

## NEW TOOLS FOR BAYESIAN INFERENCE: THE VARIATIONAL APPROXIMATION

Nikolaos P. Galatsanos  
ECE Dept. Univ. of Patras, Greece.

Dimitris Tzikas  
CS Dept. Univ. of Ioannina, Greece.

## Outline

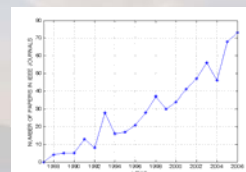
- Introduction
- Bayesian Inference Basics
- MAP Estimation
- Conjugate Distributions
- Graphical Models
- EM Algorithm
- An Alternative View of the EM
- The Variational EM framework
- Examples
  - Linear Regression
  - Blind Image Deconvolution
  - Image Restoration
    1. Constrained Variational Inference
    2. Bounded Variational Inference
  - Gaussian Mixture Modeling
- Conclusions

## Introduction

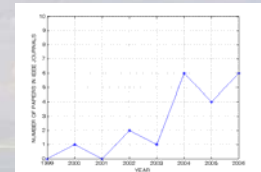


Thomas Bayes (1701-1761), left, first discovered "Bayes' theorem" in 1764. However, Bayes in his theorem used uniform priors. Pierre-Simon Laplace (1749-1827), right, unaware of Bayes' work, discovered the same theorem in more general form in a memoir he wrote at the age of 25 and showed its wide applicability.

## Introduction



(a)



(b)

- (a): # of papers /year in IEEE Journals "EM Algorithm"
- (b): # of papers /year in IEEE Journals "Variational Methodology"

## Applications of Variational Approximation

1. Mixture Modeling of pdfs & Clustering
2. ICA & PCA Analysis
3. Learning MRFs
4. Dynamic Systems Modeling
5. Image Recovery
6. Visual Tracking
7. Digital Communications
8. Acoustics and Speech Processing
9. Learning from Data Bases

## Sample of US Patents that Use Variational Inference

- US Patent 6879944 - Variational relevance vector machine, Issued on April 12, 2005
- US Patent 6556960 - Variational inference engine for probabilistic graphical models, Issued on April 29, 2003
- US Patent 20080025627-Removing camera shake from a single photograph
- US Patent: 6990447-Method and apparatus for denoising and deverbation using variational inference and Strong Speech Models, Issue date: Jan 24, 2006.
- US Patent 6591146- Method for learning switching linear dynamic system models from data, Issue date: Jul 8, 2003
- US Patent 6931374-Method of speech recognition using variational inference with switching state space models, Issue date Aug 16, 2005.
- US Patent 2004/0260548-Variational inference and learning for segmental switching state space models of hidden speech dynamics,
- US Patent 2004/0254903 A1-Systems and methods for tractable variational approximation for inference in decision graph Bayesian networks, Issue Feb 27, 2007
- Publication number US 2007/0083928-Data security and intrusion detection,
- Publication number: US 2005/0176057 A1-Diagnostic markers of mood disorders and methods of use thereof
- Publication number: US 2007/0233392 A1-Population sequencing
- Publication number: US 2006/0184260 A1-Player ranking with partial information
- Publication number: US 2007/0265718 A1-Team matching
- Publication number: US 2007/0098254 A1-Detecting humans via their pose

## Bayesian Inference Basics

- Estimation => **Parameter**
- Observations:  $x$
- Parameters:  $\theta$
- Likelihood Function:  $p(x; \theta)$  ( $\theta$  – parameter)
- Maximum Likelihood Estimation

$$\hat{\theta}_{ML} = \arg \max_{\theta} p(x; \theta)$$

## Bayesian Inference Basics

- Inference =>  $\theta$  **Random Variable**
- Find **Posterior**  $p(\theta | x)$
- Posterior **MORE** information than point estimate
  - $E(\theta | x)$  MMSE estimate
  - $Var(\theta | x)$  accuracy of estimates

## Bayesian Inference Basics

- “Hidden Variables”  $z$ :
  - Describe data generation mechanism (graphical model)
  - “links” between observations and parameters
  - Easy to compute  $p(x | z)$
  - Introduce priors  $p(z; \theta)$

## Bayesian Inference Basics

- Find Likelihood, **Marginalize** Hidden Variables

$$p(x; \theta) = \int p(x, z; \theta) dz = \int p(x | z; \theta) p(z; \theta) dz$$

- Find Posterior

$$p(z | x; \theta) = \frac{p(x | z; \theta) p(z; \theta)}{p(x; \theta)}$$

- In most cases of interest **Cannot Marginalize**

## Bayesian Inference Basics

- Main effort in Bayesian Inference techniques **bypass** or **approximate** marginalization integral.
  - Random sampling methods
    - Monte Carlo
  - Deterministic approximations
    - Laplace
    - Variational

## MAP Estimation

- Defined as mode of posterior :

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | x)$$

- Based on Bayes' theorem

$$p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)}$$

- Can be found as:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(x | \theta) p(\theta)$$

## MAP Estimation

- No need for  $p(x)$  **no marginalization**
- MAP Estimator
  - Much **easier** to find
  - **Mode** of posterior
  - No information about **shape** of posterior
- Uses Bayes' Theorem, however posterior **not found**
- **MAP=Poor Man's Bayesian Inference**

## Conjugate Priors

- Find prior which **allows closed form marginalization** of hidden variables

$$p(x) = \int p(x, z) dz = \int p(x | z) p(z) dz$$

## Conjugate Priors

- Example #1:  $\mu$  hidden variable

$$p(x | \mu, \sigma^2) = N(x; \mu, \sigma^2), p(\mu; \mu_0, \sigma_0^2) = N(\mu; \mu_0, \sigma_0^2)$$

$$p(x; \mu_0, \sigma^2, \sigma_0^2) = \int p(x | \mu, \sigma^2) p(\mu; \mu_0, \sigma_0^2) d\mu$$

## Conjugate Priors

- w.r.t. to  $\mu$  **both**  $p(x | \mu; \sigma^2)$  and  $p(\mu; \mu_0, \sigma_0^2)$  have the **same** form (Gaussian).

$$p(x | \mu, \sigma^2) = f(\mu) \propto \exp\left(-\frac{1}{2\sigma^2}(\mu^2 - 2\mu x)\right)$$

$$p(\mu; \mu_0, \sigma_0^2) = g(\mu) \propto \exp\left(-\frac{1}{2\sigma_0^2}(\mu^2 - 2\mu\mu_0)\right)$$

- $p(x | \mu, \sigma^2)$ , and  $p(\mu; \mu_0, \sigma_0^2)$  **conjugate**

## Conjugate Priors

- Marginalizing  $\mu$  possible (Gaussian Integral):

$$\begin{aligned} p(x; \mu_0, \sigma^2, \sigma_0^2) &= \int p(x | \mu, \sigma^2) p(\mu; \mu_0, \sigma_0^2) d\mu \\ &= N\left(x; \mu_0, \left(\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}\right)^{-1}\right) \end{aligned}$$

- Posterior:

$$p(\mu | x; \sigma^2, \mu_0, \sigma_0^2) = N\left(\mu; \frac{\sigma_0^2 x + \sigma^2 \mu_0}{\sigma^2 + \sigma_0^2}, \left(\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2}\right)^{-1}\right)$$

## Conjugate Priors

- Example #2:  $a$  "hidden variable"

$$p(x | a) = N(x; 0, a^{-1}),$$

$$p(a; b, c) = \text{Gamma}(a; b, c) = \frac{c^b a^{b-1} \exp(-ca)}{\Gamma(b)}$$

## Conjugate Priors

- w.r.t. to  $a$  both  $p(x|a)$  and  $p(a;b,c)$  have the **same** form (Gamma).

$$p(x|a) = f(a) \propto a^{1/2} \exp\left(-\frac{ax^2}{2}\right)$$

$$p(a;b,c) = g(a) \propto a^{b-1} \exp(-ac)$$

- $p(x|a)$  and  $p(a;b,c)$  **conjugate**

## Conjugate Priors

- Marginalize  $a$  possible

$$p(x;b,c) = \frac{\Gamma(b+1/2)}{\Gamma(b)} c^b \left(\frac{1}{2\pi}\right)^{1/2} \left(c + \frac{x^2}{2}\right)^{-b-1/2}$$

- Can write as Student's-t with  $\nu = 2b, \lambda = b/c$

$$p(x;\lambda,\nu) = St(x;\lambda,\nu) = \frac{\Gamma(\nu/2+1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left(1 + \frac{\lambda x^2}{\nu}\right)^{-\nu/2-1/2}$$

- Posterior

$$p(a|x;b,c) = Gamma\left(a; b + \frac{1}{2}, c + \frac{x^2}{2}\right)$$

## Conjugate Priors

Likelihood	Conjugate Prior Distribution	Posterior Distribution
$N(\mathbf{X} \boldsymbol{\mu},\boldsymbol{\Sigma})$	$N(\boldsymbol{\mu} \boldsymbol{\mu}_0,\boldsymbol{\Sigma}_0)$	$N(\boldsymbol{\mu}   (\boldsymbol{\Sigma}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} (\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 + n\boldsymbol{\Sigma}^{-1}\bar{\mathbf{x}}), (\boldsymbol{\Sigma}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1})$
$N(x \mu,\sigma^2)$	$Gamma(\sigma^{-2} a,b)$	$Gamma(\sigma^{-2}   a + n/2, b + \sum_{i=1}^n (x_i - \mu)^2 / 2)$
$N(\mathbf{X} \boldsymbol{\mu},\boldsymbol{\Sigma})$	$Wishart(\boldsymbol{\Sigma}^{-1} \nu,\mathbf{V})$	$Wishart(\boldsymbol{\Sigma}^{-1}   \nu + n, \mathbf{V} + \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})^T)$
$Multi(\mathbf{X} \boldsymbol{\pi})$	$Dir(\boldsymbol{\pi} \mathbf{a})$	$Dir(\boldsymbol{\pi}   \mathbf{a} + \sum_{i=1}^n \mathbf{X}_i)$

## Graphical Models

- Represent dependencies between rv.s in a statistical model
- Graph nodes represent rv.s and edges dependencies
- Directed and undirected graphs.
- Undirected *Markov Random Fields*
- Rest of presentation: directed, no cycles graphs

## Graphical Models

- $X_s$  rv associated node  $s$ ,  $\pi(s)$  parents of  $s$

- $p(x_s | x_{\pi(s)})$  conditional pdf

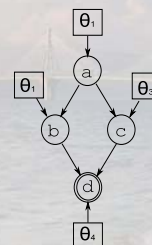
- Joint pdf over all variables

$$p(\mathbf{x}) = \prod_s p(x_s | x_{\pi(s)})$$

## Graphical Models

- Example:

$$p(a,b,c,d;\boldsymbol{\theta}) = p(a;\theta_1) p(b|a;\theta_2) p(c|a;\theta_3) p(d|b,c;\theta_4)$$





## EM Algorithm

- $x$ - observations,  $z$  -“hidden variables”,  $\theta$ -parameters.

