaussian statistics
- Soft sparsity, statistical sparsity, power law decay, ...

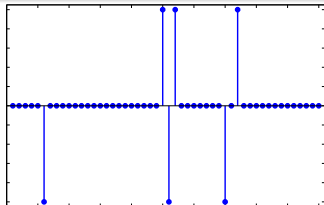




# Sparsity vs. Super-Gaussianity

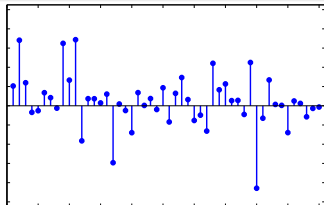
## Sparse $\mathbf{s}$

- Many/most  $s_i = 0$



## Super-Gaussian $\mathbf{s}$

- Super-Gaussian statistics
- Soft sparsity, statistical sparsity, power law decay, ...



- $P(\mathbf{s})$  super-Gaussian:  $\mathbf{s} \sim P(\mathbf{s})$  **no zeros** in general (only if  $P(\mathbf{s})$  degenerate)

# Where Are We?

- Real-world signals are not Gaussian.  
Gaussian assumptions made for convenience only.
- Super-Gaussian distributions:  
Trade-off between realistic and tractable/simple
- Latent Gaussian representations:
  - Gaussian scale mixtures
  - Super-Gaussian potentials
- Group potentials:  
Simple way to structure sparsity
- “Sparse” may mean super-Gaussian

# 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

# Variational Approximations

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

- 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 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  $\log Z$

# Variational Approximations

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

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

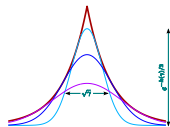
## Variational approximation

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

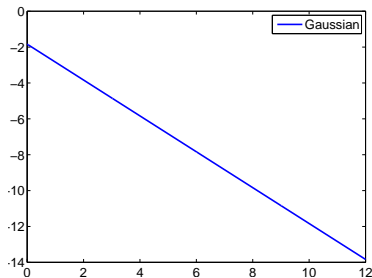
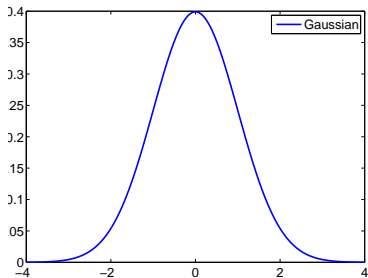
- Super-Gaussian bounding
- Expectation propagation
- Variational mean field Bayes [not here]
- Gaussian KL minimization [not here]

# Super-Gaussian Potentials

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

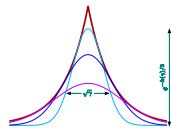


- $t(s)$  **even** and **positive**: Let's look at  $|s|^2 \mapsto 2 \log t(s)$
- What's that for a Gaussian  $t(s) = N(|s||0, \sigma^2)$ ?  
A **linear** (affine) function



# Super-Gaussian Potentials

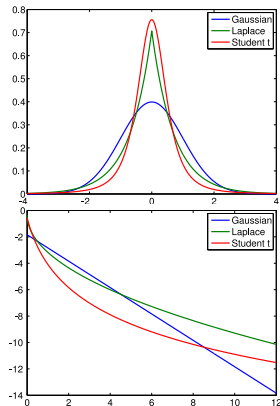
$$t(s) = \max_{\gamma \geq 0} e^{-|s|^2/(2\gamma)} e^{-h(\gamma)/2}$$



Sparsity potentials are **super-Gaussian**

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

- $t(s) = \max_{\gamma \geq 0} \dots$  Why?





# Convex (Fenchel) Duality

**Super-Gaussian:**

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

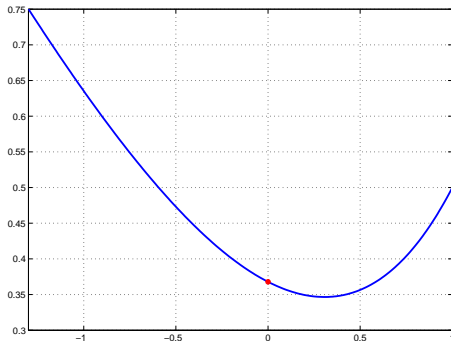
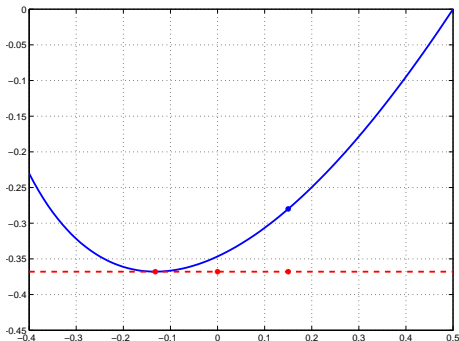
**Convex function:** Maximum of its affine lower bounds

