

In this chapter we present methods to analyze categorical attributes. Because categorical attributes have only symbolic values, many of the arithmetic operations cannot be performed directly on the symbolic values. However, we can compute the frequencies of these values and use them to analyze the attributes.

3.1 UNIVARIATE ANALYSIS

We assume that the data consists of values for a single categorical attribute, X . Let the domain of X consist of m symbolic values $dom(X) = \{a_1, a_2, \dots, a_m\}$. The data \mathbf{D} is thus an $n \times 1$ symbolic data matrix given as

$$\mathbf{D} = \begin{pmatrix} X \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

where each point $x_i \in dom(X)$.

3.1.1 Bernoulli Variable

Let us first consider the case when the categorical attribute X has domain $\{a_1, a_2\}$, with $m = 2$. We can model X as a Bernoulli random variable, which takes on two distinct values, 1 and 0, according to the mapping

$$X(v) = \begin{cases} 1 & \text{if } v = a_1 \\ 0 & \text{if } v = a_2 \end{cases}$$

The probability mass function (PMF) of X is given as

$$P(X = x) = f(x) = \begin{cases} p_1 & \text{if } x = 1 \\ p_0 & \text{if } x = 0 \end{cases}$$

where p_1 and p_0 are the parameters of the distribution, which must satisfy the condition

$$p_1 + p_0 = 1$$

Because there is only one free parameter, it is customary to denote $p_1 = p$, from which it follows that $p_0 = 1 - p$. The PMF of Bernoulli random variable X can then be written compactly as

$$P(X = x) = f(x) = p^x (1 - p)^{1-x}$$

We can see that $P(X = 1) = p^1 (1 - p)^0 = p$ and $P(X = 0) = p^0 (1 - p)^1 = 1 - p$, as desired.

Mean and Variance

The expected value of X is given as

$$\mu = E[X] = 1 \cdot p + 0 \cdot (1 - p) = p$$

and the variance of X is given as

$$\begin{aligned} \sigma^2 &= \text{var}(X) = E[X^2] - (E[X])^2 \\ &= (1^2 \cdot p + 0^2 \cdot (1 - p)) - p^2 = p - p^2 = p(1 - p) \end{aligned} \quad (3.1)$$

Sample Mean and Variance

To estimate the parameters of the Bernoulli variable X , we assume that each symbolic point has been mapped to its binary value. Thus, the set $\{x_1, x_2, \dots, x_n\}$ is assumed to be a random sample drawn from X (i.e., each x_i is IID with X).

The sample mean is given as

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{n_1}{n} = \hat{p} \quad (3.2)$$

where n_1 is the number of points with $x_i = 1$ in the random sample (equal to the number of occurrences of symbol a_1).

Let $n_0 = n - n_1$ denote the number of points with $x_i = 0$ in the random sample. The sample variance is given as

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \\ &= \frac{n_1}{n} (1 - \hat{p})^2 + \frac{n - n_1}{n} (-\hat{p})^2 \\ &= \hat{p} (1 - \hat{p})^2 + (1 - \hat{p}) \hat{p}^2 \\ &= \hat{p} (1 - \hat{p}) (1 - \hat{p} + \hat{p}) \\ &= \hat{p} (1 - \hat{p}) \end{aligned}$$

The sample variance could also have been obtained directly from Eq.(3.1), by substituting \hat{p} for p .

Example 3.1. Consider the sepal length attribute (X_1) for the Iris dataset in Table 1.1. Let us define an Iris flower as Long if its sepal length is in the range $[7, \infty]$, and Short if its sepal length is in the range $[-\infty, 7)$. Then X_1 can be treated as a categorical attribute with domain {Long, Short}. From the observed sample of size $n = 150$, we find 13 long Irises. The sample mean of X_1 is

$$\hat{\mu} = \hat{p} = 13/150 = 0.087$$

and its variance is

$$\hat{\sigma}^2 = \hat{p}(1 - \hat{p}) = 0.087(1 - 0.087) = 0.087 \cdot 0.913 = 0.079$$

Binomial Distribution: Number of Occurrences

Given the Bernoulli variable X , let $\{x_1, x_2, \dots, x_n\}$ denote a random sample of size n drawn from X . Let N be the random variable denoting the number of occurrences of the symbol a_1 (value $X = 1$) in the sample. N has a binomial distribution, given as

$$f(N = n_1 | n, p) = \binom{n}{n_1} p^{n_1} (1 - p)^{n - n_1} \quad (3.3)$$

In fact, N is the sum of the n independent Bernoulli random variables x_i IID with X , that is, $N = \sum_{i=1}^n x_i$. By linearity of expectation, the mean or expected number of occurrences of symbol a_1 is given as

$$\mu_N = E[N] = E\left[\sum_{i=1}^n x_i\right] = \sum_{i=1}^n E[x_i] = \sum_{i=1}^n p = np$$

Because x_i are all independent, the variance of N is given as

$$\sigma_N^2 = \text{var}(N) = \sum_{i=1}^n \text{var}(x_i) = \sum_{i=1}^n p(1 - p) = np(1 - p)$$

Example 3.2. Continuing with Example 3.1, we can use the estimated parameter $\hat{p} = 0.087$ to compute the expected number of occurrences N of Long sepal length Irises via the binomial distribution:

$$E[N] = n\hat{p} = 150 \cdot 0.087 = 13$$

In this case, because p is estimated from the sample via \hat{p} , it is not surprising that the expected number of occurrences of long Irises coincides with the actual occurrences. However, what is more interesting is that we can compute the variance in the number of occurrences:

$$\text{var}(N) = n\hat{p}(1 - \hat{p}) = 150 \cdot 0.079 = 11.9$$

As the sample size increases, the binomial distribution given in Eq. 3.3 tends to a normal distribution with $\mu = 13$ and $\sigma = \sqrt{11.9} = 3.45$ for our example. Thus, with confidence greater than 95% we can claim that the number of occurrences of a_1 will lie in the range $\mu \pm 2\sigma = [9.55, 16.45]$, which follows from the fact that for a normal distribution 95.45% of the probability mass lies within two standard deviations from the mean (see Section 2.5.1).

3.1.2 Multivariate Bernoulli Variable

We now consider the general case when X is a categorical attribute with domain $\{a_1, a_2, \dots, a_m\}$. We can model X as an m -dimensional Bernoulli random variable $\mathbf{X} = (A_1, A_2, \dots, A_m)^T$, where each A_i is a Bernoulli variable with parameter p_i denoting the probability of observing symbol a_i . However, because X can assume only one of the symbolic values at any one time, if $X = a_i$, then $A_i = 1$, and $A_j = 0$ for all $j \neq i$. The range of the random variable \mathbf{X} is thus the set $\{0, 1\}^m$, with the further restriction that if $X = a_i$, then $\mathbf{X} = \mathbf{e}_i$, where \mathbf{e}_i is the i th standard basis vector $\mathbf{e}_i \in \mathbb{R}^m$ given as

$$\mathbf{e}_i = (\overbrace{0, \dots, 0}^{i-1}, 1, \overbrace{0, \dots, 0}^{m-i})^T$$

In \mathbf{e}_i , only the i th element is 1 ($e_{ii} = 1$), whereas all other elements are zero ($e_{ij} = 0, \forall j \neq i$).

This is precisely the definition of a *multivariate Bernoulli variable*, which is a generalization of a Bernoulli variable from two outcomes to m outcomes. We thus model the categorical attribute X as a multivariate Bernoulli variable \mathbf{X} defined as

$$\mathbf{X}(v) = \mathbf{e}_i \text{ if } v = a_i$$

The range of \mathbf{X} consists of m distinct vector values $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m\}$, with the PMF of \mathbf{X} given as

$$P(\mathbf{X} = \mathbf{e}_i) = f(\mathbf{e}_i) = p_i$$

where p_i is the probability of observing value a_i . These parameters must satisfy the condition

$$\sum_{i=1}^m p_i = 1$$

The PMF can be written compactly as follows:

$$P(\mathbf{X} = \mathbf{e}_i) = f(\mathbf{e}_i) = \prod_{j=1}^m p_j^{e_{ij}} \quad (3.4)$$

Because $e_{ii} = 1$, and $e_{ij} = 0$ for $j \neq i$, we can see that, as expected, we have

$$f(\mathbf{e}_i) = \prod_{j=1}^m p_j^{e_{ij}} = p_1^{e_{i0}} \times \dots \times p_i^{e_{ii}} \dots \times p_m^{e_{im}} = p_1^0 \times \dots \times p_i^1 \dots \times p_m^0 = p_i$$

Table 3.1. Discretized sepal length attribute

Bins	Domain	Counts
[4.3, 5.2]	Very Short (a_1)	$n_1 = 45$
(5.2, 6.1]	Short (a_2)	$n_2 = 50$
(6.1, 7.0]	Long (a_3)	$n_3 = 43$
(7.0, 7.9]	Very Long (a_4)	$n_4 = 12$

Example 3.3. Let us consider the sepal length attribute (X_1) for the Iris dataset shown in Table 1.2. We divide the sepal length into four equal-width intervals, and give each interval a name as shown in Table 3.1. We consider X_1 as a categorical attribute with domain

$$\{a_1 = \text{VeryShort}, a_2 = \text{Short}, a_3 = \text{Long}, a_4 = \text{VeryLong}\}$$

We model the categorical attribute X_1 as a multivariate Bernoulli variable \mathbf{X} , defined as

$$\mathbf{X}(v) = \begin{cases} \mathbf{e}_1 = (1, 0, 0, 0) & \text{if } v = a_1 \\ \mathbf{e}_2 = (0, 1, 0, 0) & \text{if } v = a_2 \\ \mathbf{e}_3 = (0, 0, 1, 0) & \text{if } v = a_3 \\ \mathbf{e}_4 = (0, 0, 0, 1) & \text{if } v = a_4 \end{cases}$$

For example, the symbolic point $x_1 = \text{Short} = a_2$ is represented as the vector $(0, 1, 0, 0)^T = \mathbf{e}_2$.