- Define:

$$Q(\theta, \theta^{old}) = E[\ln p(x, z; \theta)]_{p(z|x; \theta^{old})}$$

$$= \int \ln p(x, z; \theta) p(z|x; \theta^{old}) dz$$

## EM Algorithm

1. Initial selection  $\theta^{old}$
2. **E-step:** Evaluate  $p(z|x; \theta^{old})$
3. **M-step:** Evaluate  $\theta^{new}$ 

$$\theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$$
4. Check for convergence parameters of log-likelihood, if not satisfied, *go to 2.*

## An Alternative View of the EM Algorithm

- Can write:

$$\ln p(x; \theta) = F(q, \theta) + KL(q \| p)$$

$$F(q, \theta) = \int q(z) \ln \left( \frac{p(x, z; \theta)}{q(z)} \right) dz$$

$$KL(q \| p) = - \int q(z) \ln \left( \frac{p(z|x; \theta)}{q(z)} \right) dz$$

- $q(z)$  any pdf,
- $KL(q \| p)$  Kullback-Leibler Divergence

## An Alternative View of the EM Algorithm \*\*\*

- $KL(q \| p) \geq 0$
- $\ln p(x; \theta) \geq F(q, \theta)$
- $\ln p(x; \theta) = F(q, \theta)$  when  $KL(q \| p) = 0$   
or  $q(z) = p(z|x; \theta)$

## An Alternative View of the EM Algorithm \*\*\*

- Substitute  $q(z) = p(z|x; \theta^{OLD})$
- Can write:

$$F(q, \theta) = \int p(z|x; \theta^{OLD}) \ln p(x, z; \theta) dz - \int p(z|x; \theta^{OLD}) \ln p(x|z; \theta^{OLD}) dz$$

$$= Q(\theta, \theta^{OLD}) + \text{constant}$$

- Same as before:

$$Q(\theta, \theta^{OLD}) = E[\ln p(x, z; \theta)]_{p(z|x; \theta^{OLD})}$$

## The Variational EM framework \*\*\*

- Assume  $p(z|x; \theta)$  **unknown**
- $F(q, \theta)$  **functional** in terms of  $q(z)$
- Variational EM
  - **Variational E-step:**  $q^{NEW}(z) = \max F(q, \theta^{OLD})$
  - **Variational M-step:**  $\theta^{NEW} = \arg \max_{\theta} F(q^{NEW}, \theta)$

## The Variational EM framework \*\*\*

- Key issue: **maximize**  $F(q, \theta)$  w.r.t.  $q(z)$  ?
  - Assume parametric form for  $q(z)$
  - $q(z)$  **approximates** unknown posterior  $p(z|x)$
- $\ln p(x; \theta) = F(q, \theta) + KL(q||p)$
- **max**  $F(q, \theta) \Rightarrow$  **min**  $KL(q||p)$

## Mean Field Approximation\*\*\*

- Assumption:  $q(z)$  factorizes

$$q(z) = \prod_{i=1}^M q_i(z_i)$$

- “Mean Field” approximation, **statistical physics**

## Mean Field Approximation\*\*\*

- Then optimal factor  $q_j(z_j)$  is:

$$q_j^*(z_j) = \frac{\exp\left(\langle \ln p(x, z; \theta) \rangle_{i \neq j}\right)}{\int \exp\left(\langle \ln p(x, z; \theta) \rangle_{i \neq j}\right) dz_j}$$

- with

$$\langle \ln p(x, z; \theta) \rangle_{i \neq j} = \int \ln p(x, z; \theta) \prod_{i \neq j} q_i dz_i$$

## Conjugate-Exponential models

- Prior distributions belong to the exponential family

$$p(X | Y) = \exp\left[\phi(Y)^T u(X) + f(X) + g(Y)\right]$$

- Graphical model with conjugate priors at each level

- hidden  $z$ , parents  $\pi(z)$
- $p(z | \pi(z))$  conjugate to  $p(\pi(z) | \pi(\pi(z)))$

$$p(z | \pi(\pi(z))) = \int p(z | \pi(z)) p(\pi(z) | \pi(\pi(z))) d\pi(z)$$

## Conjugate-Exponential models

- $p(x|z_1), p(z_1|z_2), p(z_2)$  exponential distributions
- $p(z_2)$  conjugate to  $p(z_1|z_2)$

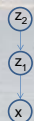
$$p(x | z_2) = \int p(x | z_1) p(z_1 | z_2) dz_1$$

- $P(x|z_1)$  conjugate to  $p(z_1|z_2)$

$$p(z_1) = \int p(z_1 | z_2) p(z_2) dz_2$$

- Cannot evaluate marginal

$$p(x) = \int p(x | z_1) p(z_1 | z_2) p(z_2) dz_1 dz_2$$



## Conjugate-Exponential models

- Tractable Variational computations

$$q(z_i) \propto \exp\left(\langle \ln p(x, z; \theta) \rangle_{i \neq j}\right)$$

- $\ln q(z_1) = \ln p(x | z_1; \theta) + \langle \ln p(z_1 | z_2; \theta) \rangle_{q(z_2)} + const$

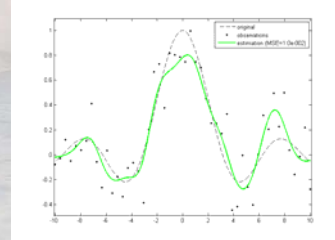
- $\ln q(z_2) = \langle \ln p(z_1 | z_2; \theta) \rangle_{q(z_1)} + \ln p(z_2; \theta) + const$

## Examples Linear Regression

## Linear Regression

- Observations at  $t_n$  find  $y(x)$

$$t_n = y(x_n) + \varepsilon_n, n = 1, \dots, N$$



## Linear Regression

- Signal  $y(x)$  modeled by

$$y(\mathbf{x}; \mathbf{w}) = \sum_{m=1}^M w_m \phi_m(\mathbf{x})$$

- Observation model

$$t_n = y(\mathbf{x}_n; \mathbf{w}) + \varepsilon_n, n = 1 \dots N$$

$$\mathbf{t} = \Phi \mathbf{w} + \boldsymbol{\varepsilon}$$

- $\Phi$   $N \times M$  design matrix

$$\Phi = (\phi_1, \dots, \phi_M), \phi_m = (\phi_m(x_1), \dots, \phi_m(x_N))^T$$

## Linear Regression

- Gaussian additive noise

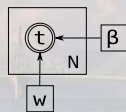
$$p(\boldsymbol{\varepsilon}) = N(\boldsymbol{\varepsilon} | 0, \beta^{-1} \mathbf{I})$$

- Likelihood of observations

$$p(\mathbf{t}; \mathbf{w}, \beta) = N(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I})$$

## Linear Regression: Least Squares ( $\mathbf{w}$ -parameters)

- Graphical Model



- Maximum Likelihood Estimation

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(\mathbf{t}; \mathbf{w}, \beta) = \arg \max_{\mathbf{w}} N(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I})$$

## Linear Regression: Parameter Estimation

- Minimize mean square error

$$E_{LS}(\mathbf{w}) = \|\mathbf{t} - \Phi \mathbf{w}\|^2 = \sum_{n=1}^N [t_n - y(\mathbf{x}_n; \mathbf{w})]^2$$

- Solution

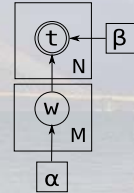
$$\mathbf{w}_{LS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad \beta^{-1} = \frac{1}{N} \|\mathbf{t} - \Phi \mathbf{w}\|^2$$

## Linear Regression: Problems with Least Squares ( $w$ -parameters)

- Maximum Likelihood Limitations:
  - $N$  observations  $M$  parameters to estimate,  $N \gg M$
  - Otherwise  $[\Phi^T \Phi]^{-1}$  Ill-conditioned
  - ML estimates large variance
- Assume constraints on the parameters  $w$ 
  - Bayesian: Use prior distribution  $p(w)$

## Bayesian Linear Regression: Model-2 ( $w$ -Gaussian iid distributed)

- Gaussian weight prior
 
$$p(w; \alpha) = \prod_{m=1}^M N(w_m | 0, \alpha^{-1})$$
- Why Gaussian?
  - Conjugate to likelihood



## Bayesian Linear Regression: Inference

- Posterior given by Bayes's law
 
$$p(w | t; \alpha, \beta) = \frac{p(t | w; \beta) p(w; \alpha)}{p(t; \alpha, \beta)}$$
- Can be found in closed form
 
$$p(w | t; \alpha, \beta) = N(w | \mu, \Sigma)$$

$$\mu = \beta \Sigma \Phi^T t$$

$$\Sigma = (\beta \Phi^T \Phi + \alpha \mathbf{I})^{-1}$$

## Bayesian Linear Regression: Parameter Estimation

- Maximum Likelihood
 
$$(\alpha_{ML}, \beta_{ML}) = \underset{\alpha, \beta}{\operatorname{argmin}} p(t; \alpha, \beta) = \int p(t | w; \beta) p(w; \alpha) dw$$

$$= \underset{\alpha, \beta}{\operatorname{argmin}} \left\{ \log |\beta^{-1} \mathbf{I} + \alpha^{-1} \Phi \Phi^T| + t^T (\beta^{-1} \mathbf{I} + \alpha^{-1} \Phi \Phi^T)^{-1} t \right\}$$
- Not straight forward constrained optimization because  $\alpha > 0, \beta > 0$ .
- Resort to EM algorithm!

## Bayesian Linear Regression: EM Algorithm

- Observations  $t$ , parameters  $\alpha, \beta$ ; **hidden variables**  $w$ .
- Key for the application of EM:  $p(w|t)$  **explicitly known**
- E-step:  $p(w|t; \alpha^0, \beta^0)$  obtained, inference of hidden variables
- M-step: ML estimates of parameters

## Bayesian Linear Regression: EM algorithm

- $Q^{(i)}(t, w; \alpha, \beta) = \langle \ln p(t, w; \alpha, \beta) \rangle_{p(w|t, \alpha^{(i)}, \beta^{(i)})} = \langle \ln p(t | w; \alpha, \beta) p(w; \alpha, \beta) \rangle_{p(w|t, \alpha^{(i)}, \beta^{(i)})}$ 

$$Q^{(i)}(t, w; \alpha, \beta) = \left\langle \frac{N}{2} \ln \beta - \frac{\beta}{2} (\|t - \Phi w\|^2) + \frac{M}{2} \ln \alpha - \frac{\alpha}{2} (\|w\|^2) \right\rangle + \text{const}$$

$$= \frac{N}{2} \ln \beta - \frac{\beta}{2} \langle \|t - \Phi w\|^2 \rangle + \frac{M}{2} \ln \alpha - \frac{\alpha}{2} \langle \|w\|^2 \rangle + \text{const}$$

$$Q^{(i)}(t, w; \alpha, \beta) = \frac{N}{2} \ln \beta - \frac{\beta}{2} (\|t - \Phi \mu^{(i)}\|^2 + \text{tr}[\Phi^T \Sigma^{(i)} \Phi])$$

$$+ \frac{M}{2} \ln \alpha - \frac{\alpha}{2} (\|\mu^{(i)}\|^2 + \text{tr}[\Sigma^{(i)}]) + \text{const}$$



