

Let (x_1, x_2, \dots, x_n) be a sequence of vectors:

$$x_i = \begin{bmatrix} x_i^1 \\ \vdots \\ x_i^m \end{bmatrix}, \quad i = 1, \dots, n.$$

In statistics one often has to compute the *sample mean* vector

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

and the *sample covariance* (or *variance-covariance*) matrix

$$\bar{V} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

where x^T is the transpose of x .

1. What *canonical* form of information would you suggest to represent the sequence (x_1, x_2, \dots, x_n) in order to compute the sample mean vector and the sample covariance matrix?

Verify that all of the “desirable” properties of canonical information are satisfied.

Solution

Consider the scalar-vector-matrix triple (n, s, T) where

$$n = \sum_{i=1}^n 1, \quad s = \sum_{i=1}^n x_i, \quad \text{and} \quad T = \sum_{i=1}^n x_i x_i^T. \quad \checkmark$$

(a) **Uniqueness:** Note that (n, s, T) is uniquely determined as each is a function of well-defined vector operations.

(i) **Elementary** canonical information: A single observation has the representation

$$x \mapsto (1, x, x x^T).$$

(ii) **Empty** canonical information: Empty information has the representation

$$\{\} \mapsto (0, 0, 0),$$

where each 0 is the respective scalar, vector, and matrix additive identity. ✓

(b) **Composition** operation: Let $(n, s, T) \oplus (n', s', T') := (n + n', s + s', T + T')$. Commutativity and associativity are inherited directly from the respective additions for scalars, vectors, and matrices. Moreover, the neutral element consists of the triple of additive identities $(0, 0, 0)$. ✓

(c) **Update** observation: Given x_{n+1} , we can update via the following schematic:

$$(S, T, n) \xrightarrow{x_{n+1}} \oplus \longrightarrow (n+1, s + x_{n+1}, T + x_{n+1} x_{n+1}^T)$$

Note that this is exactly the same operation one obtains by composing with the elementary element. ✓

- (d) **Completeness:** Recovering \bar{x} is immediately given by $\bar{x} = \frac{s}{n}$. Some matrix algebra reveals

$$\begin{aligned}\bar{V} &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{n-1} \left\{ \sum_{i=1}^n x_i x_i^T - \bar{x} \left(\sum_{i=1}^n x_i \right)^T - \left(\sum_{i=1}^n x_i \right) \bar{x}^T + n \bar{x} \bar{x}^T \right\} \\ &= \frac{1}{n-1} \left\{ T - \frac{s}{n} s^T - s \left(\frac{s}{n} \right)^T + n \frac{s}{n} \left(\frac{s}{n} \right)^T \right\} \\ &= \frac{1}{n-1} \left\{ T - \frac{1}{n} s s^T \right\}.\end{aligned}$$

Hence, we can recover each statistic using only the canonical information.

- (e) Note that we require $n \geq 1$ to compute \bar{x} and $n \geq 2$ in order to compute \bar{V} . Thus the minimum number of observations to compute (\bar{x}, \bar{V}) is 2.

10

□

2. * What *explicit* form of information would you suggest to represent the sequence (x_1, x_2, \dots, x_n) ? It should contain \bar{x} and \bar{V} and, perhaps, something else.

Solution

In the spirit of minimizing the number of quantities to keep track of, we can use the explicit variables (\bar{x}, \bar{V}, n) to form an information system.

(a) **Uniqueness** follows from the fact that both computations are unique with respect to any representation (in particular, permutation of the coordinates). However, we lack both a well-defined (i) **Elementary** and (ii) **Empty** element as both \bar{x} and \bar{V} are undefined when $n = 0$ and \bar{V} is undefined when $n = 1$.

