# 02477 – Bayesian Machine Learning: Lecture 11

## Michael Riis Andersen

Technical University of Denmark,
DTU Compute, Department of Applied Math and Computer Science

# Outline

The ELBO objective

# Variational inference: big picture

- Our goal is to approximate a posterior distribution of interest

$$p \equiv p(\boldsymbol{w}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{w})p(\boldsymbol{w})}{p(\boldsymbol{y})}$$

- Variational inference in three steps

  1. Define collection of "simple" approximate probability distributions $\mathcal{Q}$ (*the variational family*)

  2. Define a measure of "distance" between probability distributions $\mathbb{D}\left[q||p\right]$

  3. Search for the distribution $q \in \mathcal{Q}$ that resembles the exact posterior $p$ as close as possible as measured by $\mathbb{D}\left[q||p\right]$

- The variational approximation $q$ for the target distribution $p \approx q$ is defined as

$$q_* = \arg \min_{q \in \mathcal{Q}} \mathbb{D}\left[q||p\right]$$

## Understanding the ELBO objective

- Suppose our model of interest is

$$p(\boldsymbol{y}, \boldsymbol{w}) = \prod_{n=1}^{N} p(y_n|\boldsymbol{w})p(\boldsymbol{w})$$

- We optimize the ELBO $\mathcal{L}$ to minimize the KL divergence (*between approximation and target*)

$$\mathcal{L}[q] = \mathbb{E}_q[\ln p(\boldsymbol{y}, \boldsymbol{w})] - \mathbb{E}_q[\ln q(\boldsymbol{w})]$$

- Re-writing the lower bound

$$\mathcal{L}[q] \equiv \mathbb{E}_q\left[\ln \prod_{n=1}^{N} p(y_n|\boldsymbol{w})p(\boldsymbol{w})\right] - \mathbb{E}_q[\ln q(\boldsymbol{w})]$$

$$= \sum_{n=1}^{N} \mathbb{E}_q[\ln p(y_n|\boldsymbol{w})] + \mathbb{E}_q[\ln p(\boldsymbol{w})] - \mathbb{E}_q[\ln q(\boldsymbol{w})]$$

$$= \sum_{n=1}^{N} \mathbb{E}_q[\ln p(y_n|\boldsymbol{w})] - \mathbb{E}_q\left[\ln \frac{q(\boldsymbol{w})}{p(\boldsymbol{w})}\right]$$

$$= \sum_{n=1}^{N} \mathbb{E}_q[\ln p(y_n|\boldsymbol{w})] - \mathsf{KL}[q(\boldsymbol{w})||p(\boldsymbol{w})]$$

- The first term (expected log likelihood) is a *data-fit* term, while the KL term encourages $q$ to be close to the prior $p(\boldsymbol{w})$ (*regularization term*)

Free-form vs fixed-form variational inference

# Free-form variational inference

- *Factorized approximation* for approximating the target distribution $p \equiv p(\boldsymbol{w}|\boldsymbol{y})$

$$q(\boldsymbol{w}) = \prod_{j=1}^{J} q(\boldsymbol{w}_j), \quad \text{where} \quad \boldsymbol{w} = [\boldsymbol{w}_1, \boldsymbol{w}_2, ..., \boldsymbol{w}_J]$$

- Optimality condition for mean-field families

$$\ln q^*(\boldsymbol{w}_k) = \mathbb{E}_{i \neq k} [\ln p(\boldsymbol{y}, \boldsymbol{w})] + K$$

- Strategy for free-form variational inference

  1. Deduce optimal form for each factor and derive parameters
  2. Derive fixed-point algorithm based on the optimal parameters

+ Optimal functional form given assumptions
+ Fast optimization

- Requires model-specific derivations
- Required integrals may be intractable
- Optimal forms may not be "known" distributions

# Fixed-form variational inference

- *Fixed-form* variational inference as an alternative to *free-form*: We give up an the optimal functional form and simply *assume* some family of distributions $q_\psi$ with parameters $\psi$

- We often refer to $\psi$ as *variational parameters*

- **Example**: Consider a model $p(\mathbf{y}, \mathbf{w})$, where $\mathbf{w} \in \mathbb{R}^D$