## Linear Regression: EM algorithm, E-step

• E-step: Evaluate  $p(\mathbf{w}|\mathbf{t}; \alpha^{(l)}, \beta^{(l)}) = N(\boldsymbol{\mu}^{(l)}, \Sigma^{(l)})$

$$\Sigma^{(l)} = (\beta^{(l)} \Phi^T \Phi + \alpha^{(l)} \mathbf{I})^{-1}$$

$$\boldsymbol{\mu}^{(l)} = \beta^{(l)} \Sigma^{(l)} \Phi^T \mathbf{t}$$

## Bayesian Linear Regression: Parameter Estimation: M-step

• M-step

$$(\alpha^{(l+1)}, \beta^{(l+1)}) = \operatorname{argmax}_{(\alpha, \beta)} Q^{(l)}(\mathbf{t}, \mathbf{w}; \alpha, \beta)$$

$$\frac{\partial Q^{(l)}(\mathbf{t}, \mathbf{w}; \alpha, \beta)}{\partial \alpha} = \frac{M}{2\alpha} - \frac{1}{2} (\|\boldsymbol{\mu}^{(l)}\|^2 + \operatorname{tr}[\Sigma^{(l)}])$$

$$\frac{\partial Q^{(l)}(\mathbf{t}, \mathbf{w}; \alpha, \beta)}{\partial \beta} = \frac{N}{2\beta} - \frac{1}{2} (\|\mathbf{t} - \Phi \boldsymbol{\mu}^{(l)}\|^2 + \operatorname{tr}[\Phi^T \Sigma^{(l)} \Phi])$$

$$\alpha^{(l+1)} = \frac{M}{\|\boldsymbol{\mu}^{(l)}\|^2 + \operatorname{tr}[\Sigma^{(l)}]} \quad \beta^{(l+1)} = \frac{N}{\|\mathbf{t} - \Phi \boldsymbol{\mu}^{(l)}\|^2 + \operatorname{tr}[\Phi^T \Sigma^{(l)} \Phi]}$$

## Sparse Bayesian Linear Regression

- Limited model:  $w_i$  iid
- How to select basis functions?
- Sparse Linear Model
  - Consider many basis functions
  - Estimations use only few basis functions
- Advantages
  - Small variance (good generalization)
  - Fast Evaluation of estimation

## Sparse Bayesian Linear Regression: Prior Distribution

• New Prior:  $w_i$  not identically distributed

$$p(\mathbf{w} | \boldsymbol{\alpha}) = \prod_{m=1}^M N(w_m | 0, \alpha_m^{-1})$$

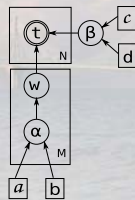
- $N$ -parameters to estimate,  $N$ -observations
- Use Conjugate Hyperpriors

$$p(\boldsymbol{\alpha}; a, b) = \prod_{m=1}^M \operatorname{Gamma}(\alpha_m | a, b)$$

$$p(\beta; c, d) = \operatorname{Gamma}(\beta | c, d)$$

## Sparse Bayesian Linear Regression: Prior Distribution

• Graphical Model:  $\mathbf{w}$ ,  $\boldsymbol{\alpha}$  hidden variables,  $a, b, c, d$ , parameters



## Sparse Bayesian Linear Regression: Prior Distribution

• "True" weight prior is **Student's-t**

$$p(\mathbf{w}; \mathbf{a}, b) = \int p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}; a, b) d\boldsymbol{\alpha}$$

$$= \int \prod_{m=1}^M N(w_m | 0, \alpha_m^{-1}) \operatorname{Gamma}(\alpha_m | a, b) d\alpha_m$$

$$= \prod_{m=1}^M \operatorname{St}(w_m | \lambda, \nu)$$

### Sparse Bayesian Linear Regression: Variational Inference

- Posterior:

$$p(\mathbf{w}, \boldsymbol{\alpha}, \beta | \mathbf{t}) = \frac{p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) p(\beta)}{p(\mathbf{t})}$$

- Cannot compute

$$p(\mathbf{t}) = \int p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) p(\beta) d\mathbf{w} d\boldsymbol{\alpha} d\beta$$

- Variational Mean Field Approximation

$$p(\mathbf{w}, \boldsymbol{\alpha}, \beta | \mathbf{t}) \approx q(\mathbf{w}, \boldsymbol{\alpha}, \beta) = q(\mathbf{w}) q(\boldsymbol{\alpha}) q(\beta)$$

### Sparse Bayesian Linear Regression: Variational Inference

$$\begin{aligned} \ln q(\mathbf{w}) &= \langle \ln p(\mathbf{t}, \mathbf{w}, \boldsymbol{\alpha}, \beta) \rangle_{q(\boldsymbol{\alpha})q(\beta)} + \text{const} \\ &= \langle \ln p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \boldsymbol{\alpha}) \rangle_{q(\boldsymbol{\alpha})q(\beta)} + \text{const} \\ &= \langle \ln p(\mathbf{t} | \mathbf{w}, \beta) + \ln p(\mathbf{w} | \boldsymbol{\alpha}) \rangle_{q(\boldsymbol{\alpha})q(\beta)} + \text{const} \\ &= \left\langle -\frac{\beta}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) - \frac{1}{2} \sum_{m=1}^M \alpha_m w_m^2 \right\rangle_{q(\boldsymbol{\alpha})q(\beta)} + \text{const} \\ &= -\frac{\langle \beta \rangle}{2} [\mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \Phi \mathbf{w} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}] - \frac{1}{2} \sum_{m=1}^M \langle \alpha_m \rangle w_m^2 + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^T (\langle \beta \rangle \Phi^T \Phi + \langle \Lambda \rangle) \mathbf{w} - \langle \beta \rangle \mathbf{w}^T \Phi^T \mathbf{t} + \text{const} \\ &= -\frac{1}{2} \mathbf{w}^T \Sigma^{-1} \mathbf{w} - \mathbf{w}^T \boldsymbol{\mu} + \text{const} \\ q(\mathbf{w}) &= N(\mathbf{w} | \boldsymbol{\mu}, \Sigma) \quad \Sigma = (\langle \beta \rangle \Phi^T \Phi + \langle \Lambda \rangle)^{-1} \quad \boldsymbol{\mu} = \langle \beta \rangle \Sigma \Phi^T \mathbf{t} \end{aligned}$$

### Sparse Bayesian Linear Regression: Variational Inference

$$\begin{aligned} \ln q(\boldsymbol{\alpha}) &= \langle \ln p(\mathbf{t}, \mathbf{w}, \boldsymbol{\alpha}, \beta) \rangle_{q(\mathbf{w})q(\beta)} + \text{const} \\ &= \langle \ln p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) \rangle_{q(\mathbf{w})} + \text{const} \\ &= \frac{1}{2} \sum_{m=1}^M \ln \alpha_m - \sum_{m=1}^M \alpha_m \langle w_m^2 \rangle + (a-1) \sum_{m=1}^M \ln \alpha_m - b \sum_{m=1}^M \alpha_m + \text{const} \\ &= \left( a - \frac{1}{2} \right) \sum_{m=1}^M \ln \alpha_m - \sum_{m=1}^M \left( \frac{1}{2} \langle w_m^2 \rangle + b \right) \alpha_m + \text{const} \\ &= \tilde{a} \sum_{m=1}^M \ln \alpha_m - \sum_{m=1}^M \tilde{b}_m \alpha_m + \text{const} \\ q(\boldsymbol{\alpha}) &= \prod_{m=1}^M \text{Gamma}(\alpha_m | \tilde{c}, \tilde{b}_m) \quad \tilde{a} = a + 1/2 \quad \tilde{b}_m = b + \frac{1}{2} \langle w_m^2 \rangle \end{aligned}$$

### Sparse Bayesian Linear Regression: Variational Inference

$$\begin{aligned} \ln q(\beta) &= \langle \ln p(\mathbf{t}, \mathbf{w}, \boldsymbol{\alpha}, \beta) \rangle_{q(\mathbf{w})q(\boldsymbol{\alpha})} + \text{const} \\ &= \langle \ln p(\mathbf{t} | \mathbf{w}) p(\beta) \rangle_{q(\mathbf{w})q(\boldsymbol{\alpha})} + \text{const} \\ &= \frac{N}{2} \ln \beta - \frac{1}{2} \beta \langle \|\mathbf{t} - \Phi \mathbf{w}\|^2 \rangle + (c-1) \ln \beta - d \beta + \text{const} \\ &= \left( c + \frac{N}{2} - 1 \right) \ln \beta - \left( \frac{1}{2} \langle \|\mathbf{t} - \Phi \mathbf{w}\|^2 \rangle + d \right) \beta + \text{const} \\ &= \tilde{c} \ln \beta - \tilde{d} \beta + \text{const} \\ q(\beta) &= \text{Gamma}(\beta | \tilde{c}, \tilde{d}) \quad \tilde{c} = c + N/2 \quad \tilde{d} = d + \frac{1}{2} \langle \|\mathbf{t} - \Phi \mathbf{w}\|^2 \rangle \end{aligned}$$

### Sparse Bayesian Linear Regression: Variational Inference

- Finding the required expectations

$$q(\beta) = \text{Gamma}(\beta | \tilde{c}, \tilde{d}), \quad \langle \beta \rangle = \tilde{c} / \tilde{d}$$

$$q(\alpha_m) = \text{Gamma}(\alpha_m | \tilde{a}, \tilde{b}_m), \quad \langle \alpha_m \rangle = \tilde{a} / \tilde{b}_m$$

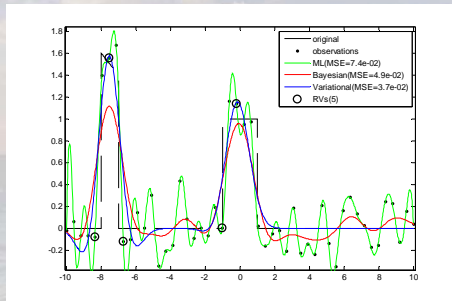
$$q(\mathbf{w}) = N(\mathbf{w} | \boldsymbol{\mu}, \Sigma), \quad \langle w_m^2 \rangle = \boldsymbol{\mu}_m + \Sigma_{mm}$$

$$\langle \|\mathbf{t} - \Phi \mathbf{w}\|^2 \rangle = \|\mathbf{t}\|^2 + \text{tr}\{\Phi \Sigma \Phi^T\} + \boldsymbol{\mu}^T \Phi^T \Phi \boldsymbol{\mu} - 2\mathbf{t}^T \Phi \boldsymbol{\mu}$$

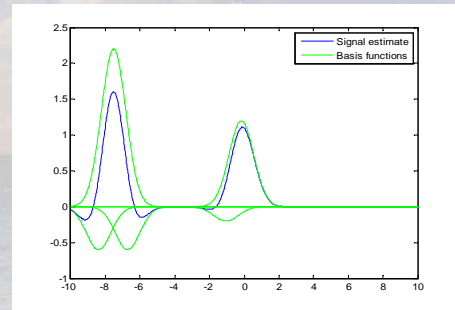
### Sparse Bayesian Linear Regression: Parameter Estimation