**Super-Gaussian function:** Maximum of its Gaussian lower bounds

# Convex (Fenchel) Duality

Super-Gaussian:

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



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

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

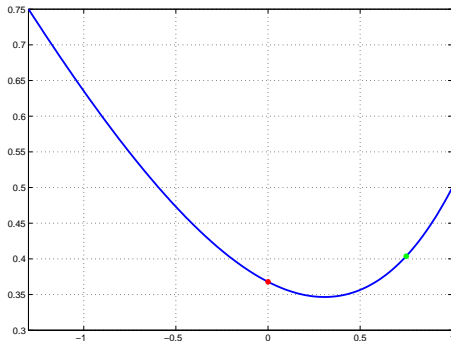
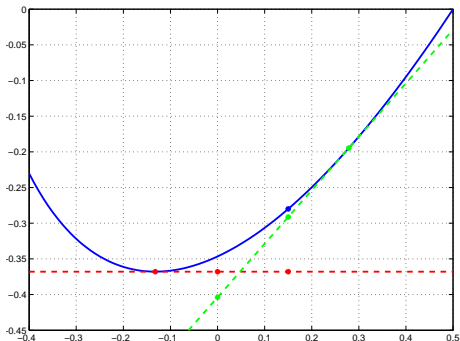
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

# Convex (Fenchel) Duality

Super-Gaussian:

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



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

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

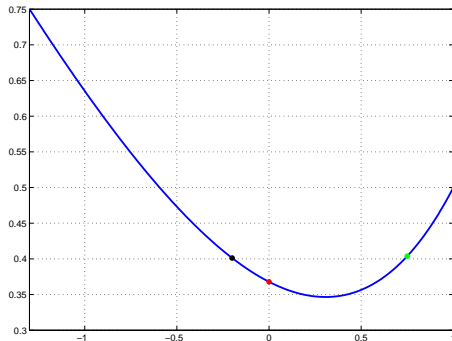
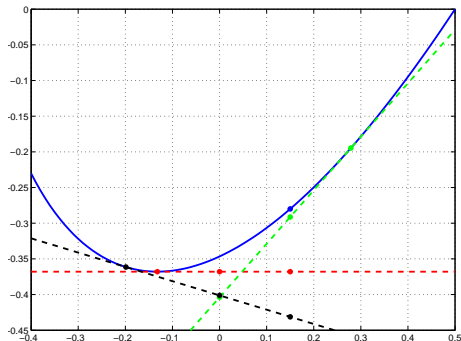
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

# Convex (Fenchel) Duality

Super-Gaussian:

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



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

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

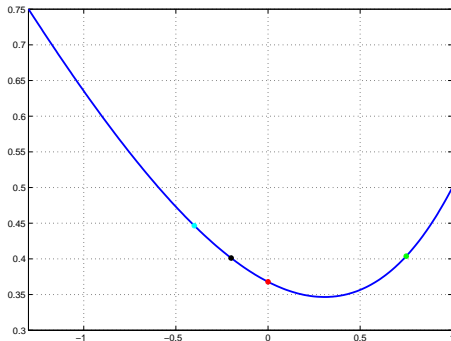
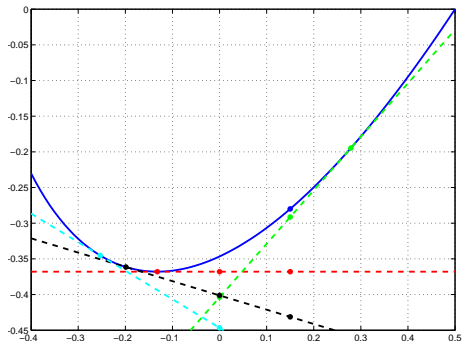
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

# Convex (Fenchel) Duality

Super-Gaussian:

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



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

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

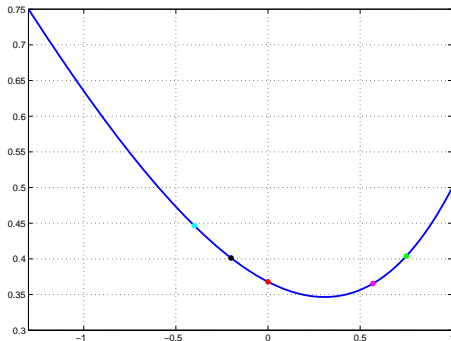
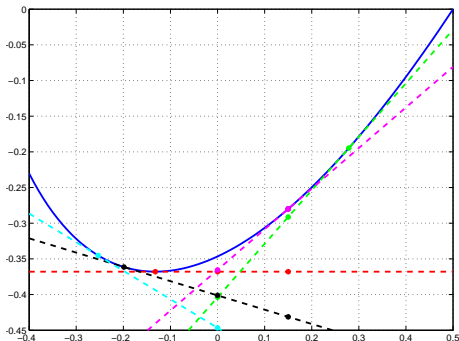
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

