

Get started

Open in app



towards  
data science

Follow

577K Followers



You have 1 free member-only story left this month. [Sign up for Medium and get an extra one](#)

## MLE, MAP and Bayesian Inference

Grasp the idea of Bayesian inference by focusing on the difference from MLE and MAP



Shota Horii Sep 21, 2019 · 7 min read ★



Photo by [Aysha Begum](#) on [Unsplash](#)

Get started

Open in app



## MLE/MAP and Bayesian inference.

In this article, I'm going to introduce Bayesian inference by focusing on the difference between MLE/MAP and Bayesian inference.

**Note:** Preliminary knowledge of MLE and MAP is assumed in this article. If you're not familiar with those methods, please refer to the following article.

### A Gentle Introduction to Maximum Likelihood Estimation and Maximum A Posteriori Estimation

Getting intuition of MLE and MAP with a football example

[towardsdatascience.com](https://towardsdatascience.com)



. . .

## The difference between MLE/MAP and Bayesian inference

Let's start from the recap of MLE and MAP.

Given the observed data  $D$ , estimations of a probabilistic model's parameter  $\theta$  by MLE and MAP are the following.

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} \{P(D|\theta)\}$$

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} \{P(\theta|D)\} = \operatorname{argmax}_{\theta} \left\{ \frac{P(D|\theta)P(\theta)}{P(D)} \right\} = \operatorname{argmax}_{\theta} \{P(D|\theta)P(\theta)\}$$

MLE gives you the value which maximises the Likelihood  $P(D|\theta)$ . And MAP gives you the value which maximises the posterior probability  $P(\theta|D)$ . As both methods give you a single fixed value, they're considered as **point estimators**.

On the other hand, Bayesian inference fully calculates the posterior probability distribution, as below formula. Hence the output is not a single value but a probability density function (when  $\theta$  is a continuous variable) or a probability mass function (when  $\theta$  is a discrete variable).

Get started

Open in app



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

This is the difference between MLE/MAP and Bayesian inference. MLE and MAP returns a single fixed value, but Bayesian inference returns probability density (or mass) function.

But why we even need to fully calculate the distribution, when we have MLE and MAP to determine the value of  $\theta$ ? To answer this question, let's see the case when MAP (and other point estimators) doesn't work well.

...

### A case when MAP (or point estimators in general) doesn't work well

Assume you're in a casino with full of slot machines with 50% winning probability. After playing for a while, you heard the rumour that there's one special slot machine with 67% winning probability.

Now, you're observing people playing 2 suspicious slot machines (you're sure that one of those is the special slot machine!) and got the following data.

*Machine A: 3 wins out of 4 plays*

*Machine B: 81 wins out of 121 plays*

By intuition, you would think *machine B* is the special one. Because 3 wins out of 4 plays on *machine A* could just happen by chance. But *machine B*'s data doesn't look like happening by chance.

But just in case, you decided to estimate those 2 machines' winning probabilities by MAP with hyperparameters  $\alpha=\beta=2$ . (Assuming that the results ( $k$  wins out of  $n$  plays) follow binomial distribution with the slot machine's winning probability  $\theta$  as its parameter.)

The formula and results are below.

$$\hat{\theta}_{MAP} = \frac{k + \alpha - 1}{n + \alpha + \beta - 2}$$

[Get started](#)
[Open in app](#)


*Machine A*:  $(3+2-1)/(4+2+2-2) = 4/6 = 66.7\%$

*Machine B*:  $(81+2-1)/(121+2+2-2) = 82/123 = 66.7\%$

Unlike your intuition, estimated winning probability  $\theta$  by MAP for the 2 machines are exactly same. Hence, by MAP, you cannot determine which one is the special slot machine.

But really? Isn't it looking obvious that *Machine B* is more likely to be the special one?

...

