

Advanced probabilistic methods

Lecture 3

Pekka Marttinen

Aalto University

February, 2019

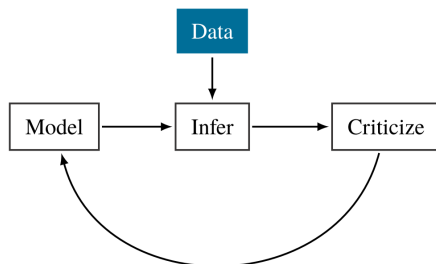
Lecture 3 overview¹

- Gaussian distribution
 - Bayesian parameter learning
- Multivariate Gaussian distribution
 - Characterization
 - Basic properties
- (Other important distributions)
- Ch. 8 in Barber's book

¹These slides build upon the book *Bayesian Reasoning and Machine Learning* and the associated teaching materials. The book and the demos can be downloaded from www.cs.ucl.ac.uk/staff/D.Barber/brml.

Recall from lecture 1

- Tools for probabilistic modeling
 - **Models:** Bayesian networks, sparse Bayesian linear regression, Gaussian mixture models, latent linear models
 - **Methods for inference:** maximum likelihood, maximum a posteriori (MAP), analytical, Laplace approximation, expectation maximization (EM), Variational Bayes (VB), stochastic variational inference (SVI)
 - **Ways to select between models**



Box's loop (Blei, 2014)

What is a model?

- A model specifies a probability distribution for a random variable Y , and it is often affected by some parameter θ . The model can be denoted as $p(y|\theta)$.
- Fitting the model (i.e. inference) corresponds to learning the value (or the distribution) of θ , after some data y have been observed.

- *Bayes' rule* tells us how to update our prior beliefs about variable θ in light of the data y to a posterior belief:

$$\underbrace{p(\theta|y)}_{\text{posterior}} = \frac{\underbrace{p(y|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}}{\underbrace{p(y)}_{\text{evidence}}}.$$

The evidence is also called the marginal likelihood.

- $p(y|\theta)$ is the probability that the model generates the observed data y when using parameter θ
 - $L(\theta) \equiv p(y|\theta)$, with y held fixed, is called the *likelihood*
 - $f(y) \equiv p(y|\theta)$, with θ held fixed, is called the *observation model*
- "*Methods for inference*" = Bayes' rule + some algorithm to do the actual computations (on this course)

Point estimates for parameters

- The *Maximum A Posteriori (MAP)* parameter value, which maximizes the posterior

$$\theta_* = \arg \max_{\theta} p(\theta|y)$$

- The Maximum likelihood assignment (ML)

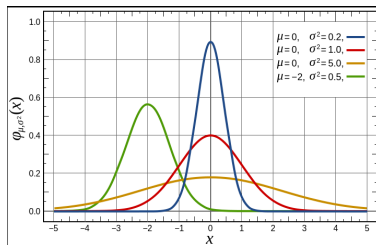
$$\theta_* = \arg \max_{\theta} p(y|\theta)$$

- The full posterior distribution $p(\theta|y)$ tells also of the uncertainty related to the value of θ .

Gaussian distribution

- $X \sim N(\mu, \sigma^2)$
- Parameters: μ : mean, σ^2 : variance
- Inverse of the variance, $\lambda = 1/\sigma^2$, is called the precision
- Standard deviation σ
- 95% credible interval equals approximately $[\mu - 2\sigma, \mu + 2\sigma]$
- PDF:

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



Gaussian (or normal) distribution (wikip.)

Bayesian estimation of the mean of a Gaussian (1/2)

- Suppose we have observations $x = (x_1, \dots, x_n)$ from $N(\mu, \sigma^2)$, where σ^2 is known.
- To learn μ , we specify a prior

$$\mu \sim N(\mu_0, \tau_0^2)$$

- Posterior

$$\begin{aligned} p(\mu|x) &= \frac{p(x|\mu)p(\mu)}{p(x)} \propto p(\mu)p(x|\mu) \\ &= \frac{1}{\sqrt{2\pi\tau_0}} e^{-\frac{1}{2\tau_0^2}(\mu-\mu_0)^2} \times \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2} \\ &\propto e^{-\frac{1}{2\tau_0^2}(\mu-\mu_0) - \frac{1}{2\sigma^2} \sum_i (x_i-\mu)^2} \\ &= \dots \text{(details in BDA course)} \end{aligned}$$

Bayesian estimation of the mean of a Gaussian (2/2)