# Convex (Fenchel) Duality

Super-Gaussian:

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



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

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

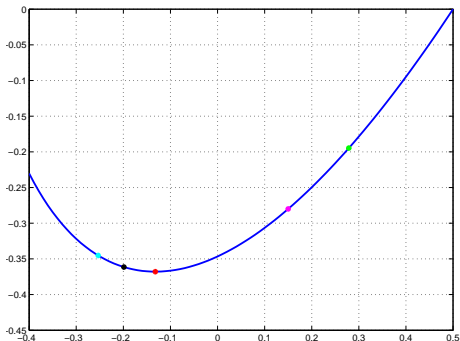
$$f^*(\pi) = \max_x \pi x - f(x)$$

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

# Convex (Fenchel) Duality

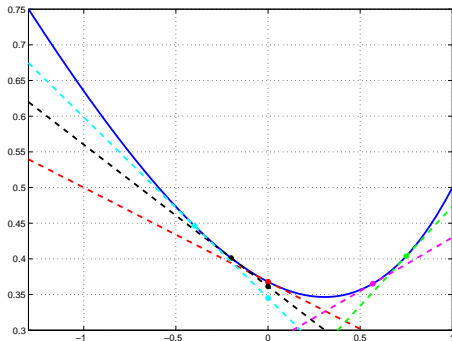
Super-Gaussian:

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



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

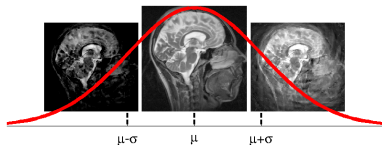
$$t(s) = \max_{\gamma} e^{(-|s|^2/\gamma - h(\gamma))/2}$$



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

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

# Super-Gaussian Potentials



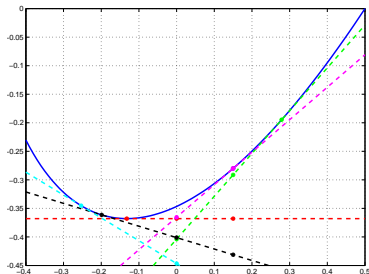
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Sparsity potentials are **super-Gaussian**

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

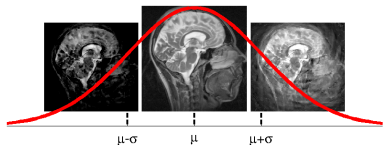
Convex (Fenchel) duality

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





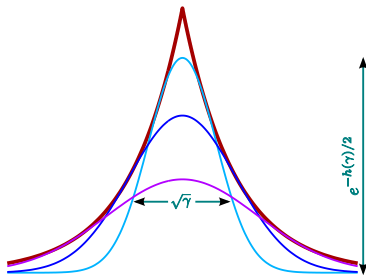
# Super-Gaussian Potentials



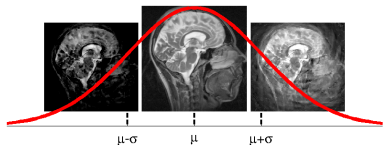
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

Sparsity potentials are **super-Gaussian**

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



# Super-Gaussian Bounding

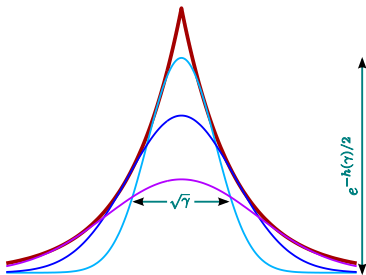


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

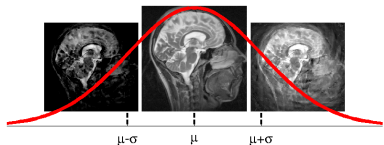
Sparsity potentials are **super-Gaussian**

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

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



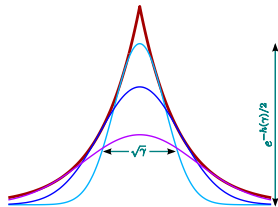
# Super-Gaussian Bounding



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

Exact representation

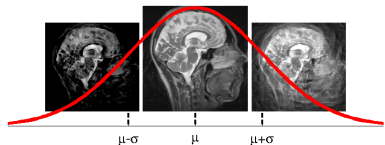
$$\begin{aligned} & \log Z \\ = & \log \int P(\mathbf{y}|\mathbf{u}) \max_{\gamma} e^{-(\mathbf{s}^H \mathbf{\Gamma}^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \end{aligned}$$



