Maximum Likelihood and Bayesian Methods

Shyue Ping Ong

University of California, San Diego

NANO281

Overview

- Preliminaries
- Maximum Likelihood Inference
- Bayesian Methods
- 4 Expectation Maximization

Preliminaries

- Thus far, we have discussed ML models in the context of minimizing a loss function, typically a sum of squares.
- These approaches have a probabilistic interpretation and are instances of the maximum likelihood (ML) approach to fitting.

Maximum Likelihood Inference

• Consider a probability distribution/mass function for the observations:

$$z_i \sim g_{\theta}(z)$$

- Also, known as a parametric model for Z, with θ being unknown parameter that govern the distribution.
- Let us assume that Z has a normal distribution $(\theta = (\mu, \sigma^2))$:

$$g_{ heta}=rac{1}{\sqrt{2\pi}\sigma}e^{-rac{(z-\mu)^2}{2\sigma^2}}$$

• The likelihood of the observed data under model g_{θ} is then given by the *likelihood* function:

$$L(\theta; \mathsf{Z}) = \prod_{i=1}^{N} g_{\theta}(z_i)$$

Maximum Likelihood Inference, cont.

• Consider the log of *L*:

$$\ell(\theta;\mathsf{Z}) = \sum_{i=1}^{N} \log g_{\theta}(z_i)$$

- Score function $\dot{\ell}(\theta; Z) = 0$ at maximum (dot means derivative).
- Information matrix:

$$\mathsf{I}(\theta) = -\sum_{i=1}^{N} \frac{\partial^{2} \ell(\theta; z_{i})}{\partial \theta \partial \theta^{T}}$$

- Observed information = $I(\hat{\theta})$ ($\hat{\theta}$ is the estimated parameter)
- Fisher information $i(\theta) = E_{\theta}[I(\theta)]$, widely used in optimal experimental design.

Linear Regression Revisited as MLE

- $Y = \beta_0 + \beta_1 X + \varepsilon$
- $\varepsilon \sim N(0, \sigma^2)$, i.e., independent Gaussian noise.

$$L(\beta_{0}, \beta_{1}, \sigma; y) = \prod_{i=1}^{N} P(y_{i}|x_{i}; \hat{\beta}_{0}, \hat{\beta}_{1}, \hat{\sigma})$$

$$= \prod_{i=1}^{N} \frac{1}{\hat{\sigma}\sqrt{2\pi}} e^{-\frac{(y_{i}-\hat{\beta}_{0}-\hat{\beta}_{1}x_{i})^{2}}{2\hat{\sigma}^{2}}}$$

$$\ell(\beta_{0}, \beta_{1}, \sigma; y) = \sum_{i=1}^{N} \log \frac{1}{\hat{\sigma}\sqrt{2\pi}} e^{-\frac{(y_{i}-\hat{\beta}_{0}-\hat{\beta}_{1}x_{i})^{2}}{2\hat{\sigma}^{2}}}$$

$$= -\frac{N}{2} \log \hat{\sigma}^{2} 2\pi - \sum_{i=1}^{N} \frac{(y_{i}-\hat{\beta}_{0}-\hat{\beta}_{1}x_{i})^{2}}{2\hat{\sigma}^{2}}$$

Clearly, MLE in this case is equivalent to minimizing least squares.

Bayesian Methods

• Posterior distribution given by:

$$P(\theta|\mathsf{Z}) = \frac{P(\mathsf{Z}|\theta)P(\theta)}{\int P(\mathsf{Z}|\theta)P(\theta)d\theta}$$

ullet Updated knowledge about heta after seeing data.

Expectation Maximization

- Popular approach for solving MLE problems.
- Algorithm (example for two-component Gaussian mixture):
 - **1** Start with initial guesses for parameters (e.g, two random y_i).
 - Expectation step: Compute the responsibilities.

$$\hat{\gamma_i} = rac{\hat{\pi}\phi_{\hat{ heta_2}}(y_i)}{(1-\hat{\pi})\phi_{\hat{ heta_1}}(y_i) + \hat{\pi}\phi_{\hat{ heta_2}}(y_i)}, i = 1, 2, ..., N$$

Maximization step: Compute weighted means and variances.

$$\hat{\mu_{1}} = \frac{\sum_{i=1}^{N} (1 - \hat{\gamma_{i}}) y_{i}}{\sum_{i=1}^{N} (1 - \hat{\gamma_{i}})}, \hat{\sigma_{1}^{2}} = \frac{\sum_{i=1}^{N} (1 - \hat{\gamma_{i}}) (y_{i} - \hat{\mu_{1}})^{2}}{\sum_{i=1}^{N} (1 - \hat{\gamma_{i}})}$$

$$\hat{\mu_{2}} = \frac{\sum_{i=1}^{N} \hat{\gamma_{i}} y_{i}}{\sum_{i=1}^{N} \hat{\gamma_{i}}}, \hat{\sigma_{2}^{2}} = \frac{\sum_{i=1}^{N} \hat{\gamma_{i}} (y_{i} - \hat{\mu_{2}})^{2}}{\sum_{i=1}^{N} \hat{\gamma_{i}}}$$

with mixing probability $\hat{\pi} = \sum_{i=1}^{N} \hat{\gamma}_i / N$.

Iterate until convergence.

Gaussian mixture using EM in scikit-learn

The End