

1 Probability and independence

a) Decomposition

We aim to validate

$$(X \perp Y, W \mid Z) \implies (X \perp Y \mid Z) \quad (1.1)$$

Proof. We suppose the statement $(X \perp Y, W \mid Z)$ is true. It follows from the definition of the conditional independence that $p(x, y, w \mid z) = p(x \mid z)p(y, w \mid z)$ for all $x \in \Omega_x$, $(y, w) \in \Omega_y \times \Omega_w$ and $z \in \Omega_z$. We then consider the marginalize $p(x, y, w \mid z)$:

$$\begin{aligned} p(x, y \mid z) &= \sum_{w \in \Omega_w} p(x, y, w \mid z) \\ &= \sum_{w \in \Omega_w} p(x \mid z)p(y, w \mid z) \\ &= p(x \mid z) \sum_{w \in \Omega_w} p(y, w \mid z) \\ &= p(x \mid z)p(y \mid z) \end{aligned}$$

from which we conclude that $(X \perp Y \mid Z)$ \square . By symmetry of the argument, we can show that $(X \perp W \mid Z)$ is true as well.

b)

We aim to validate

$$(X \perp Y \mid Z) \text{ and } (X, Y \perp W \mid Z) \implies (X \perp W \mid Z) \quad (1.2)$$

Proof. Suppose $(X, Y \perp W \mid Z)$ and $(X \perp Y \mid Z)$ are true. We know from the symmetry and decomposition properties of the conditional independence that $(X, Y \perp W \mid Z) \implies (W \perp X, Y \mid Z) \implies (X \perp W \mid Z)$. Therefore $(X \perp W \mid Z)$ is true \square .

c)

We aim to validate

$$(X \perp Y, W \mid Z) \text{ and } (Y \perp W \mid Z) \implies (X, W \perp Y \mid Z) \quad (1.3)$$

Proof. Suppose $(X \perp Y, W \mid Z)$ is true. Then it follows from the definition of conditional independence that

$$p(x, y, w \mid z) = p(x \mid z)p(y, w \mid z)$$

Then assume $(Y \perp W \mid Z)$ is true. The second factor can be factorized

$$p(x, y, w \mid z) = p(x \mid z)p(y \mid z)p(w \mid z)$$

From the decomposition property, we know $(X \perp W \mid Z)$ is true. Thus

$$p(x, y, w \mid z) = p(x, w \mid z)p(y \mid z)$$

From which we conclude $(X, W \perp Y \mid Z)$ is true \square .

d)

We aim to validate

$$(X \perp Y \mid Z) \text{ and } (X \perp Y \mid W) \implies (X \perp Y \mid Z, W) \quad (1.4)$$

Counter example. We consider the following R.V.

1. X: Person A arrive late for diner;
2. Y: Person B arrive late for diner;
3. W: They come from the same city;
4. Z: They work in the same city.

For this situation, we see that X and Y are conditionally independent when given either W or Z . If we know they are from the same city, then they might work in different cities and take different route home. Thus knowing person A was late doesn't inform us on the probability of person B to arrive late.

A similar argument can be made for $(X \perp Y \mid Z)$.

Thus the LHS of the proposition is true, yet the RHS is clearly false in our case. Assuming we were given that W and Z are true, then we are given the geolocalisation of person A and B . If we were given that person A would be late for dinner, then we'd be able to make a good guess that person B would be late as well (they would both be impacted by the same traffic jam or whatnot). Thus the proposition is false.

2 Bayesian inference and MAP

Let $\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\pi} \stackrel{\text{iid}}{\sim} \text{Multinomial}(1, \boldsymbol{\pi})$ on k element. The values are sampled from a set of cardinality 2, that is $x_j^{(i)} \in \{0, 1\}$. Each R.V. has only one non-zero entry for a given trial, that is $\sum_{j=1}^k x_j^{(i)} = 1$.

We assume a Dirichlet prior $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$ with a PDF

$$p(\boldsymbol{\pi} \mid \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{j=1}^k \alpha_j)}{\prod_{j=1}^k \Gamma(\alpha_j)} \prod_{j=1}^k \pi_j^{\alpha_j - 1}$$

a)

Since the data is IID, they are mutually independent by definition. Being given the parameters of their Multinomial distribution (or a subset for that matter) does not change the independence of the \mathbf{X} 's. Thus,

$$(\mathbf{X}_i \perp \mathbf{X}_j \mid \boldsymbol{\pi}) \forall (i, j) \in \{1, \dots, k\} \times \{1, \dots, k\}$$

Of course, none of the vector can be mutually nor conditionally independent to $\boldsymbol{\pi}$ since it contains information about the distribution of the one hot vectors \mathbf{X}_i . In this case $\boldsymbol{\pi}$ are the probabilities of one of the k entry to be equal to one. Even giving one of these away is enough to impact the posterior distribution of the conditional $p(x_i \mid x_\ell, \pi_j)$ for example.