  - Full-rank Gaussians: $\mathbf{m} \in \mathbb{R}^D$, $\mathbf{V} \in \mathbb{R}^{D \times D}$
    $$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{V}),$$

  - Low-rank Gaussians: $\mathbf{m} \in \mathbb{R}^D$, $\mathbf{B} \in \mathbb{R}^{D \times K}$, and $\mathbf{C} \in \mathbb{R}^{D \times D}$ is diagonal.
    $$q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \mathbf{B}\mathbf{B}^T + \mathbf{C}^2),$$

  - Mean-field Gaussians: $\mathbf{m} = \begin{bmatrix} m_1, \ldots, m_D \end{bmatrix} \in \mathbb{R}^D$ and $\mathbf{v} = \begin{bmatrix} v_1, \ldots, v_D \end{bmatrix} = \mathbb{R}^D$
    $$q(\mathbf{w}) = \prod_{i=1}^{D} \mathcal{N}(w_i|m_i, v_i),$$

- Non-restricted to Gaussians: Gamma, Gauss, Beta, Dirichlet, Bernouilli etc

# Fixed-form variational inference II

- *Fixed-form* variational inference as an alternative to *free-form*

- **Example**: Consider a model $p(\boldsymbol{y}, \boldsymbol{w})$, where $\boldsymbol{w} \in \mathbb{R}^D$

- Suppose we choose a full-rank Gaussian family

$$q(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{w}|\boldsymbol{m}, \boldsymbol{V}),$$

  such that the variational family $\mathcal{Q}$ consists of all multivariate Gaussian distributions

- $q_{\psi}$ is now *parametrized* by the *variational parameters* $\psi = \{\boldsymbol{m}, \boldsymbol{V}\}$

$$\begin{aligned}
\mathcal{L}\left[q_{\psi}\right] &= \mathbb{E}_{q_{\psi}}\left[\ln p(\boldsymbol{y}, \boldsymbol{w})\right] - \mathbb{E}_{q_{\psi}}\left[\ln q_{\psi}(\boldsymbol{w})\right] \\
&= \mathbb{E}_{\mathcal{N}(\boldsymbol{w}|\boldsymbol{m}, \boldsymbol{V})}\left[\ln p(\boldsymbol{y}, \boldsymbol{w})\right] - \mathbb{E}_{\mathcal{N}(\boldsymbol{w}|\boldsymbol{m}, \boldsymbol{V})}\left[\ln \mathcal{N}(\boldsymbol{w}|\boldsymbol{m}, \boldsymbol{V})\right]
\end{aligned}$$

- Fitting the approximation using gradient-based methods

$$q^* = \arg\max_{q \in \mathcal{Q}} \mathcal{L}\left[q\right] \quad \Longleftrightarrow \quad \psi^* = \arg\max_{\psi} \mathcal{L}\left[q_{\psi}\right] \quad \Longleftrightarrow \quad \boldsymbol{m}^*, \boldsymbol{V}^* = \arg\max_{\boldsymbol{m}, \boldsymbol{V}} \mathcal{L}\left[q_{\psi}\right]$$

Hyperparameter estimation in VI

## Dealing with hyperparameters for variational inference

- Consider a model with data $\boldsymbol{y}$, parameters $\boldsymbol{\theta}$, and hyperparameters $\boldsymbol{\xi}$

$$p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{\xi}) = \frac{p(\boldsymbol{y}|\boldsymbol{\theta}, \boldsymbol{\xi})p(\boldsymbol{\theta}|\boldsymbol{\xi})}{p(\boldsymbol{y}|\boldsymbol{\xi})}$$

- **Examples**

  1. Linear regression
     - $\boldsymbol{\theta} = \boldsymbol{w}$ would be the regression weights
     - $\boldsymbol{\xi} = \{\alpha, \beta\}$ would be prior and noise precision.