- Parameters  $a, b, c, d = ?$ 
  - Use fixed values that define **uninformative** priors
  - Estimate parameters in Variational M-step
- For fixed  $a, b, c, d$  iterate only between VE-step for  $q(\boldsymbol{\alpha})$ ,  $q(\beta)$  and  $q(\mathbf{w})$

## Linear Regression Example



## Linear Regression Example

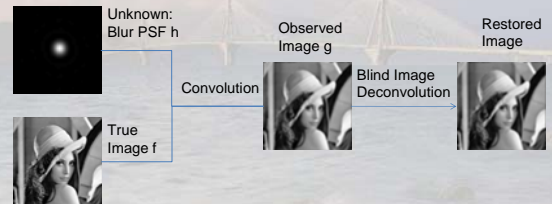


## Examples Blind Image Deconvolution

## Blind Image Deconvolution

$$\mathbf{g} = \mathbf{h} * \mathbf{f} + \mathbf{n}$$

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \mathbf{n} = \mathbf{F}\mathbf{h} + \mathbf{n}$$



## Blind Image Deconvolution

- Unknown quantities (f,h) twice than known (g)
- Properties that model should impose:
  - PSF: smooth, limited support
  - Image: smooth, preserve edges
  - Noise: robustness

## Blind Image Deconvolution: Noise Model

- Robust noise model
  - Student's t pdf

$$p(\mathbf{n} | \boldsymbol{\beta}) = \prod_{i=1}^N \mathcal{N}(n_i | 0, \beta_i^{-1}) \quad p(\boldsymbol{\beta}) = \prod_{i=1}^N \text{Gamma}(\beta_i | a^\beta, b^\beta)$$

$$p(n_i) = \int p(n_i | \beta_i) p(\beta_i) d\beta_i = \text{Student's t}$$

## Blind Image Deconvolution: PSF Model

- PSF: Sparse Linear model
- Basis functions are Gaussian kernels

$$h(x) = \sum_{i=1}^N w_i \phi_i(x) \quad \mathbf{h} = \Phi \mathbf{w}$$

$$\phi_i(x) = K(x, x_i) \quad \Phi_{i,j} = K(x_i, x_j)$$

- Sparseness weight prior

$$p(\mathbf{w} | \boldsymbol{\alpha}) = \prod_{i=1}^N N(w_i | 0, \alpha_i^{-1}), p(\alpha_i) = \text{Gamma}(\alpha_i; a^\alpha, b^\alpha)$$

$$p(\mathbf{w}) = \int p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) d\boldsymbol{\alpha} = \text{Student's t}$$

## Blind Image Deconvolution: Image Model

- Directional image differences

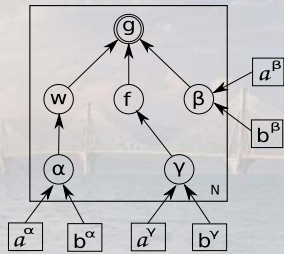
$$\varepsilon^1(x, y) = f(x, y) - f(x+1, y)$$

$$\varepsilon^2(x, y) = f(x, y) - f(x, y+1)$$

$$p(\boldsymbol{\varepsilon}^k | \boldsymbol{\gamma}^k) = \prod_{i=1}^N N(\varepsilon_i^k | 0, (\gamma_i^k)^{-1}) \quad p(\boldsymbol{\gamma}^k) = \prod_{i=1}^N \text{Gamma}(\gamma_i^k | a^\gamma, b^\gamma)$$

$$p(\mathbf{f} | \tilde{\boldsymbol{\gamma}}) = N(\mathbf{f} | 0, (\tilde{\mathbf{Q}}^T \tilde{\boldsymbol{\Gamma}} \tilde{\mathbf{Q}})^{-1})$$

## Blind Image Deconvolution: Graphical Model



- Joint pdf

$$p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}; \boldsymbol{\theta}) = p(\mathbf{g} | \mathbf{f}, \mathbf{w}, \boldsymbol{\beta}) p(\mathbf{f} | \boldsymbol{\gamma}) p(\mathbf{w} | \boldsymbol{\alpha}) p(\boldsymbol{\beta}) p(\boldsymbol{\gamma}) p(\boldsymbol{\alpha})$$

## Blind Image Deconvolution: Variational Bound

- Mean field approximation

$$q(\mathbf{f}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}) = q(\mathbf{f}) q(\mathbf{w}) q(\boldsymbol{\alpha}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})$$

- Maximization of Variational Bound results in:

$$\log q(\mathbf{f}) = \left\langle \log p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}; \boldsymbol{\theta}) \right\rangle_{q(\mathbf{w}) q(\boldsymbol{\alpha}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})}$$

$$\log q(\mathbf{w}) = \left\langle \log p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}; \boldsymbol{\theta}) \right\rangle_{q(\mathbf{f}) q(\boldsymbol{\alpha}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})}$$

$$\log q(\tilde{\boldsymbol{\alpha}}) = \left\langle \log p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}; \boldsymbol{\theta}) \right\rangle_{q(\mathbf{w}) q(\mathbf{f}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})}$$

$$\log q(\tilde{\boldsymbol{\beta}}) = \left\langle \log p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}; \boldsymbol{\theta}) \right\rangle_{q(\mathbf{w}) q(\boldsymbol{\alpha}) q(\mathbf{f}) q(\boldsymbol{\gamma})}$$

$$\log q(\tilde{\boldsymbol{\gamma}}) = \left\langle \log p(\mathbf{g}, \mathbf{f}, \mathbf{w}, \tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}}, \tilde{\boldsymbol{\gamma}}; \boldsymbol{\theta}) \right\rangle_{q(\mathbf{w}) q(\boldsymbol{\alpha}) q(\mathbf{f}) q(\boldsymbol{\beta})}$$

## Blind Image Deconvolution: Find approximate posterior $q(\mathbf{f})$

- Mean Field Optimization

$$\log q(\mathbf{f}) = \left\langle -\frac{1}{2} (\Phi \mathbf{w} \mathbf{f} - \mathbf{g})^T \mathbf{B} (\Phi \mathbf{w} \mathbf{f} - \mathbf{g}) - \frac{1}{2} \mathbf{f}^T \tilde{\mathbf{Q}}^T \tilde{\boldsymbol{\Gamma}} \tilde{\mathbf{Q}} \mathbf{f} \right\rangle_{q(\mathbf{w}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})} + \text{const}$$

$$= \left\langle -\frac{1}{2} \mathbf{f}^T (\Phi^T \mathbf{w}^T \mathbf{B} \Phi + \tilde{\mathbf{Q}}^T \tilde{\boldsymbol{\Gamma}} \tilde{\mathbf{Q}}) \mathbf{f} + \mathbf{f}^T \Phi^T \mathbf{w}^T \mathbf{B} \mathbf{g} \right\rangle_{q(\mathbf{w}) q(\boldsymbol{\beta}) q(\boldsymbol{\gamma})} + \text{const}$$

$$= -\frac{1}{2} \mathbf{f}^T (\Phi^T \langle \mathbf{w}^T \mathbf{B} \mathbf{w} \rangle \Phi + \tilde{\mathbf{Q}}^T \langle \tilde{\boldsymbol{\Gamma}} \rangle \tilde{\mathbf{Q}}) \mathbf{f} + \mathbf{f}^T \Phi^T \langle \mathbf{w}^T \mathbf{B} \rangle \mathbf{g}$$

- Completing Square results in

$$q(\mathbf{f}) = N(\mathbf{f} | \boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f)$$

$$\boldsymbol{\mu}_f = \boldsymbol{\Sigma}_f \Phi^T \langle \mathbf{w}^T \mathbf{B} \rangle \mathbf{g}$$

$$\boldsymbol{\Sigma}_f = (\Phi^T \langle \mathbf{w}^T \mathbf{B} \mathbf{w} \rangle \Phi + \tilde{\mathbf{Q}}^T \langle \tilde{\boldsymbol{\Gamma}} \rangle \tilde{\mathbf{Q}})^{-1}$$

## Blind Image Deconvolution: Approximate posterior $q(\mathbf{w})$

- Similar calculations

$$\log q(\mathbf{w}) = \left\langle -\frac{1}{2} (\mathbf{F} \Phi \mathbf{w} - \mathbf{g})^T \mathbf{B} (\mathbf{F} \Phi \mathbf{w} - \mathbf{g}) - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} \right\rangle_{q(\mathbf{f}) q(\boldsymbol{\beta}) q(\boldsymbol{\alpha})} + \text{const}$$

$$= -\frac{1}{2} \mathbf{w}^T (\Phi^T \langle \mathbf{F}^T \mathbf{B} \mathbf{F} \rangle \Phi + \langle \mathbf{A} \rangle) \mathbf{w} + \mathbf{w}^T \Phi^T \langle \mathbf{F}^T \mathbf{B} \rangle \mathbf{g} + \text{const}$$

$$q(\mathbf{w}) = N(\mathbf{w} | \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$$

$$\boldsymbol{\mu}_w = \boldsymbol{\Sigma}_w \Phi^T \langle \mathbf{F}^T \mathbf{B} \rangle \mathbf{g}$$

$$\boldsymbol{\Sigma}_w = (\Phi^T \langle \mathbf{F}^T \mathbf{B} \mathbf{F} \rangle \Phi + \langle \mathbf{A} \rangle)^{-1}$$



## Blind Image Deconvolution: Approximate posterior $q(\mathbf{a})$

- Based on the mean field approximation

$$\begin{aligned} \log q(\mathbf{a}) &= \langle \log p(\mathbf{w} | \mathbf{a}) + \log p(\mathbf{a}) \rangle_{q(\mathbf{w})} + \text{const} \\ &= \left\langle \frac{1}{2} \sum_{i=1}^N \log \alpha_i - \frac{1}{2} \sum_{i=1}^N \alpha_i w_i^2 + (a^\alpha - 1) \sum_{i=1}^N \log \alpha_i - b^\alpha \sum_{i=1}^N \alpha_i \right\rangle_{q(\mathbf{w})} + \text{const} \\ &= \sum_{i=1}^N (a^\alpha - 1 + \frac{1}{2}) \log \alpha_i - \sum_{i=1}^N (b^\alpha + \langle w_i^2 \rangle) \alpha_i + \text{const} \end{aligned}$$

- Which implies

$$q(\mathbf{a}) = \prod_{i=1}^N \text{Gamma}(\alpha_i | \tilde{a}^\alpha, \tilde{b}_i^\alpha) \quad \begin{aligned} \tilde{a}^\alpha &= a^\alpha + 1/2 \\ \tilde{b}_i^\alpha &= b^\alpha + \frac{1}{2} \langle w_i^2 \rangle \end{aligned}$$

