

Population, Sample, Parameter

1. In statistics, the **population** is any subject (not necessarily a group) which a researcher try to analyze.
2. The **sample** is a collection of data related to the population. In a typical situation, the sample is assumed to be randomly and independently extracted from the population.
3. The **parameter** represents a property of the population to be analyzed. The parameter is unknown to the researcher.
4. The goal of statistics is to obtain useful insights about the parameter of the population with the sample extracted from it.

Population Distribution

We regard the population as a probability distribution and call it the **population distribution**. Then we can interpret the sample as a set of random variables following the population distribution, and the parameters as variables which determine the “shape” of the population distribution.

Let $\mathbf{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denote the sample, and θ denote the parameter of the population distribution. To indicate that the shape of the population distribution depends on θ , the population p.f. or p.d.f. is denoted by $p(\mathbf{x}_i|\theta)$ where each \mathbf{x}_i ($i = 1, \dots, n$) is called an **observation** and supposed to be a realized value of the random variable following the population distribution. n is often called the **sample size**.

Likelihood

Suppose the sample $\mathbf{D} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ are taken from a population distribution where the parameter θ is a set of unknown parameters. The joint p.f. or the joint p.d.f of \mathbf{D} is denoted by

$$p(\mathbf{D}|\theta) = p(\mathbf{x}_1, \dots, \mathbf{x}_n|\theta).$$

In particular, if observations are independent of each other,

$$p(\mathbf{D}|\theta) = p(\mathbf{x}_1|\theta) \times \dots \times p(\mathbf{x}_n|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta).$$

When we regard $p(\mathbf{D}|\theta)$ as a function of θ , it is called the **likelihood** or **likelihood function**.

Example: Bernoulli Distribution i

Let us define a random variable X_i ($i = 1, \dots, n$) corresponding to tossing a coin such that

$$X_i = \begin{cases} 1, & \text{Head is obtained;} \\ 0, & \text{Tail is obtained,} \end{cases}$$

and

$$\Pr(X_i = 1) = \theta, \quad \Pr(X_i = 0) = 1 - \theta.$$

Then X_i follows the **Bernoulli distribution** and its p.f. is given by

$$p(x_i|\theta) = \theta^{x_i}(1 - \theta)^{1-x_i}, \quad x_i = 0, 1.$$

Example: Bernoulli Distribution ii

Then the joint p.f. of $D = (x_1, \dots, x_n)$ is

$$\begin{aligned} p(D|\theta) &= \prod_{i=1}^n p(x_i|\theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \\ &= \theta^y (1 - \theta)^{n-y}, \quad y = \sum_{i=1}^n x_i. \end{aligned}$$

Suppose we have $(x_1, x_2, x_3, x_4, x_5) = (1, 0, 1, 1, 1)$. The value of $p(D|\theta)$ depends on the value of θ .

θ	0.1000	0.2000	0.3000	0.4000	0.5000
$p(D \theta)$	0.0001	0.0013	0.0057	0.0154	0.0312
θ	0.6000	0.7000	0.8000	0.9000	
$p(D \theta)$	0.0518	0.0720	0.0819	0.0656	

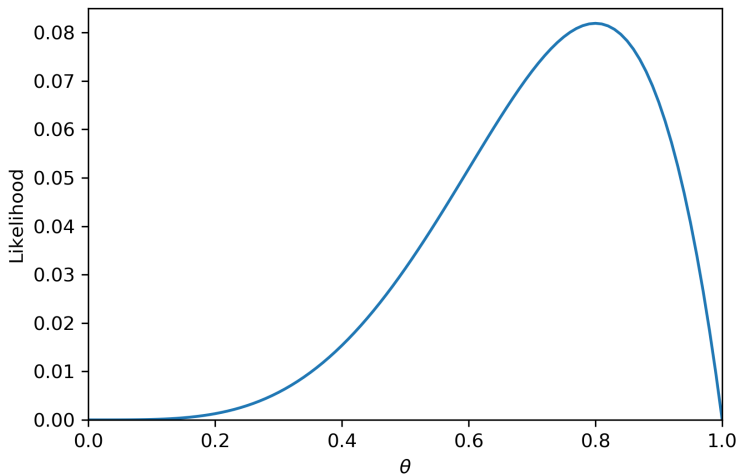


Figure 1: The likelihood of θ in the Bernoulli distribution

Interpretation Of The Likelihood

Given the sample D , the likelihood $p(D|\theta)$ is regarded as a kind of “plausibility” of a specific value of θ .

For example, the likelihood of $\theta = 0.9$ is **0.656** while that of $\theta = 0.4$ in the previous example is **0.0154**. We may say that **0.9** is about 4 times more plausible than **0.4** as the true value of θ .