  2. Gaussian process regression
     - $\boldsymbol{\theta} = \boldsymbol{f}$ would be the latent function values
     - $\boldsymbol{\xi} = \{\sigma^2, \kappa, \ell\}$ would be noise variance and kernel hyperparameters

- Hyperparameter estimation via the marginal likelihood (MLII/MAPII)

$$\hat{\boldsymbol{\xi}} = \arg \max_{\boldsymbol{\xi}} \log p(\boldsymbol{y}|\boldsymbol{\xi})$$

- But what to do for non-conjugate model, where we cannot compute the marginal likelihood?

$$p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{\xi}) \approx q(\boldsymbol{\theta})$$

# Dealing with hyperparameters for variational inference

- Consider a model with data $\boldsymbol{y}$, parameters $\boldsymbol{\theta}$, and hyperparameters $\boldsymbol{\xi}$

$$p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{\xi}) = \frac{p(\boldsymbol{y}|\boldsymbol{\theta}, \boldsymbol{\xi})p(\boldsymbol{\theta}|\boldsymbol{\xi})}{p(\boldsymbol{y}|\boldsymbol{\xi})}$$

- Variational approximation $p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{\xi}) \approx q(\boldsymbol{\theta})$ by optimizing the ELBO

$$\mathcal{L}_{\boldsymbol{\xi}}[q] = \mathbb{E}_q[\ln p(\boldsymbol{y}, \boldsymbol{\theta}|\boldsymbol{\xi})] - \mathbb{E}_q[\ln q(\boldsymbol{\theta})]$$

- The ELBO is a lowerbound on the log marginal likelihood

$$\log p(\boldsymbol{y}|\boldsymbol{\xi}) \geq \mathcal{L}_{\boldsymbol{\xi}}[q]$$

- and hence, we can do hyperparameter estimation via ELBO

$$\hat{\boldsymbol{\xi}} = \arg\max_{\boldsymbol{\xi}} p(\boldsymbol{y}|\boldsymbol{\xi}) \approx \arg\max_{\boldsymbol{\xi}} \mathcal{L}_{\boldsymbol{\xi}}[q]$$

- **Key take-away**: We can optimize the ELBO wrt. *variational parameters* $\psi$ and *hyperparameters* $\boldsymbol{\xi}$ simultaneously

$$\hat{\psi}^*, \hat{\boldsymbol{\xi}}^* = \arg\max_{\boldsymbol{\xi}} \mathcal{L}_{\boldsymbol{\xi}}[q_{\psi}]$$

Scaling Gaussian processes using variational inference
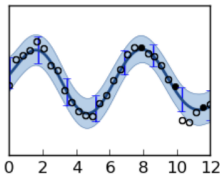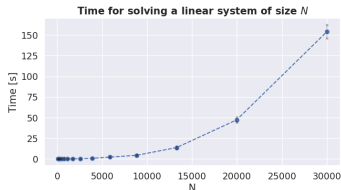
# Speeding up Gaussian process inference

- Gaussian process priors, $f \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$ can be extremely powerful, but scales poorly: $\mathcal{O}(N^3)$ in computational complexity and $\mathcal{O}(N^2)$ in memory footprint

$$p(y^*|\mathbf{y}) = \mathcal{N}\left(y^*|\mu_{y^*|\mathbf{y}}, \sigma^2_{y^*|\mathbf{y}}\right)$$

$$\mu_{y^*|\mathbf{y}} = \mathbf{k}\left(\mathbf{K} + \beta^{-1}\mathbf{I}\right)^{-1}\mathbf{y}$$

$$\sigma^2_{y^*|\mathbf{y}} = c - \mathbf{k}\left(\mathbf{K} + \beta^{-1}\mathbf{I}\right)^{-1}\mathbf{k}^T$$

- Obviously, using a subset of data would be faster, but can we do something more clever? Yes!