b)

The posterior distribution $p(\boldsymbol{\pi} \mid x_1, \dots, x_n)$ is computed via the Bayes rule

$$p(\boldsymbol{\pi} \mid \mathbf{x}_{1:n}) = \frac{p(\mathbf{x}_{1:n} \mid \boldsymbol{\pi})p(\boldsymbol{\pi})}{p(\mathbf{x}_{1:n})}$$

where $p(\boldsymbol{\pi}) = p(\boldsymbol{\pi} \mid \boldsymbol{\alpha})$ is the prior for $\boldsymbol{\pi}$ defined above. For the sake of determining the posterior distribution, we can postpone the derivation of the marginal likelihood. The likelihood $p(\mathbf{x}_{1:n} \mid \boldsymbol{\pi})$ is the probability mass function corresponding to n trials of a k -sided die throw. We define the vector $\boldsymbol{\chi} \equiv \sum_{i=1}^n \mathbf{x}_i$ with the property

$$\sum_{j=1}^k \chi_j = n$$

It becomes clear that the likelihood follows the Multinomial($n, \boldsymbol{\pi}$) distribution. The PMF is given by

$$p(\mathbf{x}_{1:n} \mid \boldsymbol{\pi}) = \binom{n}{\chi_1, \dots, \chi_k} \prod_{j=1}^k \pi_j^{\chi_j} \propto \prod_{i=1}^n \prod_{j=1}^k \pi_j^{x_j^{(i)}}$$

Where it is agreed that $\chi_j = \sum_{i=1}^n x_j^{(i)}$. Therefore, the posterior must be

$$p(\boldsymbol{\pi} \mid \mathbf{x}_{1:n}) \propto \prod_{i=1}^n \prod_{j=1}^k \pi_j^{x_j^{(i)}} \prod_{\ell=1}^k \pi_\ell^{\alpha_\ell - 1}$$

We use the fact that we can swap around product operator for real numbers.

$$p(\boldsymbol{\pi} \mid \mathbf{x}_{1:n}) \propto \prod_{i=1}^n \prod_{j=1}^k \prod_{\ell=1}^k \pi_\ell^{\alpha_\ell - 1} \pi_j^{x_j^{(i)}} = \prod_{j=1}^k \pi_j^{\sum_{i=1}^n x_j^{(i)} + \alpha_j - 1}$$

We can readily see that the resulting distribution will be a Dirichlet with updated α_ℓ 's.

The posterior distribution is a Dirichlet distribution with parameters $\alpha'_j = \alpha_j + \sum_{i=1}^n x_j^{(i)}$.

c) Marginal Likelihood

The marginal likelihood $p(\mathbf{x}_{1:n})$ is a normalizing constant defined as the integral of the numerator over all instantiation of $\boldsymbol{\pi}$

$$p(\mathbf{x}_{1:n}) = \int_{\Delta_k} p(\mathbf{x}_{1:n} | \boldsymbol{\pi}) p(\boldsymbol{\pi}) d^{(k)} \boldsymbol{\pi}$$

where Δ_k is the probability simplex. In term of the quantities defined above, this is

$$p(\mathbf{x}_{1:n}) = \int_{\Delta_k} d^{(k)} \boldsymbol{\pi} \binom{n}{\chi_1, \dots, \chi_k} \prod_{j=1}^k \pi_j^{\chi_j} \left(\frac{\Gamma(\sum_{\ell=1}^k \alpha_{\ell})}{\prod_{\ell=1}^k \Gamma(\alpha_{\ell})} \prod_{\ell=1}^k \pi_{\ell}^{\alpha_{\ell}-1} \right)$$

The π_j 's are independent variables since the simplex Δ_k is crucially defined as an affine plane in an Euclidian space which is supported by a set of orthonormal vectors. Thus, our task is to evaluate k identical integrals of the form

$$\int_0^1 d\pi_j \pi_j^{\sum_{i=1}^n x_j^{(i)} + \alpha_j - 1} = \left(\sum_{i=1}^n x_j^{(i)} + \alpha_j \right)^{-1}, \quad \{\alpha_j > 0\}$$

Thus

$$p(\mathbf{x}_{1:n}) = \binom{n}{\chi_1, \dots, \chi_k} \frac{\Gamma(\sum_{\ell=1}^k \alpha_{\ell})}{\prod_{\ell=1}^k \Gamma(\alpha_{\ell})} \prod_{j=1}^k \left(\sum_{i=1}^n x_j^{(i)} + \alpha_j \right)^{-1}$$

d) $\hat{\boldsymbol{\pi}}_{\text{MAP}}$

The maximum *a posteriori* of the Multinomial distribution is

$$\hat{\boldsymbol{\pi}}_{\text{MAP}} \equiv \underset{\boldsymbol{\pi} \in \Delta_k}{\operatorname{argmax}} p(\boldsymbol{\pi} | \mathbf{x}_{1:n})$$

