

# Chapter 14

## Bayesian Learning

supplementary slides to  
*Machine Learning Fundamentals*  
© **Hui Jiang 2020**  
published by Cambridge University Press

August 2020



# Outline

- 1 Formulation of Bayesian Learning
- 2 Conjugate Priors
- 3 Approximate Inference
- 4 Gaussian Processes

# Bayesian Learning (I)

- *frequentist vs. Bayesian* views in machine learning
  - *frequentist*: model parameters as unknown but fixed quantities
  - *Bayesian*: model parameters as **random variables**
- Bayesians use probability distributions of model parameters
- Bayes' theorem:

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

- $p(\boldsymbol{\theta})$ : prior distribution of model parameters  $\boldsymbol{\theta}$
- $p(\boldsymbol{\theta} | \mathbf{x})$ : the posterior distribution of  $\boldsymbol{\theta}$  given data  $\mathbf{x}$
- $p(\mathbf{x} | \boldsymbol{\theta})$ : the likelihood function of the model
- Bayesian learning rule: *posterior*  $\propto$  *prior*  $\times$  *likelihood*

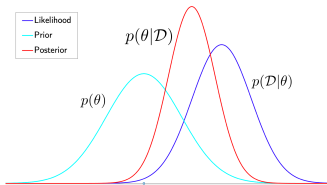
$$p(\boldsymbol{\theta} | \mathbf{x}) \propto p(\boldsymbol{\theta}) p(\mathbf{x} | \boldsymbol{\theta})$$

# Bayesian Learning (II)

- prior specification:  $p(\boldsymbol{\theta})$ 
  - use a prior distribution to describe prior knowledge on models
- Bayesian learning
  - optimally combine prior knowledge with data
  - given a training set:  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
  - Bayesian learning rule:  $p(\boldsymbol{\theta}) \xrightarrow{\mathcal{D}} p(\boldsymbol{\theta}|\mathcal{D})$

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta}) p(\mathcal{D}|\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta})$$

- Bayesian inference
  - make a decision based on  $p(\boldsymbol{\theta}|\mathcal{D})$



$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta}) p(\mathcal{D}|\boldsymbol{\theta})$$

posterior  $\propto$  prior  $\times$  likelihood

# Bayesian Inference for Classification

- given posterior  $p(\boldsymbol{\theta}|\mathcal{D})$  and likelihood  $p(\mathbf{x}|\boldsymbol{\theta})$
- define *predictive distribution* as

$$p(\mathbf{x}|\mathcal{D}) = \int_{\boldsymbol{\theta}} p(\mathbf{x}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}$$

- Bayesian classification:

- $K$  classes:  $\{\omega_1, \omega_2, \dots, \omega_K\}$
- choose prior  $p(\theta_k)$  and a training set  $\mathcal{D}_k$  for each class  $\omega_k$
- Bayesian learning:

$$p(\theta_k|\mathcal{D}_k) = \frac{p(\theta_k) p(\mathcal{D}_k|\omega_k, \theta_k)}{p(\mathcal{D}_k)} \propto p(\theta_k) p(\mathcal{D}_k|\omega_k, \theta_k)$$

- Bayesian inference:

$$\begin{aligned} g(\mathbf{x}) &= \arg \max_{k=1}^K p(\mathbf{x}|\mathcal{D}_k) \\ &= \arg \max_{k=1}^K \Pr(\omega_k) \int_{\theta_k} p(\mathbf{x}|\omega_k, \theta_k) p(\theta_k|\mathcal{D}_k) d\theta_k \end{aligned}$$

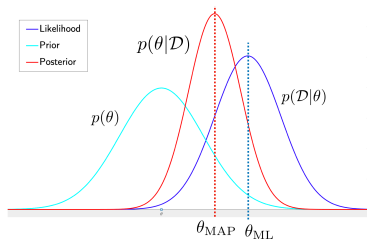
# Maximum a Posteriori (MAP) Estimation

- not easy to use a distribution  $p(\theta|\mathcal{D})$  to describe models
- point estimation: only use a point to estimate a distribution  $p(\theta|\mathcal{D})$
- maximum a posteriori (MAP) estimation:

$$\begin{aligned}\theta_{\text{MAP}} &= \arg \max_{\theta} p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} p(\theta) p(\mathcal{D} | \theta)\end{aligned}$$

- MAP estimation vs. ML estimation
  - ML solely relies on training data
  - MAP optimally combines prior knowledge with data

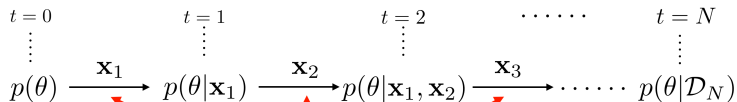
$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{D} | \theta)$$



$\theta_{\text{MAP}}$  vs.  $\theta_{\text{ML}}$

# Sequential Bayesian Learning

- Bayesian learning is an excellent tool for on-line learning, where training data come one by one
- sequential Bayesian learning
  - use the Bayesian learning to update models after each sample
  - track a slowly-changing environment