## Blind Image Deconvolution: Approximate posteriors $q(\boldsymbol{\beta})$ and $q(\boldsymbol{\gamma})$

- Similarly we get

$$q(\boldsymbol{\beta}) = \prod_{i=1}^N \text{Gamma}(\beta_i | \tilde{a}^\beta, \tilde{b}_i^\beta) \quad \begin{aligned} \tilde{a}^\beta &= a^\beta + 1/2 \\ \tilde{b}_i^\beta &= b^\beta + \frac{1}{2} ((\mathbf{Q}^k)^T \langle \mathbf{f} \mathbf{f}^T \rangle \mathbf{Q}^k)_i \end{aligned}$$

$$q(\boldsymbol{\gamma}) = \prod_{k=1}^K \prod_{i=1}^N \text{Gamma}(\gamma_i^k | \tilde{a}^\gamma, \tilde{b}_i^{\gamma^k}) \quad \begin{aligned} \tilde{a}^\beta &= a^\beta + 1/2 \\ \tilde{b}_i^\beta &= b^\beta + \frac{1}{2} \langle \mathbf{n} \mathbf{n}^T \rangle_{ii} \\ \mathbf{n} &= \mathbf{F} \Phi \mathbf{w} - \mathbf{g} \end{aligned}$$

## Blind Image Deconvolution: Statistics of Approximate Posteriors

$$\begin{aligned} \langle w \rangle &= \mu_w, \\ \langle w_i^2 \rangle &= \mu_w^2 + \Sigma_{w_{ii}}, \\ \langle f \rangle &= \mu_f, \\ \langle f f^T \rangle &= \mu_f \mu_f^T + \Sigma_f, \\ \langle \alpha_i \rangle &= \tilde{a}^\alpha / \tilde{b}_i^\alpha, \\ \langle \beta_i \rangle &= \tilde{a}^\beta / \tilde{b}_i^\beta, \\ \langle \gamma_i^k \rangle &= \tilde{a}^\gamma / \tilde{b}_i^{\gamma^k}, \\ \langle \mathbf{n} \mathbf{n}^T \rangle &= \mathbf{g} \mathbf{g}^T - 2\Phi(\mathbf{F} \mathbf{w}) \mathbf{g}^T + \Phi(\mathbf{F} \mathbf{w} \mathbf{w}^T \mathbf{F}^T) \Phi^T. \end{aligned}$$

## Blind Image Deconvolution: Statistics of Approximate Posteriors

- Make diagonal and circulant approximations

$$\begin{aligned} \Sigma_w &= (\text{diag}\{\Phi^T \langle \mathbf{F}^T \mathbf{B} \mathbf{F} \rangle \Phi\} + \langle \mathbf{A} \rangle)^{-1}, \\ \tilde{\Sigma}_f &= (\langle \tilde{\beta} \rangle \Phi^T \langle \mathbf{W}^T \mathbf{W} \rangle \Phi + \langle \tilde{\gamma} \rangle \tilde{\mathbf{Q}}^T \tilde{\mathbf{Q}})^{-1}, \\ \tilde{\gamma} &= \frac{1}{MK} \sum_{k=1}^M \sum_{i=1}^M \tilde{\gamma}_i^k, \quad \tilde{\beta} = \frac{1}{M} \sum_{i=1}^M \tilde{\beta}_i \end{aligned}$$

- This results in

$$\begin{aligned} \langle \mathbf{W}^T \mathbf{W} \rangle &= \langle \mathbf{W}^T \rangle \langle \mathbf{W} \rangle + I \sum_{i=1}^M \langle \Sigma_{w_{ii}} \rangle, \\ \langle \mathbf{F}^T \mathbf{B} \mathbf{F} \rangle &= \langle \mathbf{F}^T \rangle \langle \mathbf{B} \rangle \langle \mathbf{F} \rangle + \tilde{\Sigma}_f \sum_{i=1}^M \langle \tilde{\beta}_i \rangle. \end{aligned}$$

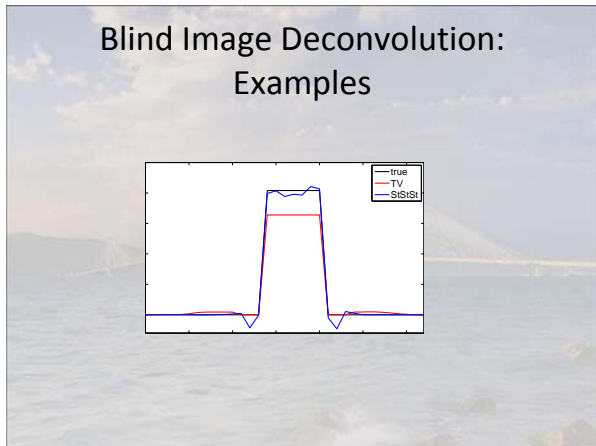
## Blind Image Deconvolution: Overall V-EM Algorithm

- Fix parameters  $\{a^\alpha, b^\alpha, a^\beta, b^\beta, a^\gamma, b^\gamma\}$  to yield **uninformative hyperpriors** (no Variational M-step).
- Iterate between estimates of statistics of  $q(\mathbf{f})$ ,  $q(\mathbf{w})$ ,  $q(\mathbf{a})$ ,  $q(\boldsymbol{\beta})$  and  $q(\boldsymbol{\gamma})$  (only Variational E-step).

## Blind Image Deconvolution: Example

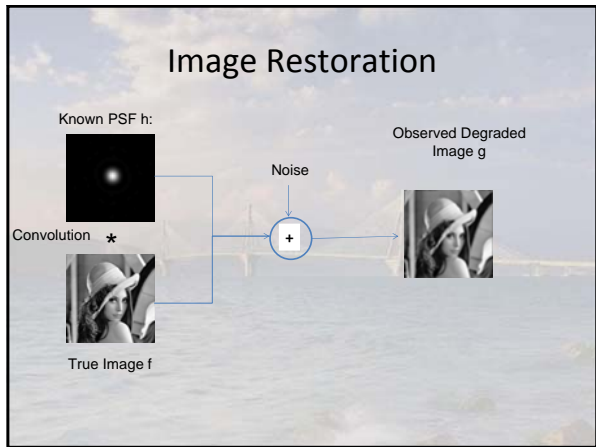
- PSF: Square 7x7, 40db noise





### Example: Image Restoration (Constrained Variational Inference\*)

\*G. Chantas, N. Galatsanos, A. Likas and M. Saunders, "Variational Bayesian Image Restoration Based on a Product of T-Distributions Image Prior", *IEEE Trans. on Image Processing*, to appear. (available on line)



### Image Restoration: Problem Definition

- Imaging Model ( $N$ -pixels)
 
$$\mathbf{g} = \mathbf{h} * \mathbf{f} + \mathbf{n} = \mathbf{H}\mathbf{f} + \mathbf{n}$$
- $\mathbf{g}: (N \times 1)$  : Degraded (*Observations*)
- $\mathbf{h}: (N \times 1)$ ,  $\mathbf{H}: (N \times N)$  : Point Spread Function (*known*)
- $\mathbf{f}: (N \times 1)$  : Original Image (*unknown*)
- $\mathbf{n} \sim N(\mathbf{0}, \beta^{-1}\mathbf{I})$  : Noise (*unknown*)

### Image Restoration: $\mathbf{f}$ Parameter

- Likelihood of Observations
 
$$p(\mathbf{g}; \mathbf{f}, \beta) = N(\mathbf{H}\mathbf{f}, \beta^{-1}\mathbf{I})$$
- $$\hat{\mathbf{f}} = \arg \max \{p(\mathbf{g}; \mathbf{f}, \beta)\} = \arg \min \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2$$
- $$\hat{\mathbf{f}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{g}$$
- Problematic **too many** parameters

Graphical Model

Double circled observed r.v.

### Image Restoration: $\mathbf{f}$ Parameter-Example

$$SNR = 20 \log_{10} \frac{\|\mathbf{f}\|}{\|\mathbf{n}\|}$$

SNR~80dB      SNR~120dB

## Image Restoration: Bayesian Inference ( $\mathbf{f}$ : hidden r.v.),

- Error from Local Linear Predictor

$$f(i) - \frac{1}{2}(f(i-1) + f(i+1)) = \varepsilon(i)$$

- Assume Gaussian i.i.d. Prediction Errors

$$p(\varepsilon(i)) = N(0, \alpha^{-1}) \quad p(\boldsymbol{\varepsilon}) = \prod_{i=1}^N p(\varepsilon(i))$$

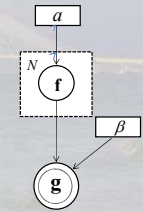
## Image Restoration: Bayesian Inference ( $\mathbf{f}$ : hidden r.v.)

- Simultaneously Autoregressive (SAR) Prior

Graphical Model

$\mathbf{Q}\mathbf{f} = \boldsymbol{\varepsilon}$ ,  $\mathbf{Q}$  Operator that describes  $\boldsymbol{\varepsilon} = \mathbf{f} - \hat{\mathbf{f}}$

$$p(\mathbf{f}; a) \propto \alpha^{\frac{N-1}{2}} \exp\left(-\alpha \frac{\mathbf{f}^T \mathbf{Q}^T \mathbf{Q} \mathbf{f}}{2}\right)$$



Double circled observed

## Image Restoration: Bayesian Inference ( $\mathbf{f}$ : hidden r.v.)

- Observation Likelihood (*computed analytically*)

$$p(\mathbf{g}; \beta, a) = \int p(\mathbf{g} | \mathbf{f}; \beta) p(\mathbf{f}; a) d\mathbf{f}$$

- Posterior of Hidden (*computed analytically*)

$$p(\mathbf{f} | \mathbf{g}; a, \beta) = N(\boldsymbol{\mu}_{f|g}, \boldsymbol{\Sigma}_{f|g})$$

- Bayesian Inference via **EM**

## Image Restoration: Bayesian Inference ( $\mathbf{f}$ hidden r.v.)

“Small”  $a/\beta$   
Amplifies Noise



“Large”  $a/\beta$   
Smooths Edges



## Image Restoration: Spatially-Varying Bayesian Model

- **Spatially Varying**  $p(\varepsilon(i)) = N(0, a_i^{-1})$
- $\alpha_i$  “hidden variables” Bayesian Inference
- Conjugate pdf  $p(\alpha_i) = \text{Gamma}(a_i; \alpha, \beta)$

## Image Restoration: Spatially-Varying, Bayesian Inference

- Use **many**  $\boldsymbol{\varepsilon}_k = \mathbf{Q}_k \mathbf{f}$ ,  $k = 1, 2, \dots, P$  in prior

- “Product” Prior

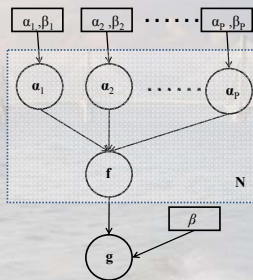
$$p(\mathbf{f} | \tilde{\mathbf{a}}) = \frac{1}{Z(\tilde{\mathbf{a}})} \prod_{k=1}^P p_k(\mathbf{f} | \mathbf{a}_k), \quad p(\mathbf{f} | \mathbf{a}_k) \propto \exp\left(-\frac{1}{2} \mathbf{f}^T \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k \mathbf{f}\right)$$

$$\mathbf{A}_k = \text{diag}(a_k(1), a_k(2), \dots, a_k(N))$$

- Prior: Enforces many properties **simultaneously**

## Image Restoration: Spatially Varying Bayesian Inference

Graphical Model



## Image Restoration: Spatially Varying, Bayesian Inference

- Difficulty: Compute normalization of prior

$$Z(\tilde{\mathbf{a}}) \propto \det \left\{ \sum_{k=1}^K \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k \right\} \mathbf{Q}_k, \mathbf{A}_k (N \times N), N=10^5-10^6$$

