# Maximum Likelihood Estimation

Given a collection of data, what is the underlying distribution?

$\longrightarrow$ Assuming a distribution, what are the likeliest parameters for that distribution.

Suppose we've conducted an experiment and our collected data is a list of values $x_1, x_2, x_3, \ldots, x_n$.

Suppose there is a function $f$ for the probability of these data to be measured. This function is parametrized by $\Theta$:

$$Prob(data) = f(x_1, x_2, x_3, \ldots, x_n \mid \Theta) \quad \leftarrow \text{ joint prob of all the data simultaneously given a (set of) parameters } \Theta.$$

Now let's make things easier: Assume each data point can occur independently from all others and each follows the same distribution (independent, identically distributed or "iid").

$$\hookrightarrow \quad f(x_1, x_2, \ldots x_n \mid \Theta) = f(x_1 \mid \Theta) \cdot f(x_2 \mid \Theta) \cdots f(x_n \mid \Theta)$$

$$= \prod_{i=1}^{n} f(x_i \mid \Theta)$$

OK, our goal is to find $\Theta$, we're given the data. We need to flip this around: what is the likelihood of $\Theta$, given the data? No problem, flip!

$$\mathcal{L}(\Theta \mid x_1, x_2, \ldots x_n) = f(x_1, x_2, \ldots x_n \mid \Theta) = \prod_{i=1}^{n} f(x_i \mid \Theta)$$

(whoa are you allowed to just do that?)

Moving forward, we want to find the $\Theta$ that maximizes $\mathcal{L}(\Theta \mid data)$. (called the maximum likelihood estimator, $\hat{\Theta}$)

Often it's easier to maximize the log-likelihood $\ell$ which is maximized at the same place:

$$\ell(\Theta \mid data) = \sum_{i=1}^{n} \ln f(x_i \mid \Theta)$$

1

# MLE cont

Some times we can find the $\hat{\theta}$ easily with calculus, other times we need a numerical method.

Using calculus, compute the derivative $\frac{\partial \ell}{\partial \theta}$, set equal to 0, solve.

Let's do a specific example:

## Poisson distribution:

Prob for a certain <u>number</u> of events to occur if the average rate of events is known and the prob. for the next event to occur does not depend on the time since the previous event.

$$P(X=k;\lambda) = P(k;\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

Let's assume the $x_i$'s are drawn from a poisson. Then we need to find $\hat{\lambda}$:

$$f(x_i|\theta) \Rightarrow f(x_i|\lambda) = \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

Compute the (log)-likelihood:

$$\mathcal{L}(\theta|data) = \prod_{i=1}^{n} f(x_i|\lambda) = \prod_{i=1}^{n} \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} = e^{-n\lambda} \prod_{i=1}^{n} \frac{\lambda^{x_i}}{x_i!}$$

$$\ell(\theta|data) = \ln \mathcal{L}(\theta|data) = \ln e^{-n\lambda} + \sum_{i=1}^{n} \ln \lambda^{x_i} - \sum_{i=1}^{n} \ln(x_i!)$$

$$= -n\lambda + \sum_{i=1}^{n} x_i \ln\lambda - \left( \searrow \right)$$

$$= -n\lambda + \ln\lambda \sum_{i=1}^{n} x_i - \left( \downarrow \right)$$

Now we find an equation for the maximum as a function of $\lambda$, and solve it for $\hat{\lambda}$:

$$\frac{\partial \ell}{\partial \lambda} = \frac{\partial}{\partial \lambda}\left( -n\lambda + \ln\lambda \sum x_i - (\ ) \right) = -n + \ln\lambda \underbrace{\frac{\partial}{\partial \lambda}\sum x_i}_{0} + \sum x_i \cdot \frac{1}{\lambda} - 0$$

Solving this for $\hat{\lambda}$

$$-n + \frac{1}{\lambda}\sum_{i=1}^{n} x_i = 0 \quad \Rightarrow \quad \hat{\lambda} = \frac{1}{n}\sum_{i=1}^{n} x_i$$

The MLE $\hat{\lambda}$ for poisson is the <u>average</u> of the data.

So if we are given a bunch of integer x-data we can compute the mean, and draw $P(X=k, \lambda)$ on top of the histogram and see if its close.

---

Recall I did something sneaky: $\mathcal{L}(\theta | data) = f(data | \theta)$

If these are both probabilities you can't just do that, you need to use <u>Bayes Theorem</u>

<u>Bayes</u>:

• The joint prob of two events A and B both occurring is commutative:

$$P(A \cap B) = P(B \cap A) \quad \text{(b/c the sets } A \cap B \text{ and } B \cap A \text{ are the same)}.$$

○ We can write the joint prob: $P(A \cap B) = P(A|B) \cdot P(B)$  (intuitive...)

Combining those and rearranging gives us Baye's thm:

$$P(A \cap B) = P(B \cap A)$$

$$P(A|B)P(B) = P(B|A)P(A)$$

$$\boxed{P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}} \quad \text{Baye's}$$

This is how we can "flip around" a prob.

Now what does this have to do w/ MLE?

$\downarrow$

3

Instead of $\mathcal{L}(\theta \mid data) = f(data \mid \theta)$ we use Bayes:

$$P(\theta \mid data) = \frac{P(data \mid \theta) P(\theta)}{P(data)}$$

$$= \frac{f(x_1, x_2, \dots x_n \mid \theta) P(\theta)}{P(x_1, x_2, \dots, x_n)}$$

Now the denominator is indep. of $\theta$, meaning it's just a constant that won't change the maximum.

So the difference really is the $P(\theta)$. Now, if we assume that is also a constant, what happens? It won't affect the location of the max either, and

$$\mathcal{L}(\theta \mid data) = P(\theta \mid data) \quad \text{except for a mult. const.}$$

Assuming $P(\theta) = const$ means all $\theta$'s are equally likely.

And that is the assumption behind MLE.

"The ML estimator is equal to the Bayesian estimator given a uniform prior distribution of the parameters."

$\rightarrow$ MLE is a special case.

4