**Learning Rule:** posterior  $\propto$  prior  $\times$  likelihood

$$p(\theta | \mathbf{x}_1) \propto p(\theta)p(\mathbf{x}_1 | \theta)$$

$$p(\theta | \mathbf{x}_1, \mathbf{x}_2) \propto p(\theta | \mathbf{x}_1) p(\mathbf{x}_2 | \theta)$$

# Example: Sequential Bayesian Learning

- a univariate Gaussian model with known variance:

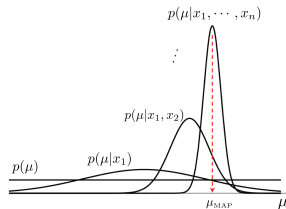
$$p(x | \mu) = \mathcal{N}(x | \mu, \sigma_0^2) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x-\mu)^2}{2\sigma_0^2}}$$

- choose a prior distribution:

$$p(\mu) = \mathcal{N}(\mu | \nu_0, \tau_0^2) = \frac{1}{\sqrt{2\pi\tau_0^2}} e^{-\frac{(\mu-\nu_0)^2}{2\tau_0^2}}$$

- $p(\mu|x_1) \propto p(\mu)p(x_1|\mu) \implies p(\mu|x_1) = \mathcal{N}(\mu|\nu_1, \tau_1^2)$   
with  $\nu_1 = \frac{\sigma_0^2}{\tau_0^2 + \sigma_0^2} \nu_0 + \frac{\tau_0^2}{\tau_0^2 + \sigma_0^2} x_1$  and  $\tau_1^2 = \frac{\tau_0^2 \sigma_0^2}{\tau_0^2 + \sigma_0^2}$

- $p(\mu | x_1, \dots, x_n) = \mathcal{N}(\mu|\nu_n, \tau_n^2)$  with  
 $\nu_n = \frac{n\tau_0^2}{n\tau_0^2 + \sigma_0^2} \bar{x}_n + \frac{\sigma_0^2}{n\tau_0^2 + \sigma_0^2} \nu_0$  and  $\tau_n^2 = \frac{\tau_0^2 \sigma_0^2}{n\tau_0^2 + \sigma_0^2}$



as  $n \rightarrow \infty$ , we have

- $\tau_n \rightarrow 0$
- $\nu_n \rightarrow \bar{x}_n$
- $\mu_{\text{MAP}} \rightarrow \mu_{\text{ML}}$



# Conjugate Priors

- *conjugate priors*: a prior is chosen to ensure its posterior has the same functional form as the prior
- conjugate to the likelihood function of the underlying model, i.e. both have the same function form
- choice of conjugate priors leads to computational convenience in Bayesian learning
- not every model has a conjugate prior, e.g. mixture models
- all e-family models have conjugate priors

# Examples of Conjugate Priors

model $p(\mathbf{x} \boldsymbol{\theta})$	conjugate prior $p(\boldsymbol{\theta})$
1-D Gaussian (known variance) $\mathcal{N}(x   \mu, \sigma_0^2)$	1-D Gaussian $\mathcal{N}(\mu   \nu, \tau^2)$
1-D Gaussian (known mean) $\mathcal{N}(x   \mu_0, \sigma^2)$	inverse-gamma $\text{gamma}^{-1}(\sigma^2   \alpha, \beta)$
Gaussian (known covariance) $\mathcal{N}(\mathbf{x}   \boldsymbol{\mu}, \boldsymbol{\Sigma}_0)$	Gaussian $\mathcal{N}(\boldsymbol{\mu}   \boldsymbol{\nu}, \Phi)$
Gaussian (known mean) $\mathcal{N}(\mathbf{x}   \boldsymbol{\mu}_0, \boldsymbol{\Sigma})$	inverse-Wishart $\mathcal{W}^{-1}(\boldsymbol{\Sigma}   \Phi, \nu)$
multivariate Gaussian $\mathcal{N}(\mathbf{x}   \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian-inverse-Wishart $\text{GIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma}   \boldsymbol{\nu}, \Phi, \lambda, \nu) =$ $\mathcal{N}(\boldsymbol{\mu}   \boldsymbol{\nu}, \frac{1}{\lambda} \boldsymbol{\Sigma}) \mathcal{W}^{-1}(\boldsymbol{\Sigma}   \Phi, \nu)$
multinomial Mult( $\mathbf{r}   \mathbf{w}$ ) = $C(\mathbf{r}) \cdot \prod_{i=1}^M w_i^{r_i}$ with $C(\mathbf{r}) = \frac{(r_1 + \dots + r_M)!}{r_1! \dots r_M!}$	Dirichlet Dir( $\mathbf{w}   \boldsymbol{\alpha}$ ) = $B(\boldsymbol{\alpha}) \cdot \prod_{i=1}^M w_i^{\alpha_i - 1}$ with $B(\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_1 + \dots + \alpha_M)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_M)}$