- We introduce *inducing points* $\mathbf{z}_i$ for $i = 1, \ldots, M$ for $M \ll N$

- Combining *variational inference* with *inducing points* allows us to approximate the posterior much faster, i.e. $\mathcal{O}(NM^2)$, and even faster, $\mathcal{O}(M^2)$, using *mini-batching* (next week)



**Time for solving a linear system of size N**



Hensman et al, 2019: Gaussian Processes for Big Data

# Introducing inducing points

- Suppose we have dataset $\mathcal{D} = \{x_n, y_n\}$, where $x_n \in \mathbb{R}^D$ and $y_n \in \mathbb{R}$
$$y_n = f(x_n) + e_n$$

- Assuming Gaussian noise and a GP prior, i.e. $f(x) \sim \mathcal{GP}(0, k(x, x'))$ yields the joint distribution
$$p(y, f) = p(y|f)p(f) = \mathcal{N}(y|f, \sigma^2 I)\mathcal{N}(f|0, K_{ff})$$

- **Goal**: fast way compute $p(f|y)$
$$f_n = f(x_n)$$

- We introduce $M$ *inducing points* $z_i \in \mathbb{R}^D$ for $m = 1, \ldots, M$ such that
$$u_i = f(z_i)$$

- Extended joint distribution for data $y$ and latent function values for both $f$ and $u$
$$p(y, f, u) = p(y|f)p(f|u)p(u) = \mathcal{N}(y|f, \sigma^2 I)\mathcal{N}(f|K_{fu}K_{uu}^{-1}u, K_{ff} - K_{fu}K_{uu}^{-1}K_{fu})\mathcal{N}(u|0, K_{uu})$$

- We can always marginalize to get the original model back
$$p(y, f) = \int p(y, f, u)\mathrm{d}u$$

## Intuition

- Our extended joint distribution

$$p(\mathbf{y}, \mathbf{f}, \mathbf{u}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})\mathcal{N}(\mathbf{f}|\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{fu})\mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{K}_{uu})$$
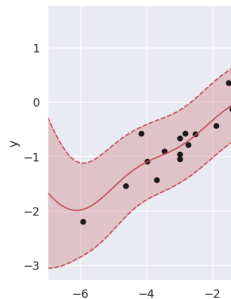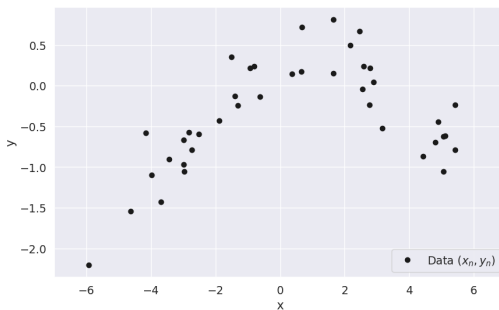


- Can we find a set of *inducing points* $\{\mathbf{z}_i\}_{i=1}^{M}$ and associated function values $u_i = f(\mathbf{z}_i)$ such that we can absorb all the information from the full data set?

# Setting up the approximation and choosing a variational family

- Our extended joint distribution

$$p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{u}) = p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{u})p(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{f}, \sigma^2\boldsymbol{I})\mathcal{N}(\boldsymbol{f}|\boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{u}, \boldsymbol{K}_{ff} - \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{fu})\mathcal{N}(\boldsymbol{u}|\boldsymbol{0}, \boldsymbol{K}_{uu})$$

- We will make a variational approximation: $p(\boldsymbol{f}, \boldsymbol{u}|\boldsymbol{y}) \approx q(\boldsymbol{f}, \boldsymbol{u})$

- A clever choice for the variational family

$$q(\boldsymbol{f}, \boldsymbol{u}) = p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{f}|\boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{u}, \boldsymbol{K}_{ff} - \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{fu})\mathcal{N}(\boldsymbol{u}|\boldsymbol{m}_u, \boldsymbol{S}_u)$$

  where

  - $\boldsymbol{m}_u \in \mathbb{R}^M$ and $\boldsymbol{S}_u \in \mathbb{R}^{M \times M}$ are variational parameters to be estimated
  - $\boldsymbol{K}_{ff}$ is the prior covariance matrix for $\boldsymbol{f}$
  - $\boldsymbol{K}_{uu}$ is the prior covariance matrix for $\boldsymbol{u}$
  - $\boldsymbol{K}_{fu}$ is the prior covariance between $\boldsymbol{f}$ and $\boldsymbol{u}$

- Recall $p(\boldsymbol{f}|\boldsymbol{u})$ is the prior conditional distribution derived from the prior $p(\boldsymbol{f}, \boldsymbol{u})$.