Mean

The mean or expected value of \mathbf{X} can be obtained as

$$\boldsymbol{\mu} = E[\mathbf{X}] = \sum_{i=1}^m \mathbf{e}_i f(\mathbf{e}_i) = \sum_{i=1}^m \mathbf{e}_i p_i = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} p_1 + \cdots + \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} p_m = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \\ p_m \end{pmatrix} = \mathbf{p} \quad (3.5)$$

Sample Mean

Assume that each symbolic point $x_i \in \mathbf{D}$ is mapped to the variable $\mathbf{x}_i = \mathbf{X}(x_i)$. The mapped dataset $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ is then assumed to be a random sample IID with \mathbf{X} . We can compute the sample mean by placing a probability mass of $\frac{1}{n}$ at each point

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \sum_{i=1}^m \frac{n_i}{n} \mathbf{e}_i = \begin{pmatrix} n_1/n \\ n_2/n \\ \vdots \\ n_m/n \end{pmatrix} = \begin{pmatrix} \hat{p}_1 \\ \hat{p}_2 \\ \vdots \\ \hat{p}_m \end{pmatrix} = \hat{\mathbf{p}} \quad (3.6)$$

where n_i is the number of occurrences of the vector value \mathbf{e}_i in the sample, which is equivalent to the number of occurrences of the symbol a_i . Furthermore, we have

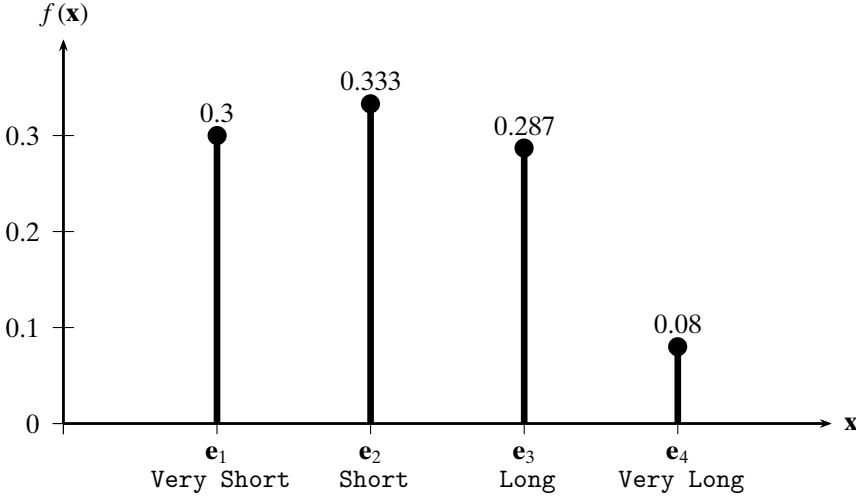


Figure 3.1. Probability mass function: sepal length.

$\sum_{i=1}^m n_i = n$, which follows from the fact that \mathbf{X} can take on only m distinct values \mathbf{e}_i , and the counts for each value must add up to the sample size n .

Example 3.4 (Sample Mean). Consider the observed counts n_i for each of the values a_i (\mathbf{e}_i) of the discretized sepal length attribute, shown in Table 3.1. Because the total sample size is $n = 150$, from these we can obtain the estimates \hat{p}_i as follows:

$$\hat{p}_1 = 45/150 = 0.3$$

$$\hat{p}_2 = 50/150 = 0.333$$

$$\hat{p}_3 = 43/150 = 0.287$$

$$\hat{p}_4 = 12/150 = 0.08$$

The PMF for \mathbf{X} is plotted in Figure 3.1, and the sample mean for \mathbf{X} is given as

$$\hat{\boldsymbol{\mu}} = \hat{\mathbf{p}} = \begin{pmatrix} 0.3 \\ 0.333 \\ 0.287 \\ 0.08 \end{pmatrix}$$

Covariance Matrix

Recall that an m -dimensional multivariate Bernoulli variable is simply a vector of m Bernoulli variables. For instance, $\mathbf{X} = (A_1, A_2, \dots, A_m)^T$, where A_i is the Bernoulli variable corresponding to symbol a_i . The variance–covariance information between the constituent Bernoulli variables yields a covariance matrix for \mathbf{X} .

Let us first consider the variance along each Bernoulli variable A_i . By Eq. (3.1), we immediately have

$$\sigma_i^2 = \text{var}(A_i) = p_i(1 - p_i)$$

Next consider the covariance between A_i and A_j . Utilizing the identity in Eq. (2.21), we have

$$\sigma_{ij} = E[A_i A_j] - E[A_i] \cdot E[A_j] = 0 - p_i p_j = -p_i p_j$$

which follows from the fact that $E[A_i A_j] = 0$, as A_i and A_j cannot both be 1 at the same time, and thus their product $A_i A_j = 0$. This same fact leads to the negative relationship between A_i and A_j . What is interesting is that the degree of negative association is proportional to the product of the mean values for A_i and A_j .

From the preceding expressions for variance and covariance, the $m \times m$ covariance matrix for \mathbf{X} is given as

$$\mathbf{\Sigma} = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{12} & \sigma_2^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \dots & \sigma_m^2 \end{pmatrix} = \begin{pmatrix} p_1(1 - p_1) & -p_1 p_2 & \dots & -p_1 p_m \\ -p_1 p_2 & p_2(1 - p_2) & \dots & -p_2 p_m \\ \vdots & \vdots & \ddots & \vdots \\ -p_1 p_m & -p_2 p_m & \dots & p_m(1 - p_m) \end{pmatrix}$$

Notice how each row in $\mathbf{\Sigma}$ sums to zero. For example, for row i , we have

$$-p_i p_1 - p_i p_2 - \dots + p_i(1 - p_i) - \dots - p_i p_m = p_i - p_i \sum_{j=1}^m p_j = p_i - p_i = 0 \quad (3.7)$$

Because $\mathbf{\Sigma}$ is symmetric, it follows that each column also sums to zero.

Define \mathbf{P} as the $m \times m$ diagonal matrix:

$$\mathbf{P} = \text{diag}(\mathbf{p}) = \text{diag}(p_1, p_2, \dots, p_m) = \begin{pmatrix} p_1 & 0 & \dots & 0 \\ 0 & p_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & p_m \end{pmatrix}$$

We can compactly write the covariance matrix of \mathbf{X} as

$$\mathbf{\Sigma} = \mathbf{P} - \mathbf{p} \cdot \mathbf{p}^T \quad (3.8)$$

Sample Covariance Matrix

The sample covariance matrix can be obtained from Eq. (3.8) in a straightforward manner:

$$\widehat{\mathbf{\Sigma}} = \widehat{\mathbf{P}} - \hat{\mathbf{p}} \cdot \hat{\mathbf{p}}^T \quad (3.9)$$

where $\widehat{\mathbf{P}} = \text{diag}(\hat{\mathbf{p}})$, and $\hat{\mathbf{p}} = \hat{\boldsymbol{\mu}} = (\hat{p}_1, \hat{p}_2, \dots, \hat{p}_m)^T$ denotes the empirical probability mass function for \mathbf{X} .

Example 3.5. Returning to the discretized `sepal length` attribute in Example 3.4, we have $\hat{\boldsymbol{\mu}} = \hat{\mathbf{p}} = (0.3, 0.333, 0.287, 0.08)^T$. The sample covariance matrix is given as

$$\begin{aligned}\hat{\boldsymbol{\Sigma}} &= \hat{\mathbf{P}} - \hat{\mathbf{p}} \cdot \hat{\mathbf{p}}^T \\ &= \begin{pmatrix} 0.3 & 0 & 0 & 0 \\ 0 & 0.333 & 0 & 0 \\ 0 & 0 & 0.287 & 0 \\ 0 & 0 & 0 & 0.08 \end{pmatrix} - \begin{pmatrix} 0.3 \\ 0.333 \\ 0.287 \\ 0.08 \end{pmatrix} (0.3 \quad 0.333 \quad 0.287 \quad 0.08) \\ &= \begin{pmatrix} 0.3 & 0 & 0 & 0 \\ 0 & 0.333 & 0 & 0 \\ 0 & 0 & 0.287 & 0 \\ 0 & 0 & 0 & 0.08 \end{pmatrix} - \begin{pmatrix} 0.09 & 0.1 & 0.086 & 0.024 \\ 0.1 & 0.111 & 0.096 & 0.027 \\ 0.086 & 0.096 & 0.082 & 0.023 \\ 0.024 & 0.027 & 0.023 & 0.006 \end{pmatrix} \\ &= \begin{pmatrix} 0.21 & -0.1 & -0.086 & -0.024 \\ -0.1 & 0.222 & -0.096 & -0.027 \\ -0.086 & -0.096 & 0.204 & -0.023 \\ -0.024 & -0.027 & -0.023 & 0.074 \end{pmatrix}\end{aligned}$$

One can verify that each row (and column) in $\hat{\boldsymbol{\Sigma}}$ sums to zero.

It is worth emphasizing that whereas the modeling of categorical attribute X as a multivariate Bernoulli variable, $\mathbf{X} = (A_1, A_2, \dots, A_m)^T$, makes the structure of the mean and covariance matrix explicit, the same results would be obtained if we simply treat the mapped values $\mathbf{X}(x_i)$ as a new $n \times m$ binary data matrix, and apply the standard definitions of the mean and covariance matrix from multivariate numeric attribute analysis (see Section 2.3). In essence, the mapping from symbols a_i to binary vectors \mathbf{e}_i is the key idea in categorical attribute analysis.

Example 3.6. Consider the sample \mathbf{D} of size $n = 5$ for the `sepal length` attribute X_1 in the Iris dataset, shown in Table 3.2a. As in Example 3.1, we assume that X_1 has only two categorical values {Long, Short}. We model X_1 as the multivariate Bernoulli variable \mathbf{X}_1 defined as

$$\mathbf{X}_1(v) = \begin{cases} \mathbf{e}_1 = (1, 0)^T & \text{if } v = \text{Long}(a_1) \\ \mathbf{e}_2 = (0, 1)^T & \text{if } v = \text{Short}(a_2) \end{cases}$$

The sample mean [Eq. (3.6)] is

$$\hat{\boldsymbol{\mu}} = \hat{\mathbf{p}} = (2/5, 3/5)^T = (0.4, 0.6)^T$$

and the sample covariance matrix [Eq. (3.9)] is

$$\begin{aligned}\hat{\boldsymbol{\Sigma}} &= \hat{\mathbf{P}} - \hat{\mathbf{p}}\hat{\mathbf{p}}^T = \begin{pmatrix} 0.4 & 0 \\ 0 & 0.6 \end{pmatrix} - \begin{pmatrix} 0.4 \\ 0.6 \end{pmatrix} (0.4 \quad 0.6) \\ &= \begin{pmatrix} 0.4 & 0 \\ 0 & 0.6 \end{pmatrix} - \begin{pmatrix} 0.16 & 0.24 \\ 0.24 & 0.36 \end{pmatrix} = \begin{pmatrix} 0.24 & -0.24 \\ -0.24 & 0.24 \end{pmatrix}\end{aligned}$$

Table 3.2. (a) Categorical dataset. (b) Mapped binary dataset. (c) Centered dataset.

(a)		(b)			(c)		
	X		A_1	A_2		Z_1	Z_2
x_1	Short	\mathbf{x}_1	0	1	\mathbf{z}_1	-0.4	0.4
x_2	Short	\mathbf{x}_2	0	1	\mathbf{z}_2	-0.4	0.4
x_3	Long	\mathbf{x}_3	1	0	\mathbf{z}_3	0.6	-0.6
x_4	Short	\mathbf{x}_4	0	1	\mathbf{z}_4	-0.4	0.4
x_5	Long	\mathbf{x}_5	1	0	\mathbf{z}_5	0.6	-0.6

To show that the same result would be obtained via standard numeric analysis, we map the categorical attribute X to the two Bernoulli attributes A_1 and A_2 corresponding to symbols Long and Short, respectively. The mapped dataset is shown in Table 3.2b. The sample mean is simply

$$\hat{\boldsymbol{\mu}} = \frac{1}{5} \sum_{i=1}^5 \mathbf{x}_i = \frac{1}{5} (2, 3)^T = (0.4, 0.6)^T$$

Next, we center the dataset by subtracting the mean value from each attribute. After centering, the mapped dataset is as shown in Table 3.2c, with attribute Z_i as the centered attribute A_i . We can compute the covariance matrix using the inner-product form [Eq. (2.30)] on the centered column vectors. We have

$$\sigma_1^2 = \frac{1}{5} \mathbf{Z}_1^T \mathbf{Z}_1 = 1.2/5 = 0.24$$

$$\sigma_2^2 = \frac{1}{5} \mathbf{Z}_2^T \mathbf{Z}_2 = 1.2/5 = 0.24$$

$$\sigma_{12} = \frac{1}{5} \mathbf{Z}_1^T \mathbf{Z}_2 = -1.2/5 = -0.24$$

Thus, the sample covariance matrix is given as

$$\hat{\boldsymbol{\Sigma}} = \begin{pmatrix} 0.24 & -0.24 \\ -0.24 & 0.24 \end{pmatrix}$$

which matches the result obtained by using the multivariate Bernoulli modeling approach.