# Conjugate Priors for Bayesian Learning: Multinomials

- a sample of some counts:  $\mathbf{r} = [r_1 r_2 \cdots r_M]$
- multinomial models:  $p(\mathbf{r} | \mathbf{w}) = \text{Mult}(\mathbf{r} | \mathbf{w}) = C(\mathbf{r}) \cdot \prod_{i=1}^M w_i^{r_i}$
- the conjugate prior is Dirichlet:

$$p(\mathbf{w}) = \text{Dir}(\mathbf{w} | \boldsymbol{\alpha}^{(0)}) = B(\boldsymbol{\alpha}^{(0)}) \cdot \prod_{i=1}^M w_i^{\alpha_i^{(0)} - 1}$$

- Bayesian learning:

$$p(\mathbf{w} | \mathbf{r}) \propto p(\mathbf{w}) p(\mathbf{r} | \mathbf{w}) \propto \prod_{i=1}^M w_i^{\alpha_i^{(0)} + r_i - 1}$$

- the posterior is also Dirichlet:

$$p(\mathbf{w} | \mathbf{r}) = \text{Dir}(\mathbf{w} | \boldsymbol{\alpha}^{(1)}) = B(\boldsymbol{\alpha}^{(1)}) \cdot \prod_{i=1}^M w_i^{\alpha_i^{(1)} - 1}$$

- MAP estimation:

$$\begin{aligned} \mathbf{w}^{(\text{MAP})} &= \arg \max_{\mathbf{w}} p(\mathbf{w} | \mathbf{r}) \quad \text{subject to} \quad \sum_{i=1}^M w_i = 1 \\ \implies w_i^{(\text{MAP})} &= \frac{\alpha_i^{(1)} - 1}{\sum_{i=1}^M \alpha_i^{(1)} - M} = \frac{r_i + \alpha_i^{(0)} - 1}{\sum_{i=1}^M (r_i + \alpha_i^{(0)}) - M} \quad \forall i = 1, 2, \dots, M \end{aligned}$$

# Conjugate Priors for Bayesian Learning: Gaussians (1)

- Gaussian models:  $p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$
- the conjugate prior is a Gaussian-inverse-Wishart (GIW) distribution:  

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \text{GIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\nu}_0, \Phi_0, \lambda_0, \nu_0) = \mathcal{N}\left(\boldsymbol{\mu} | \boldsymbol{\nu}_0, \frac{1}{\lambda_0} \boldsymbol{\Sigma}\right) \mathcal{W}^{-1}\left(\boldsymbol{\Sigma} | \Phi_0, \nu_0\right)$$

$$= c_0 |\boldsymbol{\Sigma}^{-1}|^{\frac{\nu_0 + d + 2}{2}} \exp\left[-\frac{1}{2} \lambda_0 (\boldsymbol{\mu} - \boldsymbol{\nu}_0)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\nu}_0) - \frac{1}{2} \text{tr}(\Phi_0 \boldsymbol{\Sigma}^{-1})\right]$$
- the likelihood function of a training set  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ :

$$p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^N p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \frac{|\boldsymbol{\Sigma}^{-1}|^{\frac{N}{2}}}{(2\pi)^{Nd/2}} \exp\left[-\frac{1}{2} \text{tr}(N\boldsymbol{\Sigma}^{-1}) - \frac{N}{2} (\boldsymbol{\mu} - \bar{\mathbf{x}})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \bar{\mathbf{x}})\right]$$

- Bayesian learning:

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathcal{D}) \propto \text{GIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \boldsymbol{\nu}_0, \Phi_0, \lambda_0, \nu_0) \cdot p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

## Conjugate Priors for Bayesian Learning: Gaussians (2)

- the posterior is another GIW distribution:

$$\begin{aligned}
 p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathcal{D}_N) &= \text{GIW}(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \boldsymbol{\nu}_1, \Phi_1, \lambda_1, \nu_1) \\
 &= c_1 |\boldsymbol{\Sigma}^{-1}|^{\frac{\nu_1+d+2}{2}} \exp \left[ -\frac{1}{2} \lambda_1 (\boldsymbol{\mu} - \boldsymbol{\nu}_1)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\nu}_1) - \frac{1}{2} \text{tr}(\Phi_1 \boldsymbol{\Sigma}^{-1}) \right]
 \end{aligned}$$

- $\lambda_1 = \lambda_0 + N$  and  $\nu_1 = \nu_0 + N$
  - $\boldsymbol{\nu}_1 = \frac{\lambda_0 \boldsymbol{\nu}_0 + N \bar{\mathbf{x}}}{\lambda_0 + N}$
  - $\Phi_1 = \Phi_0 + N \mathbf{S} + \frac{\lambda_0 N}{\lambda_0 + N} (\bar{\mathbf{x}} - \boldsymbol{\nu}_0)(\bar{\mathbf{x}} - \boldsymbol{\nu}_0)^\top$
- MAP estimation:  $\{\boldsymbol{\mu}_{\text{MAP}}, \boldsymbol{\Sigma}_{\text{MAP}}\} = \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} p(\boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathcal{D}_N)$

$$\boldsymbol{\mu}_{\text{MAP}} = \boldsymbol{\nu}_1 = \frac{\lambda_0 \boldsymbol{\nu}_0 + N \bar{\mathbf{x}}}{\lambda_0 + N}$$

$$\boldsymbol{\Sigma}_{\text{MAP}} = \frac{\Phi_1}{\nu_1 + d + 1} = \frac{\Phi_0 + N \mathbf{S} + \frac{\lambda_0 N}{\lambda_0 + N} (\bar{\mathbf{x}} - \boldsymbol{\nu}_0)(\bar{\mathbf{x}} - \boldsymbol{\nu}_0)^\top}{\nu_0 + N + d + 1}$$