## Image Restoration: Spatially-Varying, Bayesian Inference

- Change Observations Domain

$$\mathbf{Q}_k \mathbf{g} = \mathbf{Q}_k \mathbf{H} \mathbf{f} + \mathbf{Q}_k \mathbf{n},$$

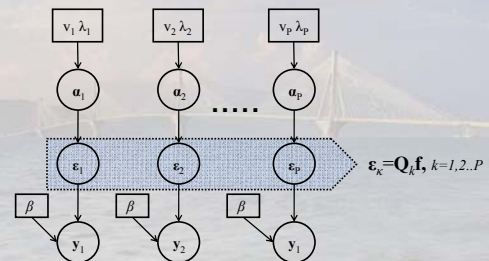
$$\mathbf{y}_k = \mathbf{H} \boldsymbol{\varepsilon}_k + \mathbf{n}_k, k = 1, \dots, P$$

- Hidden Variables:

$$\tilde{\boldsymbol{\varepsilon}} = [\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2, \dots, \boldsymbol{\varepsilon}_p], \tilde{\mathbf{a}} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p]$$

$$\boldsymbol{\varepsilon}_i = [\varepsilon_i(1), \varepsilon_i(2), \dots, \varepsilon_i(N)], \mathbf{a}_i = [a_i(1), a_i(2), \dots, a_i(N)]$$

## Image Restoration: Spatially-varying, Bayesian Inference



## Image Restoration: Spatially Varying Bayesian Inference

- Prior on  $\tilde{\boldsymbol{\varepsilon}}$  **NO difficulty** normalizing

$$p(\tilde{\boldsymbol{\varepsilon}} | \tilde{\mathbf{a}}) = \prod_{k=1}^P \prod_{i=1}^N p(\varepsilon_k(i) | a_k(i)),$$

$$p(\varepsilon_k(i) | a_k(i); \lambda_k) = N(\varepsilon_k(i); 0, \lambda_k a_k(i)^{-1})$$

$$p(a_k(i); v_k) = \text{Gamma}\left(a_k(i); \frac{v_k}{2}, \frac{v_k}{2}\right)$$

- Cannot marginalize hidden variables
- Resort to Variational Methodology

## Image Restoration: Spatially-Varying, Bayesian Inference

- "Posteriors" Mean Field Approximation

$$q(\boldsymbol{\varepsilon}_k, \mathbf{a}_k) = q(\boldsymbol{\varepsilon}_k) q(\mathbf{a}_k), k = 1, \dots, P$$

- Maximize Var. Bound

$$q(\boldsymbol{\varepsilon}_k) = N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

$$\boldsymbol{\mu}_k = \beta \boldsymbol{\Sigma}_k \mathbf{H}^T \mathbf{g},$$

$$\boldsymbol{\Sigma}_k = (\beta \mathbf{H}^T \mathbf{H} + \lambda_k \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k)^{-1}$$

- New problem: **Different**  $\mathbf{f}$  for each  $\boldsymbol{\mu}_k$



## Image Restoration: Constrained Variational Inference

- Define “**Constrained Posterior**”

$$q(\boldsymbol{\varepsilon}_k) = N(\mathbf{Q}_k \mathbf{m}, \mathbf{Q}_k \mathbf{R} \mathbf{Q}_k^T)$$

- Consistent with

$$\boldsymbol{\varepsilon}_k = \mathbf{Q}_k \mathbf{f}, \mathbf{m} = E(\mathbf{f}), \mathbf{R} = E((\mathbf{f} - \mathbf{m})(\mathbf{f} - \mathbf{m})^T)$$

- Maximize the **Variational Bound** w.r.t. **m** and **R**

## Image Restoration: Constrained Variational Inference

- VE-step:  $[q^{(t+1)}(\tilde{\mathbf{a}}), \theta_1^{(t+1)}] = \arg \max_{[q(\tilde{\mathbf{a}}), \theta_1]} (F(q(\tilde{\mathbf{a}}), \theta_1, \theta_2^{(t)}))$
- VM-step:  $\theta_2^{(t+1)} = \arg \max_{\theta_2} (F(q^{(t+1)}(\tilde{\mathbf{a}}), \theta_1^{(t+1)}, \theta_2))$

$$\theta_1 = [\mathbf{R}, \mathbf{m}], \quad \theta_2 = [\beta, \lambda_1, \dots, \lambda_p, \nu_1, \dots, \nu_p]^T$$

## Image Restoration: Constrained Variational Inference

- **VE-Step**

$$F(q, \theta) = \int \prod_{k=1}^p q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k) \log p(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\varepsilon}}, \tilde{\mathbf{a}}; \theta_2) d\tilde{\boldsymbol{\varepsilon}} d\tilde{\mathbf{a}} - \int \prod_{k=1}^p q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k) \log \prod_{k=1}^p q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k) d\tilde{\boldsymbol{\varepsilon}} d\tilde{\mathbf{a}},$$

$$F \propto \tilde{F}(\theta_1)$$

$$\tilde{F}(\theta_1) = \sum_{k=1}^p \int q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k) \log p(\mathbf{y}_k | \boldsymbol{\varepsilon}_k; \theta_2) p(\boldsymbol{\varepsilon}_k | \mathbf{a}_k; \theta_1) d\boldsymbol{\varepsilon}_k d\mathbf{a}_k - \sum_{k=1}^p \int q(\boldsymbol{\varepsilon}_k; \theta_1) \log q(\boldsymbol{\varepsilon}_k; \theta_1) d\boldsymbol{\varepsilon}_k.$$

## Image Restoration: Constrained Variational Inference

$$\begin{aligned} & \sum_{k=1}^p \int q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k) \log p(\mathbf{y}_k | \boldsymbol{\varepsilon}_k; \theta_2) p(\boldsymbol{\varepsilon}_k | \mathbf{a}_k; \theta_1) d\boldsymbol{\varepsilon}_k d\mathbf{a}_k \\ & \propto \sum_{k=1}^p \left\langle -\beta (\mathbf{H}\boldsymbol{\varepsilon}_k - \mathbf{y}_k)^T \mathbf{Q}_k^{-T} \mathbf{Q}_k^{-1} (\mathbf{H}\boldsymbol{\varepsilon}_k - \mathbf{y}_k) - \lambda_k \boldsymbol{\varepsilon}_k^T \mathbf{A}_k \boldsymbol{\varepsilon}_k \right\rangle_{q(\boldsymbol{\varepsilon}_k; \theta_1) q(\mathbf{a}_k)} = \\ & -\beta P \|\mathbf{H}\mathbf{m} - \mathbf{g}\|_2^2 - \sum_{k=1}^p \lambda_k \mathbf{m}^T \mathbf{Q}_k^T (\mathbf{A}_k) \mathbf{Q}_k \mathbf{m} - \text{trace} \left\{ \left( \beta P \mathbf{H}^T \mathbf{H} + \sum_{k=1}^p \lambda_k \mathbf{Q}_k^T (\mathbf{A}_k) \mathbf{Q}_k \right) \mathbf{R} \right\} \\ & \int q(\boldsymbol{\varepsilon}_k; \theta_1) \log q(\boldsymbol{\varepsilon}_k; \theta_1) d\boldsymbol{\varepsilon}_k \propto \frac{1}{2} \log \det |\mathbf{R}| \end{aligned}$$

## Image Restoration: Constrained Variational Inference

$$\frac{\partial \tilde{F}(\theta_1)}{\partial \mathbf{R}} = 0 \Rightarrow \frac{\partial \text{trace} \left\{ \beta P \mathbf{H}^T \mathbf{H} \mathbf{R} + \sum_{k=1}^p \lambda_k \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k \mathbf{R} \right\} - P \partial \log \det |\mathbf{R}|}{\partial \mathbf{R}} = 0$$

$$\Rightarrow \beta P \mathbf{H}^T \mathbf{H} + \sum_{k=1}^p \lambda_k \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k - P \mathbf{R}^{-1} = 0 \Rightarrow \mathbf{R} = \left( \beta P \mathbf{H}^T \mathbf{H} + \sum_{k=1}^p \lambda_k \mathbf{Q}_k^T \mathbf{A}_k \mathbf{Q}_k \right)^{-1}$$

$$\frac{\partial \tilde{F}(\theta_1)}{\partial \mathbf{m}} = 0 \Rightarrow \mathbf{m} = \beta \mathbf{R} \mathbf{H}^T \mathbf{g}$$

## Image Restoration: Constrained Variational Inference

$$\begin{aligned} q(\tilde{\mathbf{a}}) &= \exp(\log p(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\varepsilon}}, \tilde{\mathbf{a}}))_{q(\boldsymbol{\varepsilon})} \propto \\ & \prod_{k=1}^p \prod_{i=1}^N (a_k(i))^{\frac{\nu_k}{2} - 1} \exp \left\{ -\frac{\nu_k}{2} a_k(i) - \frac{1}{2} \lambda_k \left( (\mathbf{m}_k(i))^2 + C_k(i, i) \right) a_k(i) \right\} \\ q(a_k(i)) &= \text{Gamma} \left( a_k(i); \frac{\nu_k}{2} + \frac{1}{2}, \frac{\nu_k}{2} + \frac{1}{2} \lambda_k \left( (\mathbf{m}_k(i))^2 + C_k(i, i) \right) \right) \\ \mathbf{m}_k &= \mathbf{Q}_k \mathbf{m}, \text{ and } C_k = \mathbf{Q}_k \mathbf{R} \mathbf{Q}_k^T \end{aligned}$$

## Image Restoration: Constrained Variational Inference

### •VM-step

$$\frac{dF(q^{(t+1)}(\bar{\mathbf{a}}), \theta_1^{(t+1)}, \theta_2)}{d\theta_2} = \frac{d \langle \log p(\tilde{\mathbf{y}}, \tilde{\mathbf{z}}, \bar{\mathbf{a}}, \theta_2) \rangle_{q(\tilde{\mathbf{z}}, \theta_1^{(t+1)})}^{q^{(t+1)}(\bar{\mathbf{a}})} }{d\theta_2} \Big|_0$$

⇒

$$\beta^{(t+1)} = N \left( \left\| \mathbf{H}\mathbf{m}^{(t+1)} - \mathbf{g} \right\|_2^2 + \text{trace} \left\{ \mathbf{H}^T \mathbf{H} \mathbf{R} \right\} \right)^{-1}$$

$$\lambda_k^{(t+1)} = N \left( \sum_{i=1}^N \left( \left( \mathbf{m}_k^{(t+1)}(i) \right)^2 + \mathbf{C}_k^{(t+1)}(i, i) \right) \langle a_k(i) \rangle_{q^{(t+1)}(a_k(i))} \right)^{-1}$$

$$\frac{1}{N} \left( \sum_{i=1}^N \log \langle a_k(i) \rangle_{q^{(t+1)}(a_k(i))} - \sum_{i=1}^N \langle a_k(i) \rangle_{q^{(t+1)}(a_k(i))} \right) + \psi \left( \nu_k^{(t+1)} \frac{1}{2} + \frac{1}{2} \right) -$$

$$\log \left( \nu_k^{(t+1)} \frac{1}{2} + \frac{1}{2} \right) - \psi \left( \frac{\nu_k}{2} \right) + \log \left( \frac{\nu_k}{2} \right) + 1 = 0 \quad \psi(x) = \frac{d}{dx} \log \Gamma(x)$$

## Image Restoration: Constrained Variational Inference

1. Initialize,  $\mathbf{m}$  stationary model.
2. Repeat until convergence:

### VE-step:

Update,  $\mathbf{m}$  and  $\mathbf{R}$ , calculate  $\mathbf{m}_k$  and  $\mathbf{C}_k$ . Calculate expected value w.r.t.  $q(a_k(i))$ , needed for VM-step and the next VE-step.

