

Modelling lymphatic progression in head and neck cancer

Roman Ludwig

July 1, 2022

Abstract

Abstract goes here...

Dedication

To mum and dad

Declaration

I declare that..

Acknowledgements

I want to thank...

Contents

1	Previous Work	2
1.1	Bayesian Network	2
1.2	Limitations	4
2	Unilateral hidden Markov model	5
2.1	Formulating lymphatic progression as HMM	5
2.2	Parametrization of the transition matrix	6
2.3	Marginalization	7
2.4	Inference of model parameters	8
2.5	Incorporation of T-stage	9
2.6	Sampling	10
2.7	Risk assessment of microscopic involvement	11
2.8	Inference and risk assessment for incomplete diagnoses	11
2.9	Multiple diagnostic modalities	13
2.10	Combining modalities and data	13
3	Bilateral hidden Markov model	15
3.1	Expanding the unilateral model	15
3.2	Parameter symmetries and mid-line extension	16
4	Possible further extension	19
4.1	Trinary hidden random variables	19
5	Data on patterns of lymphatic progression	21
5.1	Detailed reporting of involvement in OPSCC	22
5.1.1	Abstract	22
5.1.2	Introduction	22
5.1.3	Material & methods	24
5.1.4	Data base	24
5.1.5	Graphical user interface	25

Chapter 1

Previous Work

In this chapter I will introduce models that were previously developed to capture and predict lymphatic progression in head and neck squamous cell carcinoma (HNSCC). This also includes a Bayesian network (BN) model by [18], which served as a starting point to this work. This brief recap will also introduce the notation and formalism that will be used throughout the rest of the thesis.

After that, unrelated previous works will be mentioned and lastly, we will discuss the limitations and reasons to develop a new model and formalism.

1.1 Bayesian Network

We model the state of each lymph node level (LNL) as a hidden or unobserved binary random variable, which indicates via values 0 or 1 if an LNL is healthy or involved, respectively. This state indicates if there is truly tumor present in an LNL, including the presence of occult metastases for the involved state – motivating the term hidden or unobserved state. Every LNL can be diagnosed using one or multiple modalities. Most used for diagnosis are imaging techniques like positron emission tomography (PET), computed tomography (CT) and magnetic resonance imaging (MRI), but palpation or fine needle aspiration (FNA) are also used. The diagnosis too, is modelled as binary random variable – this time an observed one – taking on 0 for negative and 1 for positive.

For notational convenience, we collect the hidden and observed random variables in a random vector each:

$$\begin{aligned} \text{hidden} \quad \mathbf{X} &= (X_v) \rightarrow \{0, 1\}^V \\ \text{observed} \quad \mathbf{Z} &= (Z_v) \rightarrow \{0, 1\}^V \end{aligned} \tag{1.1}$$

where V is the number of LNLs $v \in \{1, 2, \dots, V\}$ in the graph. The conditional probabilities that link the hidden state to the observations can be written as follows:

$$\begin{aligned} P_{BN}(Z_v = z_v \mid X_v = x_v) &= (z_v + (-1)^{z_v} \cdot s_P)(1 - x_v) \\ &\quad + ((1 - z_v) + (-1)^{1-z_v} \cdot s_N) x_v \end{aligned} \tag{1.2}$$

with s_N and s_P being the sensitivity and specificity of the used diagnostic method. For example, for the probability of a false negative observation, i.e. diagnostic