# Approximate Inference

- when conjugate priors do not exist, Bayesian learning may lead to very complicated posterior distributions
- *approximate inference*: approximate the true posterior distribution with a simple distribution for Bayesian inference
- popular approximate inference methods:
  - 1 Laplace's method
  - 2 variational Bayesian (VB) method

# Laplace's Method

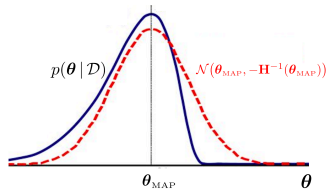
- use a Gaussian centered at  $\theta_{\text{MAP}}$  to approximate the true posterior  $p(\theta | \mathcal{D})$
- Taylor's expansion of  $f(\theta) = \ln p(\theta | \mathcal{D})$  at  $\theta_{\text{MAP}}$ :

$$f(\theta) = f(\theta_{\text{MAP}}) + \nabla(\theta_{\text{MAP}})(\theta - \theta_{\text{MAP}}) + \frac{1}{2!}(\theta - \theta_{\text{MAP}})^T \mathbf{H}(\theta_{\text{MAP}})(\theta - \theta_{\text{MAP}}) + \dots$$

- 2nd-order approximation:

$$f(\theta) \approx f(\theta_{\text{MAP}}) + \frac{1}{2}(\theta - \theta_{\text{MAP}})^T \mathbf{H}(\theta_{\text{MAP}})(\theta - \theta_{\text{MAP}})$$

$$p(\theta | \mathcal{D}) \approx C \cdot \underbrace{\exp\left(\frac{1}{2}(\theta - \theta_{\text{MAP}})^T \mathbf{H}(\theta_{\text{MAP}})(\theta - \theta_{\text{MAP}})\right)}_{\mathcal{N}(\theta_{\text{MAP}}, -\mathbf{H}^{-1}(\theta_{\text{MAP}}))}$$



# Bayesian Learning of Logistic Regression

- a training set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{0, 1\}$
- likelihood function of logistic regression:

$$p(\mathcal{D} | \mathbf{w}) = \prod_{i=1}^N \left( l(\mathbf{w}^\top \mathbf{x}_i) \right)^{y_i} \left( 1 - l(\mathbf{w}^\top \mathbf{x}_i) \right)^{1-y_i}$$

- choose a Gaussian prior:  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \Sigma_0)$
- Bayesian learning:  $p(\mathbf{w} | \mathcal{D}) \propto p(\mathbf{w}) p(\mathcal{D} | \mathbf{w})$
- the posterior  $p(\mathbf{w} | \mathcal{D})$  is not Gaussian anymore
- use **Laplace's method** to approximate the true posterior

- use a gradient descent to find  $\mathbf{w}_{\text{MAP}}$

$$\nabla(\mathbf{w}) = \nabla \ln p(\mathbf{w} | \mathcal{D}) = -\Sigma_0^{-1}(\mathbf{w} - \mathbf{w}_0) + \sum_{i=1}^N (y_i - l(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i$$

- use a Gaussian approximation:

$$p(\mathbf{w} | \mathcal{D}) \approx \mathcal{N}\left(\mathbf{w} | \mathbf{w}_{\text{MAP}}, -\mathbf{H}^{-1}(\mathbf{w}_{\text{MAP}})\right)$$

$$\text{with } \mathbf{H}(\mathbf{w}) = -\Sigma_0^{-1} - \sum_{i=1}^N l(\mathbf{w}^\top \mathbf{x}_i)(1 - l(\mathbf{w}^\top \mathbf{x}_i)) \mathbf{x}_i \mathbf{x}_i^\top$$



# Variational Bayesian Methods (I)

- variational Bayesian (VB): use a simpler *variational distribution*  $q(\boldsymbol{\theta})$  to approximate the true posterior  $p(\boldsymbol{\theta} | \mathcal{D})$ :

$$q^*(\boldsymbol{\theta}) = \arg \min_q \text{KL}\left(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} | \mathcal{D})\right)$$

- $\text{KL}\left(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} | \mathcal{D})\right) = \ln p(\mathcal{D}) - \underbrace{\int_{\boldsymbol{\theta}} q(\boldsymbol{\theta}) \ln \frac{p(\mathcal{D}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}}_{L(q)}$

- $\min_q \text{KL}\left(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta} | \mathcal{D})\right) \iff \max_q L(q)$

- assume  $q(\boldsymbol{\theta}) = q_1(\boldsymbol{\theta}_1) q_2(\boldsymbol{\theta}_2) \cdots q_I(\boldsymbol{\theta}_I)$  can be factorized over some disjoint subsets  $\boldsymbol{\theta} = \boldsymbol{\theta}_1 \cup \boldsymbol{\theta}_2 \cup \cdots \cup \boldsymbol{\theta}_I$

