

Maximum Likelihood Estimation

Wilker Aziz

February 12, 2018

Notation We use capital Roman letters (e.g. X) for random variables (rvs) and lowercase letters for assignments (e.g. x). We use X_1^n as a shorthand for X_1, \dots, X_n and similarly with x_1^n . We write P_X for probability distributions, and $P_X(X = x)$ for probability values—where we sometimes omit one or both occurrences of X , e.g. $P_X(x)$, $P(X = x)$, or $P(x)$, if no ambiguity is possible. We denote a probability mass function (pmf) by $p(x; \alpha)$, where α is a deterministic set of parameters. Throughout, we also assume that argmax returns a single point.

Assume we have a dataset of n iid observations $\mathcal{D} = \{x_1, \dots, x_n\}$ of an rv $X \sim P_X$, i.e. $(X_i \sim P_X)_{i=1}^n$. First of all, from independence, we know that

$$P_{X_1^n}(x_1, \dots, x_n) = \prod_{i=1}^n P_{X_i}(x_i) = \prod_{i=1}^n P_X(x_i) \quad (1)$$

and we then model the probability $P_X(x)$ with a member $p(x; \alpha)$ of a parametric family and proceed to derive a maximum likelihood estimate of α . In the following sections we use $\mathcal{L}(\alpha|\mathcal{D})$ for the log-likelihood function

$$\mathcal{L}(\alpha|\mathcal{D}) = \sum_{i=1}^n \log p(x_i; \alpha) \quad (2)$$

and we often omit the dependency on data writing simply $\mathcal{L}(\alpha)$. Our objective is then

$$\theta^* = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}(\theta) \quad (3)$$

where the space of valid parameters Θ is possibly subject to constraints.

1 Bernoulli

Suppose X takes on values in the set $\{0, 1\}$, then we say X is Bernoulli-distributed

$$X \sim \text{Bern}(\theta) \quad (4)$$

where $0 < \theta < 1$ the Bernoulli parameter. The Bernoulli pmf is

$$p(x; \theta) = \text{Bern}(X = x | \theta) = \theta^x (1 - \theta)^{(1-x)} \quad (5)$$

and therefore the Bernoulli parameter corresponds to the probability of the positive class—i.e. $P_X(X = 1) = \theta$.

We now derive the maximum likelihood estimate of the parameter θ . We start by rewriting the objective (3) in terms for the Bernoulli pmf (4)

$$\theta^* = \underset{\theta \in [0,1]}{\text{argmax}} \mathcal{L}(\theta) \quad (6a)$$

$$= \underset{\theta \in (0,1)}{\text{argmax}} \sum_{i=1}^n \log p(x_i; \theta) \quad (6b)$$

$$= \underset{\theta \in (0,1)}{\text{argmax}} \sum_{i=1}^n x_i \log \theta + (1 - x_i) \log(1 - \theta) \quad (6c)$$

$$= \underset{\theta \in (0,1)}{\text{argmax}} \log \theta \underbrace{\left(\sum_{i=1}^n x_i \right)}_{n_1} + \log(1 - \theta) \underbrace{\left(\sum_{i=1}^n 1 - x_i \right)}_{n_0} \quad (6d)$$

$$= \underset{\theta \in (0,1)}{\text{argmax}} n_1 \log \theta + n_0 \log(1 - \theta) \quad (6e)$$

where we use n_1 for the number of positive observations and n_0 for the number of negative observations—and note that $n = n_1 + n_0$ is the total number of observations.

Now we find the first derivative of $\mathcal{L}(\theta)$ with respect to θ :

$$\frac{d\mathcal{L}(\theta)}{d\theta} = \frac{d}{d\theta} [n_1 \log \theta + n_0 \log(1 - \theta)] \quad (7a)$$

$$= n_1 \frac{d}{d\theta} \log \theta + n_0 \frac{d}{d\theta} \log(1 - \theta) \quad (7b)$$

$$= n_1 \frac{d}{d\theta} \log \theta + n_0 \underbrace{\frac{d}{d(1 - \theta)} \log(1 - \theta) \frac{d(1 - \theta)}{d\theta}}_{\text{chain rule of derivatives}} \quad (7c)$$

$$= \frac{n_1}{\theta} + \frac{n_0}{1 - \theta} (-1) \quad (7d)$$

Setting the derivative to 0 and solving for θ gives us the MLE solution.

$$0 = \frac{n_1}{\theta} + \frac{n_0}{1-\theta}(-1) \quad (8a)$$

$$= \frac{n_1(1-\theta) - n_0\theta}{\theta(1-\theta)} \quad (8b)$$

$$= n_1(1-\theta) - n_0\theta \quad (8c)$$

$$= n_1 - n_1\theta - n_0\theta \quad (8d)$$

$$n_1 = \theta(n_1 + n_0) \quad (8e)$$

$$\theta = \frac{n_1}{n_1 + n_0} \quad (8f)$$

$$= \frac{n_1}{n} \quad (8g)$$

Note that in (8b) we have to assume that the denominator $\theta(1-\theta) \neq 0$, which excludes $\theta = 0$ and $\theta = 1$. That explains why we define the Bernoulli parameter in the open interval $(0, 1)$. But note that this is not a bad thing, as $\theta = 0$ or $\theta = 1$ would lead to X being deterministic.