Multinomial Distribution: Number of Occurrences

Given a multivariate Bernoulli variable \mathbf{X} and a random sample $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ drawn from \mathbf{X} . Let N_i be the random variable corresponding to the number of occurrences of symbol a_i in the sample, and let $\mathbf{N} = (N_1, N_2, \dots, N_m)^T$ denote the vector random variable corresponding to the joint distribution of the number of occurrences over all the symbols. Then \mathbf{N} has a multinomial distribution, given as

$$f(\mathbf{N} = (n_1, n_2, \dots, n_m) \mid \mathbf{p}) = \binom{n}{n_1 n_2 \dots n_m} \prod_{i=1}^m p_i^{n_i}$$

We can see that this is a direct generalization of the binomial distribution in Eq. (3.3). The term

$$\binom{n}{n_1 n_2 \dots n_m} = \frac{n!}{n_1! n_2! \dots n_m!}$$

denotes the number of ways of choosing n_i occurrences of each symbol a_i from a sample of size n , with $\sum_{i=1}^m n_i = n$.

The mean and covariance matrix of \mathbf{N} are given as n times the mean and covariance matrix of \mathbf{X} . That is, the mean of \mathbf{N} is given as

$$\boldsymbol{\mu}_{\mathbf{N}} = E[\mathbf{N}] = nE[\mathbf{X}] = n \cdot \boldsymbol{\mu} = n \cdot \mathbf{p} = \begin{pmatrix} np_1 \\ \vdots \\ np_m \end{pmatrix}$$

and its covariance matrix is given as

$$\boldsymbol{\Sigma}_{\mathbf{N}} = n \cdot (\mathbf{P} - \mathbf{p}\mathbf{p}^T) = \begin{pmatrix} np_1(1-p_1) & -np_1p_2 & \cdots & -np_1p_m \\ -np_1p_2 & np_2(1-p_2) & \cdots & -np_2p_m \\ \vdots & \vdots & \ddots & \vdots \\ -np_1p_m & -np_2p_m & \cdots & np_m(1-p_m) \end{pmatrix}$$

Likewise the sample mean and covariance matrix for \mathbf{N} are given as

$$\hat{\boldsymbol{\mu}}_{\mathbf{N}} = n\hat{\mathbf{p}} \quad \hat{\boldsymbol{\Sigma}}_{\mathbf{N}} = n(\hat{\mathbf{P}} - \hat{\mathbf{p}}\hat{\mathbf{p}}^T)$$

3.2 BIVARIATE ANALYSIS

Assume that the data comprises two categorical attributes, X_1 and X_2 , with

$$\text{dom}(X_1) = \{a_{11}, a_{12}, \dots, a_{1m_1}\}$$

$$\text{dom}(X_2) = \{a_{21}, a_{22}, \dots, a_{2m_2}\}$$

We are given n categorical points of the form $\mathbf{x}_i = (x_{i1}, x_{i2})^T$ with $x_{i1} \in \text{dom}(X_1)$ and $x_{i2} \in \text{dom}(X_2)$. The dataset is thus an $n \times 2$ symbolic data matrix:

$$\mathbf{D} = \begin{pmatrix} X_1 & X_2 \\ x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{pmatrix}$$

We can model X_1 and X_2 as multivariate Bernoulli variables \mathbf{X}_1 and \mathbf{X}_2 with dimensions m_1 and m_2 , respectively. The probability mass functions for \mathbf{X}_1 and \mathbf{X}_2 are

given according to Eq. (3.4):

$$P(\mathbf{X}_1 = \mathbf{e}_{1i}) = f_1(\mathbf{e}_{1i}) = p_i^1 = \prod_{k=1}^{m_1} (p_i^1)^{e_{ik}^1}$$

$$P(\mathbf{X}_2 = \mathbf{e}_{2j}) = f_2(\mathbf{e}_{2j}) = p_j^2 = \prod_{k=1}^{m_2} (p_j^2)^{e_{jk}^2}$$

where \mathbf{e}_{1i} is the i th standard basis vector in \mathbb{R}^{m_1} (for attribute X_1) whose k th component is e_{ik}^1 , and \mathbf{e}_{2j} is the j th standard basis vector in \mathbb{R}^{m_2} (for attribute X_2) whose k th component is e_{jk}^2 . Further, the parameter p_i^1 denotes the probability of observing symbol a_{1i} , and p_j^2 denotes the probability of observing symbol a_{2j} . Together they must satisfy the conditions: $\sum_{i=1}^{m_1} p_i^1 = 1$ and $\sum_{j=1}^{m_2} p_j^2 = 1$.

The joint distribution of \mathbf{X}_1 and \mathbf{X}_2 is modeled as the $d' = m_1 + m_2$ dimensional vector variable $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$, specified by the mapping

$$\mathbf{X}((v_1, v_2)^T) = \begin{pmatrix} \mathbf{X}_1(v_1) \\ \mathbf{X}_2(v_2) \end{pmatrix} = \begin{pmatrix} \mathbf{e}_{1i} \\ \mathbf{e}_{2j} \end{pmatrix}$$

provided that $v_1 = a_{1i}$ and $v_2 = a_{2j}$. The range of \mathbf{X} thus consists of $m_1 \times m_2$ distinct pairs of vector values $\{(\mathbf{e}_{1i}, \mathbf{e}_{2j})^T\}$, with $1 \leq i \leq m_1$ and $1 \leq j \leq m_2$. The joint PMF of \mathbf{X} is given as

$$P(\mathbf{X} = (\mathbf{e}_{1i}, \mathbf{e}_{2j})^T) = f(\mathbf{e}_{1i}, \mathbf{e}_{2j}) = p_{ij} = \prod_{r=1}^{m_1} \prod_{s=1}^{m_2} p_{ij}^{e_{ir}^1 \cdot e_{js}^2}$$

where p_{ij} the probability of observing the symbol pair (a_{1i}, a_{2j}) . These probability parameters must satisfy the condition $\sum_{i=1}^{m_1} \sum_{j=1}^{m_2} p_{ij} = 1$. The joint PMF for \mathbf{X} can be expressed as the $m_1 \times m_2$ matrix

$$\mathbf{P}_{12} = \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1m_2} \\ p_{21} & p_{22} & \cdots & p_{2m_2} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m_1 1} & p_{m_1 2} & \cdots & p_{m_1 m_2} \end{pmatrix} \quad (3.10)$$

Example 3.7. Consider the discretized `sepal length` attribute (X_1) in Table 3.1. We also discretize the `sepal width` attribute (X_2) into three values as shown in Table 3.3. We thus have

$$\text{dom}(X_1) = \{a_{11} = \text{VeryShort}, a_{12} = \text{Short}, a_{13} = \text{Long}, a_{14} = \text{VeryLong}\}$$

$$\text{dom}(X_2) = \{a_{21} = \text{Short}, a_{22} = \text{Medium}, a_{23} = \text{Long}\}$$

The symbolic point $\mathbf{x} = (\text{Short}, \text{Long}) = (a_{12}, a_{23})$, is mapped to the vector

$$\mathbf{X}(\mathbf{x}) = \begin{pmatrix} \mathbf{e}_{12} \\ \mathbf{e}_{23} \end{pmatrix} = (0, 1, 0, 0 \mid 0, 0, 1)^T \in \mathbb{R}^7$$

Table 3.3. Discretized sepal width attribute

Bins	Domain	Counts
[2.0, 2.8]	Short (a_1)	47
(2.8, 3.6]	Medium (a_2)	88
(3.6, 4.4]	Long (a_3)	15

where we use $|$ to demarcate the two subvectors $\mathbf{e}_{12} = (0, 1, 0, 0)^T \in \mathbb{R}^4$ and $\mathbf{e}_{23} = (0, 0, 1)^T \in \mathbb{R}^3$, corresponding to symbolic attributes `sepal length` and `sepal width`, respectively. Note that \mathbf{e}_{12} is the second standard basis vector in \mathbb{R}^4 for \mathbf{X}_1 , and \mathbf{e}_{23} is the third standard basis vector in \mathbb{R}^3 for \mathbf{X}_2 .

Mean

The bivariate mean can easily be generalized from Eq. (3.5), as follows:

$$\boldsymbol{\mu} = E[\mathbf{X}] = E\left[\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}\right] = \begin{pmatrix} E[\mathbf{X}_1] \\ E[\mathbf{X}_2] \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \end{pmatrix}$$

where $\boldsymbol{\mu}_1 = \mathbf{p}_1 = (p_1^1, \dots, p_{m_1}^1)^T$ and $\boldsymbol{\mu}_2 = \mathbf{p}_2 = (p_1^2, \dots, p_{m_2}^2)^T$ are the mean vectors for \mathbf{X}_1 and \mathbf{X}_2 . The vectors \mathbf{p}_1 and \mathbf{p}_2 also represent the probability mass functions for \mathbf{X}_1 and \mathbf{X}_2 , respectively.

Sample Mean

The sample mean can also be generalized from Eq. (3.6), by placing a probability mass of $\frac{1}{n}$ at each point:

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} \begin{pmatrix} \sum_{i=1}^{m_1} n_i^1 \mathbf{e}_{1i} \\ \sum_{j=1}^{m_2} n_j^2 \mathbf{e}_{2j} \end{pmatrix} = \frac{1}{n} \begin{pmatrix} n_1^1 \\ \vdots \\ n_{m_1}^1 \\ n_1^2 \\ \vdots \\ n_{m_2}^2 \end{pmatrix} = \begin{pmatrix} \hat{p}_1^1 \\ \vdots \\ \hat{p}_{m_1}^1 \\ \hat{p}_1^2 \\ \vdots \\ \hat{p}_{m_2}^2 \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{p}}_1 \\ \hat{\mathbf{p}}_2 \end{pmatrix} = \begin{pmatrix} \hat{\boldsymbol{\mu}}_1 \\ \hat{\boldsymbol{\mu}}_2 \end{pmatrix}$$

where n_j^i is the observed frequency of symbol a_{ij} in the sample of size n , and $\hat{\boldsymbol{\mu}}_i = \hat{\mathbf{p}}_i = (p_1^i, p_2^i, \dots, p_{m_i}^i)^T$ is the sample mean vector for \mathbf{X}_i , which is also the empirical PMF for attribute \mathbf{X}_i .

