Estimating Probabilities and Naive Bayes

Shuhao Xia

Shanghaitech University

March 30, 2020

Outline

Estimating Probabilities

Bayes rule Maximum Likelihood Estimate (MLE) Maximum a Posterior (MAP)

Naive Bayes

Conditional Independence Naive Bayes for Discrete Inputs Naive Bayes for Continuous Inputs

Bayes rule

Given observations D, our goal is to estimate the parameter θ . Through Bayes rule, we have the following identity,

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

where we call $P(\theta)$ the prior, $P(\theta|D)$ the posterior and $P(D|\theta)$ the likelihood.

MLE

One approach to estimate probabilities is to maximize the likelihood as follows,

$$\hat{\theta}^{MLE} = \arg \max_{\theta} P(D|\theta),$$

which is the general definition of MLE.

Intuition

We observe training data D, we should choose the value of θ that makes D most probable.

An example

- X be a random variable for a coin, 1 or 0,
- ▶ θ is the probability of X taking 1, e.g., $P(X = 1) = \theta$, and unknown,
- ▶ D is the observations produced by flip a coin X $N = \alpha_1 + \alpha_0$ times where α_1 the number of X = 1,
- Assuming I.I.D.

Likelihood is defined as $L(\theta) = P(D|\theta)$. With the conditions claimed before, we have the following formula,

$$L(\theta) = P(D|\theta) = \theta^{\alpha_1} (1 - \theta)^{\alpha_0}.$$

The MLE is to choose θ to maximize $P(D|\theta)$. For convenient, we take the log of $L(\theta)$,

$$I(\theta) = \ln L(\theta) = \alpha_1 \ln \theta + \alpha_0 \ln(1 - \theta),$$

where $I(\theta)$ is called as log-likelihood. Since $I(\theta)$ is a concave function of θ , we just calculate the derivative of $I(\theta)$ with respect to θ ,

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{\partial \ln P(D|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$
$$= \frac{\partial \ln \left[\boldsymbol{\theta}^{\alpha_1} (1 - \boldsymbol{\theta})^{\alpha_0}\right]}{\partial \boldsymbol{\theta}}$$

$$= \frac{\partial \left[\alpha_{1} \ln \theta + \alpha_{0} \ln(1-\theta)\right]}{\partial \theta}$$

$$= \alpha_{1} \frac{\partial \ln \theta}{\partial \theta} + \alpha_{0} \frac{\partial \ln(1-\theta)}{\partial \theta}$$

$$= \alpha_{1} \frac{\partial \ln \theta}{\partial \theta} + \alpha_{0} \frac{\partial \ln(1-\theta)}{\partial (1-\theta)} \cdot \frac{\partial (1-\theta)}{\partial \theta}$$

$$\implies \frac{\partial \ell(\theta)}{\partial \theta} = \alpha_{1} \frac{1}{\theta} + \alpha_{0} \frac{1}{(1-\theta)} \cdot (-1)$$

$$\implies \theta = \frac{\alpha_{1}}{\alpha_{1} + \alpha_{0}}$$

Thus,

$$\hat{\theta}^{MLE} = \arg\max_{\theta} P(D|\theta) = \arg\max_{\theta} \ln P(D|\theta) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

Given the observed data D and the prior $P(\theta)$, we want to maximize the posterior probability. By using Bayes rule, we have

$$\hat{\theta}^{MAP} = \arg\max_{\theta} P(\theta|D) = \arg\max_{\theta} \frac{P(D|\theta)P(\theta)}{P(D)}$$

Comparing the MLE algorithm, the only difference is the extra $P(\theta)$.

Intuition

Given new evidence, update the prior knowledge.

