

Section 8: Machine Learning

1. Why Boba Cares About MAP:

You don't understand why there's no boba place within walking distance around campus, so you decide to start one. In order to estimate the amount of ingredients needed and the time you will spend in the business (you still need to study), you want to estimate how many orders you will receive per hour. After taking CS109, you are pretty confident that incoming orders can be considered as independent events and the process can be modeled with a Poisson.

Now the question is - what is the λ parameter of the Poisson? In the first hour of your soft opening, you are visited by 4 curious students, each of whom made an order. You have a prior belief that $f(\Lambda = \lambda) = K \cdot \lambda \cdot e^{-\frac{\lambda}{2}}$. What is the MLE estimate? What is inference of λ given the observation? What is the Maximum-A-Posteriori (MAP) estimate of λ ? Through your process try to identify what is a point-estimate, and what is a distribution.

To find the MLE, we start from finding the likelihood function (i.e. joint probability of observed events) and find the λ that maximizes the likelihood function.

$$L(\lambda) = \frac{\lambda^4 \cdot e^{-\lambda}}{4!}$$

$$LL(\lambda) = 4 \log(\lambda) - \lambda$$

$$\frac{\partial LL}{\partial \lambda} = \frac{4}{\lambda} - 1$$

Set $\frac{\partial LL}{\partial \lambda}$ to 0 and solve for λ .

$$\lambda = 4$$

Inference of λ given the observation:

$$f(\lambda|X=4) = \frac{P(X=4|\lambda) \cdot f(\lambda)}{P(X=4)}$$

MAP estimate of λ : we find the λ that maximizes the inference given the observation, i.e. we want to solve:

$$\begin{aligned} \arg \max_{\lambda} f(\lambda|X=4) &= \arg \max_{\lambda} \frac{P(X=4|\lambda) \cdot f(\lambda)}{P(X=4)} \\ &= \arg \max_{\lambda} P(X=4|\lambda) \cdot f(\lambda) \\ &= \arg \max_{\lambda} \frac{\lambda^4 \cdot e^{-\lambda}}{4!} \cdot K \cdot \lambda \cdot e^{-\frac{\lambda}{2}} \end{aligned}$$

Take log.

$$\log\left(\frac{\lambda^4 \cdot e^{-\lambda}}{4!} \cdot K \cdot \lambda \cdot e^{-\frac{\lambda}{2}}\right) = 4\log(\lambda) - \lambda + 1 + \log(K) + \log(\lambda) - \frac{\lambda}{2}$$

Differentiate with respect to λ , set to 0 and solve.

$$\begin{aligned} \frac{5}{\lambda} - 1 - \frac{1}{2} &= 0 \\ \lambda &= \frac{10}{3} \end{aligned}$$

2. Vision Test MLE

You decide that the vision tests given by eye doctors would be more precise if we used an approach inspired by logistic regression. In a vision test a user looks at a letter with a particular font size and either correctly guesses the letter or incorrectly guesses the letter.

You assume that the probability that a particular patient is able to guess a letter correctly is:

$$p = \sigma(\theta + f)$$

Where θ is the user's vision score and f is the font size of the letter. This formula uses the sigmoid function:

$$\begin{aligned} \sigma(z) &= \frac{1}{1 + e^{-z}} \\ \frac{\partial \sigma(z)}{\partial z} &= \sigma(z)[1 - \sigma(z)] \end{aligned}$$

Explain how you could estimate a user's vision score (θ) based on their 20 responses $(f^{(1)}, y^{(1)}) \dots (f^{(20)}, y^{(20)})$, where $y^{(i)}$ is an indicator variable for whether the user correctly identified the i th letter and $f^{(i)}$ is the font size of the i th letter. Solve for any and all partial derivatives required by the approach you describe in your answer.

We are going to solve this problem by finding the MLE estimate of θ . To find the MLE estimate, we are going to find the argmax of the log likelihood function. To calculate argmax we are going to use gradient ascent, which requires that we know the partial derivative of the log likelihood function with respect to theta.