Covariance Matrix

The covariance matrix for \mathbf{X} is the $d' \times d' = (m_1 + m_2) \times (m_1 + m_2)$ matrix given as

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{pmatrix} \quad (3.11)$$

where $\boldsymbol{\Sigma}_{11}$ is the $m_1 \times m_1$ covariance matrix for \mathbf{X}_1 , and $\boldsymbol{\Sigma}_{22}$ is the $m_2 \times m_2$ covariance matrix for \mathbf{X}_2 , which can be computed using Eq. (3.8). That is,

$$\begin{aligned} \boldsymbol{\Sigma}_{11} &= \mathbf{P}_1 - \mathbf{p}_1 \mathbf{p}_1^T \\ \boldsymbol{\Sigma}_{22} &= \mathbf{P}_2 - \mathbf{p}_2 \mathbf{p}_2^T \end{aligned}$$

where $\mathbf{P}_1 = \text{diag}(\mathbf{p}_1)$ and $\mathbf{P}_2 = \text{diag}(\mathbf{p}_2)$. Further, $\boldsymbol{\Sigma}_{12}$ is the $m_1 \times m_2$ covariance matrix between variables \mathbf{X}_1 and \mathbf{X}_2 , given as

$$\begin{aligned}\boldsymbol{\Sigma}_{12} &= E[(\mathbf{X}_1 - \boldsymbol{\mu}_1)(\mathbf{X}_2 - \boldsymbol{\mu}_2)^T] \\ &= E[\mathbf{X}_1 \mathbf{X}_2^T] - E[\mathbf{X}_1]E[\mathbf{X}_2]^T \\ &= \mathbf{P}_{12} - \boldsymbol{\mu}_1 \boldsymbol{\mu}_2^T \\ &= \mathbf{P}_{12} - \mathbf{p}_1 \mathbf{p}_2^T \\ &= \begin{pmatrix} p_{11} - p_1^1 p_1^2 & p_{12} - p_1^1 p_2^2 & \cdots & p_{1m_2} - p_1^1 p_{m_2}^2 \\ p_{21} - p_2^1 p_1^2 & p_{22} - p_2^1 p_2^2 & \cdots & p_{2m_2} - p_2^1 p_{m_2}^2 \\ \vdots & \vdots & \ddots & \vdots \\ p_{m_1 1} - p_{m_1}^1 p_1^2 & p_{m_1 2} - p_{m_1}^1 p_2^2 & \cdots & p_{m_1 m_2} - p_{m_1}^1 p_{m_2}^2 \end{pmatrix}\end{aligned}$$

where \mathbf{P}_{12} represents the joint PMF for \mathbf{X} given in Eq. (3.10).

Incidentally, each row and each column of $\boldsymbol{\Sigma}_{12}$ sums to zero. For example, consider row i and column j :

$$\begin{aligned}\sum_{k=1}^{m_2} (p_{ik} - p_i^1 p_k^2) &= \left(\sum_{k=1}^{m_2} p_{ik} \right) - p_i^1 = p_i^1 - p_i^1 = 0 \\ \sum_{k=1}^{m_1} (p_{kj} - p_k^1 p_j^2) &= \left(\sum_{k=1}^{m_1} p_{kj} \right) - p_j^2 = p_j^2 - p_j^2 = 0\end{aligned}$$

which follows from the fact that summing the joint mass function over all values of \mathbf{X}_2 , yields the marginal distribution of \mathbf{X}_1 , and summing it over all values of \mathbf{X}_1 yields the marginal distribution for \mathbf{X}_2 . Note that p_j^2 is the probability of observing symbol a_{2j} ; it should not be confused with the square of p_j . Combined with the fact that $\boldsymbol{\Sigma}_{11}$ and $\boldsymbol{\Sigma}_{22}$ also have row and column sums equal to zero via Eq. (3.7), the full covariance matrix $\boldsymbol{\Sigma}$ has rows and columns that sum up to zero.

Sample Covariance Matrix

The sample covariance matrix is given as

$$\hat{\boldsymbol{\Sigma}} = \begin{pmatrix} \hat{\boldsymbol{\Sigma}}_{11} & \hat{\boldsymbol{\Sigma}}_{12} \\ \hat{\boldsymbol{\Sigma}}_{12}^T & \hat{\boldsymbol{\Sigma}}_{22} \end{pmatrix} \quad (3.12)$$

where

$$\begin{aligned}\hat{\boldsymbol{\Sigma}}_{11} &= \hat{\mathbf{P}}_1 - \hat{\mathbf{p}}_1 \hat{\mathbf{p}}_1^T \\ \hat{\boldsymbol{\Sigma}}_{22} &= \hat{\mathbf{P}}_2 - \hat{\mathbf{p}}_2 \hat{\mathbf{p}}_2^T \\ \hat{\boldsymbol{\Sigma}}_{12} &= \hat{\mathbf{P}}_{12} - \hat{\mathbf{p}}_1 \hat{\mathbf{p}}_2^T\end{aligned}$$

Here $\hat{\mathbf{P}}_1 = \text{diag}(\hat{\mathbf{p}}_1)$ and $\hat{\mathbf{P}}_2 = \text{diag}(\hat{\mathbf{p}}_2)$, and $\hat{\mathbf{p}}_1$ and $\hat{\mathbf{p}}_2$ specify the empirical probability mass functions for \mathbf{X}_1 , and \mathbf{X}_2 , respectively. Further, $\hat{\mathbf{P}}_{12}$ specifies the empirical joint PMF for \mathbf{X}_1 and \mathbf{X}_2 , given as

$$\hat{\mathbf{P}}_{12}(i, j) = \hat{f}(\mathbf{e}_{1i}, \mathbf{e}_{2j}) = \frac{1}{n} \sum_{k=1}^n I_{ij}(\mathbf{x}_k) = \frac{n_{ij}}{n} = \hat{p}_{ij} \quad (3.13)$$

where I_{ij} is the indicator variable

$$I_{ij}(\mathbf{x}_k) = \begin{cases} 1 & \text{if } \mathbf{x}_{k1} = \mathbf{e}_{1i} \text{ and } \mathbf{x}_{k2} = \mathbf{e}_{2j} \\ 0 & \text{otherwise} \end{cases}$$

Taking the sum of $I_{ij}(\mathbf{x}_k)$ over all the n points in the sample yields the number of occurrences, n_{ij} , of the symbol pair (a_{1i}, a_{2j}) in the sample. One issue with the cross-attribute covariance matrix $\widehat{\Sigma}_{12}$ is the need to estimate a quadratic number of parameters. That is, we need to obtain reliable counts n_{ij} to estimate the parameters p_{ij} , for a total of $O(m_1 \times m_2)$ parameters that have to be estimated, which can be a problem if the categorical attributes have many symbols. On the other hand, estimating $\widehat{\Sigma}_{11}$ and $\widehat{\Sigma}_{22}$ requires that we estimate m_1 and m_2 parameters, corresponding to p_i^1 and p_j^2 , respectively. In total, computing Σ requires the estimation of $m_1 m_2 + m_1 + m_2$ parameters.

Example 3.8. We continue with the bivariate categorical attributes X_1 and X_2 in Example 3.7. From Example 3.4, and from the occurrence counts for each of the values of sepal width in Table 3.3, we have

$$\hat{\mu}_1 = \hat{\mathbf{p}}_1 = \begin{pmatrix} 0.3 \\ 0.333 \\ 0.287 \\ 0.08 \end{pmatrix} \quad \hat{\mu}_2 = \hat{\mathbf{p}}_2 = \frac{1}{150} \begin{pmatrix} 47 \\ 88 \\ 15 \end{pmatrix} = \begin{pmatrix} 0.313 \\ 0.587 \\ 0.1 \end{pmatrix}$$

Thus, the mean for $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$ is given as

$$\hat{\mu} = \begin{pmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{p}}_1 \\ \hat{\mathbf{p}}_2 \end{pmatrix} = (0.3, 0.333, 0.287, 0.08 \mid 0.313, 0.587, 0.1)^T$$

From Example 3.5 we have

$$\widehat{\Sigma}_{11} = \begin{pmatrix} 0.21 & -0.1 & -0.086 & -0.024 \\ -0.1 & 0.222 & -0.096 & -0.027 \\ -0.086 & -0.096 & 0.204 & -0.023 \\ -0.024 & -0.027 & -0.023 & 0.074 \end{pmatrix}$$

In a similar manner we can obtain

$$\widehat{\Sigma}_{22} = \begin{pmatrix} 0.215 & -0.184 & -0.031 \\ -0.184 & 0.242 & -0.059 \\ -0.031 & -0.059 & 0.09 \end{pmatrix}$$

Next, we use the observed counts in Table 3.4 to obtain the empirical joint PMF for \mathbf{X}_1 and \mathbf{X}_2 using Eq. (3.13), as plotted in Figure 3.2. From these probabilities we get

$$E[\mathbf{X}_1 \mathbf{X}_2^T] = \widehat{\mathbf{P}}_{12} = \frac{1}{150} \begin{pmatrix} 7 & 33 & 5 \\ 24 & 18 & 8 \\ 13 & 30 & 0 \\ 3 & 7 & 2 \end{pmatrix} = \begin{pmatrix} 0.047 & 0.22 & 0.033 \\ 0.16 & 0.12 & 0.053 \\ 0.087 & 0.2 & 0 \\ 0.02 & 0.047 & 0.013 \end{pmatrix}$$

Table 3.4. Observed Counts (n_{ij}): sepal length and sepal width

		X_2		
		Short (\mathbf{e}_{21})	Medium (\mathbf{e}_{22})	Long (\mathbf{e}_{23})
X_1	Very Short (\mathbf{e}_{11})	7	33	5
	Short (\mathbf{e}_{22})	24	18	8
	Long (\mathbf{e}_{13})	13	30	0
	Very Long (\mathbf{e}_{14})	3	7	2

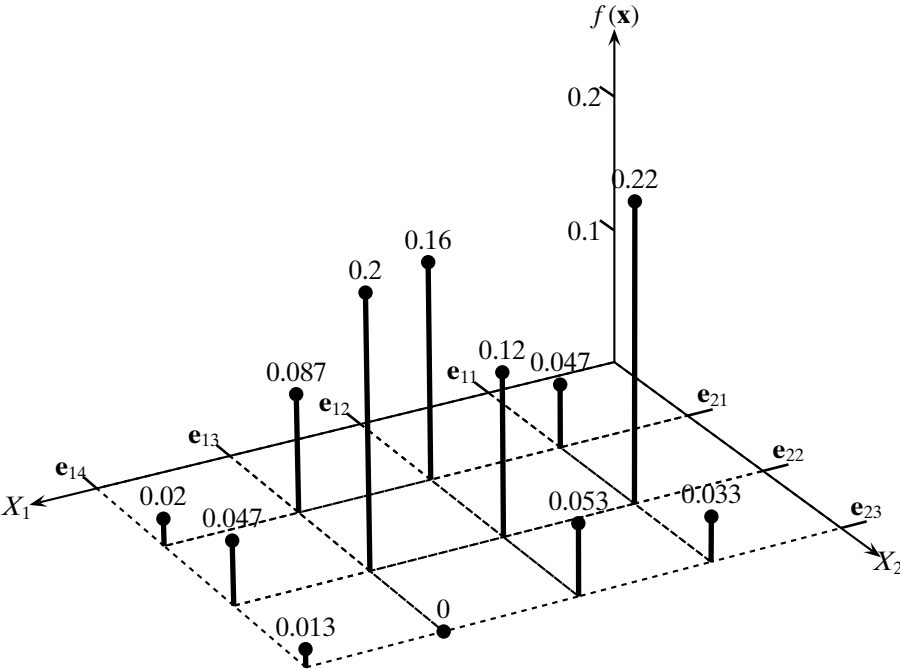


Figure 3.2. Empirical joint probability mass function: sepal length and sepal width.