- $L(q) = \int_{\boldsymbol{\theta}} \prod_{i=1}^I q_i(\boldsymbol{\theta}_i) \ln p(\mathcal{D}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \sum_{i=1}^I \int_{\boldsymbol{\theta}_i} q_i(\boldsymbol{\theta}_i) \ln q_i(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i$

## Variational Bayesian Methods (II)

- maximize  $L(q)$  w.r.t. each  $q_i(\boldsymbol{\theta}_i)$  separately

$$\max_{q_i} \int_{\boldsymbol{\theta}_i} q_i(\boldsymbol{\theta}_i) \left[ \underbrace{\int_{\boldsymbol{\theta}_{j \neq i}} \prod_{j \neq i} q_j(\boldsymbol{\theta}_j) \ln p(\mathcal{D}, \boldsymbol{\theta}) d\boldsymbol{\theta}_{j \neq i}}_{\mathbb{E}_{j \neq i} [\ln p(\mathcal{D}, \boldsymbol{\theta})]} \right] d\boldsymbol{\theta}_i - \int_{\boldsymbol{\theta}_i} q_i(\boldsymbol{\theta}_i) \ln q_i(\boldsymbol{\theta}_i) d\boldsymbol{\theta}_i$$

- define a new distribution:  $\tilde{p}(\boldsymbol{\theta}_i; \mathcal{D}) \propto \exp\left(\mathbb{E}_{j \neq i} [\ln p(\mathcal{D}, \boldsymbol{\theta})]\right)$

- we have  $q_i^*(\boldsymbol{\theta}_i) = \arg \max_{q_i} \int_{\boldsymbol{\theta}_i} q_i(\boldsymbol{\theta}_i) \ln \frac{\tilde{p}(\boldsymbol{\theta}_i; \mathcal{D})}{q_i(\boldsymbol{\theta}_i)} d\boldsymbol{\theta}_i$

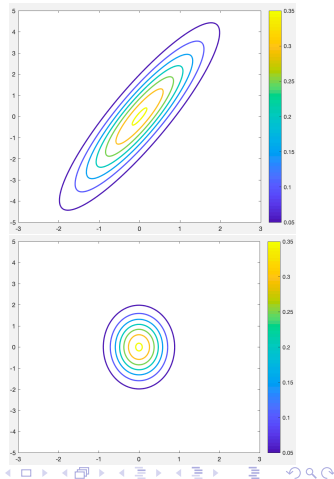
$$\implies q_i^*(\boldsymbol{\theta}_i) = \arg \min_{q_i} \text{KL}\left(q_i(\boldsymbol{\theta}_i) \parallel \tilde{p}(\boldsymbol{\theta}_i; \mathcal{D})\right)$$

- derive  $q_i^*(\boldsymbol{\theta}_i) = \tilde{p}(\boldsymbol{\theta}_i; \mathcal{D}) \propto \exp\left(\mathbb{E}_{j \neq i} [\ln p(\mathcal{D}, \boldsymbol{\theta})]\right)$  or

$$\ln q_i^*(\boldsymbol{\theta}_i) = \mathbb{E}_{j \neq i} [\ln p(\mathcal{D}, \boldsymbol{\theta})] + C$$

# Variational Bayesian Methods (III)

- mean field theory: use a factorizable variational distribution to approximate a true posterior distribution
- a 2-D Gaussian with  $\Sigma = \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix}$
- approximate with a variational distribution  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$
- best-fit is found by minimizing the KL-divergence



# Variational Bayesian Learning of GMMs (I)

- a Gaussian mixture model (GMM):

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{m=1}^M w_m \cdot \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

where model parameters  $\boldsymbol{\theta} = \{w_m, \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m | m = 1, 2, \dots, M\}$

- no conjugate prior exists for GMMs
- choose a prior distribution as

$$p(\boldsymbol{\theta}) = p(w_1, \dots, w_M) \prod_{m=1}^M p(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

with

$$p(w_1, \dots, w_M) = \text{Dir}(w_1, \dots, w_M | \alpha_1^{(0)}, \dots, \alpha_M^{(0)})$$

$$p(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = \text{GIW}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m | \boldsymbol{\nu}_m^{(0)}, \Phi_m^{(0)}, \lambda_m^{(0)}, \nu_m^{(0)})$$

## Variational Bayesian Learning of GMMs (II)

- introduce 1-of- $M$  latent variable  $\mathbf{z} = [z_1 z_2 \cdots z_M]$  for GMMs:

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = \prod_{m=1}^M (w_m)^{z_m} \left( \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \right)^{z_m}$$

- use the variational Bayesian method to approximate the posterior distribution  $p(\mathbf{z}, \boldsymbol{\theta} | \mathbf{x})$
- introduce a variational distribution factorized as:

$$q(\mathbf{z}, \boldsymbol{\theta}) = q(\mathbf{z})q(\boldsymbol{\theta}) = q(\mathbf{z}) q(w_1, \cdots, w_M) \prod_{m=1}^M q(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

- derive the best-fit variational distribution  $q^*(\mathbf{z}, \boldsymbol{\theta})$

