

0.1 Inference and risk assessment for incomplete diagnoses

A diagnosis is often not complete, meaning that not all lymph node levels (LNLs) might have been checked with a diagnostic modality. E.g., fine needle aspiration (FNA) is usually only performed in a subset of LNLs. Hence, we must be able to deal with “incomplete” observations for some LNLs. To do so, we first introduce a new observation variable

$$d_v \in \{0, 1, \emptyset\} \quad (1)$$

where \emptyset indicates *unobserved*. Furthermore, we define a *match function*

$$\text{match}(\mathbf{d}, \mathbf{z}) := \begin{cases} \text{true} & \text{if } d_v = z_v \vee d_v = \emptyset; \forall v \\ \text{false} & \text{else} \end{cases} \quad (2)$$

which returns *true* if a - potentially incomplete - diagnosis \mathbf{d} is consistent with a complete observation \mathbf{z} . We will use this function for conveniently marginalizing over the missing observations. In analogy to ??, we can compute the risk for an incomplete observation as

$$\begin{aligned} R(X_v = 1 \mid \mathbf{d}, \theta) &= \frac{P(\mathbf{d} \mid X_v = 1, \theta) P(X_v = 1 \mid \theta)}{P(\mathbf{d} \mid \theta)} \\ &= \sum_{i: \xi_{iv}=1} \frac{P(\mathbf{d} \mid \xi_i, \theta) P(\xi_i \mid \theta)}{P(\mathbf{d} \mid \theta)} \end{aligned} \quad (3)$$

where the enumerator of the second line can now be rewritten using the match function:

$$\begin{aligned} P(\mathbf{d} \mid \xi_i, \theta) P(\xi_i \mid \theta) &= \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} P(\zeta_j \mid \xi_i, \theta) P(\xi_i \mid \theta) \\ &= \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} B_{ij} \left[p_T(\mathbf{t}) \cdot \Lambda \right]_i \end{aligned} \quad (4)$$

In this case B_{ij} denotes the element of the observation matrix that corresponds to state ξ_i and observation ζ_j . Again, the indices $\{i : \xi_{iv} = 1\}$ in eq. (3) correspond to all possible states with a positive involvement in lymph node level X_v . Essentially, the whole term is the likelihood of an observation \mathbf{d} where we have removed all entries that correspond to states with $X_v \neq 1$ both from the observation matrix and the resulting probability vector of the evolution. It can therefore be easily computed algebraically, too.

The evidence in the denominator of eq. (3) becomes a marginalization over all possible diagnoses that are not available to us or that we deem unimportant

$$P(\mathbf{d} \mid \theta) = \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} \left[p_T(\mathbf{t}) \cdot \Lambda \right]_j \quad (5)$$

We can make this summation a bit more elegant using a column vector $\mathbf{c}^{\mathbf{d}}$ that has entries corresponding to the match-function

$$c_i^{\mathbf{d}} = \text{match}(\mathbf{d}, \zeta_i) \quad (6)$$

where every *true* corresponds to a 1 and every *false* to a 0. This way we can rewrite eq. (5) in the following way:

$$P(\mathbf{d} \mid \theta) = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \cdot \mathbf{B} \cdot \mathbf{c}^{\mathbf{d}} \quad (7)$$

Using this notation for marginalizing over unknown or incomplete observations also allows us to encode entire datasets $\mathcal{D} = (\mathbf{d}_1 \ \mathbf{d}_2 \ \cdots \ \mathbf{d}_N)$ of (potentially incomplete) observations in the form of a matrix

$$\mathbf{C} = (\mathbf{c}^{\mathbf{d}_1} \ \mathbf{c}^{\mathbf{d}_2} \ \cdots \ \mathbf{c}^{\mathbf{d}_N}) \quad (8)$$

so that the row-vector of likelihoods reads as

$$P(\mathcal{D} \mid \theta) = (P(\mathbf{d}_n \mid \theta))_{n \in [1, N]} = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \cdot \mathbf{B} \cdot \mathbf{C} \quad (9)$$