- Posterior

$$\begin{aligned} p(\mu|x) &\propto e^{-\frac{1}{2\tau_n^2}(\mu-\mu_n)^2} \\ &\propto N(\mu|\mu_n, \tau_n^2) \end{aligned}$$

where

$$\mu_n = \frac{\frac{1}{\tau_0^2}\mu_0 + \frac{n}{\sigma^2}\bar{x}}{\frac{n}{\sigma^2} + \frac{1}{\tau_0^2}} \quad \text{and} \quad \frac{1}{\tau_n^2} = \frac{n}{\sigma^2} + \frac{1}{\tau_0^2}.$$

- Posterior precision $1/\tau_n^2$: sum of prior precision $1/\tau_0^2$ and data precision n/σ^2
- Posterior mean μ_n : precision weighted average of prior mean μ_0 and data mean \bar{x} .

Conjugate prior distributions (1/2)

- In the previous example

$$\text{Prior: } \mu \sim N(\mu_0, \tau_0^2)$$

$$\text{Posterior: } \mu \sim N(\mu_n, \tau_n^2).$$

If the prior and posterior belong to the same family of distributions, we say that the prior is conjugate to the likelihood used.

- For example, normal prior $\mu \sim N(\mu_0, \tau_0^2)$ is conjugate to the normal likelihood $N(x|\mu, \sigma^2)$.
- Conjugacy is useful, because it makes computations easy.

Conjugate prior distributions (2/2)

- With conjugate prior, the posterior is available in a closed form

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

- Drop all terms not depending on θ
- Recognize the result as a density function belonging to the same family of distributions as the prior $p(\theta)$, but with different parameters.
- Examples (likelihood - conjugate prior):
 - Likelihood for normal mean - Normal prior
 - Likelihood for normal variance - Inverse-Gamma prior
 - Bernoulli - Beta
 - Binomial - Beta
 - Exponential - Gamma
 - Poisson - Gamma

Conjugate prior example (1/2)

- Suppose we have observations $x = (x_1, \dots, x_n)$ from $N(\mu, \lambda^{-1})$, where μ is known.
- To learn the precision λ , we specify a prior

$$\lambda \sim \text{Gam}(a, b)$$

Gamma distribution

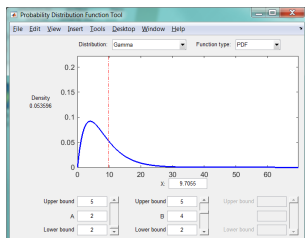
- Distribution for positive real values.

$$\lambda \sim \text{Gam}(a, b), \quad a > 0 : \text{shape}, \quad b > 0 : \text{rate}$$

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} e^{-b\lambda}$$

- Alternative parameterization uses $\lambda \sim \text{Gam}(a, \theta)$, $\theta = 1/b$ is called the **scale**

$$\text{Gam}(\lambda|a, \theta) = \frac{1}{\Gamma(a)\theta^a} \lambda^{a-1} e^{-\lambda/\theta}$$



disttool in Matlab

Conjugate prior example (2/2)

- Observations $x = (x_1, \dots, x_n)$ from $N(\mu, \lambda^{-1})$, where μ is known; $\lambda \sim \text{Gam}(a, b)$.

$$\begin{aligned} p(\lambda|x) &\propto p(x|\lambda)p(\lambda) \\ &= \prod_{i=1}^n \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2}(x_i-\mu)^2} \times \frac{1}{\Gamma(a)} b^a \lambda^{a-1} e^{-b\lambda} \\ &\propto \lambda^{\frac{n}{2}} e^{-\frac{\lambda}{2} \sum_i (x_i-\mu)^2} \times \lambda^{a-1} e^{-b\lambda} \\ &= \lambda^{\frac{n}{2}+a-1} e^{-\lambda[\frac{1}{2} \sum_i (x_i-\mu)^2 + b]} \\ &\propto \text{Gam}(\lambda|a_n, b_n), \end{aligned}$$

with

$$\begin{aligned} a_n &= a + \frac{n}{2} \\ b_n &= b + \frac{1}{2} \sum_i (x_i - \mu)^2 \end{aligned}$$

Gaussian distribution, unknown mean and precision (1/2)

- Suppose we have observations $x = (x_1, \dots, x_n)$ from $N(\mu, \lambda^{-1})$, where both the mean μ and the precision λ are unknown.
- The conjugate prior distribution is the normal-gamma distribution

$$\begin{aligned} p(\mu, \lambda | \mu_0, \beta, a, b) &= N(\mu | \mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda | a, b) \\ &\equiv \text{Normal-Gamma}(\mu, \lambda | \mu_0, \beta, a, b) \end{aligned}$$

Note the dependency of the prior of μ on the value of λ .