# Variational Bayesian Learning of GMMs (III)

$$\begin{aligned}
 \mathbf{1} \quad \ln q^*(\mathbf{z}) &= \mathbb{E}_{\boldsymbol{\theta}} [\ln p(\mathbf{x}, \mathbf{z}, \boldsymbol{\theta})] + C = \mathbb{E}_{\boldsymbol{\theta}} [\ln p(\boldsymbol{\theta}) + \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})] + C \\
 &\Rightarrow \ln q^*(\mathbf{z}) = C' + \\
 &\sum_{m=1}^M z_m \underbrace{\left( \mathbb{E}[\ln w_m] - \mathbb{E}\left[\frac{\ln |\boldsymbol{\Sigma}_m|}{2}\right] - \mathbb{E}\left[\frac{(\mathbf{x} - \boldsymbol{\mu}_m)^\top \boldsymbol{\Sigma}_m^{-1} (\mathbf{x} - \boldsymbol{\mu}_m)}{2}\right] \right)}_{\ln \rho_m}
 \end{aligned}$$

- $q^*(\mathbf{z})$  is a multinomial:  $q^*(\mathbf{z}) \propto \prod_{m=1}^M (\rho_m)^{z_m} \propto \prod_{m=1}^M (r_m)^{z_m}$ ,  
where  $r_m = \frac{\rho_m}{\sum_{m=1}^M \rho_m}$  for all  $m$

$$\begin{aligned}
 \mathbf{2} \quad \ln q^*(w_1, \dots, w_M) &= \mathbb{E}_{\mathbf{z}, \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m} [\ln p(\boldsymbol{\theta}) + \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})] \\
 &= \sum_{m=1}^M (\alpha_m^{(0)} - 1) \ln w_m + \sum_{m=1}^M r_m \ln w_m + C \\
 &\quad \circ \quad q^*(w_1, \dots, w_M) \text{ is a Dirichlet distribution:}
 \end{aligned}$$

$$q^*(w_1, \dots, w_M) = \text{Dir}(w_1, \dots, w_M \mid \alpha_1^{(1)}, \dots, \alpha_M^{(1)})$$

where  $\alpha_m^{(1)} = \alpha_m^{(0)} + r_m$  for all  $m = 1, 2, \dots, M$

# Variational Bayesian Learning of GMMs (IV)

- 3  $\ln q^*(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = \mathbb{E}_{\mathbf{z}, w_m} [\ln p(\boldsymbol{\theta}) + \ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})] + C$   
 $= \ln p(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) + \mathbb{E}[z_m] \ln \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) + C'$
- $q^*(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$  is also a GIW distribution:

$$q^*(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = \text{GIW}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m | \boldsymbol{\nu}_m^{(1)}, \Phi_m^{(1)}, \lambda_m^{(1)}, \nu_m^{(1)})$$

where

$$\lambda_m^{(1)} = \lambda_m^{(0)} + r_m$$

$$\nu_m^{(1)} = \nu_m^{(0)} + r_m \mathbf{x}$$

$$\boldsymbol{\nu}_m^{(1)} = \frac{\lambda_m^{(0)} \boldsymbol{\nu}_m^{(0)} + r_m \mathbf{x}}{\lambda_m^{(0)} + r_m}$$

$$\Phi_m^{(1)} = \Phi_m^{(0)} + \frac{\lambda^{(0)} r_m}{\lambda_m^{(0)} + r_m} (\mathbf{x} - \boldsymbol{\nu}_m^{(0)}) (\mathbf{x} - \boldsymbol{\nu}_m^{(0)})^\top$$

# Variational Bayesian Learning of GMMs (V)

- based on the above distributions, we have

$$\ln \pi_m \triangleq \mathbb{E}[\ln w_k] = \psi(\alpha_m^{(1)}) - \psi\left(\sum_{m=1}^M \alpha_m^{(1)}\right)$$

$$\ln B_m \triangleq \mathbb{E}[\ln |\Sigma_m|] = \sum_{i=1}^d \psi\left(\frac{\lambda_m + 1 - i}{2}\right) - \ln |\Phi_m^{(1)}|$$

$$\mathbb{E}\left[(\mathbf{x} - \boldsymbol{\mu}_m)^\top \Sigma_m^{-1} (\mathbf{x} - \boldsymbol{\mu}_m)\right] = \frac{d}{\nu_m^{(1)}} + \lambda_m^{(1)} (\mathbf{x} - \boldsymbol{\nu}_m^{(1)})^\top (\Phi_m^{(1)})^{-1} (\mathbf{x} - \boldsymbol{\nu}_m^{(1)})$$

to compute  $\rho_m$  as well as  $r_m$  ( $\forall m = 1, 2, \dots, M$ )

- derive an EM-like algorithm to solve mutual dependency



# Variational Bayesian Learning of GMMs (VI)

## Variational Bayesian GMMs

**Input:**  $\{\alpha_m^{(0)}, \nu_m^{(0)}, \Phi_m^{(0)}, \lambda_m^{(0)}, \nu_m^{(0)} \mid m = 1, 2, \dots, M\}$

set  $n = 0$

**while** not converge **do**

**E-step:** collect statistics:

$$\{\alpha_m^{(n)}, \nu_m^{(n)}, \Phi_m^{(n)}, \lambda_m^{(n)}, \nu_m^{(n)}\} + \mathbf{x} \longrightarrow \{r_m\}$$

