

Reading Notes for ch2 Probability Distributions

Xiang Pan

September 10, 2021

1 Priors

1.1 Parametric Priors

Conjugate priors: lead to posterior distributions having the same functional form as the prior.

- Multinomial distribution: Dirichlet distribution
- Gaussian distribution: Gaussian distribution

For the Exponential Family, conjugate prior,

$$p(\boldsymbol{\eta} \mid \boldsymbol{\chi}, \nu) = f(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp \{ \nu \boldsymbol{\eta}^\top \boldsymbol{\chi} \}$$

conjugate posterior,

$$p(\boldsymbol{\eta} \mid \mathbf{X}, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+N} \exp \left\{ \boldsymbol{\eta}^\top \left(\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + \nu \boldsymbol{\chi} \right) \right\}$$

For the parameter estimation, posterior mean of θ , averaged over the distribution generating the data, is equal to the prior mean of θ .

$$\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\theta}] = \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\theta} \mid \mathcal{D}]]$$

$$\text{var}_{\boldsymbol{\theta}}[\boldsymbol{\theta}] = \mathbb{E}_{\mathcal{D}} [\text{var}_{\boldsymbol{\theta}}[\boldsymbol{\theta} \mid \mathcal{D}]] + \text{var}_{\mathcal{D}} [\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\theta} \mid \mathcal{D}]]$$

1.2 Nonparametric Density Estimation

Notation Region R , total number K of points that lie inside R . V is the volume of R , N is the dataset number.

We can decide the K and V from the data by fixing V and determining the value of K from the data (kernel density estimator) or fixed value of K and use the data to find an appropriate value for V (Nearest-neighbour methods).

kernel density estimator (e.g. Gaussian kernel density estimator): Gaussian density model is obtained by placing a Gaussian over each data point and then adding up the contributions over the whole data set, and then dividing by N .

Nearest-Neighbour methods ¹

2 Gaussian Distribution

Multivariate Gaussian distribution is that if two sets of variables are jointly Gaussian, then the conditional distribution of one set conditioned on the other is again Gaussian. Similarly, the marginal distribution of either set is also Gaussian.

¹Limited to the page limit, you can get the full version note at https://github.com/Xiang-Pan/NYU_Baysian_Machine_Learning/blob/master/reading_notes/ch2/note2.pdf

2.1 Conditional Gaussian

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1})$$
$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

2.2 Marginal Gaussian

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

2.3 Mixtures of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
$$\ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

The maximum likelihood solution for the parameters no longer has a closed-form analytical solution.

3 Continuous Estimation

Robbins-Monro procedure for parameter estimation:

$$\theta^{(N)} = \theta^{(N-1)} + a_{N-1} \frac{\partial}{\partial \theta^{(N-1)}} \ln p(x_N | \theta^{(N-1)})$$