Where the probability simplex is defined as

$$\Delta_k = \left\{ \boldsymbol{\pi} \in \mathbb{R}^k \mid \pi_j \in [0, 1] \text{ and } \sum_{j=1}^k \pi_j = 1 \right\}$$

To satisfy the constraint $g(\boldsymbol{\pi}) = 1 - \sum_{j=1}^k \pi_j$, we use the Lagrange multiplier λ s.t. the optimisation of the log posterior becomes

$$\hat{\boldsymbol{\pi}}_{\text{MAP}} = \underset{(\boldsymbol{\pi}, \lambda) \in \mathbb{R}^{k+1}}{\operatorname{argmax}} \sum_{j=1}^k \left(\sum_{i=1}^n x_j^{(i)} + \alpha_j - 1 \right) \log \pi_j + \lambda g(\boldsymbol{\pi})$$

Here we ignore the normalizing constants which become additive constants in the log posterior optimization problem. The solution is found where

$$\begin{aligned} \nabla_{\boldsymbol{\pi}} \log p(\boldsymbol{\pi} | \mathbf{x}_{1:n}) + \lambda g(\boldsymbol{\pi}) &= 0 \\ g(\boldsymbol{\pi}) &= 0 \end{aligned}$$

The first condition yields

$$[\nabla_{\boldsymbol{\pi}} \log p(\boldsymbol{\pi} | \mathbf{x}_{1:n}) + \lambda g(\boldsymbol{\pi})]_{\ell} \Big|_{\substack{\pi_{\ell} = \pi_{\ell}^* \\ \lambda = \lambda^*}} = 0 \implies \frac{\sum_{i=1}^n x_{\ell}^{(i)} + \alpha_{\ell} - 1}{\pi_{\ell}^*} = \lambda^*$$

Replacing this result in the second condition, we get

$$1 - \sum_{j=1}^k \frac{\sum_{i=1}^n x_j^{(i)} + \alpha_j - 1}{\lambda^*} = 0 \implies \lambda^* = n + \sum_{j=1}^k \alpha_j - 1$$

Where we swapped the sum over the $x_j^{(i)}$ and used the fact that \mathbf{x}_j are one hot vectors. Thus

$$\boxed{(\hat{\boldsymbol{\pi}}_{\text{MAP}})_j = \pi_j^* = \frac{\sum_{i=1}^n x_j^{(i)} + \alpha_j - 1}{n + \alpha_j - 1} \in [0, 1]}$$

3 Properties of estimators

a) Poisson

Let n trials $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$ where $\lambda = \mathbb{E}_x[x]$. The pmf of the Poisson is

$$p(x | \lambda) = \frac{\lambda^x}{x!} e^{-\lambda} \quad \forall x \in \mathbb{N}$$

Such that the pmf of n trials should be

$$p(x_{1:n} | \lambda) \propto \prod_{j=1}^n p(x_j | \lambda)$$

I MLE

Using the log likelihood, we define the MLE estimation of λ as

$$\hat{\lambda}_{\text{MLE}} = \underset{\lambda \in \mathbb{R}_{>0}}{\text{argmax}} \sum_{j=1}^n (x_j \log \lambda - \lambda)$$

Which is found where

$$\left. \nabla_{\lambda} \log p(x_{1:n} | \lambda) \right|_{\lambda=\lambda^*} = 0$$

That is

$$\nabla_{\lambda} \log p(x_{1:n} | \lambda) = \lambda^{-1} \sum_{j=1}^n (x_j - 1)$$

Thus

$$\boxed{\hat{\lambda}_{\text{MLE}} = \frac{1}{n} \sum_{j=1}^n x_j}$$

II Bias

The bias is defined as

$$\text{Bias}(\lambda, \hat{\lambda}_{\text{MLE}}) \equiv \mathbb{E}_x[\hat{\lambda}_{\text{MLE}}] - \lambda$$

The expectation value of the MLE estimator is

$$\begin{aligned} \mathbb{E}_x[\hat{\lambda}_{\text{MLE}}] &= \mathbb{E}_x \left[\frac{1}{n} \sum_{j=1}^n x'_j \right] \\ &= \frac{1}{n} \sum_{j=1}^n \mathbb{E}_x[x_j] \\ &= \lambda \end{aligned}$$

Therefore the **MLE estimator of a Poisson distribution is an unbiased estimator.**

III Variance

The variance of the estimator is

$$\text{Var}(\hat{\lambda}_{\text{MLE}}) \equiv \mathbb{E}_X[\hat{\lambda}_{\text{MLE}}^2] - \mathbb{E}_x[\hat{\lambda}_{\text{MLE}}]^2$$