First we write the log likelihood:

$$L(\theta) = \prod_{i=1}^{20} p^{y_i} (1-p)^{[1-y_i]}$$

$$LL(\theta) = \sum_{i=1}^{20} (y_i \log(p) + (1-y_i) \log(1-p))$$

Then we find the derivative of log likelihood with respect to θ . We first do this for one data point:

$$\frac{\partial LL}{\partial \theta} = \frac{\partial LL}{\partial p} \cdot \frac{\partial p}{\partial \theta}$$

We can calculate both the smaller partial derivatives independently:

$$\frac{\partial LL}{\partial p} = \frac{y_i}{p} - \frac{1-y_i}{1-p}$$

$$\frac{\partial p}{\partial \theta} = p[1-p]$$

Putting it all together for one letter:

$$\begin{aligned} \frac{\partial LL}{\partial \theta} &= \frac{\partial LL}{\partial p} \cdot \frac{\partial p}{\partial \theta} \\ &= \left[\frac{y_i}{p} - \frac{1-y_i}{1-p} \right] p[1-p] \\ &= y_i(1-p) - p(1-y_i) \\ &= y_i - p \\ &= y_i - \sigma(\theta - f) \end{aligned}$$

For all twenty examples:

$$\frac{\partial LL}{\partial \theta} = \sum_{i=1}^{20} y_i - \sigma(\theta + f^{(i)})$$

3. Multiclass Bayes

Note: we don't expect folks to get to this problem in section. It is here just for students who would like some extra review!

In this problem we are going to explore how to write Naive Bayes for multiple output classes. We want to predict a single output variable Y which represents how a user feels about a book. Unlike in your homework, the output variable Y can take on one of the *four* values in the

set {Like, Love, Haha, Sad}. We will base our predictions off of three binary feature variables $X_1, X_2,$ and X_3 which are indicators of the user's taste. All values $X_i \in \{0, 1\}$.

We have access to a dataset with 10,000 users. Each user in the dataset has a value for X_1, X_2, X_3 and Y . You can use a special query method **count** that returns the number of users in the dataset with the given *equality* constraints (and only equality constraints). Here are some example usages of **count**:

- count**($X_1 = 1, Y = \text{Haha}$) returns the number of users where $X_1 = 1$ and $Y = \text{Haha}$.
- count**($Y = \text{Love}$) returns the number of users where $Y = \text{Love}$.
- count**($X_1 = 0, X_3 = 0$) returns the number of users where $X_1 = 0,$ and $X_3 = 0$.

You are given a new user with $X_1 = 1, X_2 = 1, X_3 = 0$. What is the best prediction for how the user will feel about the book (Y)? You may leave your answer in terms of an argmax function. You should explain how you would calculate all probabilities used in your expression. Use **Laplace estimation** when calculating probabilities.

We can make the Naive Bayes assumption of independence and simplify argmax of $P(Y|\mathbf{X})$ to get an expression for \hat{Y} , the predicted output value, and evaluate it using the provided **count** function.

$$\begin{aligned} \hat{Y} &= \arg \max_y \frac{P(X_1 = 1, X_2 = 1, X_3 = 0|Y = y)P(Y = y)}{P(X_1 = 1, X_2 = 1, X_3 = 0)} \\ &= \arg \max_y P(X_1 = 1, X_2 = 1, X_3 = 0|Y = y)P(Y = y) \\ &= \arg \max_y P(X_1 = 1|Y = y)P(X_2 = 1|Y = y)P(X_3 = 0|Y = y)P(Y = y), \text{ where:} \end{aligned}$$

$$\begin{aligned} P(X_1 = 1|Y = y) &= [\text{count}(X_1 = 1, Y = y) + 1]/\text{count}(Y = y) + 2 \\ P(X_2 = 1|Y = y) &= [\text{count}(X_2 = 1, Y = y) + 1]/\text{count}(Y = y) + 2 \\ P(X_3 = 1|Y = y) &= [\text{count}(X_3 = 1, Y = y) + 1]/\text{count}(Y = y) + 2 \\ P(X_1 = 0|Y = y) &= [\text{count}(X_1 = 0, Y = y) + 1]/\text{count}(Y = y) + 2 \\ P(X_2 = 0|Y = y) &= [\text{count}(X_2 = 0, Y = y) + 1]/\text{count}(Y = y) + 2 \\ P(X_3 = 0|Y = y) &= [\text{count}(X_3 = 0, Y = y) + 1]/\text{count}(Y = y) + 2 \end{aligned}$$

and you don't need to use MAP to estimate Y :

$$P(Y = y) = \text{count}(Y = y)/10,000$$