# Adaptive Importance Sampling

$$\int P(x) = \int \frac{P(x)}{q(x)} q(x) \approx \frac{1}{N} \sum_{x_i=1}^{N} \frac{P(x_i)}{q(x_i)} \equiv \mu^N \quad where \quad x_i \sim q(x)$$

uncertainty (squared):

$$Var(\mu^N) = \frac{1}{N} \left[ \int \frac{P(x)}{q(x)} P(x) dx - \left( \int P(x) dx \right)^2 \right]$$

- minimize the uncertainty with respect to q

  – it suffices to minimize: $\log\left( \int \frac{P(x)}{q(x)} P(x) dx \right)$

  – by Jensen's inequality: $\geq \int \left( \log \frac{P(x)}{q(x)} \right) P(x) dx$

  – which is the Kullback-Leibler divergence $KL(P\|q)$

# Adaptive Importance Sampling

Note:

- $\left(0\leq\right) Var\left(\mu^{N}\right) = 0$ if and only if $P=q$

- $\left(0\leq\right) KL\left(P\|q\right) = 0$ if and only if $P=q$

- $\left(0\leq\right) KL\left(q\|P\right) = 0$ if and only if $P=q$

Although not guaranteed, there is a good chance to decrease $Var\left(\mu^{N}\right)$ while minimizing $KL\left(P\|q\right)$ or $KL\left(q\|P\right)$ with respect to q.

# Adaptive Importance Sampling

- Conventional adaptive Importance Sampling approach:
    - restrict q to be a Gaussian or Student T mixture with a fixed number of components K:

$$q(x) = \sum_{i=1}^{K} \alpha_i q_i(x | \mu_i, \Sigma_i, \nu_i) \quad where \quad q_i \in \{\mathcal{N}, \mathcal{T}\}$$

    - minimize $KL(P \| q) \rightarrow$ EM-like parameter updates
    - no a priori information about q (i.e. flat priors for the parameters)

# Adaptive Importance Sampling

- Variational Bayes adaptive Importance Sampling approach:

  - restrict q to be a Gaussian or Student T mixture with a fixed number of components K:

  $$q(x) = \sum_{i=1}^{K} \alpha_i q_i(x|\mu_i, \Sigma_i, \nu_i) \ \ where \ \ q_i \in \{\mathcal{N}, \mathcal{T}\}$$

  - minimize $KL(q\|P) \rightarrow$ EM-like hyperparameter updates
  - include prior information about $\alpha_i, \mu_i, \Sigma_i, \nu_i$

- Very detailed in "Pattern Recognition and Machine learning" (Christopher M. Bishop)

# Variational Bayes

Definition: Hyperparameter

Consider a probabilistic model with parameters $\boldsymbol{\theta}$. Then the parameters of the prior distribution are called hyperparameters $\mathbf{h}$:

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\boldsymbol{h})$$

# Variational Bayes

Example: Hyperparameter

Be **X** data known to be Gaussian distributed, but with unknown mean and covariance:

$$P(\boldsymbol{x})=\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) \quad \boldsymbol{X} \sim g(\boldsymbol{x}) \quad \boldsymbol{\mu},\boldsymbol{\Sigma} \text{ unknown}$$

Suppose the prior can be written as follows:

$$P(\underbrace{\boldsymbol{\mu},\boldsymbol{\Sigma}}_{\text{parameters}} | \underbrace{\boldsymbol{m},\beta,\boldsymbol{W},\nu}_{\text{hyperparameters}})=\mathcal{N}(\boldsymbol{\mu}|\boldsymbol{m},(\beta\boldsymbol{\Lambda})^{-1})\,\mathcal{W}(\boldsymbol{\Lambda}|\boldsymbol{W},\nu)$$

Then m, β, W and $\nu$ are hyperparameters.

# Variational Bayes

- We end up with an EM-like algorithm for the hyperparameters.

- Can do Importance Sampling by taking the mode of $p(\boldsymbol{\alpha}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\nu} | X)$ as parameters for q.

- Can include knowledge from previous sampling runs.

- Can even correctly include Markov Chain prerun on multimodal target.

# Variational Bayes (details)

Given data **X**, we can write for the evidence of any arbitrary model:

$$\ln p(X) = \mathcal{L}(u) + KL(u \| p)$$

where we define:

$$\mathcal{L}(u) \equiv \int u(\mathbf{Z}, \boldsymbol{\theta}) \ln \left\{ \frac{p(X, \mathbf{Z}, \boldsymbol{\theta})}{u(\mathbf{Z}, \boldsymbol{\theta})} \right\} d\mathbf{Z} \, d\boldsymbol{\theta}$$

$$KL(u \| p) \equiv - \int u(\mathbf{Z}, \boldsymbol{\theta}) \ln \left\{ \frac{p(\mathbf{Z}, \boldsymbol{\theta}|X)}{u(\mathbf{Z}, \boldsymbol{\theta})} \right\} d\mathbf{Z} \, d\boldsymbol{\theta}$$

**θ**: model parameters $(\alpha_i, \mu_i, \Sigma_i, \nu_i)$
**Z**: latent variables (next slide)
u: an arbitrary proper probability distribution

# Variational Bayes (details)

- We want: $p(\mathbf{Z},\boldsymbol{\alpha},\boldsymbol{\mu},\boldsymbol{\Sigma},\boldsymbol{\nu}|X)$

- We need: analytically tractable approximation

  - assume $u(\mathbf{Z},\boldsymbol{\theta})=u(\mathbf{Z})u(\boldsymbol{\theta})$

  - maximize $\mathcal{L}(u)$ $\left(\Leftrightarrow \text{ minimize } KL(u\|p)\right)$

# Variational Bayes (details)

Definition: Latent (hidden) variables

Be $\mathbf{X} = \{x_1, \dots, x_N\}$ the (visible) data, then
$\mathbf{Z} = \{z_1, \dots, z_N\}$ is called latent if $\mathbf{Z}$ contains information which cannot uniquely be determined from $\mathbf{X}$ but inferred in a probabilistic model.

Example (Gaussian mixtures):

$$g(x) = \sum_i \alpha_i \, \mathcal{N}_i(x | \mu_i, \Sigma_i)$$
$$X \sim g(x)$$

Then the index i is a latent variable