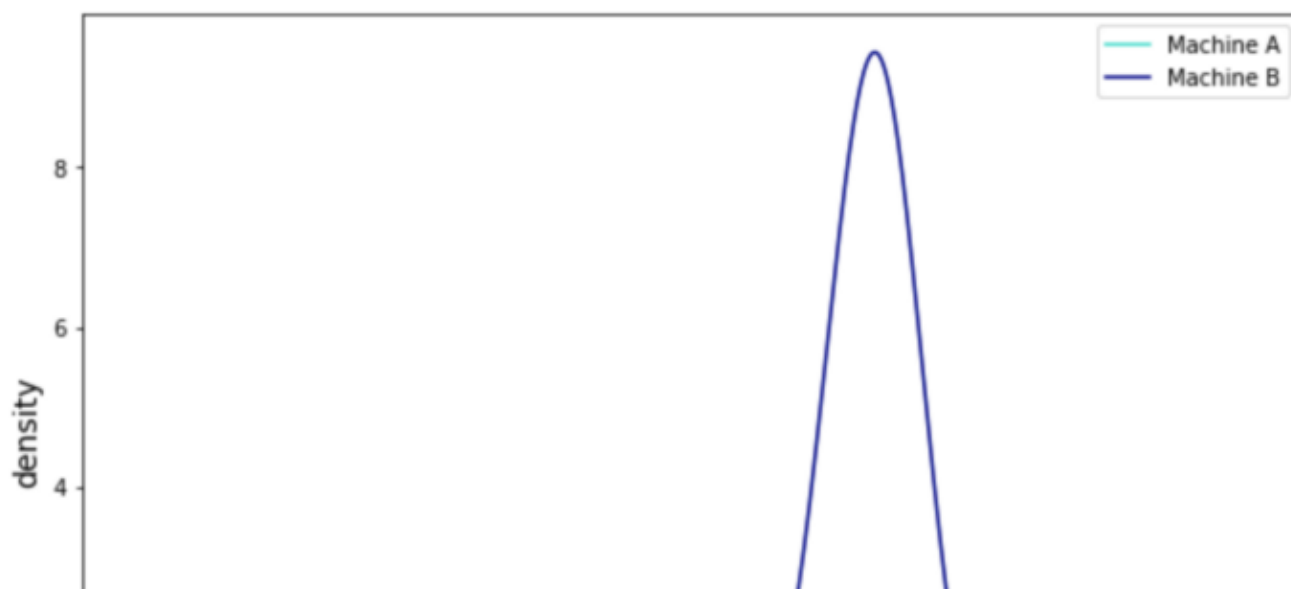
## Bayesian Inference

To see if there really be no difference between *machine A* and *machine B*, let's fully calculate the posterior probability distribution, not only MAP estimates.

In the case above, the posterior probability distribution  $P(\theta|D)$  is calculated as below. (Detailed computation will be covered in the next section.)

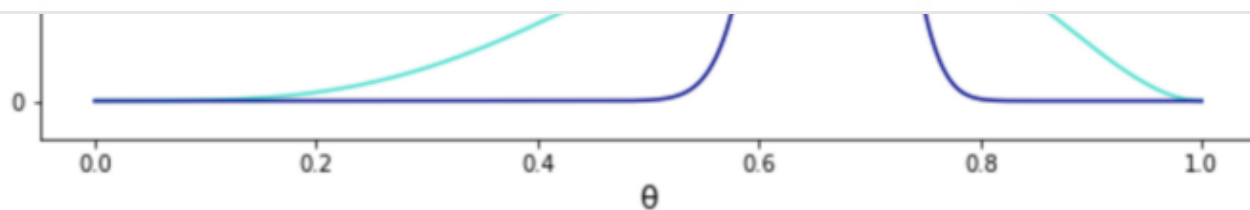
$$P(\theta|D) = \frac{\Gamma(n + \alpha + \beta)}{\Gamma(k + \alpha)\Gamma(n - k + \beta)} \theta^{k+\alpha-1} (1 - \theta)^{n-k+\beta-1}$$

And  $P(\theta|D)$  for *machine A* and *machine B* are drawn as below.