**M-step:** update all hyperparameters:

$$\begin{aligned} & \{\alpha_m^{(n)}, \nu_m^{(n)}, \Phi_m^{(n)}, \lambda_m^{(n)}, \nu_m^{(n)}\} + \{r_m\} + \mathbf{x} \\ \longrightarrow & \{\alpha_m^{(n+1)}, \nu_m^{(n+1)}, \Phi_m^{(n+1)}, \lambda_m^{(n+1)}, \nu_m^{(n+1)}\} \end{aligned}$$

$n = n + 1$

**end while**

# Non-Parametric Bayesian Methods

- Bayesian learning of parametric models: rely on prior/posterior distributions of model *parameters*
- how about Bayesian learning of non-parametric models?
- non-parametric Bayesian methods: use stochastic processes as priors for non-parametric models
  - **Gaussian processes**
  - Dirichlet processes

# Gaussian Processes: Concepts (I)

- given an arbitrary function  $f(\mathbf{x})$
- for any set of  $N$  points in  $\mathbb{R}^d$ , i.e.  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
- function values form an  $N$ -dimensional real-valued vector

$$\mathbf{f} = [f(\mathbf{x}_1) f(\mathbf{x}_2) \cdots f(\mathbf{x}_N)]^\top$$

- assume  $\mathbf{f}$  follows a multivariate Gaussian distribution

$$\mathbf{f} = [f(\mathbf{x}_1) f(\mathbf{x}_2) \cdots f(\mathbf{x}_N)]^\top \sim \mathcal{N}(\boldsymbol{\mu}_{\mathcal{D}}, \boldsymbol{\Sigma}_{\mathcal{D}})$$

where  $\boldsymbol{\mu}_{\mathcal{D}}$  and  $\boldsymbol{\Sigma}_{\mathcal{D}}$  depends on  $N$  data points in  $\mathcal{D}$

- it holds for any  $\mathcal{D}$ ,  $f(\mathbf{x})$  is a sample from a **Gaussian process**:

$$f(\mathbf{x}) \sim \text{GP}(\mathbf{m}(\mathbf{x}), \Phi(\mathbf{x}, \mathbf{x}'))$$

- $\mathbf{m}(\mathbf{x})$ : *mean function*  $\implies \boldsymbol{\mu}_{\mathcal{D}}$
- $\Phi(\mathbf{x}, \mathbf{x}')$ : *covariance function*  $\implies \boldsymbol{\Sigma}_{\mathcal{D}}$

# Gaussian Processes: Concepts (II)

- how to specify a Gaussian process?
- mean function  $\mathbf{m}(\mathbf{x}) = 0$
- covariance function: Mercer's condition

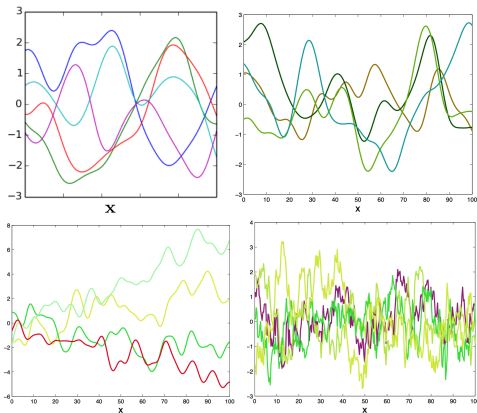
$$\Phi(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j))$$

$$\circ \Sigma_{\mathcal{D}} = \begin{bmatrix} \Phi(\mathbf{x}_i, \mathbf{x}_j) \end{bmatrix}_{N \times N}$$

- RBF kernel function

$$\Phi(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2l^2}}$$

- $\sigma$ : vertical scale
- $l$ : horizontal scale



# Gaussian Processes for Non-Parametric Bayesian Learning

- Gaussian processes as a non-parametric prior
  - randomly sample a function  $f(\cdot)$  from a Gaussian process
  - a prior can be implicitly computed with a data set

$$\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- function values  $\mathbf{f}$  follow a multivariate Gaussian distribution
- non-parametric prior:

$$p(f | \mathcal{D}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \Sigma_{\mathcal{D}})$$

- Gaussian processes for regression or classification
  - input-output pairs yield likelihood function
  - apply Bayesian learning rule:

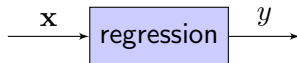
$$\text{posterior} \propto \text{prior} \times \text{likelihood}$$

# Gaussian Processes for Regression (I)

basic setting for regression:

- $f(\mathbf{x}) \sim \text{GP}(\mathbf{0}, \Phi(\mathbf{x}, \mathbf{x}'))$

- $y = f(\mathbf{x}) + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, \sigma_0^2)$



- given a training set:  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
- the corresponding outputs:  $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^\top$
- a non-parametric prior:

$$p(f | \mathcal{D}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \Sigma_{\mathcal{D}})$$

- the likelihood function due to the residual Gaussian noise  $\epsilon$ :

$$p(\mathbf{y} | f, \mathcal{D}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma_0^2 \mathbf{I})$$