Further, we have

$$\begin{aligned} E[\mathbf{X}_1]E[\mathbf{X}_2]^T &= \hat{\boldsymbol{\mu}}_1\hat{\boldsymbol{\mu}}_2^T = \hat{\mathbf{p}}_1\hat{\mathbf{p}}_2^T \\ &= \begin{pmatrix} 0.3 \\ 0.333 \\ 0.287 \\ 0.08 \end{pmatrix} \begin{pmatrix} 0.313 & 0.587 & 0.1 \end{pmatrix} \\ &= \begin{pmatrix} 0.094 & 0.176 & 0.03 \\ 0.104 & 0.196 & 0.033 \\ 0.09 & 0.168 & 0.029 \\ 0.025 & 0.047 & 0.008 \end{pmatrix} \end{aligned}$$

We can now compute the across-attribute sample covariance matrix $\widehat{\Sigma}_{12}$ for \mathbf{X}_1 and \mathbf{X}_2 using Eq. (3.11), as follows:

$$\begin{aligned}\widehat{\Sigma}_{12} &= \widehat{\mathbf{P}}_{12} - \widehat{\mathbf{p}}_1 \widehat{\mathbf{p}}_2^T \\ &= \begin{pmatrix} -0.047 & 0.044 & 0.003 \\ 0.056 & -0.076 & 0.02 \\ -0.003 & 0.032 & -0.029 \\ -0.005 & 0 & 0.005 \end{pmatrix}\end{aligned}$$

One can observe that each row and column in $\widehat{\Sigma}_{12}$ sums to zero. Putting it all together, from $\widehat{\Sigma}_{11}$, $\widehat{\Sigma}_{22}$ and $\widehat{\Sigma}_{12}$ we obtain the sample covariance matrix as follows

$$\begin{aligned}\widehat{\Sigma} &= \begin{pmatrix} \widehat{\Sigma}_{11} & \widehat{\Sigma}_{12} \\ \widehat{\Sigma}_{12}^T & \widehat{\Sigma}_{22} \end{pmatrix} \\ &= \left(\begin{array}{cccc|ccc} 0.21 & -0.1 & -0.086 & -0.024 & -0.047 & 0.044 & 0.003 \\ -0.1 & 0.222 & -0.096 & -0.027 & 0.056 & -0.076 & 0.02 \\ -0.086 & -0.096 & 0.204 & -0.023 & -0.003 & 0.032 & -0.029 \\ -0.024 & -0.027 & -0.023 & 0.074 & -0.005 & 0 & 0.005 \\ \hline -0.047 & 0.056 & -0.003 & -0.005 & 0.215 & -0.184 & -0.031 \\ 0.044 & -0.076 & 0.032 & 0 & -0.184 & 0.242 & -0.059 \\ 0.003 & 0.02 & -0.029 & 0.005 & -0.031 & -0.059 & 0.09 \end{array} \right)\end{aligned}$$

In $\widehat{\Sigma}$, each row and column also sums to zero.

3.2.1 Attribute Dependence: Contingency Analysis

Testing for the independence of the two categorical random variables X_1 and X_2 can be done via *contingency table analysis*. The main idea is to set up a hypothesis testing framework, where the null hypothesis H_0 is that \mathbf{X}_1 and \mathbf{X}_2 are independent, and the alternative hypothesis H_1 is that they are dependent. We then compute the value of the chi-square statistic χ^2 under the null hypothesis. Depending on the p -value, we either accept or reject the null hypothesis; in the latter case the attributes are considered to be dependent.

Contingency Table

A contingency table for \mathbf{X}_1 and \mathbf{X}_2 is the $m_1 \times m_2$ matrix of observed counts n_{ij} for all pairs of values $(\mathbf{e}_{1i}, \mathbf{e}_{2j})$ in the given sample of size n , defined as

$$\mathbf{N}_{12} = n \cdot \widehat{\mathbf{P}}_{12} = \begin{pmatrix} n_{11} & n_{12} & \cdots & n_{1m_2} \\ n_{21} & n_{22} & \cdots & n_{2m_2} \\ \vdots & \vdots & \ddots & \vdots \\ n_{m_11} & n_{m_12} & \cdots & n_{m_1m_2} \end{pmatrix}$$

Table 3.5. Contingency table: sepal length vs. sepal width

Sepal length (X_1)	Sepal width (X_2)				Row Counts
		Short	Medium	Long	
		a_{21}	a_{22}	a_{23}	
	Very Short (a_{11})	7	33	5	$n_1^1 = 45$
	Short (a_{12})	24	18	8	$n_2^1 = 50$
	Long (a_{13})	13	30	0	$n_3^1 = 43$
	Very Long (a_{14})	3	7	2	$n_4^1 = 12$
	Column Counts	$n_1^2 = 47$	$n_2^2 = 88$	$n_3^2 = 15$	$n = 150$

where $\hat{\mathbf{P}}_{12}$ is the empirical joint PMF for \mathbf{X}_1 and \mathbf{X}_2 , computed via Eq. (3.13). The contingency table is then augmented with row and column marginal counts, as follows:

$$\mathbf{N}_1 = n \cdot \hat{\mathbf{p}}_1 = \begin{pmatrix} n_1^1 \\ \vdots \\ n_{m_1}^1 \end{pmatrix} \quad \mathbf{N}_2 = n \cdot \hat{\mathbf{p}}_2 = \begin{pmatrix} n_1^2 \\ \vdots \\ n_{m_2}^2 \end{pmatrix}$$

Note that the marginal row and column entries and the sample size satisfy the following constraints:

$$n_i^1 = \sum_{j=1}^{m_2} n_{ij} \quad n_j^2 = \sum_{i=1}^{m_1} n_{ij} \quad n = \sum_{i=1}^{m_1} n_i^1 = \sum_{j=1}^{m_2} n_j^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} n_{ij}$$

It is worth noting that both \mathbf{N}_1 and \mathbf{N}_2 have a multinomial distribution with parameters $\mathbf{p}_1 = (p_1^1, \dots, p_{m_1}^1)$ and $\mathbf{p}_2 = (p_1^2, \dots, p_{m_2}^2)$, respectively. Further, \mathbf{N}_{12} also has a multinomial distribution with parameters $\mathbf{P}_{12} = \{p_{ij}\}$, for $1 \leq i \leq m_1$ and $1 \leq j \leq m_2$.

Example 3.9 (Contingency Table). Table 3.4 shows the observed counts for the discretized sepal length (X_1) and sepal width (X_2) attributes. Augmenting the table with the row and column marginal counts and the sample size yields the final contingency table shown in Table 3.5.

χ^2 Statistic and Hypothesis Testing

Under the null hypothesis \mathbf{X}_1 and \mathbf{X}_2 are assumed to be independent, which means that their joint probability mass function is given as

$$\hat{p}_{ij} = \hat{p}_i^1 \cdot \hat{p}_j^2$$

Under this independence assumption, the expected frequency for each pair of values is given as

$$e_{ij} = n \cdot \hat{p}_{ij} = n \cdot \hat{p}_i^1 \cdot \hat{p}_j^2 = n \cdot \frac{n_i^1}{n} \cdot \frac{n_j^2}{n} = \frac{n_i^1 n_j^2}{n} \quad (3.14)$$

However, from the sample we already have the observed frequency of each pair of values, n_{ij} . We would like to determine whether there is a significant difference in the observed and expected frequencies for each pair of values. If there is no

significant difference, then the independence assumption is valid and we accept the null hypothesis that the attributes are independent. On the other hand, if there is a significant difference, then the null hypothesis should be rejected and we conclude that the attributes are dependent.

The χ^2 statistic quantifies the difference between observed and expected counts for each pair of values; it is defined as follows:

$$\chi^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \frac{(n_{ij} - e_{ij})^2}{e_{ij}} \quad (3.15)$$

At this point, we need to determine the probability of obtaining the computed χ^2 value. In general, this can be rather difficult if we do not know the sampling distribution of a given statistic. Fortunately, for the χ^2 statistic it is known that its sampling distribution follows the *chi-squared* density function with q degrees of freedom:

$$f(x|q) = \frac{1}{2^{q/2} \Gamma(q/2)} x^{\frac{q}{2}-1} e^{-\frac{x}{2}} \quad (3.16)$$

where the gamma function Γ is defined as

$$\Gamma(k > 0) = \int_0^{\infty} x^{k-1} e^{-x} dx \quad (3.17)$$

The degrees of freedom, q , represent the number of independent parameters. In the contingency table there are $m_1 \times m_2$ observed counts n_{ij} . However, note that each row i and each column j must sum to n_i^1 and n_j^2 , respectively. Further, the sum of the row and column marginals must also add to n ; thus we have to remove $(m_1 + m_2)$ parameters from the number of independent parameters. However, doing this removes one of the parameters, say $n_{m_1 m_2}$, twice, so we have to add back one to the count. The total degrees of freedom is therefore

$$\begin{aligned} q &= |dom(X_1)| \times |dom(X_2)| - (|dom(X_1)| + |dom(X_2)|) + 1 \\ &= m_1 m_2 - m_1 - m_2 + 1 \\ &= (m_1 - 1)(m_2 - 1) \end{aligned}$$

p-value

The *p-value* of a statistic θ is defined as the probability of obtaining a value at least as extreme as the observed value, say z , under the null hypothesis, defined as

$$p\text{-value}(z) = P(\theta \geq z) = 1 - F(\theta)$$

where $F(\theta)$ is the cumulative probability distribution for the statistic.

The *p-value* gives a measure of how surprising is the observed value of the statistic. If the observed value lies in a low-probability region, then the value is more surprising. In general, the lower the *p-value*, the more surprising the observed value, and the

Table 3.6. Expected counts

		X_2		
		Short (a_{21})	Medium (a_{22})	Short (a_{23})
X_1	Very Short (a_{11})	14.1	26.4	4.5
	Short (a_{12})	15.67	29.33	5.0
	Long (a_{13})	13.47	25.23	4.3
	Very Long (a_{14})	3.76	7.04	1.2

more the grounds for rejecting the null hypothesis. The null hypothesis is rejected if the p -value is below some *significance level*, α . For example, if $\alpha = 0.01$, then we reject the null hypothesis if $p\text{-value}(z) \leq \alpha$. The significance level α corresponds to the probability of rejecting the null hypothesis when it is true. For a given significance level α , the value of the test statistic, say z , with a p -value of $p\text{-value}(z) = \alpha$, is called a *critical value*. An alternative test for rejection of the null hypothesis is to check if $\chi^2 > z$, as in that case the p -value of the observed χ^2 value is bounded by α , that is, $p\text{-value}(\chi^2) \leq p\text{-value}(z) = \alpha$. The value $1 - \alpha$ is also called the *confidence level*.

