

# 1. Binary, Beta, Multinomial, Dirichlet Distribution

## (1) Binary Variable.

for a R.V.  $X \in \{0, 1\}$   $P(X=1) = \mu$ . we can have

$$\text{Bern}(X) = \mu^X (1-\mu)^{1-X} \Leftrightarrow \text{Bern}(X) = \begin{cases} \mu & X=1 \\ 1-\mu & X=0 \\ 0 & \text{else.} \end{cases}$$

Bernulli Distribution

$$\mathbb{E}(X) = \mu \cdot 1 + (1-\mu) \cdot 0 = \mu$$

$$\begin{aligned} \text{Var}(X) &= \mu(1-\mu)^2 + (1-\mu)(0-\mu)^2 \\ &= \mu(1-\mu)(1-\mu+\mu) \\ &= \mu(1-\mu) \end{aligned}$$

If we repeat the experiment  $N$  times and collect

$D = [X_1, X_2, \dots, X_N]$ . the joint pmf would be

$$p(D) = \prod_{i=1}^N p(X_i) = \prod_{i=1}^N \mu^{X_i} (1-\mu)^{1-X_i} = \mu^{\sum X_i} (1-\mu)^{N - \sum X_i}$$

Use the traditional logarithm trick

$$\ln p(D) = \sum X_i \ln \mu + (N - \sum X_i) \ln (1-\mu)$$

$$\frac{\partial}{\partial \mu} \ln p(D) = \sum X_i \cdot \frac{1}{\mu} + (N - \sum X_i) \frac{-1}{1-\mu} = 0 \quad (\text{MLE})$$

$$\Rightarrow \mu_{ML} = \frac{1}{N} \sum_{i=1}^N X_i$$

## Problem for MLE.

We have derived the estimation formula for  $\mu$  that

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_i X_i$$

But if our experiment can not repeat for sufficiently large times, MLE may have very large bias.

E.g. We flip a coin, 1-head, 0-tail. We see 3 head among 3 experiment. According to  $\hat{\mu}_{MLE} = \frac{\# \text{ of head}}{\# \text{ of exp}} = 1$ , which is far away from our common sense. This is an instance of over-fitting.

We can use another model to describe the experiment when we repeat the flipping coin experiment  $N$  times, which is Binomial Distribution

for  $X_i \in \{0, 1\}$ .  $P(X_i=1) = \mu$ . we repeat the exp  $N$  times, we see  $X_i=1$  for  $m$  times.

$$P(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$

$$E(m) = \sum_{i=1}^N P(X_i=1) X_i = N\mu$$

$$\text{Var}(m) = \sum_{i=1}^N \text{Var}(X_i=1) = N\mu(1-\mu)$$

## (2). Beta Distribution. (Conjugate Prior)

In part 1, we know that MLE can have over-fit problem when data size is small. We can introduce Bayesian's method to improve this method. We want to turn MLE to MAP

From the perspective of Bayesian, parameter  $\mu$  also satisfies a distribution, which is called prior. For each experiment,  $\mu$  may change. All we can have is the posterior and likelihood.

$$P(\mu|x) = \frac{P(\mu) \cdot P(x|\mu)}{P(x)}$$

First, we need to choose a prior. For the likelihood function

$$P(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$

We want to find a form/model for prior which can make likelihood of the same form as posterior. In this case, likelihood/posterior is proportional to  $\mu^m (1-\mu)^{N-m}$ .

We choose Beta distribution

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

remind you of  $\binom{a+b}{a}$ ?

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

This is a continuous function though.

$\Gamma(x)$  has very interesting properties.

①  $\Gamma(x+1) = x \Gamma(x)$

②  $\Gamma(1) = 1$

③  $\Gamma(x) = (x-1)!$  if  $x$  is integer

## Proof

$$\begin{aligned} 1. \Gamma(x+1) &= \int_0^{\infty} t^x e^{-t} dt = -\int_0^{\infty} t^x de^{-t} = \underbrace{-t^x e^{-t}}_0 \Big|_0^{\infty} + \int_0^{\infty} e^{-t} dt^x \\ &= \int_0^{\infty} e^{-t} dt^x = \int_0^{\infty} e^{-t} x \cdot t^{x-1} dt = x \int_0^{\infty} t^{x-1} e^{-t} dt \\ &= x \cdot \Gamma(x) \end{aligned}$$

$$2. \Gamma(1) = \int_0^{\infty} t^0 e^{-t} dt = -e^{-t} \Big|_0^{\infty} = 1$$

$$3. \Gamma(x) = (x-1) \Gamma(x-1) = \prod_{i=1}^{x-1} (x-i) = (x-1)!$$

Let's prove that Beta distribution is a valid distribution.

