# Bayesian Inference for Gaussians (Contd)

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

# Multivariate Gaussian

- The (multivariate) Gaussian with mean $\boldsymbol{\mu}$ and cov. matrix $\boldsymbol{\Sigma}$

$$\mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left\{ -\frac{1}{2}(x - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(x - \boldsymbol{\mu}) \right\}$$

Trace notation

$$= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left\{ -\frac{1}{2}\mathrm{trace}\left[ \boldsymbol{\Sigma}^{-1}\mathbf{S} \right] \right\} \qquad \text{where } \mathbf{S} = (x - \boldsymbol{\mu})(x - \boldsymbol{\mu})^\top$$

- An alternate representation: The "information form"

Quadratic in $\boldsymbol{x}$

Linear in $\boldsymbol{x}$

$$\mathcal{N}_c(\boldsymbol{x}|\boldsymbol{\xi}, \boldsymbol{\Lambda}) = (2\pi)^{-D/2}|\boldsymbol{\Lambda}|^{1/2} \exp\left\{ -\frac{1}{2}\left( \boldsymbol{x}^\top \boldsymbol{\Lambda} \boldsymbol{x} + \boldsymbol{\xi}^\top \boldsymbol{\Lambda}^{-1}\boldsymbol{\xi} - 2\boldsymbol{x}^\top \boldsymbol{\xi} \right) \right\}$$
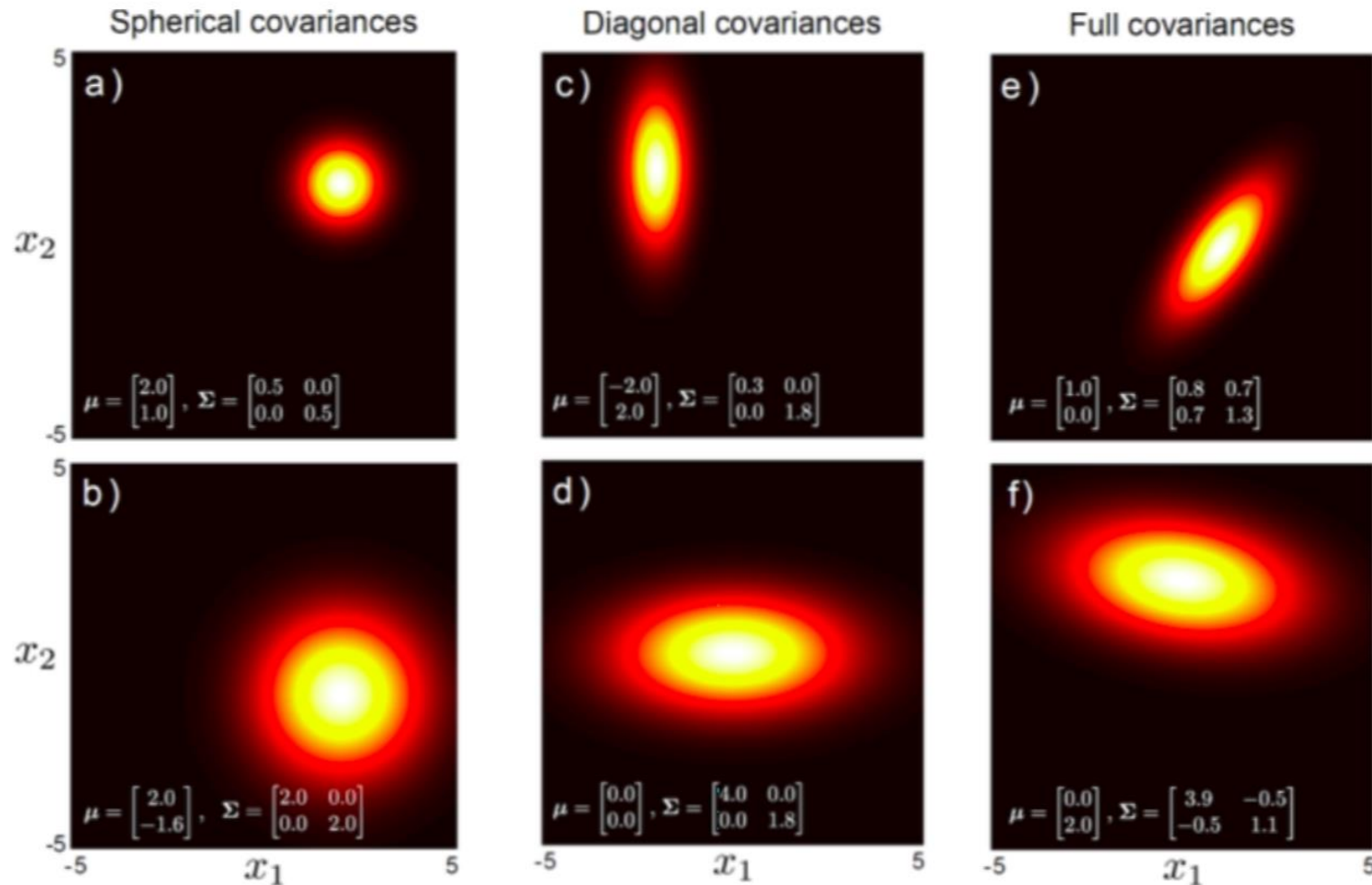
More when we discuss exp. family