Example 3.10. Consider the contingency table for sepal length and sepal width in Table 3.5. We compute the expected counts using Eq. (3.14); these counts are shown in Table 3.6. For example, we have

$$e_{11} = \frac{n_1^1 n_1^2}{n} = \frac{45 \cdot 47}{150} = \frac{2115}{150} = 14.1$$

Next we use Eq. (3.15) to compute the value of the χ^2 statistic, which is given as $\chi^2 = 21.8$.

Further, the number of degrees of freedom is given as

$$q = (m_1 - 1) \cdot (m_2 - 1) = 3 \cdot 2 = 6$$

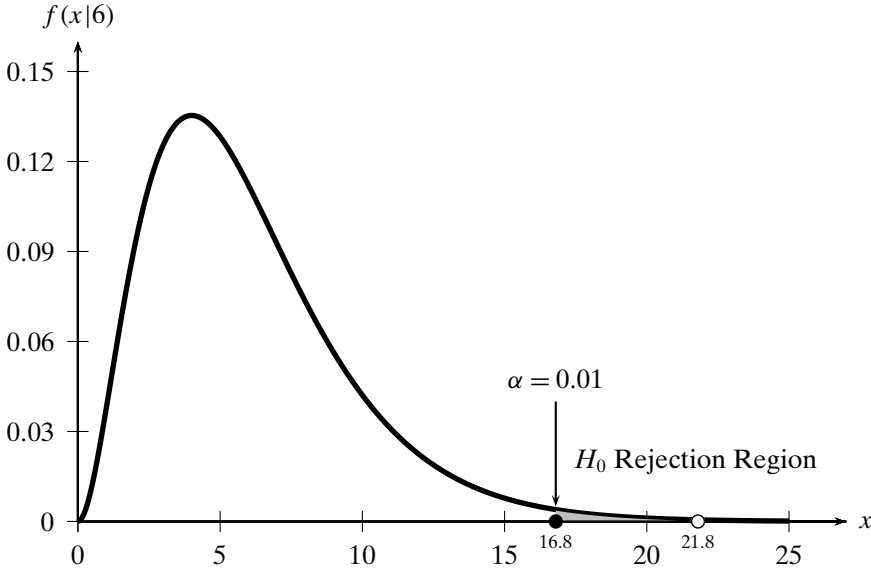
The plot of the chi-squared density function with 6 degrees of freedom is shown in Figure 3.3. From the cumulative chi-squared distribution, we obtain

$$p\text{-value}(21.8) = 1 - F(21.8|6) = 1 - 0.9987 = 0.0013$$

At a significance level of $\alpha = 0.01$, we would certainly be justified in rejecting the null hypothesis because the large value of the χ^2 statistic is indeed surprising. Further, at the 0.01 significance level, the critical value of the statistic is

$$z = F^{-1}(1 - 0.01|6) = F^{-1}(0.99|6) = 16.81$$

This critical value is also shown in Figure 3.3, and we can clearly see that the observed value of 21.8 is in the rejection region, as $21.8 > z = 16.81$. In effect, we reject the null hypothesis that sepal length and sepal width are independent, and accept the alternative hypothesis that they are dependent.

Figure 3.3. Chi-squared distribution ($q = 6$).

3.3 MULTIVARIATE ANALYSIS

Assume that the dataset comprises d categorical attributes X_j ($1 \leq j \leq d$) with $\text{dom}(X_j) = \{a_{j1}, a_{j2}, \dots, a_{jm_j}\}$. We are given n categorical points of the form $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ with $x_{ij} \in \text{dom}(X_j)$. The dataset is thus an $n \times d$ symbolic matrix

$$\mathbf{D} = \begin{pmatrix} X_1 & X_2 & \cdots & X_d \\ x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{pmatrix}$$

Each attribute X_i is modeled as an m_i -dimensional multivariate Bernoulli variable \mathbf{X}_i , and their joint distribution is modeled as a $d' = \sum_{j=1}^d m_j$ dimensional vector random variable

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_d \end{pmatrix}$$

Each categorical data point $\mathbf{v} = (v_1, v_2, \dots, v_d)^T$ is therefore represented as a d' -dimensional binary vector

$$\mathbf{X}(\mathbf{v}) = \begin{pmatrix} \mathbf{X}_1(v_1) \\ \vdots \\ \mathbf{X}_d(v_d) \end{pmatrix} = \begin{pmatrix} \mathbf{e}_{1k_1} \\ \vdots \\ \mathbf{e}_{dk_d} \end{pmatrix}$$

provided $v_i = a_{ik_i}$, the k_i th symbol of X_i . Here \mathbf{e}_{ik_i} is the k_i th standard basis vector in \mathbb{R}^{m_i} .

Mean

Generalizing from the bivariate case, the mean and sample mean for \mathbf{X} are given as

$$\boldsymbol{\mu} = E[\mathbf{X}] = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \vdots \\ \boldsymbol{\mu}_d \end{pmatrix} = \begin{pmatrix} \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_d \end{pmatrix} \quad \hat{\boldsymbol{\mu}} = \begin{pmatrix} \hat{\boldsymbol{\mu}}_1 \\ \vdots \\ \hat{\boldsymbol{\mu}}_d \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{p}}_1 \\ \vdots \\ \hat{\mathbf{p}}_d \end{pmatrix}$$

where $\mathbf{p}_i = (p_1^i, \dots, p_{m_i}^i)^T$ is the PMF for \mathbf{X}_i , and $\hat{\mathbf{p}}_i = (\hat{p}_1^i, \dots, \hat{p}_{m_i}^i)^T$ is the empirical PMF for \mathbf{X}_i .

Covariance Matrix

The covariance matrix for \mathbf{X} , and its estimate from the sample, are given as the $d' \times d'$ matrices:

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} & \cdots & \boldsymbol{\Sigma}_{1d} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} & \cdots & \boldsymbol{\Sigma}_{2d} \\ \cdots & \cdots & \ddots & \cdots \\ \boldsymbol{\Sigma}_{1d}^T & \boldsymbol{\Sigma}_{2d}^T & \cdots & \boldsymbol{\Sigma}_{dd} \end{pmatrix} \quad \hat{\boldsymbol{\Sigma}} = \begin{pmatrix} \hat{\boldsymbol{\Sigma}}_{11} & \hat{\boldsymbol{\Sigma}}_{12} & \cdots & \hat{\boldsymbol{\Sigma}}_{1d} \\ \hat{\boldsymbol{\Sigma}}_{12}^T & \hat{\boldsymbol{\Sigma}}_{22} & \cdots & \hat{\boldsymbol{\Sigma}}_{2d} \\ \cdots & \cdots & \ddots & \cdots \\ \hat{\boldsymbol{\Sigma}}_{1d}^T & \hat{\boldsymbol{\Sigma}}_{2d}^T & \cdots & \hat{\boldsymbol{\Sigma}}_{dd} \end{pmatrix}$$

where $d' = \sum_{i=1}^d m_i$, and $\boldsymbol{\Sigma}_{ij}$ (and $\hat{\boldsymbol{\Sigma}}_{ij}$) is the $m_i \times m_j$ covariance matrix (and its estimate) for attributes \mathbf{X}_i and \mathbf{X}_j :

$$\boldsymbol{\Sigma}_{ij} = \mathbf{P}_{ij} - \mathbf{p}_i \mathbf{p}_j^T \quad \hat{\boldsymbol{\Sigma}}_{ij} = \hat{\mathbf{P}}_{ij} - \hat{\mathbf{p}}_i \hat{\mathbf{p}}_j^T \quad (3.18)$$

Here \mathbf{P}_{ij} is the joint PMF and $\hat{\mathbf{P}}_{ij}$ is the empirical joint PMF for \mathbf{X}_i and \mathbf{X}_j , which can be computed using Eq. (3.13).

Example 3.11 (Multivariate Analysis). Let us consider the 3-dimensional subset of the Iris dataset, with the discretized attributes sepal length (X_1) and sepal width (X_2), and the categorical attribute class (X_3). The domains for X_1 and X_2 are given in Table 3.1 and Table 3.3, respectively, and $\text{dom}(X_3) = \{\text{iris-versicolor}, \text{iris-setosa}, \text{iris-virginica}\}$. Each value of X_3 occurs 50 times.

The categorical point $\mathbf{x} = (\text{Short}, \text{Medium}, \text{iris-versicolor})$ is modeled as the vector

$$\mathbf{X}(\mathbf{x}) = \begin{pmatrix} \mathbf{e}_{12} \\ \mathbf{e}_{22} \\ \mathbf{e}_{31} \end{pmatrix} = (0, 1, 0, 0 \mid 0, 1, 0 \mid 1, 0, 0)^T \in \mathbb{R}^{10}$$

From Example 3.8 and the fact that each value in $\text{dom}(X_3)$ occurs 50 times in a sample of $n = 150$, the sample mean is given as

$$\hat{\boldsymbol{\mu}} = \begin{pmatrix} \hat{\boldsymbol{\mu}}_1 \\ \hat{\boldsymbol{\mu}}_2 \\ \hat{\boldsymbol{\mu}}_3 \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{p}}_1 \\ \hat{\mathbf{p}}_2 \\ \hat{\mathbf{p}}_3 \end{pmatrix} = (0.3, 0.333, 0.287, 0.08 \mid 0.313, 0.587, 0.1 \mid 0.33, 0.33, 0.33)^T$$

Using $\hat{\mathbf{p}}_3 = (0.33, 0.33, 0.33)^T$ we can compute the sample covariance matrix for X_3 using Eq. (3.9):

$$\hat{\Sigma}_{33} = \begin{pmatrix} 0.222 & -0.111 & -0.111 \\ -0.111 & 0.222 & -0.111 \\ -0.111 & -0.111 & 0.222 \end{pmatrix}$$

Using Eq. (3.18) we obtain

$$\hat{\Sigma}_{13} = \begin{pmatrix} -0.067 & 0.16 & -0.093 \\ 0.082 & -0.038 & -0.044 \\ 0.011 & -0.096 & 0.084 \\ -0.027 & -0.027 & 0.053 \end{pmatrix}$$

$$\hat{\Sigma}_{23} = \begin{pmatrix} 0.076 & -0.098 & 0.022 \\ -0.042 & 0.044 & -0.002 \\ -0.033 & 0.053 & -0.02 \end{pmatrix}$$

Combined with $\hat{\Sigma}_{11}$, $\hat{\Sigma}_{22}$ and $\hat{\Sigma}_{12}$ from Example 3.8, the final sample covariance matrix is the 10×10 symmetric matrix given as

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{11} & \hat{\Sigma}_{12} & \hat{\Sigma}_{13} \\ \hat{\Sigma}_{12}^T & \hat{\Sigma}_{22} & \hat{\Sigma}_{23} \\ \hat{\Sigma}_{13}^T & \hat{\Sigma}_{23}^T & \hat{\Sigma}_{33} \end{pmatrix}$$