$$\text{Target: } \int_0^1 \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} d\mu = 1 \text{ for } \forall a, b \in (0, \infty)$$

Instead of proving brutally, we turn to prove.

$$\Gamma(a+b) \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \Gamma(a)\Gamma(b).$$

Proof:

$$\begin{aligned} \Gamma(a+b) \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu &= \int_0^{\infty} \int_0^1 \gamma^{a+b-1} e^{-\gamma} d\gamma \int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu \\ &= \int_0^{\infty} \int_0^1 \gamma^{a+b-1} e^{-\gamma} \mu^{a-1} (1-\mu)^{b-1} d\gamma d\mu \\ &= \int_0^{\infty} \int_0^1 (\gamma\mu)^{a-1} [\gamma(1-\mu)]^{b-1} \gamma e^{-\gamma} d\gamma d\mu \end{aligned}$$

Now, use variable substitution

$$\text{let } x = \delta\mu, \quad y = \delta(1-\mu)$$

we have:

$$x \in (0, \infty) \quad y \in (0, \infty) \quad \delta = x+y \quad \mu = \frac{x}{x+y}$$

Jacobian matrix is

$$\begin{bmatrix} \frac{\partial \delta}{\partial x} & \frac{\partial \mu}{\partial x} \\ \frac{\partial \delta}{\partial y} & \frac{\partial \mu}{\partial y} \end{bmatrix} = \begin{bmatrix} 1 & \frac{y}{(x+y)^2} \\ 1 & \frac{-x}{(x+y)^2} \end{bmatrix},$$

The absolute value of determinant is  $\frac{1}{x+y}$

So, our transformed integration turns to be

$$\begin{aligned} & \int_0^{\infty} \int_0^{\infty} x^{a-1} y^{b-1} (x+y) e^{-(x+y)} \cdot \frac{1}{x+y} dx dy \\ &= \int_0^{\infty} \int_0^{\infty} x^{a-1} y^{b-1} e^{-x-y} dx dy \\ &= \int_0^{\infty} x^{a-1} e^{-x} dx \int_0^{\infty} y^{b-1} e^{-y} dy \\ &= \Gamma(a) \Gamma(b) \end{aligned}$$

proof end.

---

Let's compute the expectation and variance of Beta. R.V.

The result is

$$\mathbb{E}(\mu) = \frac{a}{a+b}$$

$$\text{Var}(\mu) = \frac{ab}{(a+b)^2(a+b+1)}$$

Proof.

We will prove a much more general formula, the moment of  $\mu$ .

$$E(\mu^k) = \int_0^1 \mu^k \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} d\mu$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_0^1 \mu^{a+k-1} (1-\mu)^{b-1} d\mu$$

Turn to page 4

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+k)\Gamma(b)}{\Gamma(a+b+k)} = \frac{(a+b-1)! (a+k-1)!}{(a-1)! (a+b+k-1)!}$$

$$= \frac{\text{Beta}(\mu, a+k, b)}{\text{Beta}(\mu, a, b)}$$

So.

$$E(\mu) = \frac{(a+b-1)! a!}{(a-1)! (a+b)!} = \frac{a}{a+b}$$

$$\begin{aligned} \text{Var}(\mu) &= E(\mu^2) - E(\mu)^2 = \frac{(a+b-1)! (a+1)!}{(a-1)! (a+b+1)!} - \frac{a^2}{(a+b)^2} \\ &= \frac{(a+1) \cdot a}{(a+b)(a+b+1)} - \frac{a^2}{(a+b)^2} = \frac{a}{a+b} \cdot \frac{a^2 + ab + a + b - a^2 - ab - a}{(a+b)(a+b+1)} \\ &= \frac{ab}{(a+b)^2 (a+b+1)} \end{aligned}$$

---

$a, b$  in Beta distribution is called hyperparam, because it controls the distribution of  $\mu$ . We can randomly or manually choose at first. As the data accumulates, the effects of  $a$  &  $b$  will decrease to  $\rightarrow 0$ .