Gaussian distribution, unknown mean and precision (2/2)

- The conjugate prior distribution is the normal-gamma distribution

$$p(\mu, \lambda | \mu_0, \beta, a, b) = \text{Normal-Gamma}(\mu, \lambda | \mu_0, \beta, a, b)$$

- Posterior

$$p(\mu, \lambda | x) = \text{Normal-Gamma}(\mu, \lambda | \mu_n, \beta_n, a_n, b_n),$$

with

$$\mu_n = \frac{\beta \mu_0 + n \bar{x}}{\beta + n}$$

$$\beta_n = \beta + n$$

$$a_n = a + \frac{n}{2}$$

$$b_n = b + \frac{1}{2} \left(ns + \frac{\beta n (\bar{x} - \mu_0)^2}{\beta + n} \right)$$

Gaussian distribution, unknown mean and precision, example (1/2)

- Simulate samples from $N(\mu = 2, \sigma^2 = 0.25)$
 - precision $\lambda = 4$
- Try to learn μ and λ
- Specify prior

$$p(\mu, \lambda | \mu_0, \beta, a, b) = \text{Normal-Gamma}(\mu, \lambda | \mu_0, \beta, a, b)$$

with

$$\mu_0 = 0, \quad \beta = 0.001, \quad a = 0.01, \quad b = 0.01$$

- See: *normal_example.m*

Gaussian distribution, unknown mean and precision, example (2/2)

- When μ and λ have distribution

$$\text{Normal-Gamma}(\mu, \lambda | \mu_n, \beta_n, a_n, b_n) = N(\mu | \mu_n, (\beta_n \lambda)^{-1}) \text{Gam}(\lambda | a_n, b_n),$$

marginal distribution of λ can be plotted using the PDF of $\text{Gam}(\lambda | a_n, b_n)$

- To plot the marginal distribution of μ , we need to take the dependence on λ into account.
 - we compute the marginal distribution of μ by averaging over $N(\mu | \mu_n, (\beta_n \lambda_i)^{-1})$, for multiple λ_i simulated from $\text{Gam}(\lambda | a_n, b_n)$
 - (could also be done analytically...)

- If $p(x|\theta_t)$ is the true data generating mechanism, and A is a neighborhood of θ_t , then

$$p(\theta \in A|x) \xrightarrow{n \rightarrow \infty} 1.$$

- The posterior distribution concentrates around the true value (if such a value exists!). See the *normal_example.m*
- It follows that

$$\bar{\theta}_{MAP} \xrightarrow{n \rightarrow \infty} \theta_t \quad \text{and} \quad \bar{\theta}_{ML} \xrightarrow{n \rightarrow \infty} \theta_t$$

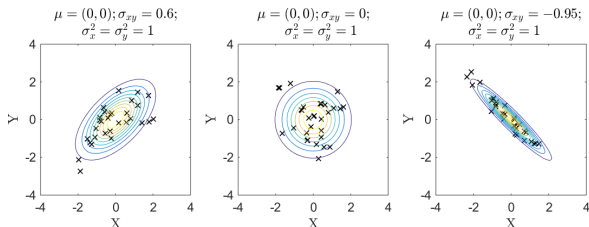
Multivariate Gaussian distribution

$$N_D(x|\mu, \Sigma) \equiv (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

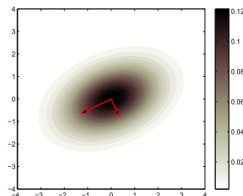
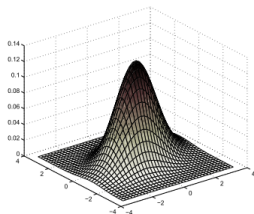
- D : dimension, μ : mean, Σ : covariance matrix. With $D = 2$:

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$$

- $\sigma_{12} = \sigma_{21}$: covariance between x_1 and x_2 . (tells direction of dependency)
- $\rho_{12} = \sigma_{12} / (\sigma_1 \sigma_2)$: correlation between x_1 and x_2 . (direction and strength)



Multivariate Gaussian - characterization (1/2)



- Eigendecomposition

$$\Sigma = E\Lambda E^T,$$

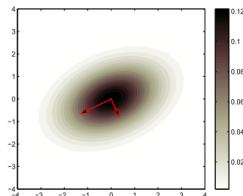
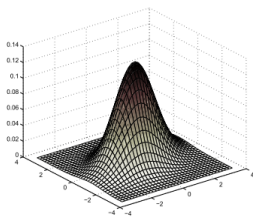
where $E^T E = I$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_D)$.