modality misses the presence of tumor, we get

$$P_{BN}(Z_v = 0 \mid X_v = 1) = 1 - s_N \quad (1.3)$$

Spread of the tumor through the lymphatic network is represented in this model by directed arcs to and between LNLs as illustrated in ???. We introduce an additional vertex to the graph representing the primary tumor, which we assume to be the only one. Directed arcs from the primary tumor to an LNL represent direct spread of tumor cells from the primary tumor to the LNL. These arcs are associated with parameters b_v that we call base probabilities, and which indicate the probability that the tumor spreads directly to LNL v . When LNL s receives efferent lymphatics from LNL r , this too is represented by a directed arc from LNL r to s , and $r = \text{pa}(s)$ which is called a parent node of s . These arcs are associated with a transition probability t_{rs} from r to s . The resulting directed acyclic graph (DAG) is shown in ??, comprising ipsilateral levels I, II, II, through IV, and will be used throughout this work. However, when more data of detailed LNL involvement including additional levels becomes available and/or contralateral involvement, the model can be extended. The parameters b_v and t_{rs} associated with the directed arcs represent conditional probabilities, i.e. b_v answers the question given that all parent nodes are healthy, how likely is it that the primary tumor spreads to node v ? t_{rs} on the other hand, can answer the question assuming no efferent spread from the primary tumor and given that all parent nodes except r are healthy, what is the likelihood of spread to node s ? The conditional probability for involvement of LNL v given the state of its parent nodes is then given by

$$\begin{aligned} P_{BN}(X_v = x_v \mid X_{\text{pa}(v)} = x_{\text{pa}(v)}, b_v, t_{\text{pa}(v)v}) \\ = x_v + (-1)^{x_v} (1 - b_v) (1 - t_{\text{pa}(v)v})^{x_{\text{pa}(v)}} \end{aligned} \quad (1.4)$$

We note here that this parametrization assumes the independence of causal influences (ICI), thereby allowing us to describe the model using only a few interpretable parameters. Dropping this assumption, a BN can also be defined using conditional probability tables (CPT) that have columns for every possible combinations of parent states. However, with the increase of the number of parent nodes (causes) in the graph, the number of parameters in the respective CPT would grow exponentially.

For the graph in ??? we can write down the parametrized CPT in the following manner:

$$\begin{aligned} P_{BN}(X_v = 0 \mid X_{\text{pa}(v)} = 0) &= 1 - b_v \\ P_{BN}(X_v = 1 \mid X_{\text{pa}(v)} = 0) &= b_v \\ P_{BN}(X_v = 0 \mid X_{\text{pa}(v)} = 1) &= (1 - b_v) (1 - t_{\text{pa}(v)v}) \\ P_{BN}(X_v = 1 \mid X_{\text{pa}(v)} = 1) &= 1 - (1 - b_v) (1 - t_{\text{pa}(v)v}) \end{aligned} \quad (1.5)$$

In case of a more general network, in which some LNLs receive efferent lymphatics from multiple other LNLs, eq. (1.5) can be generalised and the conditional probability of the hidden state becomes

$$\begin{aligned} P_{BN}(X_v = x_v \mid \{X_r = x_r, t_{rv}\}_{r \in \text{pa}(v)}, b_v) \\ = x_v + (-1)^{x_v} (1 - b_v) \prod_{r \in \text{pa}(v)} (1 - t_{rv})^{x_r} \end{aligned} \quad (1.6)$$

where we marginalized over all hidden variables X . Here we have assumed that each patient’s diagnosis $\mathbf{z} = (z_1 \ z_2 \ \cdots \ z_V)$ is complete, meaning that we have a diagnosis for each LNL. The likelihood can then be used to infer the model parameters via maximum likelihood inference or sampling.

1.2 Limitations

While BNs can model the probabilistic relationship between involvement in different levels, they lack an explicit way to describe the evolution of the tumor over time. The concept of dynamic Bayesian networks (DBNs) has been developed to introduce the notion of time into probabilistic models. DBNs are generalizations of hidden Markov models (HMMs) and formally similar to what we will introduce now. The metastatic spread in the lymphatic system is a dynamic system and by modeling it with a formalism that can capture this, we obtain a more intuitive model of the problem and a framework that can incorporate T-stage into estimating the risk of LNL involvement. We can do this because tumors go through the stages T1 to T4 sequentially, meaning that – for a given tumor – it is a surrogate of time.