Now posterior.

$$P(\mu | \underline{x}) = \frac{1}{\int} P(\underline{x} | \mu) \cdot \text{Beta}(\mu | a, b)$$

Normalization

Term

$$= \frac{1}{\int} \binom{N}{m} \mu^m (1-\mu)^{N-m} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

Remain only terms involving  $\mu$ .

$$\propto \mu^{m+a-1} (1-\mu)^{N-m+b-1}$$

$$\Rightarrow \frac{\Gamma(N+a+b)}{\Gamma(m+a)\Gamma(N-m+b)} \mu^{m+a-1} (1-\mu)^{N-m+b-1}$$

$$= \text{Beta}(\mu | a+m, b+N-m)$$

So, if we have a dataset with  $m$ , " $x=1$ " cases and  $N-m$ , " $x=0$ " cases we will modify/improve our  $\hat{\mu}$  by  $\text{Beta}(\mu | a, b) \rightarrow \text{Beta}(\mu | a+m, b+N-m)$

Another important property is that no matter what prior we choose (say,  $a$  &  $b$ ), as our experiment data go to  $\infty$ , we will always converge to a satisfying result. This is called Sequential learning. We can therefore introduce an online training algorithm based on MLE.

But, if we are not given infinite number of data, all we have is dataset  $D$ . the best we can do is MAP rather than ML.

MAP.

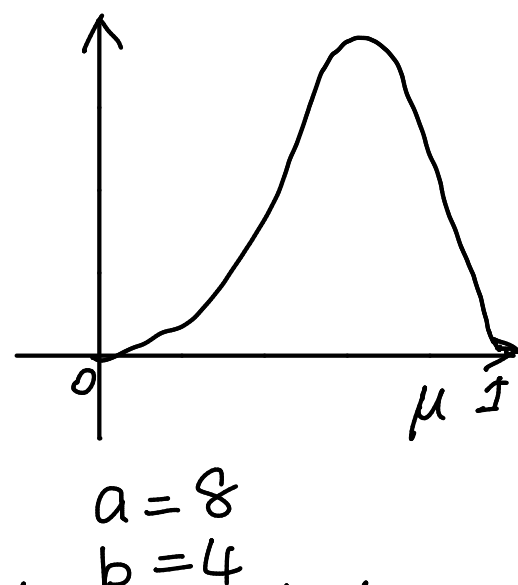
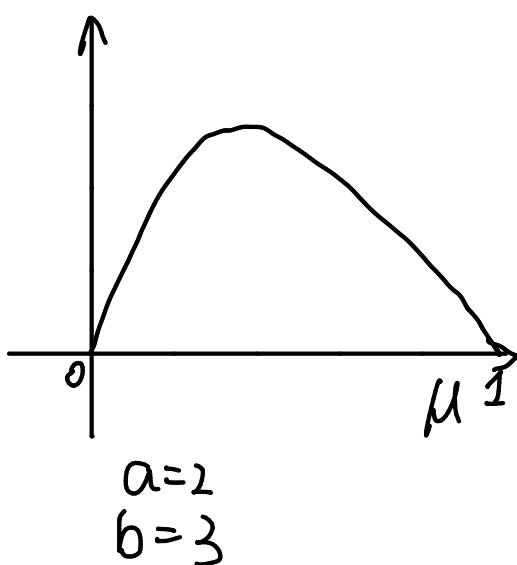
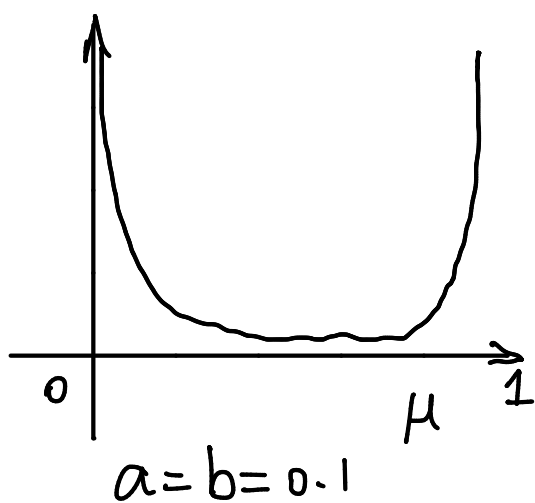
$$P(X=1|D) = \int_0^1 P(X=1|\mu) P(\mu|D) d\mu = \int_0^1 \mu P(\mu|D) d\mu = \mathbb{E}(\mu|D)$$

If there are  $m$  samples of  $X=1$ ,  $N$  samples in all.

$$P(X=1|D) = \frac{m+a}{N+a+b}$$

Also, if we have infinite number of samples, MAP will have same result as MLE.

Beta( $\mu|a, b$ )



As experiment number keeps increasing, the distribution will become sharper and narrower.

Idea:

As we get more and more data, uncertainty represented by posterior distribution decrease steadily.

This is true.