To make comparison between two competing values of θ , say θ_0 and θ_1 , we introduce the **likelihood ratio**:

$$\text{likelihood ratio} = \frac{p(D|\theta_0)}{p(D|\theta_1)}.$$

Prior Knowledge On Parameters

In practice, researchers often have information on unknown parameters before they start analysis. For example,

- θ must take a value between 0 and 1 because it is probability;
- in case of tossing a coin, θ is supposed to be 50% if the coin is fair.

In Bayesian statistics, we construct a distribution of unknown parameters that reflect our prior knowledge on their true values. This is call the **prior distribution**. Let $p(\theta)$ denote the prior distribution.

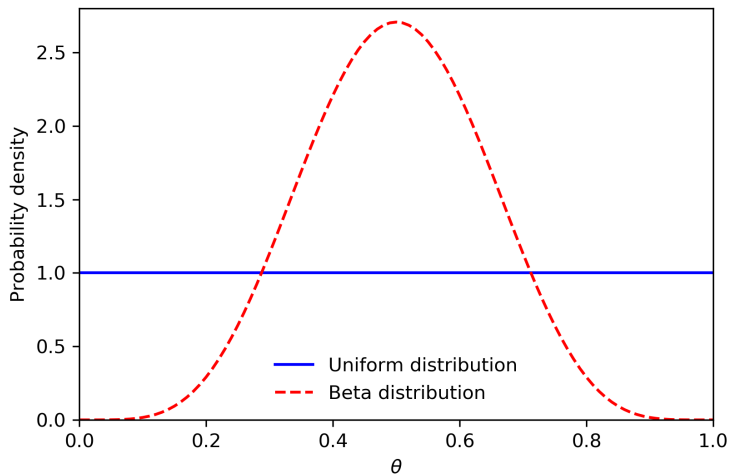


Figure 2: Prior distributions of θ in the Bernoulli distribution

The **uniform distribution** $\text{Uniform}(\mathbf{a}, \mathbf{b})$ is

$$p(\mathbf{x}|\mathbf{a}, \mathbf{b}) = \begin{cases} \frac{1}{b-a}, & (\mathbf{a} \leq \mathbf{x} \leq \mathbf{b}); \\ 0, & (\text{otherwise}). \end{cases}$$

In the above figure, we set $\mathbf{a} = \mathbf{0}$ and $\mathbf{b} = \mathbf{1}$.

The **beta distribution** $\text{Beta}(\alpha, \beta)$ is

$$p(\mathbf{x}|\alpha, \beta) = \frac{\mathbf{x}^{\alpha-1}(1-\mathbf{x})^{\beta-1}}{B(\alpha, \beta)}, \quad 0 \leq \mathbf{x} \leq 1.$$

where $B(\alpha, \beta)$ is the beta function:

$$B(\alpha, \beta) = \int_0^1 \mathbf{x}^{\alpha-1}(1-\mathbf{x})^{\beta-1} d\mathbf{x}.$$

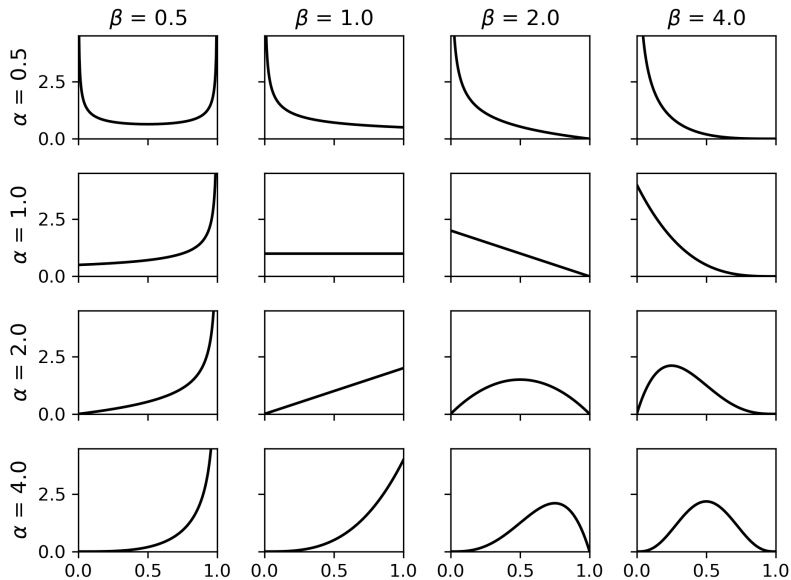


Figure 3: Beta distributions with various (α, β)

Bayes' Theorem And Posterior Distribution

Suppose $p(\theta, D)$ is the joint distribution θ and D . Using the definition of the conditional distribution, we have

$$p(\theta, D) = p(\theta|D)p(D) = p(D|\theta)p(\theta).$$

Arranging the middle and the right-hand side of the above equation, we have

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}.$$

- The above formula is called **Bayes' theorem**.
- $p(\theta|D)$ is called the **posterior distribution**.
- $p(D)$ is called the **normalizing constant**.