Get started

Open in app



Although both distributions have their **mode** on  $\theta=0.666\dots$  (that's why their MAP estimates are the same value), the shapes of the distributions are quite different. Density around the mode is much higher in the distribution of *machine B* than the one of *machine A*. This is why you want to calculate full distribution, not only the MAP estimate.

...

## Computation of Bayesian Inference

As we skipped the computation of  $P(\theta|D)$  in the previous section, let's go through the detailed calculation process in this section.

Both MAP and Bayesian inference are based on Bayes' theorem. The computational difference between Bayesian inference and MAP is that, in Bayesian inference, we need to calculate  $P(D)$  called **marginal likelihood** or **evidence**. It's the denominator of Bayes' theorem and it assures that the integrated value\* of  $P(\theta|D)$  over all possible  $\theta$  becomes 1. (\* Sum of  $P(\theta|D)$ , if  $\theta$  is a discrete variable.)

$P(D)$  is obtained by marginalisation of joint probability. When  $\theta$  is a continuous variable, the formula is as below.

$$P(D) = \int_{\theta} P(D, \theta) d\theta$$

Considering the chain rule, we obtain the following formula.

$$P(D) = \int_{\theta} P(D|\theta)P(\theta)d\theta$$

[Get started](#)[Open in app](#)

Now, put this into the original formula of the posterior probability distribution.

Calculating below is the goal of Bayesian Inference.

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

...

Let's calculate  $P(\theta|D)$  for the case above.

Beginning with  $P(D|\theta)$  — **Likelihood** — which is the probability that data  $D$  is observed when parameter  $\theta$  is given. In the case above,  $D$  is “3 wins out of 4 matches”, and parameter  $\theta$  is the winning probability of *machine A*. As we assume that the number of wins follows binomial distribution, the formula is as below, where  $n$  is the number of matches and  $k$  is the number of wins.

$$P(D|\theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

Then  $P(\theta)$  — the **prior probability distribution** of  $\theta$  — which is the probability distribution expressing our prior knowledge about  $\theta$ . Here, specific probability distributions are used corresponding to the probability distribution of Likelihood  $P(D|\theta)$ . It's called conjugate prior distribution.

Since the conjugate prior of binomial distribution is Beta distribution, we use Beta distribution to express  $P(\theta)$  here. Beta distribution is described as below, where  $\alpha$  and  $\beta$  are hyperparameters.

$$P(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

[Get started](#)[Open in app](#)

$$\begin{aligned}
 P(D|\theta)P(\theta) &= \binom{n}{k} \theta^k (1-\theta)^{n-k} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\
 &= \binom{n}{k} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{k+\alpha-1} (1-\theta)^{n-k+\beta-1}
 \end{aligned}$$

Then,  $P(D)$  — the denominator of the formula — is calculated as follows. Note that the possible range of  $\theta$  is  $0 \leq \theta \leq 1$ .

$$\begin{aligned}
 P(D) &= \int_{\theta} P(D|\theta)P(\theta)d\theta \\
 &= \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta \\
 &= \binom{n}{k} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 \theta^{k+\alpha-1} (1-\theta)^{n-k+\beta-1} d\theta
 \end{aligned}$$

With [Euler integral of the first kind](#), the above formula can be deformed to below.

$$P(D) = \binom{n}{k} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(k+\alpha)\Gamma(n-k+\beta)}{\Gamma(n+\alpha+\beta)}$$

Finally, we can obtain  $P(\theta|D)$  as below.