## The approximate posterior distribution

- Our choice of variational family:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = \mathcal{N}(\mathbf{f}|\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{fu})\mathcal{N}(\mathbf{u}|\mathbf{m}_u, \mathbf{S}_u)$$

implies a marginal variational approximation for $q(\mathbf{f})$ (linear Gaussian system)

$$\begin{aligned}
q(\mathbf{f}) &= \int p(\mathbf{f}|\mathbf{u})q(\mathbf{u})\mathrm{d}\mathbf{u} \\
&= \int \mathcal{N}(\mathbf{f}|\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{fu})\mathcal{N}(\mathbf{u}|\mathbf{m}_u, \mathbf{S}_u)\mathrm{d}\mathbf{u} \\
&= \mathcal{N}(\mathbf{f}|\underbrace{\mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{m}_u}_{\mathbf{m}_f}, \underbrace{\mathbf{K}_{ff} + \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}(\mathbf{S}_u - \mathbf{K}_{uu})\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf}}_{\mathbf{S}_f}) \\
&= \mathcal{N}(\mathbf{f}|\mathbf{m}_f, \mathbf{S}_f)
\end{aligned}$$

- ... and similarly we can make predictions for new input points to get $f^* = f(\mathbf{x}^*)$

$$\begin{aligned}
p(f^*|\mathbf{y}) &= \int p(f^*|\mathbf{u})p(\mathbf{u}|\mathbf{y})\mathrm{d}\mathbf{u} \\
&\approx \int p(f^*|\mathbf{u})q(\mathbf{u})\mathrm{d}\mathbf{u} \\
&= \mathcal{N}(f^*|\underbrace{\mathbf{K}_{f^*u}\mathbf{K}_{uu}^{-1}\mathbf{m}_u}_{\mathbf{m}_{f^*}}, \underbrace{\mathbf{K}_{f^*f^*} + \mathbf{K}_{f^*u}\mathbf{K}_{uu}^{-1}(\mathbf{S}_u - \mathbf{K}_{uu})\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf^*}}_{\mathbf{S}_{f^*s}})
\end{aligned}$$

- So all we need to do is to estimate the mean and covariance for $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}_u, \mathbf{S}_u)$

## Calculating the ELBO

We substitute our model and variational approximation into the ELBO

$$\mathcal{L}[q] = \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}, \boldsymbol{f}, \boldsymbol{u})\right] - \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log q(\boldsymbol{f}, \boldsymbol{u})\right]$$

$$= \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}|\boldsymbol{f})p(\boldsymbol{f}|\boldsymbol{u})p(\boldsymbol{u})\right] - \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})\right]$$

$$= \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}|\boldsymbol{f})\right] + \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{f}|\boldsymbol{u})\right] + \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{u})\right]$$

$$- \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{f}|\boldsymbol{u})\right] - \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log q(\boldsymbol{u})\right]$$

$$= \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}|\boldsymbol{f})\right] + \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{u})\right] - \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log q(\boldsymbol{u})\right]$$

$$= \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}|\boldsymbol{f})\right] + \mathbb{E}_{q(\boldsymbol{u})}\left[\log p(\boldsymbol{u})\right] - \mathbb{E}_{q(\boldsymbol{u})}\left[\log q(\boldsymbol{u})\right]$$

$$= \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log p(\boldsymbol{y}|\boldsymbol{f})\right] - \mathrm{KL}\left[q(\boldsymbol{u})||p(\boldsymbol{u})\right]$$

Next, recall that $p(\boldsymbol{y}|\boldsymbol{f}) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{f}, \sigma^2\boldsymbol{I}) = \prod_{n=1}^{N} \mathcal{N}(y_n|f_n, \sigma^2)$

$$\mathcal{L}[q] = \sum_{n=1}^{N} \mathbb{E}_{p(\boldsymbol{f}|\boldsymbol{u})q(\boldsymbol{u})}\left[\log \mathcal{N}(y_n|f_n, \sigma^2)\right] - \mathrm{KL}\left[q(\boldsymbol{u})||p(\boldsymbol{u})\right]$$

$$= \sum_{n=1}^{N} \mathbb{E}_{p(f_n|\boldsymbol{u})q(\boldsymbol{u})}\left[\log \mathcal{N}(y_n|f_n, \sigma^2)\right] - \mathrm{KL}\left[q(\boldsymbol{u})||p(\boldsymbol{u})\right]$$