where $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ and $\boldsymbol{\xi} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ are known as natural parameters of Gaussian

- Information form can help recognize $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ of a multivariate Gaussian when doing algebraic manipulations (e.g., when computing a posterior)

# Multivariate Gaussian

- The covariance matrix can be spherical, diagonal, or full

# Marginals and Conditionals from Gaussian Joint

- Assume $x$ having multivar Gaussian distr $\mathcal{N}(x|\mu, \Sigma)$ with $\Lambda = \Sigma^{-1}$. Suppose

$$x = \begin{bmatrix} x_a \\ x_b \end{bmatrix} \qquad \mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \qquad \Lambda = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}$$

- The marginal distribution of one block, say $x_a$, is a Gaussian

$$p(x_a) = \int p(x_a, x_b) dx_b = \mathcal{N}(x_a|\mu_a, \Sigma_{aa})$$

- The conditional distribution of $x_a$ given $x_b$, is Gaussian $\quad p(x_a|x_b) = \mathcal{N}(x_a|\mu_{a|b}, \Sigma_{a|b})$

Extremely useful results when working with Gaussian joint distributions

Note that $\Sigma_{a|b}$ is "smaller" than $\Sigma_{aa}$ (conditioning reduces variance)

$$\Sigma_{a|b} = \Lambda_{aa}^{-1} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}$$

$$\mu_{a|b} = \Sigma_{a|b}\{\Lambda_{aa}\mu_a - \Lambda_{ab}(x_b - \mu_b)\}$$

$$= \mu_a - \Lambda_{aa}^{-1}\Lambda_{ab}(x_b - \mu_b)$$

$$= \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(x_b - \mu_b)$$

# Linear Transformations of Random Variables

- Let $\boldsymbol{x} = f(\boldsymbol{z}) = \boldsymbol{A}\boldsymbol{z} + \boldsymbol{b}$ be a linear function of a vector r.v. $\boldsymbol{z}$

Need not be a Gaussian random var

- Suppose $\mathbb{E}[\boldsymbol{z}] = \boldsymbol{\mu}$ and $\text{cov}[\boldsymbol{z}] = \boldsymbol{\Sigma}$ then

$$\mathbb{E}[\boldsymbol{x}] = \mathbb{E}[\boldsymbol{A}\boldsymbol{z} + \boldsymbol{b}] = \boldsymbol{A}\boldsymbol{\mu} + \boldsymbol{b}$$

$$\text{cov}[\boldsymbol{x}] = \text{cov}[\boldsymbol{A}\boldsymbol{z} + \boldsymbol{b}] = \boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^{\top}$$

- Likewise if $x = f(\boldsymbol{z}) = \boldsymbol{a}^{\top}\boldsymbol{z} + b$ is a scalar-valued linear function of the above r.v. $\boldsymbol{z}$