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



Get started

Open in app



$$\begin{aligned}
 &= \frac{\binom{n}{k} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \theta^k (1-\theta)^{n-k}}{\binom{n}{k} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \frac{\Gamma(k+\alpha)\Gamma(n-k+\beta)}{\Gamma(n+\alpha+\beta)}} \\
 &= \frac{\Gamma(n+\alpha+\beta)}{\Gamma(k+\alpha)\Gamma(n-k+\beta)} \theta^{k+\alpha-1} (1-\theta)^{n-k+\beta-1} \\
 &\quad \dots
 \end{aligned}$$

### Expected A Posteriori (EAP)

As you may have noticed, the estimate by MAP is the *mode* of the posterior distribution. But we can also use other statistics for the point estimation, such as *expected value* of  $\theta|D$ . The estimation using the expected value of  $\theta|D$  is called **Expected A Posteriori**.

$$\begin{aligned}
 \hat{\theta}_{EAP} &= E[\theta|D] \\
 &= \int_{\theta} \theta P(\theta|D) d\theta \\
 &\quad \dots
 \end{aligned}$$

Let's estimate the winning probability of the 2 machines using EAP. From the discussion above,  $P(\theta|D)$  in this case is below.

$$P(\theta|D) = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(k+\alpha)\Gamma(n-k+\beta)} \theta^{k+\alpha-1} (1-\theta)^{n-k+\beta-1}$$



[Get started](#)[Open in app](#)

Thus the estimate is described as below.

$$\begin{aligned}
 \hat{\theta}_{EAP} &= \int_{\theta} \theta P(\theta|D) d\theta \\
 &= \int_0^1 \theta \frac{\Gamma(n + \alpha + \beta)}{\Gamma(k + \alpha)\Gamma(n - k + \beta)} \theta^{k+\alpha-1} (1 - \theta)^{n-k+\beta-1} d\theta \\
 &= \frac{\Gamma(n + \alpha + \beta)}{\Gamma(k + \alpha)\Gamma(n - k + \beta)} \int_0^1 \theta^{k+\alpha} (1 - \theta)^{n-k+\beta-1} d\theta
 \end{aligned}$$

With Euler integral of the first kind and the definition of Gamma function, above formula can be deformed to below.