We need to evaluate the first term. To do this, we first use the Multinomial theorem to expand $\hat{\lambda}_{\text{MLE}}^2$:

$$\hat{\lambda}_{\text{MLE}}^2 = \frac{1}{n^2} \sum_{k_1 + \dots + k_n = 2} \binom{2}{k_1, \dots, k_n} x_1^{k_1} \dots x_n^{k_n}$$

Then we use both the linearity of the expectation operator and the fact that the R.V. X_1, \dots, X_n are independent to factorize the expectation of a "cross" product

$$\mathbb{E}_x[X_i X_j] = \mathbb{E}_x[X_i] \mathbb{E}_x[X_j], \quad \forall i, j \in \{\text{iid}\}$$

to get

$$\mathbb{E}_x[\hat{\lambda}_{\text{MLE}}^2] = \frac{1}{n^2} \sum_{k_1 + \dots + k_n = 2} \binom{2}{k_{1:n}} \mathbb{E}_x[x_1^{k_1}] \dots \mathbb{E}_x[x_n^{k_n}]$$

The sum can be separated into quadratic and linear term, s.t.

$$\mathbb{E}_x[\hat{\lambda}_{\text{MLE}}^2] = \frac{1}{n^2} \left(n \mathbb{E}_x[x^2] + 2 \binom{n}{2} \lambda^2 \right)$$

We used the fact that $\mathbb{E}_x[1] = 1$ and $\mathbb{E}_x[x_j] = \lambda$, $\forall j$. To estimate the quadratic term, we can use a magic trick by adding zero inside the operator argument. Using its linear property

$$\mathbb{E}_x[x^2] = \mathbb{E}_x[x(x-1) + x] = \mathbb{E}_x[x(x-1)] + \lambda$$

It turns out that

$$\begin{aligned} \mathbb{E}_x[x(x-1)] &= \sum_{x=0}^{\infty} x(x-1) \frac{\lambda^x}{x!} e^{-\lambda} \\ &= e^{-\lambda} \sum_{x=2}^{\infty} \frac{\lambda^x}{(x-2)!} \\ &= e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^{x+2}}{x!} \\ &= \lambda^2 \end{aligned}$$

By noticing the sum is the Taylor series of e^λ . In the end, we get

$$\mathbb{E}_x[\hat{\lambda}_{\text{MLE}}^2] = \frac{1}{n^2} (n\lambda + n^2\lambda^2)$$

Where we expanded the binom coefficient $2 \binom{n}{2} = n(n-1)$. The variance is thus

$$\text{Var}(\hat{\lambda}_{\text{MLE}}) = \frac{\lambda}{n}$$

IV Consistency

As $n \rightarrow \infty$, the estimator give an unbiased estimate of λ with a variance that goes to 0. Thus, the **estimator is consistent**.

b) Bernoulli

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ and let $n > 10$. We consider the estimator

$$\hat{p} \equiv \frac{1}{10} \sum_{i=1}^{10} X_i$$

I Bias

We first note that the expected value of a Bernoulli is

$$\mathbb{E}_x[x] = p$$

Since $x \in \{0, 1\}$ and p is the probability that $X = 1$. Therefore,

$$\text{Bias}(p, \hat{p}) = \frac{1}{10} \sum_{j=1}^{10} \mathbb{E}_x[x_j] - p = 0$$

Therefore \hat{p} is an unbiased estimator.

II Variance

The variance is

$$\begin{aligned}\text{Var}(\hat{p}) &= \mathbb{E}_x[\hat{p}^2] - \mathbb{E}_x^2[\hat{p}] \\ &= \frac{1}{100} (10 \mathbb{E}_x[x^2] + 90 \mathbb{E}_x^2[x]) - p^2 \\ &= \frac{p}{10} + p^2 \left(\frac{90}{100} - 1 \right) \\ &= \frac{1}{10} (p - p^2)\end{aligned}$$

III Consistency

This estimator is not consistent since the variance is constant as $n \rightarrow \infty$.

c) Uniform

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Uniform}(0, \theta)$. The pdf of this distribution is

$$p(x_i \mid 0, \theta) = \begin{cases} \theta^{-1}, & x \in [0, \theta] \\ 0, & \text{otherwise} \end{cases}$$

for $\theta \in \mathbb{R}_{>0}$. Therefore the likelihood will be proportional to

$$p(x_{1:n} \mid 0, \theta) \propto \begin{cases} \theta^{-n}, & x_{1:n} \in [0, \theta]^k \\ 0, & \text{otherwise} \end{cases}$$

I MLE

Assuming $x \in [0, \theta]$, we can compute the MLE using the log likelihood

$$\hat{\theta}_{\text{MLE}} = \underset{\theta \in \mathbb{R}_{>0}}{\text{argmax}} (-1) \log \theta$$