## Calculating the first term

$$\mathbb{E}_{p(f_n|\boldsymbol{u})q(\boldsymbol{u})}\left[\log\mathcal{N}(y_n|f_n,\sigma^2)\right] = \iint p(f_n|\boldsymbol{u})q(\boldsymbol{u})\log\mathcal{N}(y_n|f_n,\sigma^2)\mathrm{d}f_n\mathrm{d}\boldsymbol{u}$$

$$= \iint p(f_n|\boldsymbol{u})q(\boldsymbol{u})\mathrm{d}\boldsymbol{u}\log\mathcal{N}(y_n|f_n,\sigma^2)\mathrm{d}f_n \qquad \text{(Recall: } q(\boldsymbol{f}) = \mathcal{N}(\boldsymbol{f}|\boldsymbol{m}_f,\boldsymbol{S}_f))$$

$$= \int \mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)\log\mathcal{N}(y_n|f_n,\sigma^2)\mathrm{d}f_n$$

$$= \mathbb{E}_{\mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)}\left[\log\mathcal{N}(y_n|f_n,\sigma^2)\right]$$

$$= \mathbb{E}_{\mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)}\left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n - f_n)^2\right]$$

$$= \mathbb{E}_{\mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)}\left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n^2 + f_n^2 - 2y_nf_n)\right]$$

$$= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n^2 + \mathbb{E}_{\mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)}\left[f_n^2\right] - 2y_n\mathbb{E}_{\mathcal{N}\left(f_n|m_{f_n},\sigma_{f_n}^2\right)}[f_n])$$

$$= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n^2 + m_{f_n}^2 + \sigma_{f_n}^2 - 2y_nm_{f_n})$$

$$= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n^2 + m_{f_n}^2 - 2y_nm_{f_n}) - \frac{1}{2\sigma}\sigma_{f_n}^2$$

$$= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma}(y_n - m_{f_n})^2 - \frac{1}{2\sigma}\sigma_{f_n}^2$$

$$= \log\mathcal{N}(y_n|m_{f_n},\sigma^2)\mathrm{d}f_n - \frac{1}{2\sigma}\sigma_{f_n}^2$$

## The KL term

- The two terms in the lowerbound

$$\mathcal{L}[q] = \sum_{n=1}^{N} \mathbb{E}_{p(f_n|\boldsymbol{u})q(\boldsymbol{u})}\left[\log\mathcal{N}(y_n|f_n, \sigma^2)\right] - \mathsf{KL}\left[q(\boldsymbol{u})||p(\boldsymbol{u})\right]$$

- The KL divergence between two multivariate Gaussians can be computed in closed-form:

$$\mathsf{KL}\left[\mathcal{N}(\boldsymbol{u}|\boldsymbol{m}_0, \boldsymbol{\Sigma}_0)||\mathcal{N}(\boldsymbol{u}|\boldsymbol{m}_1, \boldsymbol{\Sigma}_1)\right]$$
$$= \frac{1}{2}\left[\mathsf{trace}(\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\Sigma}_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T\boldsymbol{\Sigma}_1^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - D + \log\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|}\right]$$

- So for $q(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{u}|\boldsymbol{m}_u, \boldsymbol{S}_u)$ and $p(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{u}|\boldsymbol{0}, \boldsymbol{K}_{uu})$, we get

$$\mathsf{KL}\left[q(\boldsymbol{u})||p(\boldsymbol{u})\right] = \frac{1}{2}\left[\mathsf{trace}(\boldsymbol{K}_{uu}^{-1}\boldsymbol{S}_u) + \boldsymbol{m}_u^T\boldsymbol{K}_{uu}^{-1}\boldsymbol{m}_u - M + \log\frac{|\boldsymbol{K}_{uu}|}{|\boldsymbol{S}_u|}\right]$$

## Combining everything

- We derived

$$\mathcal{L}[q] = \sum_{n=1}^{N} \mathbb{E}_{p(f_n|\boldsymbol{u})q(\boldsymbol{u})}\left[\log \mathcal{N}(y_n|f_n, \sigma^2)\right] - \mathsf{KL}[q(\boldsymbol{u})||p(\boldsymbol{u})]$$