Chapter 2

Unilateral hidden Markov model

This chapter concerns itself with modelling the unilateral lymphatic spread using the formalism of HMMs that will also be introduced below.

The content of this chapter is largely based on our publication [14].

2.1 Formulating lymphatic progression as HMM

We consider discrete time-steps $t \in \{0, 1, 2, \dots, t_{\max}\}$. We will start by defining the hidden random variable for the state of the HMM at time t to be

$$\mathbf{X}[t] = (X_v[t]) \quad (2.1)$$

which represents the patient's state of LNL involvement as in the BN, but for each time-step we have an instance of it. For the diagnosis \mathbf{Z} on the other hand, we do not need to differentiate between different times, since in practice we will only ever see one diagnosis. This is illustrated in ???. The reason for this is that, if we diagnose a patient with cancer, treatment starts timely and we no longer observe the natural progression of the disease. From a modeling standpoint however, this is a problem that we will address later.

A hidden Markov model is fully described by the starting state $\mathbf{X}[0] := \boldsymbol{\pi}$ and the two conditional probability functions that govern the progression from a state $X[t]$ at time t to a state $X[t+1]$ at the following time-step

$$P_{HMM}(\mathbf{X}[t+1] \mid \mathbf{X}[t]) \quad (2.2)$$

and the probability of a diagnostic observation given the true state of the patient

$$P_{HMM}(\mathbf{Z} \mid \mathbf{X}[t]) \quad (2.3)$$

Since both our state space and our observation space are discrete and finite, it is possible to enumerate all possible states and observations and collect them in a table or matrix. This so-called *transition matrix* would then be

$$\mathbf{A} = (A_{ij}) = (P_{HMM}(\mathbf{X}[t+1] = \boldsymbol{\xi}_i \mid \mathbf{X}[t] = \boldsymbol{\xi}_j)) \quad (2.4)$$

and the *observation matrix*

$$\mathbf{B} = (B_{ij}) = (P_{HMM}(\mathbf{Z} = \zeta_i \mid \mathbf{X}[t] = \boldsymbol{\xi}_j)) \quad (2.5)$$

Here, ξ_i and ζ_j are no new variables, but just \mathbf{x} and \mathbf{z} renamed and reordered. The indices i and j for one of the possible states or observations for the entire patient, not for an individual LNL. In total, there are $S = |\{0, 1\}|^V$ different states and the same number of different possible observations per diagnostic modality. We order the hidden states from

$$\xi_1 = (0 \ 0 \ 0 \ 0) \quad (2.6)$$

to

$$\xi_{16} = (1 \ 1 \ 1 \ 1) \quad (2.7)$$

in this case of $V = 4$. The exact ordering does not matter, it is just a convenience for the notation. our ordering of the states can be seen in the axes of ???. In analogy, we order the observations ζ_j from 1 to 2^V . Note that for now we will not consider multiple diagnostic modalities and how to combine them. We will get back to that topic in ??.

In our case, the starting state corresponds to a primary tumor being present but all LNLs are still in the healthy state. The observation matrix \mathbf{B} is specified via sensitivity and specificity as described in eq. (2.7). The main task is to infer the transition matrix \mathbf{A} . Usually, it is inferred from a series of observations and there exist efficient algorithms for that, e.g. the sum-product algorithm, which is particularly efficient in chains. Unfortunately, these algorithms cannot be applied for our problem for two profound reasons:

1. We only have a single observation instead of a consecutive series of observations.
2. It is unclear how many time-steps it took from the starting state to the one observation we have at the time of diagnosis.

In the remainder of this section, we will detail the HMM step-by-step, starting with the parameterization of the transition matrix \mathbf{A} in section 2.2. Afterwards, in section 2.3, I will tackle the aforementioned problems, followed up by explaining how we perform inference on this model (section 2.4), incorporate information about a patient's T-stage (section 2.5) and assess the risk of LNL involvement in a new patient (section 2.7). Lastly, we will introduce a way to incorporate incomplete observations in section 2.8.