3.3.1 Multiway Contingency Analysis

For multiway dependence analysis, we have to first determine the empirical joint probability mass function for \mathbf{X} :

$$\hat{f}(\mathbf{e}_{1i_1}, \mathbf{e}_{2i_2}, \dots, \mathbf{e}_{di_d}) = \frac{1}{n} \sum_{k=1}^n I_{i_1 i_2 \dots i_d}(\mathbf{x}_k) = \frac{n_{i_1 i_2 \dots i_d}}{n} = \hat{p}_{i_1 i_2 \dots i_d}$$

where $I_{i_1 i_2 \dots i_d}$ is the indicator variable

$$I_{i_1 i_2 \dots i_d}(\mathbf{x}_k) = \begin{cases} 1 & \text{if } x_{k1} = \mathbf{e}_{1i_1}, x_{k2} = \mathbf{e}_{2i_2}, \dots, x_{kd} = \mathbf{e}_{di_d} \\ 0 & \text{otherwise} \end{cases}$$

The sum of $I_{i_1 i_2 \dots i_d}$ over all the n points in the sample yields the number of occurrences, $n_{i_1 i_2 \dots i_d}$, of the symbolic vector $(a_{1i_1}, a_{2i_2}, \dots, a_{di_d})$. Dividing the occurrences by the sample size results in the probability of observing those symbols. Using the notation $\mathbf{i} = (i_1, i_2, \dots, i_d)$ to denote the index tuple, we can write the joint empirical PMF as the d -dimensional matrix $\hat{\mathbf{P}}$ of size $m_1 \times m_2 \times \dots \times m_d = \prod_{i=1}^d m_i$, given as

$$\hat{\mathbf{P}}(\mathbf{i}) = \{\hat{p}_{\mathbf{i}}\} \text{ for all index tuples } \mathbf{i}, \text{ with } 1 \leq i_1 \leq m_1, \dots, 1 \leq i_d \leq m_d$$

where $\hat{p}_{\mathbf{i}} = \hat{p}_{i_1 i_2 \dots i_d}$. The d -dimensional contingency table is then given as

$$\mathbf{N} = n \times \hat{\mathbf{P}} = \{n_{\mathbf{i}}\} \text{ for all index tuples } \mathbf{i}, \text{ with } 1 \leq i_1 \leq m_1, \dots, 1 \leq i_d \leq m_d$$

where $n_i = n_{i_1 i_2 \dots i_d}$. The contingency table is augmented with the marginal count vectors \mathbf{N}_i for all d attributes \mathbf{X}_i :

$$\mathbf{N}_i = n \hat{\mathbf{p}}_i = \begin{pmatrix} n_1^i \\ \vdots \\ n_{m_i}^i \end{pmatrix}$$

where $\hat{\mathbf{p}}_i$ is the empirical PMF for \mathbf{X}_i .

χ^2 -Test

We can test for a d -way dependence between the d categorical attributes using the null hypothesis H_0 that they are d -way independent. The alternative hypothesis H_1 is that they are not d -way independent, that is, they are dependent in some way. Note that d -dimensional contingency analysis indicates whether all d attributes taken together are independent or not. In general we may have to conduct k -way contingency analysis to test if any subset of $k \leq d$ attributes are independent or not.

Under the null hypothesis, the expected number of occurrences of the symbol tuple $(a_{1i_1}, a_{2i_2}, \dots, a_{di_d})$ is given as

$$e_i = n \cdot \hat{p}_i = n \cdot \prod_{j=1}^d \hat{p}_{i_j} = \frac{n_{i_1}^1 n_{i_2}^2 \dots n_{i_d}^d}{n^{d-1}} \quad (3.19)$$

The chi-squared statistic measures the difference between the observed counts n_i and the expected counts e_i :

$$\chi^2 = \sum_i \frac{(n_i - e_i)^2}{e_i} = \sum_{i_1=1}^{m_1} \sum_{i_2=1}^{m_2} \dots \sum_{i_d=1}^{m_d} \frac{(n_{i_1, i_2, \dots, i_d} - e_{i_1, i_2, \dots, i_d})^2}{e_{i_1, i_2, \dots, i_d}} \quad (3.20)$$

The χ^2 statistic follows a chi-squared density function with q degrees of freedom. For the d -way contingency table we can compute q by noting that there are ostensibly $\prod_{i=1}^d |\text{dom}(X_i)|$ independent parameters (the counts). However, we have to remove $\sum_{i=1}^d |\text{dom}(X_i)|$ degrees of freedom because the marginal count vector along each dimension \mathbf{X}_i must equal \mathbf{N}_i . However, doing so removes one of the parameters d times, so we need to add back $d - 1$ to the free parameters count. The total number of degrees of freedom is given as

$$\begin{aligned} q &= \prod_{i=1}^d |\text{dom}(X_i)| - \sum_{i=1}^d |\text{dom}(X_i)| + (d - 1) \\ &= \left(\prod_{i=1}^d m_i \right) - \left(\sum_{i=1}^d m_i \right) + d - 1 \end{aligned} \quad (3.21)$$

To reject the null hypothesis, we have to check whether the p -value of the observed χ^2 value is smaller than the desired significance level α (say $\alpha = 0.01$) using the chi-squared density with q degrees of freedom [Eq. (3.16)].

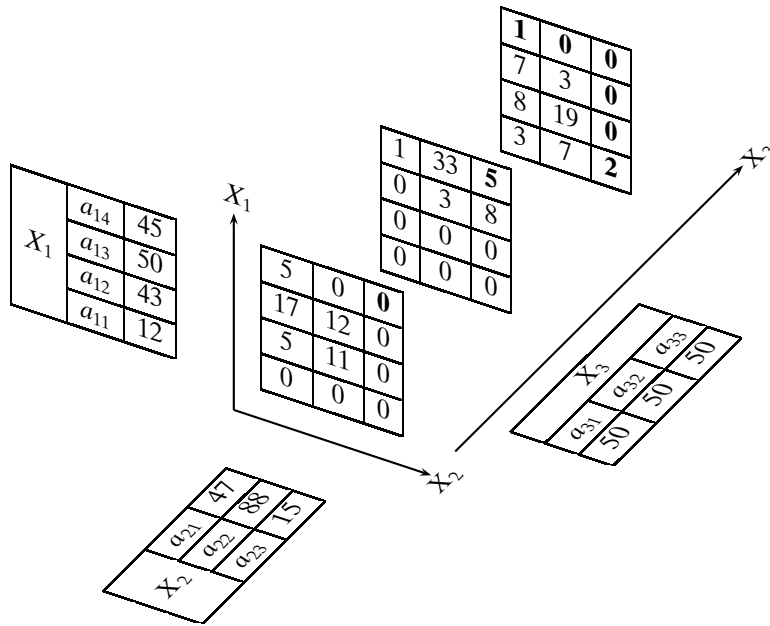


Figure 3.4. 3-Way contingency table (with marginal counts along each dimension).

Table 3.7. 3-Way expected counts

		$X_3(a_{31}/a_{32}/a_{33})$		
		X_2		
		a_{21}	a_{22}	a_{23}
X_1	a_{11}	1.25	2.35	0.40
	a_{12}	4.49	8.41	1.43
	a_{13}	5.22	9.78	1.67
	a_{14}	4.70	8.80	1.50

Example 3.12. Consider the 3-way contingency table in Figure 3.4. It shows the observed counts for each tuple of symbols (a_{1i}, a_{2j}, a_{3k}) for the three attributes sepal length (X_1), sepal width (X_2), and class (X_3). From the marginal counts for X_1 and X_2 in Table 3.5, and the fact that all three values of X_3 occur 50 times, we can compute the expected counts [Eq. (3.19)] for each cell. For instance,

$$e_{(4,1,1)} = \frac{n_4^1 \cdot n_1^2 \cdot n_1^3}{150^2} = \frac{45 \cdot 47 \cdot 50}{150 \cdot 150} = 4.7$$

The expected counts are the same for all three values of X_3 and are given in Table 3.7.

The value of the χ^2 statistic [Eq. (3.20)] is given as

$$\chi^2 = 231.06$$

Using Eq. (3.21), the number of degrees of freedom is given as

$$q = 4 \cdot 3 \cdot 3 - (4 + 3 + 3) + 2 = 36 - 10 + 2 = 28$$

In Figure 3.4 the counts in bold are the dependent parameters. All other counts are independent. In fact, any eight distinct cells could have been chosen as the dependent parameters.

For a significance level of $\alpha = 0.01$, the critical value of the chi-square distribution is $z = 48.28$. The observed value of $\chi^2 = 231.06$ is much greater than z , and it is thus extremely unlikely to happen under the null hypothesis. We conclude that the three attributes are not 3-way independent, but rather there is some dependence between them. However, this example also highlights one of the pitfalls of multiway contingency analysis. We can observe in Figure 3.4 that many of the observed counts are zero. This is due to the fact that the sample size is small, and we cannot reliably estimate all the multiway counts. Consequently, the dependence test may not be reliable as well.

3.4 DISTANCE AND ANGLE

With the modeling of categorical attributes as multivariate Bernoulli variables, it is possible to compute the distance or the angle between any two points \mathbf{x}_i and \mathbf{x}_j :

$$\mathbf{x}_i = \begin{pmatrix} \mathbf{e}_{1i_1} \\ \vdots \\ \mathbf{e}_{di_d} \end{pmatrix} \quad \mathbf{x}_j = \begin{pmatrix} \mathbf{e}_{1j_1} \\ \vdots \\ \mathbf{e}_{dj_d} \end{pmatrix}$$

The different measures of distance and similarity rely on the number of matching and mismatching values (or symbols) across the d attributes \mathbf{X}_k . For instance, we can compute the number of matching values s via the dot product:

$$s = \mathbf{x}_i^T \mathbf{x}_j = \sum_{k=1}^d (\mathbf{e}_{ki_k})^T \mathbf{e}_{kj_k}$$

On the other hand, the number of mismatches is simply $d - s$. Also useful is the norm of each point:

$$\|\mathbf{x}_i\|^2 = \mathbf{x}_i^T \mathbf{x}_i = d$$

Euclidean Distance

The Euclidean distance between \mathbf{x}_i and \mathbf{x}_j is given as

$$\delta(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\| = \sqrt{\mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{x}_j + \mathbf{x}_j^T \mathbf{x}_j} = \sqrt{2(d - s)}$$

Thus, the maximum Euclidean distance between any two points is $\sqrt{2d}$, which happens when there are no common symbols between them, that is, when $s = 0$.

Hamming Distance

The *Hamming distance* between \mathbf{x}_i and \mathbf{x}_j is defined as the number of mismatched values:

$$\delta_H(\mathbf{x}_i, \mathbf{x}_j) = d - s = \frac{1}{2} \delta(\mathbf{x}_i, \mathbf{x}_j)^2$$

Hamming distance is thus equivalent to half the squared Euclidean distance.