We know that

$$\mathbb{E}_\theta[\theta] = \mathbb{E}_D[\mathbb{E}(\theta|D)]$$

Because

$$\begin{aligned}\mathbb{E}_D[\mathbb{E}(\theta|D)] &= \int_D \int_\theta \theta \cdot P(\theta|D) d\theta \cdot P(D) dD \\ &= \int \theta P(\theta) d\theta = \mathbb{E}_\theta[\theta]\end{aligned}$$

and

$$\underbrace{\text{Var}_\theta[\theta]}_{\text{Prior's Variance}} = \underbrace{\mathbb{E}_D[\text{Var}_\theta[\theta|D]]}_{\text{average of posterior's variance } (\geq 0)} + \underbrace{\text{Var}_D[\mathbb{E}_\theta(\theta|D)]}_{\text{Variance of MAP estimator}}$$

So, on average MAP's variance  $\leq$  Prior's variance after averaging over  $D$ .

### (3) Multinomial Variable.

Multinomial Variable is a generalization of Binary Variable. We now have  $K$  possible outputs.

E.g.  $K=6$ . and we see  $X_3=1$ , then we have

$$\underline{X} = (0, 0, 1, 0, 0, 0)^T \quad \text{one-bit hot code}$$

$$P(\underline{X}|\underline{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

$$\sum_{\underline{X}} P(\underline{X}|\underline{\mu}) = \sum_{k=1}^K \mu_k = 1$$

$$\bar{E}(\underline{X}|\underline{\mu}) = \sum_{\underline{X}} P(\underline{X}|\underline{\mu}) \cdot \underline{X} = (\mu_1, \mu_2, \dots, \mu_K) = \underline{\mu}$$

The likelihood function is also similar to Binary Variable's

$$P(\mathcal{D}|\underline{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{\sum_n x_{nk}} = \prod_{k=1}^K \mu_k^{m_k}$$

$$m_k = \sum_n x_{nk} \dots \# \text{ of state where } X=k \text{ shows up.}$$

Also, we get start from MLE.

$$\ln P(\mathcal{D}|\underline{\mu}) = \sum_{k=1}^K m_k \ln \mu_k$$

Because we should have another constrain  $\sum_k \mu_k = 1$ , so we use Lagrangian.

$$\text{Solve: } \max \ln P(\mathcal{D}|\underline{\mu}) + \lambda \left( \sum_k \mu_k - 1 \right) = L(\underline{\mu}, \lambda)$$

$$\frac{\partial}{\partial \mu_k} L(\underline{\mu}, \lambda) = \frac{m_k}{\mu_k} + \lambda = 0 \Rightarrow \mu_k = -\frac{m_k}{\lambda}$$

$$\text{So } \sum_k \mu_k = \frac{-\sum m_k}{\lambda} = 1 \Rightarrow \lambda = -\sum_k m_k = -N$$

$$\boxed{\mu_k = \frac{m_k}{N}}$$

If we similarly consider multinomial distribution

$$\text{Multi}(m_1, m_2, \dots, m_k | \underline{\mu}, N) = \binom{N}{m_1, m_2, \dots, m_k} \prod_k \mu_k^{m_k}$$

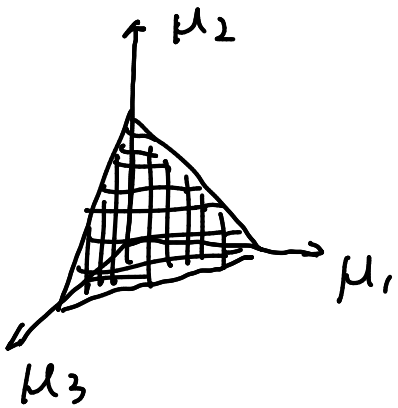
$$\binom{N}{m_1, m_2, \dots, m_k} = \frac{N!}{m_1! \dots m_k!}$$

## (4) Dirichlet Distribution

Now, we will find the prior of the Multinomial Distribution, which is Dirichlet Distribution, which is the conjugate prior to multinomial.

$$\text{Dir}(\underline{\mu} | \underline{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \Gamma(\alpha_2) \dots \Gamma(\alpha_k)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

if  $k=3$ , say  $\mu_1 + \mu_2 + \mu_3 = 1$ ,  $\mu_1, \mu_2, \mu_3 \in [0, 1]$  its distribution is called simplex.



also, the posterior is another Dirichlet dist.

$$P(\underline{\mu} | D, \underline{\alpha}) = \text{Dir}(\underline{\mu} | \underline{\alpha} + \underline{m})$$

$$\underline{m} = [m_1, m_2, \dots, m_k]^T$$