

CS 340 Lec. 13: Maximum Likelihood

AD

February 2011

- You are given data $\{\mathbf{x}^i\}_{i=1}^N$ ($\{\mathbf{x}^i, y^i\}_{i=1}^N$ in the supervised learning case).
- You have a probabilistic model for the data; i.e. typically in most learning problem

$$p(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N | \theta) = \prod_{i=1}^N p(\mathbf{x}^i | \theta)$$

- **Aim:** you want to pick the best $\theta \in \Theta$.
- Two main approaches considered here: Maximum Likelihood and Bayesian.

- The most standard approach consists of selecting

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta \in \Theta} p(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N | \theta) \\ &= \arg \max_{\theta \in \Theta} \log p(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N | \theta)\end{aligned}$$

- You select the value of $\theta \in \Theta$ that maximizes the probability of observing $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N)$.
- Example:** Assume independent $\{\mathbf{x}^i\}$ where $\mathbf{x}^i = x_1^i = x^i$ with

$$p(x | \theta) = \theta^{\mathbb{I}(x=1)} (1 - \theta)^{\mathbb{I}(x=0)}$$

then $\hat{\theta}_{MLE} = \sum_{i=1}^N \mathbb{I}(x^i = 1) / N$.

Maximum Likelihood for Poisson Data

- **Example:** Assume you have independent Poisson observations $\{x^i\}_{i=1}^N$ such that

$$p(x|\theta) = e^{-\theta} \frac{\theta^x}{x!}$$

for $\theta > 0$ and $x = 0, 1, 2, \dots$

- In this case, we have

$$\begin{aligned} l(\theta) &= \log p(x^{1:N} | \theta) \\ &= -N\theta + \log \theta \sum_{i=1}^N x^i - \sum_{i=1}^N \log x^i! \end{aligned}$$

- By setting $\frac{\partial l(\theta)}{\partial \theta} = 0$, we obtain

$$\hat{\theta}_{MLE} = \frac{\sum_{i=1}^N x^i}{N}.$$

Maximum Likelihood for Gaussian Data

- **Example:** Assume you independent observations $\{x^i\}_{i=1}^N$ such that

$$p(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

where $\theta = (\sigma^2, \mu)$.

- We have

$$\begin{aligned}l(\theta) &= \log p(x^{1:N}|\theta) \\ &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x^i - \mu)^2\end{aligned}$$

- By setting $\frac{\partial l(\theta)}{\partial \mu} = 0$ and $\frac{\partial l(\theta)}{\partial \sigma^2} = 0$, we obtain

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x^i, \quad \hat{\sigma}^2_{MLE} = \frac{1}{N} \sum_{i=1}^N (x^i - \hat{\mu}_{MLE})^2.$$

Maximum Likelihood for Multinomial Data

- **Example:** Assume you have independent observations $\{\mathbf{x}^i\}_{i=1}^N$ such that

$$p(\mathbf{x}|\theta) = \frac{P!}{\prod_{i=1}^d x_k!} \prod_{k=1}^d \theta_k^{x_k}$$

for $\theta_k > 0$, $\sum_{k=1}^d \theta_k = 1$ and $x_k = 0, 1, 2, \dots, P$ with $\sum_k x_k = P$.

- In this case, we have

$$\begin{aligned} l(\theta) &= \log p(\mathbf{x}^{1:N}|\theta) \\ &= \sum_{i=1}^N \log \left(\frac{P!}{\prod_{k=1}^d x_k^i!} \right) + \sum_{k=1}^d \left(\sum_{i=1}^N x_k^i \right) \log \theta_k \end{aligned}$$

- Be careful: It is a constrained optimization problem as $\sum_{k=1}^d \theta_k = 1$.

