

towards
data science

493K Followers · About Follow

This is your **last** free member-only story this month. [Sign up for Medium and get an extra one](#)

Conjugate Prior Explained

With examples & proofs



Aerin Kim Jan 8 · 4 min read ★

1. What is Prior?

Prior probability is **the probability of an event before we see the data**.
In Bayesian Inference, the prior is our guess about the probability based on what we know now, before new data becomes available.

2. What is Conjugate Prior?

Conjugate prior just can not be understood without knowing Bayesian inference.

Bayesian Inference — Intuition and Implementation
The art of Bayesian Inference lies in how you implement it...

For the rest of the blog, I'll assume you know the concepts of prior, sampling and posterior.

Conjugate prior in essence

For some likelihood functions, if you choose a certain prior, the posterior ends up being in the same distribution as the prior. Such a prior then is called a Conjugate Prior.

It is always best understood through examples. Below is the code to **calculate the posterior of the binomial likelihood**. θ is the probability of success and our goal is to pick the θ that maximizes the posterior probability.

```
1  import numpy as np
2  import scipy.stats as stats
3
4  success_prob = 0.3
5  data = np.random.binomial(n=1, p=success_prob, size=1000) # success is 1, failure is 0.
6
7  # Domain  $\theta$ 
8  theta_range = np.linspace(0, 1, 1000)
9
10 # Prior  $P(\theta)$ 
11 a = 2
12 b = 8
13 theta_range_e = theta_range + 0.0001
14 prior = stats.beta.cdf(x = theta_range_e, a=a, b=b) - stats.beta.cdf(x = theta_range, a
15
16 # The sampling dist. aka Likelihood  $P(X|\theta)$ 
17 likelihood = stats.binom.pmf(k = np.sum(data), n = len(data), p = theta_range)
18
19 # Posterior
20 posterior = likelihood * prior
21 normalized_posterior = posterior / np.sum(posterior)
```

conjugate_prior.py hosted with ❤ by GitHub

[view raw](#)

A question to you: Is there anything that concerns you in the code block above?

There are two things that make the posterior calculation expensive.

First, we are computing the posterior for every single θ .

Why do we have to calculate the posterior for thousands of thetas? Because you are normalizing the posterior (line 21). Even if you choose not to normalize the posterior, the end goal is to find the **maximum** of the posteriors (Maximum a posteriori). In order to find the maximum in a vanilla way, we need to consider every candidate — the likelihood $P(X|\theta)$ for every θ .

Second, if there is no closed-form formula of the posterior distribution, we have to find the maximum by numerical optimization, such as gradient descent or newtons method.

3. How does the Conjugate Prior help?

When you know that your prior is a conjugate prior, you can skip the `posterior = likelihood * prior` computation. Furthermore, if your prior distribution has a closed-form expression, you already know what the maximum posterior is going to be.

In the example above, the beta distribution is a conjugate prior to the binomial likelihood. What does this mean? It means **during the modeling phase, we already know the posterior will also be a beta distribution**. Therefore, after carrying out more experiments, **you can compute the posterior simply by adding the number of acceptances and rejections to the existing parameters α , β respectively**, instead of multiplying the likelihood with the prior distribution. This is very convenient! (Proof in the next section.)

As a data/ML scientist, your model is never complete. You have to update your model as more data come in (and that's why we use Bayesian Inference).

As you saw, the computations in Bayesian Inference can be heavy or sometimes even intractable. However, if we could use the closed-form formula of the conjugate prior, the computation becomes very light.

4. Proof — Why is a Beta distribution a conjugate prior to Binomial

likelihood?

When we use the Beta distribution as a prior, a posterior of binomial likelihood will also follow the beta distribution.

Show Beta begets Beta.

What do the PDFs of Binomial and Beta look like?

	PDF
Binomial	$f(x) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$ ↓ the function of x
Beta	$g(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$ ↓ the function of θ

Let's plug them into the famous Bayes formula.

θ is the probability of success.

x is the number of successes.

n is the total number of trials, therefore $n-x$ is the number of failures.

$$\begin{aligned}
 \overset{\text{posterior}}{P(\theta|X)} &= \frac{\overset{\text{sampling}}{P(X|\theta)} \cdot \overset{\text{prior}}{P(\theta)}}{\int_{\theta} P(X|\theta) \cdot P(\theta) d\theta \rightarrow \text{normalizing constant}} \\
 &\quad \downarrow \text{data} \quad \downarrow \text{for every possible } \theta \\
 &= \frac{\binom{n}{x} \theta^x (1-\theta)^{n-x} \cdot \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int_0^1 \binom{n}{x} \theta^x (1-\theta)^{n-x} \cdot \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta} \\
 &= \frac{\frac{n!}{x!(n-x)!} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}{\frac{n!}{x!(n-x)!} \underbrace{\int_0^1 \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta}_{B(x+\alpha, n-x+\beta)}} \\
 &= \text{Beta}(x+\alpha, n-x+\beta) \quad \square
 \end{aligned}$$

Why does the last integral become $B(x+\alpha, n-x+\beta)$? → <https://bit.ly/2t1i2KT>

The prior distribution $P(\theta)$ was **Beta(α, β)** and after getting x successes and $n-x$ failures from the experiments, the posterior also becomes a Beta distribution with parameters **($x+\alpha, n-x+\beta$)**.

What's nice here is you will know this analytically without doing the computation.

5. Conjugate Prior Distributions

The Beta distribution is a conjugate prior for the **Bernoulli**, **binomial**, **negative binomial** and **geometric** distributions (seems like those are the distributions that involve success & failure).

<Beta posterior>

Beta prior * **Bernoulli** likelihood → Beta posterior
 Beta prior * **Binomial** likelihood → Beta posterior
 Beta prior * **Negative Binomial** likelihood → Beta posterior
 Beta prior * **Geometric** likelihood → Beta posterior

<Gamma posterior>

Gamma prior * **Poisson** likelihood → Gamma posterior
 Gamma prior * **Exponential** likelihood → Gamma posterior

<Normal posterior>

Normal prior * Normal likelihood (mean) → Normal posterior
 Conjugate prior $P(\theta)$ in an equation:

$P(\theta)$ such that $P(\theta|D) = P(\theta)$

An interesting way to put this is that even if you do all those experiments and multiply your likelihood to the prior, your initial choice of the prior distribution was **so good** that the final distribution is the same as the prior.

A few things to note:

1. When we use the conjugate prior, sequential estimation (updating the counts after each observation) gives the same result as a batch estimation.
2. In order to **find the maximum posterior**, you don't have to **normalize** the multiplication of likelihood (sampling) and the prior (the integration for every possible θ in the denominator).

$$\begin{array}{c}
 \nearrow \text{posterior} \\
 P(\theta|x) = \frac{\overset{\nearrow \text{sampling}}{P(x|\theta)} \cdot \overset{\nearrow \text{prior}}{P(\theta)}}{\int_{\theta} P(x|\theta) \cdot P(\theta) d\theta \rightarrow \text{normalizing constant}} \\
 \downarrow \text{data} \qquad \downarrow \theta \\
 \text{for every possible } \theta
 \end{array}$$

You can still find the maximum without normalizing. However, if you want to compare posteriors from different models, or calculate the point estimates, you need to normalize.

Sign up for The Daily Pick

By Towards Data Science

Hands-on real-world examples, research, tutorials, and cutting-edge techniques delivered Monday to Thursday. Make learning your daily ritual. [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.

Data Science

Machine Learning

Artificial Intelligence

Bayesian Statistics

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

Get the Medium app



Download on the
App Store



GET IT ON
Google Play