- ... and then showed that the first term simplifies to

$$\mathbb{E}_{p(f_n|\boldsymbol{u})q(\boldsymbol{u})}\left[\log \mathcal{N}(y_n|f_n, \sigma^2)\right] = \log \mathcal{N}(y_n|m_{f_n}, \sigma^2)\mathsf{d}f_n - \frac{1}{2\sigma}\sigma_{f_n}^2$$

- Combining yields

$$\mathcal{L}[q] = \sum_{n=1}^{N}\left[\log \mathcal{N}(y_n|m_{f_n}, \sigma^2)\mathsf{d}f_n - \frac{1}{2\sigma}\sigma_{f_n}^2\right] - \mathsf{KL}[q(\boldsymbol{u})||p(\boldsymbol{u})]$$

$$= \sum_{n=1}^{N}\log \mathcal{N}(y_n|m_{f_n}, \sigma^2) - \frac{1}{2\sigma}\sum_{n=1}^{N}\sigma_{f_n}^2 - \mathsf{KL}[q(\boldsymbol{u})||p(\boldsymbol{u})]$$

$$= \log \mathcal{N}(\boldsymbol{y}|\boldsymbol{m}_f, \sigma^2 \boldsymbol{I}) - \frac{1}{2\sigma}\mathsf{trace}(\boldsymbol{S}_f) - \mathsf{KL}[q(\boldsymbol{u})||p(\boldsymbol{u})]$$

where

$$\boldsymbol{m}_f = \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{m}_u$$
$$\boldsymbol{S}_f = \boldsymbol{K}_{ff} + \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}(\boldsymbol{S}_u - \boldsymbol{K}_{uu})\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{uf}$$

- We can optimize this bound wrt. the *variational parameters* $\boldsymbol{m}_u$ and $\boldsymbol{S}_u$ and wrt. the *hyperparameters* simultaneously!

## One step more step

- It turns out we can optimize the bound wrt. $\boldsymbol{m}_u$ and $\boldsymbol{S}_u$ analytically

$$
\begin{aligned}
\mathcal{L}[q] &= \sum_{n=1}^{N} \left[ \log \mathcal{N}(y_n | m_{f_n}, \sigma^2) \mathrm{d}f_n - \frac{1}{2\sigma} \sigma_{f_n}^2 \right] - \mathrm{KL}\left[ q(\boldsymbol{u}) || p(\boldsymbol{u}) \right] \\
&= \sum_{n=1}^{N} \log \mathcal{N}(y_n | m_{f_n}, \sigma^2) - \frac{1}{2\sigma} \sum_{n=1}^{N} \sigma_{f_n}^2 - \mathrm{KL}\left[ q(\boldsymbol{u}) || p(\boldsymbol{u}) \right] \\
&= \log \mathcal{N}(\boldsymbol{y} | \boldsymbol{m}_f, \sigma^2 \boldsymbol{I}) - \frac{1}{2\sigma} \mathrm{trace}(\boldsymbol{S}_f) - \mathrm{KL}\left[ q(\boldsymbol{u}) || p(\boldsymbol{u}) \right]
\end{aligned}
$$

- ... to get

$$
\boldsymbol{S}_u^{-1} = \frac{1}{\sigma^2} \boldsymbol{K}_{uu}^{-1} \boldsymbol{K}_{uf} \boldsymbol{K}_{fu} \boldsymbol{K}_{uu}^{-1} + \boldsymbol{K}_{uu}^{-1}
$$

$$
\boldsymbol{m}_u = \frac{1}{\sigma^2} \boldsymbol{S}_u \boldsymbol{K}_{uu}^{-1} \boldsymbol{K}_{uf} \boldsymbol{y}
$$

- ... which leads to the *collapsed lowerbound*

$$
\mathcal{L}[q] = \log \mathcal{N}(\boldsymbol{y} | \boldsymbol{0}, \boldsymbol{K}_{fu} \boldsymbol{K}_{uu} \boldsymbol{K}_{uf} + \sigma^2 \boldsymbol{I}) - \frac{1}{2\sigma} \mathrm{trace}(\boldsymbol{K}_{ff} - \boldsymbol{K}_{fu} \boldsymbol{K}_{uu}^{-1} \boldsymbol{K}_{uf})
$$