2 Categorical

Suppose X takes on values in the discrete interval $[1, k]$, then we say X is Categorical-distributed

$$X \sim \text{Cat}(\theta_1, \dots, \theta_k) \quad (9)$$

where $\theta_1, \dots, \theta_k$ are the Categorical parameters subject to $\theta_x > 0$ and $\sum_{x=1}^k \theta_x = 1$. The Categorical pmf is

$$p(a; \theta_1^k) = \text{Cat}(X = a | \theta_1, \dots, \theta_k) = \prod_{x=1}^k \theta_x^{\delta_{xa}} \quad (10)$$

where δ_{ij} is the [Kronecker delta](#) and therefore each Categorical parameter corresponds to the probability of the respective category—i.e. $P_X(X = x) = \theta_x$.

We now derive the maximum likelihood estimate of the parameters θ_1^k . We

start by rewriting the objective (3) in terms for the Categorical pmf (4)

$$\theta^* = \operatorname{argmax}_{\theta_1^k \in \Delta} \mathcal{L}(\theta_1^k) \quad (11a)$$

$$= \operatorname{argmax}_{\theta_1^k \in \Delta} \sum_{i=1}^n \log p(x_i; \theta_1^k) \quad (11b)$$

$$= \operatorname{argmax}_{\theta_1^k \in \Delta} \sum_{i=1}^n \log \prod_{x=1}^k \theta_x^{\delta_{xx_i}} \quad (11c)$$

$$= \operatorname{argmax}_{\theta_1^k \in \Delta} \sum_{i=1}^n \sum_{x=1}^k \delta_{xx_i} \log \theta_x \quad (11d)$$

$$= \operatorname{argmax}_{\theta_1^k \in \Delta} \sum_{x=1}^k \log \theta_x \underbrace{\sum_{i=1}^n \delta_{xx_i}}_{n_x} \quad (11e)$$

$$= \operatorname{argmax}_{\theta_1^k \in \Delta} \sum_{x=1}^k n_x \log \theta_x \quad (11f)$$

where we denote the number of observations of class $x \in [1, k]$ by n_x .

To avoid optimising the constrained objective, where $\sum_{x=1}^k \theta_x = 1$, we employ a [Lagrange multiplier](#) λ such that the new objective is

$$\theta^* = \operatorname{argmax}_{\theta_1^k \in \mathbb{R}^k} \underbrace{\mathcal{L}(\theta_1^k) - \lambda \left[\left(\sum_{x=1}^k \theta_x \right) - 1 \right]}_{\mathcal{L}(\theta_1^k, \lambda)} \quad (12a)$$

$$\text{s.t. } \theta_x > 0 \text{ for } x \in [1, k]$$

Now we find the first partial derivative of $\mathcal{L}(\theta_1^k, \lambda)$ with respect to λ

$$\frac{\partial}{\partial \lambda} \mathcal{L}(\theta_1^k, \lambda) = \frac{\partial}{\partial \lambda} \left\{ \sum_{x=1}^k n_x \log \theta_x - \lambda \left[\left(\sum_{x=1}^k \theta_x \right) - 1 \right] \right\} \quad (13a)$$

$$= \frac{\partial}{\partial \lambda} \left\{ \sum_{x=1}^k n_x \log \theta_x \right\} - \left[\left(\sum_{x=1}^k \theta_x \right) - 1 \right] \quad (13b)$$

$$= 1 - \sum_{x=1}^k \theta_x \quad (13c)$$

where setting the derivative to zero yields

$$\sum_{x=1}^k \theta_x = 1 \quad . \quad (14)$$

We now turn to the first partial derivative of $L(\theta_1^k, \lambda)$ with respect to θ_j for $j \in [1, k]$

$$\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta_1^k, \lambda) = \frac{\partial}{\partial \theta_j} \left\{ \sum_{x=1}^k n_x \log \theta_x - \lambda \left[\left(\sum_{x=1}^k \theta_x \right) - 1 \right] \right\} \quad (15a)$$

$$= \sum_{x=1}^k n_x \frac{\partial}{\partial \theta_j} \log \theta_x - \lambda \sum_{x=1}^k \frac{\partial}{\partial \theta_j} \theta_x \quad (15b)$$

$$= \sum_{x=1}^k n_x \frac{\delta_{jx}}{\theta_x} - \lambda \sum_{x=1}^k \delta_{jx} \quad (15c)$$

$$= \frac{n_j}{\theta_j} - \lambda \quad (15d)$$

where setting the derivative to zero yields

$$0 = \frac{n_j}{\theta_j} - \lambda \quad (16a)$$

$$\lambda = \frac{n_j}{\theta_j} \quad (16b)$$

$$\theta_j = \frac{n_j}{\lambda} \quad (16c)$$

Now substituting (16c) into (14) we have

$$\sum_{x=1}^k \theta_x = \sum_{x=1}^k \frac{n_x}{\lambda} = \frac{1}{\lambda} \sum_{x=1}^k n_x = \frac{1}{\lambda} n = 1 \quad (17a)$$

and therefore $\lambda = \frac{1}{n}$. And finally, substituting λ in (16c)

$$\theta_j = \frac{n_j}{n} \quad (18)$$

yields the maximum likelihood estimate. Note that in (16b) we have to assume θ_j is not 0, and that's why we defined the categorical pmf for $\theta_j \in \mathbb{R}_{>0}$. Also note that the solution in (18) is strictly positive as long as $n_j > 0$.