Cosine Similarity

The cosine of the angle between \mathbf{x}_i and \mathbf{x}_j is given as

$$\cos \theta = \frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \cdot \|\mathbf{x}_j\|} = \frac{s}{d}$$

Jaccard Coefficient

The *Jaccard Coefficient* is a commonly used similarity measure between two categorical points. It is defined as the ratio of the number of matching values to the number of distinct values that appear in both \mathbf{x}_i and \mathbf{x}_j , across the d attributes:

$$J(\mathbf{x}_i, \mathbf{x}_j) = \frac{s}{2(d-s) + s} = \frac{s}{2d-s}$$

where we utilize the observation that when the two points do not match for dimension k , they contribute 2 to the distinct symbol count; otherwise, if they match, the number of distinct symbols increases by 1. Over the $d-s$ mismatches and s matches, the number of distinct symbols is $2(d-s) + s$.

Example 3.13. Consider the 3-dimensional categorical data from Example 3.11. The symbolic point (Short, Medium, iris-versicolor) is modeled as the vector

$$\mathbf{x}_1 = \begin{pmatrix} \mathbf{e}_{12} \\ \mathbf{e}_{22} \\ \mathbf{e}_{31} \end{pmatrix} = (0, 1, 0, 0 \mid 0, 1, 0 \mid 1, 0, 0)^T \in \mathbb{R}^{10}$$

and the symbolic point (VeryShort, Medium, iris-setosa) is modeled as

$$\mathbf{x}_2 = \begin{pmatrix} \mathbf{e}_{11} \\ \mathbf{e}_{22} \\ \mathbf{e}_{32} \end{pmatrix} = (1, 0, 0, 0 \mid 0, 1, 0 \mid 0, 1, 0)^T \in \mathbb{R}^{10}$$

The number of matching symbols is given as

$$\begin{aligned} s &= \mathbf{x}_1^T \mathbf{x}_2 = (\mathbf{e}_{12})^T \mathbf{e}_{11} + (\mathbf{e}_{22})^T \mathbf{e}_{22} + (\mathbf{e}_{31})^T \mathbf{e}_{32} \\ &= (0 \ 1 \ 0 \ 0) \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} + (0 \ 1 \ 0) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + (1 \ 0 \ 0) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ &= 0 + 1 + 0 = 1 \end{aligned}$$

The Euclidean and Hamming distances are given as

$$\delta(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{2(d-s)} = \sqrt{2 \cdot 2} = \sqrt{4} = 2$$

$$\delta_H(\mathbf{x}_1, \mathbf{x}_2) = d - s = 3 - 1 = 2$$

The cosine and Jaccard similarity are given as

$$\cos \theta = \frac{s}{d} = \frac{1}{3} = 0.333$$

$$J(\mathbf{x}_1, \mathbf{x}_2) = \frac{s}{2d-s} = \frac{1}{5} = 0.2$$

3.5 DISCRETIZATION

Discretization, also called *binning*, converts numeric attributes into categorical ones. It is usually applied for data mining methods that cannot handle numeric attributes. It can also help in reducing the number of values for an attribute, especially if there is noise in the numeric measurements; discretization allows one to ignore small and irrelevant differences in the values.

Formally, given a numeric attribute X , and a random sample $\{x_i\}_{i=1}^n$ of size n drawn from X , the discretization task is to divide the value range of X into k consecutive intervals, also called *bins*, by finding $k-1$ boundary values v_1, v_2, \dots, v_{k-1} that yield the k intervals:

$$[x_{\min}, v_1], (v_1, v_2], \dots, (v_{k-1}, x_{\max}]$$

where the extremes of the range of X are given as

$$x_{\min} = \min_i \{x_i\} \qquad x_{\max} = \max_i \{x_i\}$$

The resulting k intervals or bins, which span the entire range of X , are usually mapped to symbolic values that comprise the domain for the new categorical attribute X .

Equal-Width Intervals

The simplest binning approach is to partition the range of X into k *equal-width* intervals. The interval width is simply the range of X divided by k :

$$w = \frac{x_{\max} - x_{\min}}{k}$$

Thus, the i th interval boundary is given as

$$v_i = x_{\min} + iw, \text{ for } i = 1, \dots, k-1$$

Equal-Frequency Intervals

In *equal-frequency* binning we divide the range of X into intervals that contain (approximately) equal number of points; equal frequency may not be possible due to repeated values. The intervals can be computed from the empirical quantile or

inverse cumulative distribution function $\hat{F}^{-1}(q)$ for X [Eq. (2.2)]. Recall that $\hat{F}^{-1}(q) = \min\{x \mid P(X \leq x) \geq q\}$, for $q \in [0, 1]$. In particular, we require that each interval contain $1/k$ of the probability mass; therefore, the interval boundaries are given as follows:

$$v_i = \hat{F}^{-1}(i/k) \text{ for } i = 1, \dots, k-1$$

Example 3.14. Consider the `sepal length` attribute in the Iris dataset. Its minimum and maximum values are

$$x_{\min} = 4.3 \qquad x_{\max} = 7.9$$

We discretize it into $k = 4$ bins using equal-width binning. The width of an interval is given as

$$w = \frac{7.9 - 4.3}{4} = \frac{3.6}{4} = 0.9$$

and therefore the interval boundaries are

$$v_1 = 4.3 + 0.9 = 5.2 \qquad v_2 = 4.3 + 2 \cdot 0.9 = 6.1 \qquad v_3 = 4.3 + 3 \cdot 0.9 = 7.0$$

The four resulting bins for `sepal length` are shown in Table 3.1, which also shows the number of points n_i in each bin, which are not balanced among the bins.

For equal-frequency discretization, consider the empirical inverse cumulative distribution function (CDF) for `sepal length` shown in Figure 3.5. With $k = 4$ bins, the bin boundaries are the quartile values (which are shown as dashed lines):

$$v_1 = \hat{F}^{-1}(0.25) = 5.1 \qquad v_2 = \hat{F}^{-1}(0.50) = 5.8 \qquad v_3 = \hat{F}^{-1}(0.75) = 6.4$$

The resulting intervals are shown in Table 3.8. We can see that although the interval widths vary, they contain a more balanced number of points. We do not get identical

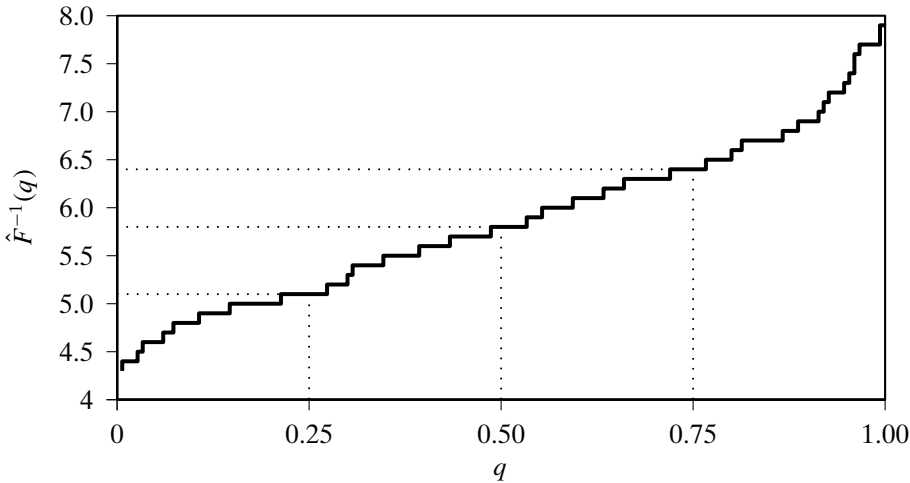


Figure 3.5. Empirical inverse CDF: `sepal length`.

Table 3.8. Equal-frequency discretization: sepal length

Bin	Width	Count
[4.3, 5.1]	0.8	$n_1 = 41$
(5.1, 5.8]	0.7	$n_2 = 39$
(5.8, 6.4]	0.6	$n_3 = 35$
(6.4, 7.9]	1.5	$n_4 = 35$

counts for all the bins because many values are repeated; for instance, there are nine points with value 5.1 and there are seven points with value 5.8.

3.6 FURTHER READING

For a comprehensive introduction to categorical data analysis see Agresti (2012). Some aspects also appear in Wasserman (2004). For an entropy-based supervised discretization method that takes the class attribute into account see Fayyad and Irani (1993).

Agresti, A. (2012). *Categorical Data Analysis*, 3rd ed. Hoboken, NJ: John Wiley & Sons.

Fayyad, U. M. and Irani, K. B. (1993). Multi-interval Discretization of Continuous-valued Attributes for Classification Learning. *In Proceedings of the 13th International Joint Conference on Artificial Intelligence*. Morgan-Kaufmann, pp. 1022–1027.

Wasserman, L. (2004). *All of Statistics: A Concise Course in Statistical Inference*. New York: Springer Science + Business Media.

3.7 EXERCISES

- Q1.** Show that for categorical points, the cosine similarity between any two vectors in lies in the range $\cos \theta \in [0, 1]$, and consequently $\theta \in [0^\circ, 90^\circ]$.
- Q2.** Prove that $E[(\mathbf{X}_1 - \mu_1)(\mathbf{X}_2 - \mu_2)^T] = E[\mathbf{X}_1 \mathbf{X}_2^T] - E[\mathbf{X}_1]E[\mathbf{X}_2]^T$.

Table 3.9. Contingency table for Q3

	$Z = f$		$Z = g$	
	$Y = d$	$Y = e$	$Y = d$	$Y = e$
$X = a$	5	10	10	5
$X = b$	15	5	5	20
$X = c$	20	10	25	10

Table 3.10. χ^2 Critical values for different p -values for different degrees of freedom (q): For example, for $q = 5$ degrees of freedom, the critical value of $\chi^2 = 11.070$ has p -value = 0.05.

q	0.995	0.99	0.975	0.95	0.90	0.10	0.05	0.025	0.01	0.005
1	—	—	0.001	0.004	0.016	2.706	3.841	5.024	6.635	7.879
2	0.010	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210	10.597
3	0.072	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345	12.838
4	0.207	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277	14.860
5	0.412	0.554	0.831	1.145	1.610	9.236	11.070	12.833	15.086	16.750
6	0.676	0.872	1.237	1.635	2.204	10.645	12.592	14.449	16.812	18.548

- Q3.** Consider the 3-way contingency table for attributes X, Y, Z shown in Table 3.9. Compute the χ^2 metric for the correlation between Y and Z . Are they dependent or independent at the 95% confidence level? See Table 3.10 for χ^2 values.
- Q4.** Consider the “mixed” data given in Table 3.11. Here X_1 is a numeric attribute and X_2 is a categorical one. Assume that the domain of X_2 is given as $dom(X_2) = \{a, b\}$. Answer the following questions.
- (a) What is the mean vector for this dataset?
- (b) What is the covariance matrix?
- Q5.** In Table 3.11, assuming that X_1 is discretized into three bins, as follows:

$$c_1 = (-2, -0.5]$$

$$c_2 = (-0.5, 0.5]$$

$$c_3 = (0.5, 2]$$

Answer the following questions:

- (a) Construct the contingency table between the discretized X_1 and X_2 attributes. Include the marginal counts.
- (b) Compute the χ^2 statistic between them.
- (c) Determine whether they are dependent or not at the 5% significance level. Use the χ^2 critical values from Table 3.10.

Table 3.11. Dataset for Q4 and Q5

X_1	X_2
0.3	a
-0.3	b
0.44	a
-0.60	a
0.40	a
1.20	b
-0.12	a
-1.60	b
1.60	b
-1.32	a