### VM-step:

Update  $\beta, \lambda_k, \nu_k, k=1, 2, \dots, P$

3. Use  $\mathbf{m}$  as restored image estimate.

## Image Restoration: Constrained Variational Inference



## Example: Image Restoration (Bounded Variational Inference\*)

\*Babacan, S.D.; Molina, R.; Katsaggelos, A.K.; "Parameter Estimation in TV Image Restoration Using Variational Distribution Approximation", IEEE Trans. on Image Processing, Volume 17, Issue 3, March 2008 Page(s):326 - 339

## Image Restoration: Bounded Variational Inference

- Imaging Model ( $N$ -pixels)

$$\mathbf{g} = \mathbf{h} * \mathbf{f} + \mathbf{n} = \mathbf{H}\mathbf{f} + \mathbf{n}$$

- Observations Likelihood

$$p(\mathbf{g} / \mathbf{f}; \beta) \sim N(\mathbf{H}\mathbf{f}, \beta^{-1}\mathbf{I})$$

## Image Restoration: Bounded Variational Inference

- TV Based Image Prior

$$p(\mathbf{f} | a) = \alpha^{\frac{N^2}{2}} \exp(-\alpha \text{TV}(\mathbf{f}))$$

$$\text{TV}(\mathbf{f}) = \sum_{i,j=1}^N \sqrt{([\mathbf{Q}\mathbf{f}]_i)^2 + ([\mathbf{Q}\mathbf{f}]_j)^2}$$

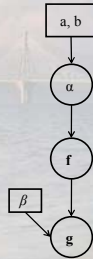
$$[\mathbf{Q}\mathbf{f}]_i = \mathbf{f}(i, j) - \mathbf{f}(i, j-1), \quad [\mathbf{Q}\mathbf{f}]_j = \mathbf{f}(i, j) - \mathbf{f}(i-1, j)$$

- Conjugate Hyperprior

$$p(\alpha; a, b) = \text{Gamma}(\alpha; a, b)$$

## Image Restoration: Bounded Variational Inference

Graphical Model



## Image Restoration: Bounded Variational Inference

- Difficulty in VE-step

$$\log q(\alpha) = \langle \log p(\mathbf{g}, \mathbf{f}, \alpha; \beta, a, b) \rangle_{q(\mathbf{f})} = \left\langle -\frac{1}{2} \beta \|H\mathbf{f} - \mathbf{g}\|^2 - \alpha \sum_{i=1}^N \sqrt{([\mathbf{Q}_1 \mathbf{f}]_i)^2 + ([\mathbf{Q}_2 \mathbf{f}]_i)^2} \right\rangle_{q(\mathbf{f})} + \text{const}$$

- Due to  $\sqrt{\cdot}$  cannot Compute Expectation

## Image Restoration: Bounded Variational Inference

- Bypass difficulty: Maximize a Lower Bound of Variational Bound

- Use Upper Bound:  $f(w) = \sqrt{w} \leq \frac{w+u}{2\sqrt{u}} = g(u, w), \forall u > 0$

- Bound gets "tight":  $f(w) = g(u^*, w), u^* = w$

## Image Restoration: Bounded Variational Inference

- Define Function

$$p(\mathbf{f} | \alpha) \geq M(\mathbf{f}, \alpha, \mathbf{u}) = \alpha^{N/2} \exp \left\{ -\alpha \sum_{i=1}^N \frac{([\mathbf{Q}_1 \mathbf{f}]_i)^2 + ([\mathbf{Q}_2 \mathbf{f}]_i)^2 + u_i}{2\sqrt{u_i}} \right\}$$

$$p(\mathbf{g}, \mathbf{f}, \alpha; \theta) \geq p(\mathbf{g} | \mathbf{f}) M(\mathbf{f}, \alpha, \mathbf{u}) p(\alpha) = \tilde{M}(\mathbf{g}, \mathbf{f}, \mathbf{u}, \alpha; \theta)$$

- Lower Bound of Variational Bound

$$F(q, \theta) \geq \int q(\mathbf{f}, \alpha) \ln \frac{\tilde{M}(\mathbf{g}, \mathbf{f}, \mathbf{u}, \alpha; \theta)}{q(\mathbf{f}, \alpha)} d\mathbf{f} d\alpha = F^b(q, \mathbf{u}, \theta)$$

## Image Restoration: Bounded Variational Inference

- VE-Step

$$[q^{(t+1)}(\mathbf{f}), q^{(t+1)}(\alpha)] = \arg \max_{q(\mathbf{f}), q(\alpha)} [F^b(q(\mathbf{f}), q(\alpha), \mathbf{u}^{(t)}, \theta^{(t)})]$$

- VM-Step

"Bound Tightening"

$$\mathbf{u}^{(t+1)} = \arg \max_{\mathbf{u}} [F^b(q^{(t+1)}(\mathbf{f}), q^{(t+1)}(\alpha), \mathbf{u}, \theta^{(t)})]$$

$$\theta^{(t+1)} = \arg \max_{\theta} [F^b(q^{(t+1)}(\mathbf{f}), q^{(t+1)}(\alpha), \mathbf{u}^{(t+1)}, \theta)]$$

## Image Restoration: Bounded Variational Inference

- "Tightest" Bound

$$f(w) = g(u^*, w), u^* = w \Rightarrow$$

$$\mathbf{u}_i^{(t+1)} = \left\langle ([\mathbf{Q}_1 \mathbf{f}]_i)^2 + ([\mathbf{Q}_2 \mathbf{f}]_i)^2 \right\rangle_{q^{(t+1)}(\mathbf{f})}$$

$$= [\mathbf{Q}_1 \boldsymbol{\mu}^{(t+1)}]_i^2 + [\mathbf{Q}_2 \boldsymbol{\mu}^{(t+1)}]_i^2 + [\mathbf{Q}_1 \boldsymbol{\Sigma}^{(t+1)} \mathbf{Q}_1^T]_{ii} + [\mathbf{Q}_2 \boldsymbol{\Sigma}^{(t+1)} \mathbf{Q}_2^T]_{ii}, i = 1, 2, \dots, N$$

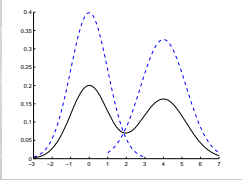
$$\text{with } q^{(t+1)}(\mathbf{f}) = N(\mathbf{f}; \boldsymbol{\mu}^{(t+1)}, \boldsymbol{\Sigma}^{(t+1)})$$

- $\mathbf{u}_i$  Captures Local Spatial Activity

## Example: Gaussian Mixture Models

### Gaussian Mixture Models

- $p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{j=1}^M \pi_j N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ 
  - Model any pdf
  - Soft clustering
- Parameters
 
$$\boldsymbol{\theta} = \{ \pi_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j \}_{j=1, \dots, M}$$
- Maximum Likelihood Estimation Difficult
 
$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \log \sum_{j=1}^M \pi_j N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$



### Gaussian Mixture Models: Data Generation Mechanism

- Introduce binary hidden variable  $\mathbf{z}$ 
  1. Select component  $\mathbf{z}$ 

$$z_j = 1, \mathbf{z} = (0, \dots, 0, \overset{1}{1}, 0, \dots, 0)$$

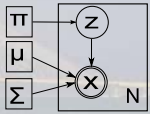
$$\mathbf{z} \sim p(\mathbf{z}; \boldsymbol{\pi}) = \prod_{j=1}^M \pi_j^{z_j} \text{ multinomial}$$
  2. Generate sample from selected component
 
$$x \sim p_j(\mathbf{x}) = N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

### Gaussian Mixture Models: Data Generation Mechanism

- Joint pdf
 
$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})$$

$$= \prod_{j=1}^M N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)^{z_j} (\pi_j)^{z_j}$$
- Marginal pdf
 
$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$$

$$= \sum_{j=1}^M \pi_j N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$



### Gaussian Mixture Models: Posterior

- Posterior (responsibility) can be computed analytically
 
$$p(z_j = 1 | \mathbf{x}) = \frac{p(\mathbf{x} | z_j = 1) p(z_j = 1)}{p(\mathbf{x})} = \frac{\pi_j N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{l=1}^M \pi_l N(\mathbf{x} | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}$$

$$p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}) = \prod_{j=1}^M \left[ \frac{\pi_j N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{l=1}^M \pi_l N(\mathbf{x} | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)} \right]^{z_j}$$

### Gaussian Mixture Models: Parameter Estimation

- Maximum likelihood
 
$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{X}; \boldsymbol{\theta})$$

$$= \arg \max_{\boldsymbol{\theta}} \log \sum_{j=1}^M \pi_j N(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$
- Use EM
  - Simplifies optimization
  - Proved convergence
  - Satisfies positivity constraints
 
$$\pi_j > 0, \sum_{j=1}^M \pi_j = 1$$



## Gaussian Mixture Models: Parameter Estimation EM

- $p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{j=1}^M N(\mathbf{x}_n; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)^{\tau_{jn}} (\pi_j)^{\tau_{jn}}$
- $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(l)}) = \langle \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) \rangle_{p(\mathbf{z}, \mathbf{x}; \boldsymbol{\theta}^{(l)})}$   
 $= \sum_{n=1}^N \sum_{j=1}^M \langle \tau_{jn}^{(l)} \rangle \log \pi_j + \sum_{n=1}^N \sum_{j=1}^M \langle \tau_{jn}^{(l)} \rangle \log N(\mathbf{x}_n; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$

## Gaussian Mixture Models: Parameter Estimation EM

### E-step

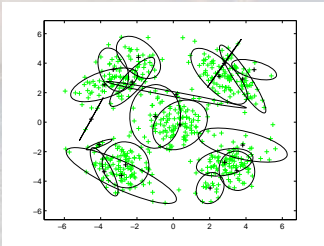
$$\langle \tau_{jn}^{(l)} \rangle = \sum_{j=1}^M \left[ \frac{\pi_j^{(l)} N(\mathbf{x}_n; \boldsymbol{\mu}_j^{(l)}, \boldsymbol{\Sigma}_j^{(l)})}{\sum_{j=1}^M \pi_j^{(l)} N(\mathbf{x}_n; \boldsymbol{\mu}_j^{(l)}, \boldsymbol{\Sigma}_j^{(l)})} \right]^{\tau_{jn}^{(l)}} = \frac{\pi_j^{(l)} N(\mathbf{x}_n; \boldsymbol{\mu}_j^{(l)}, \boldsymbol{\Sigma}_j^{(l)})}{\sum_{j=1}^M \pi_j^{(l)} N(\mathbf{x}_n; \boldsymbol{\mu}_j^{(l)}, \boldsymbol{\Sigma}_j^{(l)})}$$

### M-step

$$\pi_j^{(l+1)} = \frac{1}{N} \sum_{n=1}^N \langle \tau_{jn}^{(l)} \rangle \quad \boldsymbol{\mu}_j^{(l+1)} = \frac{\sum_{n=1}^N \langle \tau_{jn}^{(l)} \rangle \mathbf{x}_n}{\sum_{n=1}^N \langle \tau_{jn}^{(l)} \rangle} \quad \boldsymbol{\Sigma}_j^{(l+1)} = \frac{\sum_{n=1}^N \langle \tau_{jn}^{(l)} \rangle (\mathbf{x}_n - \boldsymbol{\mu}_j^{(l)}) (\mathbf{x}_n - \boldsymbol{\mu}_j^{(l)})^T}{\sum_{n=1}^N \langle \tau_{jn}^{(l)} \rangle}$$

## Gaussian Mixture Models: Limitations

- How many components?
- Ill-conditioned covariance matrices

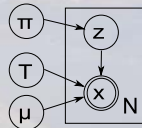


## Variational Bayesian Gaussian Mixture Models

- H. Attias, "A Variational Bayesian Framework for Graphical Models", *Proc. NIPS 12*, pp. 209-216, MIT Press, 2000.