## Gaussian Processes for Regression (II)

- Bayesian learning for the predictive distribution:

$$\begin{aligned} p(\mathbf{y} | \mathcal{D}) &= \int_f p(\mathbf{y}, f | \mathcal{D}) df = \int_f p(\mathbf{y} | f, \mathcal{D}) p(f | \mathcal{D}) df \\ &= \int_{\mathbf{f}} \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma_0^2 \mathbf{I}) \mathcal{N}(\mathbf{f} | \mathbf{0}, \Sigma_{\mathcal{D}}) d\mathbf{f} \\ &= \mathcal{N}(\mathbf{y} | \mathbf{0}, \Sigma_{\mathcal{D}} + \sigma_0^2 \mathbf{I}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{C}_N) \end{aligned}$$

- hyper-parameter learning:

$$\{\sigma^*, l^*, \sigma_0^*\} = \arg \max_{\sigma, l, \sigma_0} p(\mathbf{y} | \mathcal{D}, \sigma, l, \sigma_0) = \arg \max_{\sigma, l, \sigma_0} \ln \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{C}_N)$$

- may use a gradient descent method

# Gaussian Processes for Regression (III)

- predict output  $\tilde{y}$  for a new input  $\tilde{\mathbf{x}}$ :

$$p(\mathbf{y}, \tilde{y} | \mathcal{D}, \mathbf{x}) = \mathcal{N}(\mathbf{y}, \tilde{y} | \mathbf{0}, \mathbf{C}_{N+1})$$

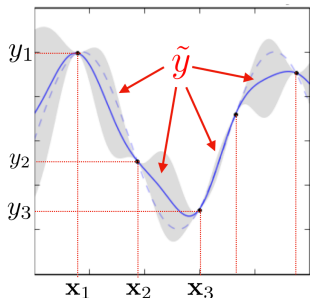
with

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^\top & \kappa^2 \end{bmatrix}$$

where  $\kappa^2 = \Phi(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}) + \sigma_0^2$  and  $\mathbf{k}_i = \Phi(\mathbf{x}_i, \tilde{\mathbf{x}})$

- the predictive distribution:

$$\begin{aligned} p(\tilde{y} | \mathcal{D}, \mathbf{y}, \tilde{\mathbf{x}}) &= \frac{p(\mathbf{y}, \tilde{y} | \mathcal{D}, \tilde{\mathbf{x}})}{p(\mathbf{y} | \mathcal{D})} \\ &= \mathcal{N}(\tilde{y} | \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{y}, \kappa^2 - \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{k}) \end{aligned}$$



point estimation (MAP or mean):

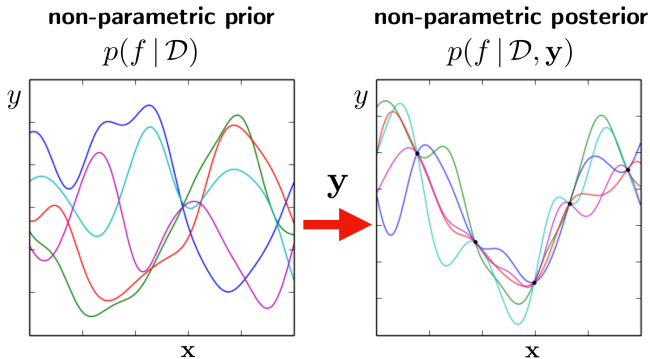
$$\begin{aligned} \mathbb{E}[\tilde{y} | \mathcal{D}, \mathbf{y}, \tilde{\mathbf{x}}] &= \tilde{y}_{\text{MAP}} \\ &= \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{y} \end{aligned}$$





# Gaussian Processes for Regression (IV)

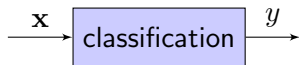
- derive a non-parametric prior from  $\mathcal{D}$  and  $\Phi(\mathbf{x}, \mathbf{x}')$
- non-parametric Bayesian learning based on  $\mathbf{y}$ :



# Gaussian Processes for Classification

basic setting for binary classification  $y \in \{0, 1\}$  :

- $f(\mathbf{x}) \sim \text{GP}(\mathbf{0}, \Phi(\mathbf{x}, \mathbf{x}'))$
- $\Pr(y = 1 | \mathbf{x}) = l(f(\mathbf{x})) = \frac{1}{1+e^{-f(\mathbf{x})}}$



- given a training set  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  and the corresponding outputs  $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^\top$
- non-parametric prior:  $p(f | \mathcal{D}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \Sigma_{\mathcal{D}})$
- likelihood:  $p(\mathbf{y} | f, \mathcal{D}) = \prod_{i=1}^N \left( l(f(\mathbf{x}_i)) \right)^{y_i} \left( 1 - l(f(\mathbf{x}_i)) \right)^{1-y_i}$
- no closed-form solution to derive the marginal and predictive distributions, i.e.  $p(\mathbf{y} | \mathcal{D})$  and  $p(\tilde{y} | \mathcal{D}, \mathbf{y}, \tilde{\mathbf{x}})$
- require approximate inference, such as Laplace's method