Maximum Likelihood for Multinomial Data

- We introduce a Lagrange multiplier λ and propose to maximize instead w.r.t θ and λ

$$l(\theta, \lambda) = l(\theta) + \lambda \left(1 - \sum_{k=1}^d \theta_k \right).$$

- Setting $\frac{\partial l(\theta, \lambda)}{\partial \lambda} = 0 \Rightarrow \sum_{k=1}^d \theta_k = 1$ and setting

$$\frac{\partial l(\theta, \lambda)}{\partial \theta_i} = 0 \Rightarrow \frac{\sum_{i=1}^N x_k^i}{\theta_k} - \lambda = 0 \Leftrightarrow \lambda \theta_k = \sum_{i=1}^N x_k^i$$

- It follows that, as $\sum_{k=1}^d \theta_k = 1$, then $\lambda = \left(\sum_{k=1}^d \sum_{i=1}^N x_k^i \right)$

$$\hat{\theta}_{k, MLE} = \frac{\sum_{i=1}^N x_k^i}{\lambda} = \frac{\sum_{i=1}^N x_k^i}{\sum_{i=1}^N \sum_{k=1}^d x_k^i}.$$

Application to Naive Bayes

- We assume that we have N data $(\mathbf{x}^i, y^i) \in \mathbb{N}^d \times \{0, 1\}^C$ and we use the model

$$p(\mathbf{x}, y = c | \theta) = \pi_c \frac{P!}{\prod_{i=1}^d x_k!} \prod_{k=1}^d \theta_{k,c}^{x_k}$$

where $\theta = (\pi_1, \dots, \pi_C, \theta_{1,1}, \dots, \theta_{d,1}, \dots, \theta_{1,C}, \dots, \theta_{d,C})$.

- We have

$$\begin{aligned} l(\theta) &= \sum_{k=1}^n \log p(\mathbf{x}^i, y^i | \theta) \\ &= \sum_{c=1}^C \left(\sum_{k=1}^N \mathbb{I}(y_k^i = c) \right) \left\{ \log \pi_c + \sum_{k=1}^d x_k^i \log \theta_{k,c} + Cste \right\} \end{aligned}$$

yields with $N_c = \sum_{k=1}^N \mathbb{I}(y_k^i = c)$

$$\hat{\pi}_{c,MLE} = \frac{N_c}{N}, \quad \hat{\theta}_{k,c,MLE} = \frac{\sum_{i=1}^N x_k^i \mathbb{I}(y_k^i = c)}{\sum_{k=1}^d \sum_{i=1}^N x_k^i \mathbb{I}(y_k^i = c)}.$$

Asymptotics of Maximum Likelihood Estimate

- Assume you have independent data $\{\mathbf{x}^i\}_{i=1}^N$ distributed according to $p(\mathbf{x}|\theta^*)$; i.e. θ^* is the true parameter. Under regularity assumptions, we have $\hat{\theta}_{MLE} \rightarrow \theta^*$ as $N \rightarrow \infty$.
- This follows from the fact that first

$$\frac{l(\theta)}{N} = \frac{1}{N} \sum_{i=1}^N \log p(\mathbf{x}^i|\theta) \xrightarrow{N \rightarrow \infty} \bar{l}(\theta) = \int p(\mathbf{x}|\theta^*) \log p(\mathbf{x}|\theta) d\mathbf{x}.$$

- Second, the average log-likelihood $\bar{l}(\theta)$ is maximized θ^* ; for any $\theta \in \Theta$ as

$$\begin{aligned} \bar{l}(\theta) - \bar{l}(\theta^*) &= \int p(\mathbf{x}|\theta^*) \log \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x}|\theta^*)} d\mathbf{x} \\ &\leq \log \left(\int p(\mathbf{x}|\theta^*) \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x}|\theta^*)} d\mathbf{x} \right) \quad (\text{Jensen's inequality}) \\ &\leq 0. \end{aligned}$$

Limitations of Maximum Likelihood

- Maximum likelihood estimation overfits!
- Suppose we are N Bernoulli data such that $\hat{\theta}_{MLE} = \sum_{i=1}^N x^i / N = 0$ then we have

$$p(x = 1 | \hat{\theta}_{MLE}) = 0.$$

- Similarly, suppose we have N multinomial data such that $\hat{\theta}_{k,MLE} = \sum_{i=1}^N x_k^i / \sum_{i=1}^N \sum_{k=1}^d x_k^i = 0$ then we have

$$p(x_1, \dots, x_{k-1}, x_k = 1, x_{k+1}, \dots, x_d | \hat{\theta}_{MLE}) = 0.$$

- Hence if we have not observed such events in our training set, we predict that we will never observe them, ever!
- Failing to predict that certain events are possible is analogous to a problem in philosophy called the black swan paradox. This is based on the ancient Western conception that all swans were white. In that context, a black swan was a metaphor for something that could not exist. (Black swans were discovered in Australia by European explorers in the 17th Century.)