## Variational Bayesian Gaussian Mixture Models: Prior Distribution

- Treat parameters as hidden variables  $\mathbf{h} = \{\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \mathbf{T}\}$
- Introduce conjugate priors



$$p(\boldsymbol{\pi} | \boldsymbol{\alpha}) = \text{Dir}(\boldsymbol{\pi} | \alpha_1, \dots, \alpha_M) = \frac{\Gamma(\sum_{j=1}^M \alpha_j)}{\prod_{j=1}^M \Gamma(\alpha_j)} \prod_{j=1}^M \pi_j^{\alpha_j - 1}$$

$$p(\boldsymbol{\mu}_j | \mathbf{T}_j) = N(\boldsymbol{\mu}_j; \boldsymbol{\mu}_0, \beta_0 \mathbf{T}_j)$$

$$p(\mathbf{T}_j | \nu, \mathbf{V}) = W(\mathbf{T}_j | \nu, \mathbf{V}) = \frac{|\mathbf{T}_j|^{(\nu-d-1)/2} \exp\{tr\{-\frac{1}{2}\mathbf{V}\mathbf{T}_j\}\}}{2^{\nu d/2} \pi^{d(d-1)/4} |\mathbf{V}|^{-\nu/2} \prod_{i=1}^d \Gamma((\nu+1-i)/2)}$$

## Variational Bayesian Gaussian Mixture Models: Mean-Field Approximation

- Exact Bayesian Inference Intractable  
– Variational Mean Field Approximation

$$q(\mathbf{h}) = q_z(\mathbf{Z}) q_\pi(\boldsymbol{\pi}) q_{\boldsymbol{\mu}}(\boldsymbol{\mu}, \mathbf{T})$$

$$q_{\boldsymbol{\mu}}(\boldsymbol{\mu}, \mathbf{T}) = \prod_{j=1}^M q_{\boldsymbol{\mu}}(\boldsymbol{\mu}_j | \mathbf{T}_j) q_{\mathbf{T}}(\mathbf{T}_j)$$

## Variational Bayesian Gaussian Mixture Models: Approximate Posteriors

$$\bullet q_z(\mathbf{Z}) = \prod_{n=1}^N \prod_{j=1}^M r_{jn}^{z_{jn}}$$

$$\log r_{jn} = \langle \log \pi_j \rangle + \langle \log |\mathbf{T}_j| \rangle - \frac{1}{2} (\mathbf{y}_n - \boldsymbol{\mu}_j)^T (\mathbf{y}_n - \boldsymbol{\mu}_j) - \frac{d}{2\beta_j} + \text{const}$$

$$\bullet q_\pi(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \{\lambda_j\})$$

$$\lambda_j = \bar{N}_j + \alpha \quad \bar{N}_j = \sum_{n=1}^N z_{jn}$$

## Variational Bayesian Gaussian Mixture Models: Approximate Posteriors

$$\bullet q_{\boldsymbol{\mu}_j}(\boldsymbol{\mu}_j | \mathbf{T}) = \prod_{j=1}^M N(\boldsymbol{\mu}_j; \mathbf{m}_j, \beta_j \mathbf{T}_j)$$

$$\mathbf{m}_j = (\bar{N}_j \bar{\boldsymbol{\mu}}_j + \beta^0 \boldsymbol{\mu}^0) / (\bar{N}_j + \beta^0)$$

$$\bar{\boldsymbol{\mu}}_j = \frac{1}{\bar{N}_j} \sum_{n=1}^N \mathbf{z}_{jn} \mathbf{y}_n$$

$$\beta_j = \bar{N}_j + \beta^0$$

$$\bullet q_{\mathbf{T}}(\mathbf{T}) = \prod_{j=1}^M W(\mathbf{T}_j; \boldsymbol{\eta}_j, \mathbf{U}_j)$$

$$\mathbf{U}_j = \bar{N}_j \boldsymbol{\Sigma}_j + \bar{N}_j \beta^0 (\boldsymbol{\mu}_j - \boldsymbol{\mu}^0)(\boldsymbol{\mu}_j - \boldsymbol{\mu}^0)^T / (\bar{N}_j + \beta^0) + \mathbf{T}^0$$

$$\boldsymbol{\eta}_j = \bar{N}_j + \nu$$

$$\bar{\boldsymbol{\Sigma}}_j = \frac{1}{\bar{N}_j} \sum_{n=1}^N \mathbf{z}_{jn} (\mathbf{y}_n - \bar{\boldsymbol{\mu}}_j)(\mathbf{y}_n - \bar{\boldsymbol{\mu}}_j)^T$$

## Variational Bayesian Gaussian Mixture Models: Discussion

- Select Parameters that define uninformative priors (e.g.  $\alpha_j = 1/M$ )
- Advantages
  - Disallow singular covariance matrices
  - Bayesian model selection
- Dirichlet distribution for mixing coefficients  $\pi$  disallows pruning of unnecessary components

## Variational Bayesian Gaussian Mixture Models: Removing the prior from the mixing weights

- A. Corduneanu and C. Bishop, "Variational Bayesian Model Selection for Mixture Distributions", *Proc. AI and Statistics Conference*, January 2001

## Variational Bayesian GMM: Removing the prior from the mixing weights

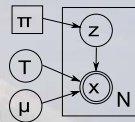
- Treat  $\pi$  as parameter
- Include M-step to update  $\pi$

$$\pi_j = \frac{\sum_{n=1}^N r_{jn}}{\sum_{k=1}^M \sum_{n=1}^N r_{kn}}$$

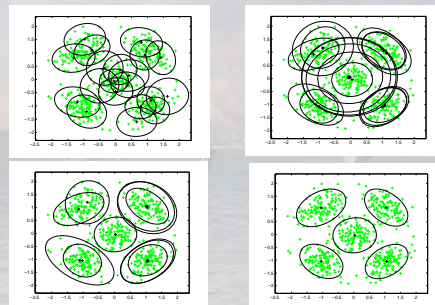
$$r_{jn} = \frac{\tilde{r}_{jn}}{\sum_{k=1}^M \tilde{r}_{kn}}$$

$$\tilde{r}_{jn} = \pi_j \exp \left\{ \frac{1}{2} \langle \log |\mathbf{T}_j| \rangle - \frac{1}{2} \text{tr} \left\{ \langle \mathbf{T}_j \rangle \left( \mathbf{x}_n \mathbf{x}_n^T - \mathbf{x}_n \langle \boldsymbol{\mu}_j \rangle^T + \langle \boldsymbol{\mu}_j \rangle \mathbf{x}_n^T + \langle \boldsymbol{\mu}_j \boldsymbol{\mu}_j^T \rangle \right) \right\} \right\}$$

- Advantage
  - Eliminates irrelevant components



## Variational Bayesian GMM: Example

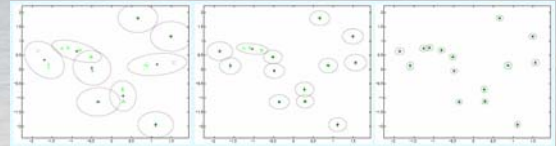


## Incremental Variational Bayesian Gaussian Mixture Models

Constantinopoulos C. and Likas, A., "Unsupervised Learning of Gaussian Mixtures Based on Variational Component Splitting", IEEE Trans. on Neural Networks, vol. 18, no. 3, pp. 745-755, 2007.

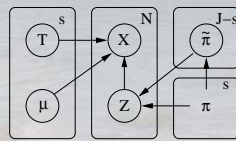
## Incremental Variational Bayesian Gaussian Mixture Models

- Solutions depend on:
  - maximum initial number of components
  - initialization of component parameters
  - specification of the scale matrix  $V$  of  $p(T_j) = \text{Wishart}(v, V)$



## Incremental Variational Bayesian Gaussian Mixture Models

- Divide components as 'fixed' and 'free'
- Restrict competition among 'free' components only
- $$p(\tilde{\pi} | \pi) = \left( \prod_{j=1}^s (1 - \pi_j) \right)^{-M^{k+1}} \frac{\Gamma(\sum_{j=s+1}^M \alpha_j)}{\prod_{j=s+1}^M \Gamma(\alpha_j)} \prod_{j=s+1}^M \left( \frac{\tilde{\pi}_j}{1 - \sum_{k=1}^s \pi_k} \right)^{\alpha_j - 1}$$



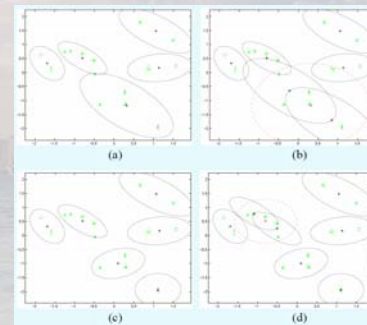
## Incremental Variational Bayesian Gaussian Mixture Models

- We start by training a GMM with two components
- At each step:
  - Select a component  $j$
  - Set  $V = d\lambda I$ , where  $\lambda$  the max eigenvalue of  $\langle T_j^{-1} \rangle$
  - Split the component in two subcomponents
  - Apply VB learning considering the two components as free

## Incremental Variational Bayesian Gaussian Mixture Models

- If the data in the region of component  $j$  suggest the existence of more than two components then the two components will be retained
- Otherwise one of the two components will be removed from the model

## Incremental Variational Bayesian Gaussian Mixture Models



## Conclusions

- Variational Approximation Pros:

1. Very Flexible Tool
2. Nice Theoretical Properties
3. Gives Tractable Algorithms
4. Applied to Many Problems\*

- Variational Approximation Cons:

1. Tightness of Bound
2. Sometimes Difficult Calculations