$p(\boldsymbol{x})$ will also be Gaussian with mean and covariance/variance given by these expressions

$$\mathbb{E}[x] = \mathbb{E}[\boldsymbol{a}^{\top}\boldsymbol{z} + b] = \boldsymbol{a}^{\top}\boldsymbol{\mu} + b$$

$$\text{var}[x] = \text{var}[\boldsymbol{a}^{\top}\boldsymbol{z} + b] = \boldsymbol{a}^{\top}\boldsymbol{\Sigma}\boldsymbol{a}$$

Especially when $p(\boldsymbol{z})$ is Gaussian

- These properties are often helpful in obtaining the marginal distribution $p(\boldsymbol{x})$ from $p(\boldsymbol{z})$

# Linear Gaussian Model

Independently added and drawn from $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \boldsymbol{L}^{-1})$

- Consider linear transf. of a r.v. $\boldsymbol{z}$ with $p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$, plus Gaussian noise $\boldsymbol{\epsilon}$

$$\boldsymbol{x} = \boldsymbol{Az} + \boldsymbol{b} + \boldsymbol{\epsilon}$$

- Easy to see that, conditioned on $\boldsymbol{z}$, $\boldsymbol{x}$ too has a Gaussian distribution

$$p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{Az} + \boldsymbol{b}, \boldsymbol{L}^{-1})$$

- A Linear Gaussian Model. Very commonly encountered in probabilistic modeling

- The following two distributions are of interest. Assuming $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \boldsymbol{A}^\top \boldsymbol{L} \boldsymbol{A})^{-1}$

$$p(\boldsymbol{z}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z})}{p(\boldsymbol{z})} = \mathcal{N}(\boldsymbol{z}|\boldsymbol{\Sigma}\{\mathbf{A}^\top \mathbf{L}(\boldsymbol{x} - \boldsymbol{b}) + \boldsymbol{\Lambda}\mu\}, \boldsymbol{\Sigma})$$

If $p(\boldsymbol{z})$ is a prior and $p(\boldsymbol{x}|\boldsymbol{z})$ is likelihood then this is the posterior

$$p(\boldsymbol{x}) = \int p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z})d\boldsymbol{z} = \mathcal{N}(\boldsymbol{x}|\mathbf{A}\mu + \boldsymbol{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top + \mathbf{L}^{-1})$$

If $p(\boldsymbol{z})$ is a prior and $p(\boldsymbol{x}|\boldsymbol{z})$ is likelihood then this is the marginal likelihood

- Exercise: Prove the above results (MLAPP Chap. 4 and PRML Chap. 2 contain proof)

# Applications of Gaussian-based Models

- Gaussians and Linear Gaussian Models widely used in probabilistic models, e.g.,
  - Probability density estimation: Given $x_1, x_2, \ldots, x_N$, estimate $p(x)$ assuming Gaussian lik./noise
  - Given $N$ sensor obs. $\{x_n\}_{n=1}^N$ with $x_n = \mu + \epsilon_n$ (zero-mean Gaussian noise $\epsilon_n$) estimate the underlying true value $\mu$ (possibly along with the variance of the estimate of $\mu$)
  - Estimating missing data: $p(x_n^{\text{miss}} | x_n^{\text{obs}})$ or $\mathbb{E}[x_n^{\text{miss}} | x_n^{\text{obs}}]$
  - Linear Regression with Gaussian Likelihood

    Training feat. mat

    The prior $p(\boldsymbol{w})$ is Gaussian

    i.i.d. Gaussian noise

    Training responses

    $$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{w} + \boldsymbol{\epsilon}$$

  - Linear latent variable models (probabilistic PCA, factor analysis, Kalman filters) and their mixtures
  - Gaussian Processes (GP) extensively use Gaussian conditioning and marginalization rules

    $\boldsymbol{y} = \boldsymbol{f} + \text{noise}$          (GP assumes $\boldsymbol{f} = [f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_N)]$ is jointly Gaussian)

  - More complex models where parts of the model use Gaussian likelihoods/priors

# Coming Up Next

- Exponential Family distributions
- Conditional Models for supervised learning