To define the composition operation, let us first denote the map that takes the canonical information to the data in **Problem 1(d) Completeness**

$$\tau(n, s, T) = (n, \bar{x}, \bar{V}).$$

Observe

$$\begin{aligned}\bar{V} &= \frac{1}{n-1} \left(T - \frac{s s^T}{n} \right) \\ \iff T &= (n-1)\bar{V} + \frac{1}{n}(n\bar{x})(n\bar{x})^T = (n-1)\bar{V} + n\bar{x}\bar{x}^T\end{aligned}$$

strictly speaking, τ is not defined for $n=0$ and $n=1$

Hence, τ is invertible by

$$\tau^{-1}(n, \bar{x}, \bar{V}) = \left(n, n\bar{x}, (n-1)\bar{V} + n\bar{x}\bar{x}^T \right).$$

Schematically, we can construct the (b) **Composition** operation

$$\begin{array}{c}
 (n, \bar{x}, \bar{V}) \xrightarrow{\tau^{-1}} (n, s, T) \\
 \tilde{\oplus} : \quad \quad \quad \searrow \quad \quad \quad \oplus \longrightarrow (\tilde{n}, \tilde{s}, \tilde{T}) \xrightarrow{\tau} (\tilde{n}, \tilde{x}, \tilde{V}) \\
 (n', \bar{x}', \bar{V}') \xrightarrow{\tau^{-1}} (n', s', T') \nearrow
 \end{array}$$

The commutative monoid properties are inherited from \oplus in Problem 1. I.e. ✓

$$\begin{aligned}
 (n, \bar{x}, \bar{V}) \tilde{\oplus} (n', \bar{x}', \bar{V}') &= \tau \left(\tau^{-1}(n, \bar{x}, \bar{V}) \oplus \tau^{-1}(n', \bar{x}', \bar{V}') \right) \\
 &= \tau \left(\tau^{-1}(n', \bar{x}', \bar{V}') \oplus \tau^{-1}(n, \bar{x}, \bar{V}) \right) \\
 &= (n', \bar{x}', \bar{V}') \tilde{\oplus} (n, \bar{x}, \bar{V})
 \end{aligned}$$

and

$$\begin{aligned}
 ((n, \bar{x}, \bar{V}) \tilde{\oplus} (n', \bar{x}', \bar{V}')) \tilde{\oplus} (n'', \bar{x}'', \bar{V}'') &= \tau \left(\tau^{-1}(n', \bar{x}', \bar{V}') \oplus \tau^{-1}(n, \bar{x}, \bar{V}) \right) \tilde{\oplus} (n'', \bar{x}'', \bar{V}'') \\
 &= \tau \left(\tau^{-1} \tau \left(\tau^{-1}(n', \bar{x}', \bar{V}') \oplus \tau^{-1}(n, \bar{x}, \bar{V}) \right) \oplus \tau^{-1}(n'', \bar{x}'', \bar{V}'') \right) \\
 &= \tau \left(\tau^{-1}(n', \bar{x}', \bar{V}') \oplus \tau^{-1} \tau \left(\tau^{-1}(n, \bar{x}, \bar{V}) \oplus \tau^{-1}(n'', \bar{x}'', \bar{V}'') \right) \right) \\
 &= \tau \left(\tau^{-1}(n', \bar{x}', \bar{V}') \oplus \tau^{-1} \left((n, \bar{x}, \bar{V}) \tilde{\oplus} (n'', \bar{x}'', \bar{V}'') \right) \right) \\
 &= (n, \bar{x}, \bar{V}) \tilde{\oplus} \left((n', \bar{x}', \bar{V}') \tilde{\oplus} (n'', \bar{x}'', \bar{V}'') \right) \quad \text{✓}
 \end{aligned}$$

and

$$\begin{aligned}
 (n, \bar{x}, \bar{V}) \tilde{\oplus} (0, 0, 0) &= \tau \left(\tau^{-1}(n, \bar{x}, \bar{V}) \oplus (0, 0 \cdot 0, (0 - 1)0 + 0) \right) \\
 &= \tau \tau^{-1}(n, \bar{x}, \bar{V}) \\
 &= (n, \bar{x}, \bar{V}). \quad \text{✓}
 \end{aligned}$$