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

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

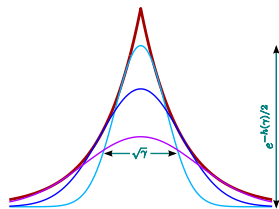
## Super-Gaussian Bounding



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

Lower bound

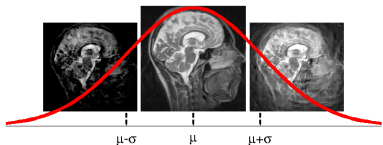
$$\begin{aligned} & \log Z \\ &= \log \int P(\mathbf{y}|\mathbf{u}) \max_{\gamma} e^{-(\mathbf{s}^H \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \\ &\geq \max_{\gamma} \log \int P(\mathbf{y}|\mathbf{u}) e^{-(\mathbf{s}^H \Gamma^{-1} \mathbf{s} + h(\gamma))/2} d\mathbf{u} \end{aligned}$$



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

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

# Super-Gaussian Bounding



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

Lower bound

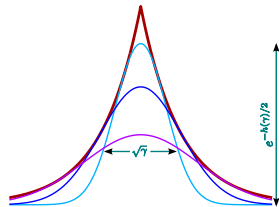
$\log Z$

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

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

Gaussian approximation

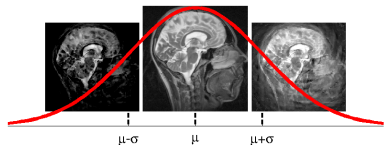
$$Q(\mathbf{u}|\mathbf{y}) = Z_Q^{-1} P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^H \mathbf{\Gamma}^{-1} \mathbf{s}/2}, \quad \mathbf{s} = \mathbf{B}\mathbf{u}$$



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

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

## Super-Gaussian Bounding



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

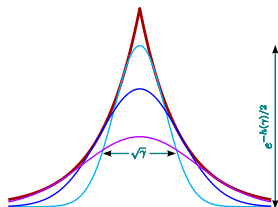
Variational problem:  $Q(\mathbf{u}|\mathbf{y}) \approx P(\mathbf{u}|\mathbf{y})$

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

Gaussian approximation

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

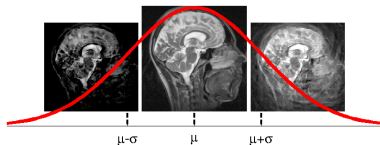
$$Z_Q = \int P(\mathbf{y}|\mathbf{u}) e^{-\mathbf{s}^H \Gamma^{-1} \mathbf{s} / 2} d\mathbf{u}$$



$$t_i(s_i) =$$

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

# Super-Gaussian Bounding

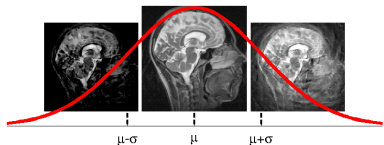


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

What did we do?

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

# Super-Gaussian Bounding



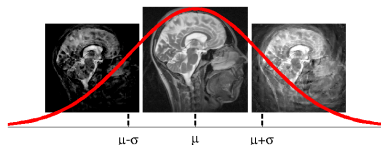
$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

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  $\log Z$



# Super-Gaussian Bounding

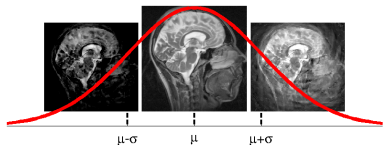


$$P(\mathbf{u}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{u}) \times P(\mathbf{u})}{P(\mathbf{y})}$$

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  $\mathcal{Q} = \{Q(\mathbf{u}|\mathbf{y})\}$

# 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**  $Q(\mathbf{u}|\gamma)$   
 $\Rightarrow$  Approximation family  $\mathcal{Q} = \{Q(\mathbf{u}|\gamma)\}$
- Divergence  $Q(\mathbf{u}|\gamma) \leftrightarrow P(\mathbf{u}|\mathbf{y})$ ? Maximize lower bound!  
 $\Rightarrow \phi(\gamma) = -2 \log Z_Q + h(\gamma)$

## MAP Estimation and Variational Inference

## MAP Estimation

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

## Bayesian Inference

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

# Coordinate Descent Algorithm

- Simple algorithm: Update **single variables**  $\gamma_j$

**repeat**

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

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

Gaussian propagation of pseudo-evidence change

**end for**

Refresh representation

**until** convergence

# Coordinate Descent Algorithm

- Simple algorithm: Update **single variables**  $\gamma_j$

**repeat**

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

Update  $\gamma_j$ , 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

Seeger, JMLR 2008

# Coordinate Descent Algorithm

- Simple algorithm: Update **single variables**  $\gamma_j$

**repeat**

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

Update  $\gamma_j$ , 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 Seeger, JMLR 2008
- Large scale problems?  
This algorithm is not scalable. Can do much better . . .