- Now the transformation

$$y = \Lambda^{-\frac{1}{2}} E^T (x - \mu)$$

can be shown to have the distribution $N_D(0, I)$ (product of D independent standard Gaussians)

Multivariate Gaussian - characterization (2/2)



- Thus, $x = E\Lambda^{\frac{1}{2}}y + \mu$ with distribution $N_D(\mu, \Sigma)$ is obtained from standard independent Gaussians y by
 - *scaling* by the square roots of eigenvalues
 - *rotating* by the eigenvectors
 - *shifting* by adding the mean

Marginalization and conditioning (1/2)

- Let $z \sim N(\mu, \Sigma)$ and consider partitioning it as:

$$z = \begin{pmatrix} x \\ y \end{pmatrix}$$

with

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}.$$

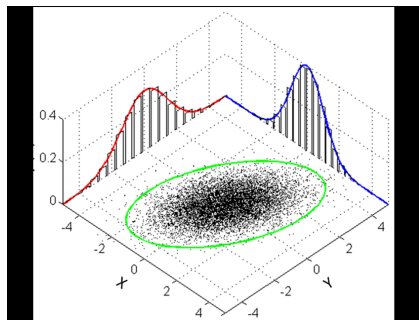
Marginalization and conditioning (2/2)

- Then

$$p(x) \sim N(\mu_x, \Sigma_{xx}) \quad (\text{marginalization})$$

$$p(x|y) = N(\mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}) \quad (\text{conditioning})$$

\implies Marginals and conditionals of M-V Gaussians are still M-V Gaussian.



- **Linear transformation:** if

$$y = Mx + \eta,$$

where $x \sim N(\mu_x, \Sigma_x)$ and $\eta \sim N(\mu, \Sigma)$, then

$$p(y) = N(y | M\mu_x + \mu, M\Sigma_x M^T + \Sigma)$$

- **Completing the square:**

$$\frac{1}{2}x^T Ax - b^T x = \frac{1}{2}(x - A^{-1}b)^T A(x - A^{-1}b) - \frac{1}{2}b^T A^{-1}b$$

From which one can derive, for example

$$\int \exp\left(-\frac{1}{2}x^T Ax + b^T x\right) dx = \sqrt{\det(2\pi A^{-1})} \exp\left(\frac{1}{2}b^T A^{-1}b\right)$$

- Let $x = (x_1, \dots, x_n)$ be from $N(\mu, \Sigma)$ with unknown μ and Σ .
Log-likelihood, assuming data are *i.i.d.*:

$$\begin{aligned} L(\mu, \Sigma) &= \sum_{i=1}^N \log p(x_i | \mu, \Sigma) \\ &= -\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{N}{2} \log \det(2\pi\Sigma) \end{aligned}$$

Multivariate Gaussian - ML fitting

- Differentiate $L(\mu, \Sigma)$ w.r.t. the vector μ :

$$\nabla_{\mu} L(\mu, \Sigma) = \sum_{i=1}^N \Sigma^{-1} (x_i - \mu)$$

Equating to zero gives

$$\sum_{i=1}^N \Sigma^{-1} x_i = N \Sigma^{-1} \mu.$$

Thus we get

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

- Similarly one can derive:

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

- Gaussian-Wishart is the conjugate prior, when $X_i \sim N(\mu, \Lambda)$ and both mean μ and precision Λ are unknown:

$$p(\mu, \Lambda | \mu_0, \beta, W, \nu) = N(\mu | \mu_0, (\beta\Lambda)^{-1}) \mathcal{W}(\Lambda | W, \nu)$$

- If X_i are scalar, this is equivalent to the Gaussian-Gamma distribution.
- Posterior

$$p(\mu, \Lambda | x) = N(\mu | \mu_n, (\beta_n\Lambda)^{-1}) \mathcal{W}(\Lambda | W_n, \nu_n)$$

- Wishart distribution is a distribution for nonnegative-definite matrix-valued random variables

$$\Lambda \sim \mathcal{W}(\Lambda|W, \nu)$$

$$E(\Lambda) = \nu W$$

$$\text{Var}(\Lambda_{ij}) = n(w_{ij}^2 + w_{ii}w_{jj})$$

- Further: exercises...

- Bayesian learning of the Gaussian distribution using conjugate priors
- Multivariate Gaussian
 - Characterization
 - Marginal & conditional distributions
 - Linear transformations & completing the square
 - ML-fitting