Explicitly, this results in the expression $(n, \bar{x}, \bar{V}) \tilde{\oplus} (n', \bar{x}', \bar{V}') = (\tilde{n}, \tilde{x}, \tilde{V})$ where

$$\tilde{n} = n + n', \quad \tilde{x} = \frac{s + s'}{n + n'} = \frac{n\bar{x} + n'\bar{x}'}{n + n'},$$

and

$$\begin{aligned}
 \tilde{V} &= \left[(T + T') \overset{\text{red}}{\oplus} \frac{1}{n + n'} (s + s')(s + s')^T \right] \overset{\text{red}}{\frac{n\bar{x} + n'\bar{x}'}{n + n'}} \\
 &= \left((n - 1)\bar{V} + n\bar{x}\bar{x}^T + (n' - 1)\bar{V}' + n'\bar{x}'\bar{x}'^T \right) + \frac{1}{n + n'} (n\bar{x} + n'\bar{x}')(n\bar{x} + n'\bar{x}')^T \\
 &= \left((n - 1)\bar{V} + n\bar{x}\bar{x}^T + (n' - 1)\bar{V}' + n'\bar{x}'\bar{x}'^T \right) + \frac{1}{n + n'} (n\bar{x}\bar{x}^T + n'n(\bar{x}'\bar{x} + \bar{x}\bar{x}'^T) + n'^2\bar{x}'\bar{x}'^T). \\
 &= \frac{1}{n + n' - 1} \left[(n - 1)\bar{V} + (n' - 1)\bar{V}' + \frac{nn'}{n + n'} (\bar{x} - \bar{x}')(\bar{x} - \bar{x}')^T \right]
 \end{aligned}$$

This suggests that if we wish to save on computation time, we could add to the canonical information $W = \bar{x}\bar{x}'$, and the expression above simplifies to

$$\dots = \left((n-1)\bar{V} + nW + (n'-1)\bar{V}' + n'W' \right) + \frac{1}{n+n'} (nW + n'n(\bar{x}'\bar{x} + \bar{x}\bar{x}'^T) + n'^2W').$$

The **(c) Update** map is very similar to the one derived for scalar mean and variance:

$$\begin{aligned}\bar{x}_{n+1} &= \frac{1}{n+1} \sum_{i=1}^{n+1} \\ &= \frac{n}{n+1} \bar{x}_n + \frac{1}{n+1} x_{n+1} \\ &= \bar{x}_n + \frac{1}{n+1} (x_{n+1} - \bar{x}_n),\end{aligned}$$

and

$$\begin{aligned}\bar{V}_{n+1} &= \frac{1}{n} \sum_{i=1}^{n+1} (x_i - \bar{x}_i)(x_i - \bar{x}_i)^T \leftarrow n+1 \\ &= \frac{n-1}{n} \bar{V}_n + \frac{1}{n} (x_{n+1} - \bar{x}_{n+1})(x_{n+1} - \bar{x}_{n+1})^T \\ &= \bar{V}_n + \frac{1}{n} \left((x_{n+1} - \bar{x}_{n+1})(x_{n+1} - \bar{x}_{n+1}) - \bar{V}_n \right).\end{aligned}$$

$\bar{V}_n = \frac{1}{n-1} \sum (x_i - \bar{x}_n)(x_i - \bar{x}_n)^T \neq \bar{V}_{n+1}$

Note that \bar{V}_{n+1} is expressed in terms of \bar{x}_{n+1} , so the computation for \bar{x}_{n+1} should precede \bar{V}_{n+1} . Also, this suggests that $(0, 0, 0)$ and $(1, x_1, 0)$ could be used for the **(ii) Empty** and **(i) Elementary** elements respectively.

Being an explicit representation, this is clearly **(e) Complete**, and as before, meaningful \bar{x} and \bar{V} are obtained only for $n \geq 2$.

10 □