An example

As in our coin flip example, the most common form of prior is a Beta distribution:

$$P(\theta) = \operatorname{Beta}(\beta_0, \beta_1) = \frac{\theta^{\beta_1 - 1}(1 - \theta)^{\beta_0 - 1}}{B(\beta_0, \beta_1)}.$$

Recall the expression for $P(D|\theta)$, we have:

$$\begin{split} \hat{\theta}^{MAP} &= \arg\max_{\theta} P(D|\theta) P(\theta) \\ &= \arg\max_{\theta} \theta^{\alpha_1} (1-\theta)^{\alpha_0} \frac{\theta^{\beta_1-1} (1-\theta)^{\beta_0-1}}{B\left(\beta_0,\beta_1\right)} \\ &= \arg\max_{\theta} \frac{\theta^{\alpha_1+\beta_1-1} (1-\theta)^{\alpha_0+\beta_0-1}}{B\left(\beta_0,\beta_1\right)} \\ &= \arg\max_{\theta} \theta^{\alpha_1+\beta_1-1} (1-\theta)^{\alpha_0+\beta_0-1}. \end{split}$$

Substitute $(\alpha_1 + \beta_1 - 1)$ for α_1 and $(\alpha_0 + \beta_0 - 0)$ for α_0 in $\hat{\theta}^{MLE}$, we have

$$\hat{\theta}^{MAP} = \arg\max_{\theta} P(D|\theta)P(\theta) = \frac{(\alpha_1 + \beta_1 - 1)}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}.$$

Outline

Estimating Probabilities

Bayes rule Maximum Likelihood Estimate (MLE) Maximum a Posterior (MAP)

Naive Bayes

Conditional Independence Naive Bayes for Discrete Inputs Naive Bayes for Continuous Inputs

Conditional Independence

Definition

$$P(Y = y_k | X_1 ... X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

The Naive Bayes classification rule:

$$Y = \arg \max_{y_k} \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$
$$= \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k).$$

Naive Bayes for Discrete Input

We want to estimate two sets of parameters,

$$\theta_{ijk} \equiv P\left(X_i = x_{ij} | Y = y_k\right)$$

and

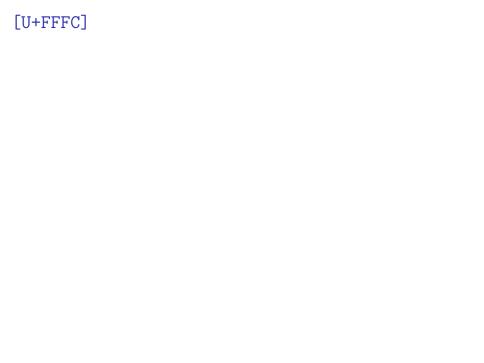
$$\pi_k \equiv P(Y = y_k)$$
.

For example, by using MLE, we have

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij} | Y = y_k) = \frac{\#D\{X_i = x_{ij} \land Y = y_k\}}{\#D\{Y = y_k\}}.$$

Adding a smoothing term, the estimate is given by

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij}|Y = y_k) = \frac{\#D\{X_i = x_{ij} \land Y = y_k\} + I}{\#D\{Y = y_k\} + IJ}.$$



For π_k , we have

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\}}{|D|}$$

and the smoothed one

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + I}{|D| + IK}.$$

Naive Bayes for Continuous Inputs

Assume that for each possible discrete value y_k , the distribution of each continuous X_i is Gaussian. In order to train such a Naive Bayes classifier we must therefore estimate the mean and standard deviation of each of these Gaussians:

$$\mu_{ik} = E[X_i|Y = y_k]$$

$$\sigma_{ik}^2 = E[(X_i - \mu_{ik})^2|Y = y_k].$$

By using MLE approach, we obtain:

$$\hat{\mu}_{ik} = \frac{1}{\sum_{j} \delta\left(Y^{j} = y_{k}\right)} \sum_{j} X_{i}^{j} \delta\left(Y^{j} = y_{k}\right)$$

$$\hat{\sigma}_{ik}^{2} = \frac{1}{\sum_{j} \delta\left(Y^{j} = y_{k}\right)} \sum_{j} \left(X_{i}^{j} - \hat{\mu}_{ik}\right)^{2} \delta\left(Y^{j} = y_{k}\right).$$