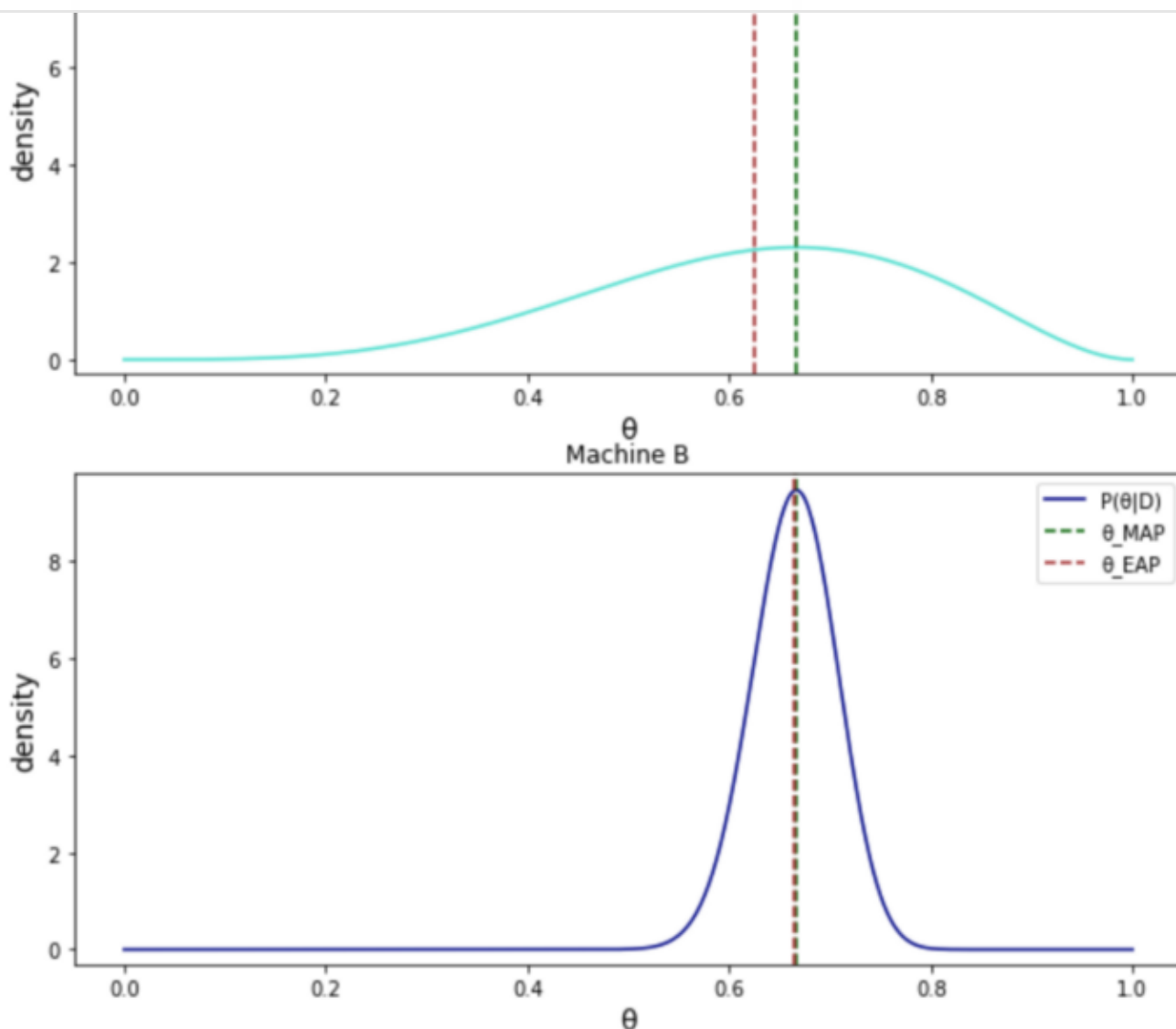
$$\begin{aligned}
 \hat{\theta}_{EAP} &= \frac{\Gamma(n + \alpha + \beta)}{\Gamma(k + \alpha)\Gamma(n - k + \beta)} \frac{\Gamma(k + \alpha + 1)\Gamma(n - k + \beta)}{\Gamma(n + \alpha + \beta + 1)} \\
 &= \frac{\Gamma(n + \alpha + \beta)}{\Gamma(k + \alpha)\Gamma(n - k + \beta)} \frac{(k + \alpha)\Gamma(k + \alpha)\Gamma(n - k + \beta)}{(n + \alpha + \beta)\Gamma(n + \alpha + \beta)} \\
 &= \frac{k + \alpha}{n + \alpha + \beta}
 \end{aligned}$$

Hence, EAP estimate of 2 machines' winning probabilities with hyperparameters  $\alpha=\beta=2$  are below.

Machine A:  $(3+2)/(4+2+2) = 5/8 = 62.5\%$

Machine B:  $(81+2)/(121+2+2) = 83/125 = 66.4\%$

Machine A

[Get started](#)
[Open in app](#)


• • •

## Conclusion

As seen above, Bayesian inference provides much more information than point estimators like MLE and MAP. However, it also has a drawback — the complexity of its integral computation. The case in this article was quite simple and solved analytically, but it's not always the case in real-world applications. We then need to use MCMC or other algorithms as a substitute for the direct integral computation. Hope this article helped you to understand Bayesian inference.

**Sign up for The Variable**

Get started

Open in app



and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Your email



Get this newsletter

By signing up, you will create a Medium account if you don't already have one. Review our [Privacy Policy](#) for more information about our privacy practices.

Machine Learning

Statistics

Bayesian Statistics



[About](#) [Help](#) [Legal](#)

Get the Medium app



Download on the  
App Store



GET IT ON  
Google Play