## The big picture and how to use it in practice

**Goal**: fast way compute $p(\boldsymbol{f}|\boldsymbol{y})$

1. Choose a set of *inducing points* $\boldsymbol{z}_i$, where $u_i = f(\boldsymbol{z}_i)$

2. Optimize the *collapsed bound* wrt. our hyperparameters $\boldsymbol{\theta}$

$$\hat{\boldsymbol{\theta}}^* = \arg\max_{\boldsymbol{\theta}} \mathcal{L}_{\boldsymbol{\theta}}\left[q\right] = \log \mathcal{N}(\boldsymbol{y}|\boldsymbol{0}, \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}\boldsymbol{K}_{uf} + \sigma^2 \boldsymbol{I}) - \frac{1}{2\sigma}\mathrm{trace}(\boldsymbol{K}_{ff} - \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{uf})$$

3. Compute the posterior mean and covariance for latent function values
   $q(\boldsymbol{u}) = \mathcal{N}(\boldsymbol{u}|\boldsymbol{m}_u, \boldsymbol{S}_u)$

$$\boldsymbol{S}_u^{-1} = \frac{1}{\sigma^2}\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{uf}\boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1} + \boldsymbol{K}_{uu}^{-1}$$

$$\boldsymbol{m}_u = \frac{1}{\sigma^2}\boldsymbol{S}_u\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{uf}\boldsymbol{y}$$

4. Compute the approximate posterior distribution for $q(\boldsymbol{f}) = \mathcal{N}(\boldsymbol{f}|\boldsymbol{m}_f, \boldsymbol{S}_f)$

$$\boldsymbol{m}_f = \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}\boldsymbol{m}_u$$

$$\boldsymbol{S}_f = \boldsymbol{K}_{ff} + \boldsymbol{K}_{fu}\boldsymbol{K}_{uu}^{-1}(\boldsymbol{S}_u - \boldsymbol{K}_{uu})\boldsymbol{K}_{uu}^{-1}\boldsymbol{K}_{uf}$$

5. Use $p(\boldsymbol{f}|\boldsymbol{y}) \approx q(\boldsymbol{f})$ to make predictions

## Airline delays dataset

- Flight arrival and departure times for every commercial flight in the USA from Jan. 2008 to April 2008

- 2 million flights: 700000 flight for training, 100000 for testing

- Target variable: Delay in minutes:

- $D = 8$ features: age of aircraft, flight distance, airtime, departure time, arrival time, day of the week, day of the month, month

- Squared exponential kernel with separate lengthscale $\ell_i > 0$ for each dimension

$$k(\mathbf{x}, \mathbf{x}') = \kappa^2 \exp\left[-\frac{1}{2} \sum_{i=1}^{D} \ell_i^{-1} |\mathbf{x}_i - \mathbf{x}_i'|^2\right] + \tau$$

- Using $M = 1000$ inducing points

Hensman et al: Gaussian Processes for Big Data (2013)
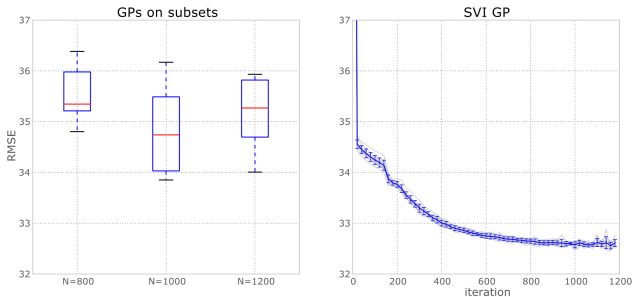
# Airline delays dataset



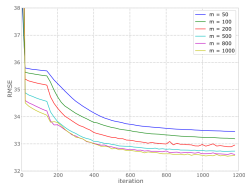Figure 7: Root mean squared errors in predicting flight delays using information about the flight.



Figure 8: Root mean square errors for models with different numbers of inducing variables.