Marginal Likelihood

By marginalizing the joint distribution $p(\theta, D)$, $p(D)$ is given as

$$p(D) = \int p(\theta, D) d\theta = \int p(D|\theta)p(\theta) d\theta.$$

Thus $p(D)$ is interpreted as “averaged likelihood” in terms of the prior $p(\theta)$. In this sense, $p(D)$ is called the **marginal likelihood**.

Then Bayes' theorem is rewritten as

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta) d\theta}.$$

Example: Bernoulli Distribution

Suppose the prior distribution is $\text{Beta}(\alpha_0, \beta_0)$.

The posterior distribution of θ is given by

$$\begin{aligned} p(\theta|D) &\propto p(D|\theta)p(\theta) \\ &\propto \theta^y(1-\theta)^{n-y} \times \theta^{\alpha_0-1}(1-\theta)^{\beta_0-1} \\ &\propto \theta^{y+\alpha_0-1}(1-\theta)^{n-y+\beta_0-1} \\ &\propto \theta^{\alpha_\star-1}(1-\theta)^{\beta_\star-1}, \\ \alpha_\star &= y + \alpha_0, \quad \beta_\star = n - y + \beta_0. \end{aligned}$$

This is the beta distribution **$\text{Beta}(\alpha_\star, \beta_\star)$** .

Bayesian Learning

Bayes' theorem is rearranged as

$$\frac{p(\theta|D)}{p(\theta)} = \frac{p(D|\theta)}{p(D)} \propto p(D|\theta).$$

The left-hand side is

$$\left\{ \begin{array}{l} \frac{p(\theta|D)}{p(\theta)} > 1, \text{ plausibility of } \theta \text{ is increased;} \\ \frac{p(\theta|D)}{p(\theta)} < 1, \text{ plausibility of } \theta \text{ is decreased.} \end{array} \right.$$

In other words, an update in belief is proportional to the likelihood $p(D|\theta)$.

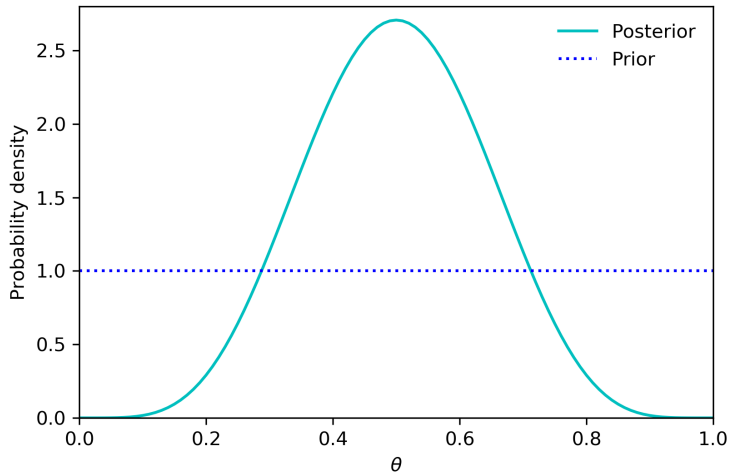


Figure 4: Posterior distributions of the probability of success

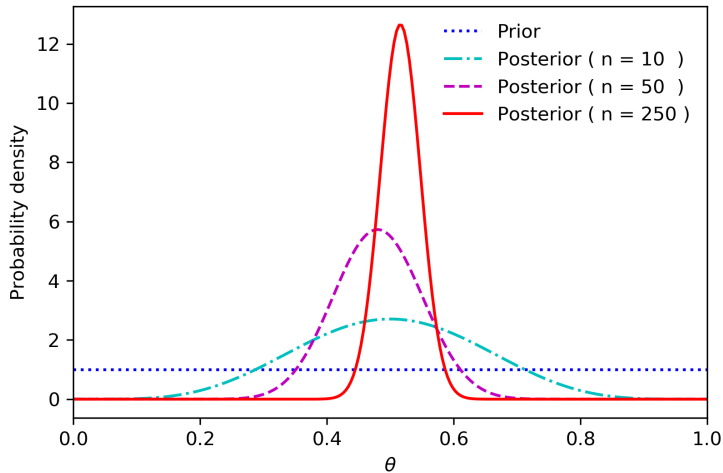


Figure 5: Sequential updating of the posterior distribution