2.2 Parametrization of the transition matrix

The square transition matrix \mathbf{A} has $S = 2^{2V}$ entries and therefore $S(S - 1) = 2^{2V} - 2^V$ degrees of freedom. Although searching the full space of viable transition matrices is possible via unparametrized sampling techniques, it is computationally challenging and hard to interpret. To achieve this reduction in degrees of freedom, and also preserve the anatomically and medically motivated structure of the Bayesian network from ??, we can represent the transition probability from one state $\mathbf{x}[t]$ to another state $\mathbf{x}[t + 1]$ using the conditional probabilities defined for the BN. The difference is that the probability of observing a certain state of LNL v now depends on the state of the patient one time-step before. Note that from

here on, we will mostly drop the probabilistically correct notation $P(X = x)$ and just write $P(x)$ for brevity

$$P_{HMM}(\mathbf{x}[t+1] \mid \mathbf{x}[t]) = \prod_{v \leq V} Q(x_v[t+1]; x_v[t]) \times \left[P_{BN}(x_v[t+1] \mid \{x_r[t], \tilde{t}_{rv}\}_{r \in \text{pa}(v)}, \tilde{b}_v) \right]^{1-x_v[t]} \quad (2.8)$$

Here we have reused the conditional probability from the BN for each LNL, but we take it to the power of one minus that node's previous value. This ensures that an involved node stays involved with probability 1. The parameters $\tilde{t}_{\text{pa}(v)v}$ and \tilde{b}_v take the same role as in the BN, but they are now probability *rates*, since they act per time-step. Lastly, the first term Q in the product formalizes the fact that a metastatic lymph node level cannot become healthy again once it was involved. This also means that several entries in the transition matrix \mathbf{A} must be zero. In a table the values of $Q(x_v[t+1]; x_v[t])$ can be written like this:

$$\begin{aligned} Q(X_v[t+1] = 0; X_v[t] = 0) &= 1 \\ Q(X_v[t+1] = 0; X_v[t] = 1) &= 0 \\ Q(X_v[t+1] = 1; X_v[t] = 0) &= 1 \\ Q(X_v[t+1] = 1; X_v[t] = 1) &= 1 \end{aligned} \quad (2.9)$$

which gives rise to a "mask" for \mathbf{A} which can be seen in ??.

To illustrate eq. (2.8), it helps to look at a specific example. E.g., the transition probability from state $\xi_5 = (0 \ 1 \ 0 \ 0)$ to state $\xi_7 = (0 \ 1 \ 1 \ 0)$, which represents starting with involvement only in LNL II and asking for the probability that LNL III becomes involved as well over the next time-step:

$$\begin{aligned} P_{HMM}(\mathbf{X}[t+1] = \xi_7 \mid \mathbf{X}[t] = \xi_5) \\ = & Q(X_1[t+1] = 0; X_1[t] = 0) P_{BN}(X_1[t+1] = 0 \mid \tilde{b}_1)^1 \\ & \times Q(X_2[t+1] = 1; X_2[t] = 1) P_{BN}(X_2[t+1] = 1 \mid X_1[t] = 0, \tilde{t}_{12}, \tilde{b}_2)^0 \\ & \times Q(X_3[t+1] = 1; X_3[t] = 0) P_{BN}(X_3[t+1] = 1 \mid X_2[t] = 1, \tilde{t}_{23}, \tilde{b}_3)^1 \\ & \times Q(X_4[t+1] = 0; X_4[t] = 0) P_{BN}(X_4[t+1] = 0 \mid X_3[t] = 0, \tilde{t}_{34}, \tilde{b}_4)^1 \\ = & (1 - \tilde{b}_1) \cdot 1 \cdot (\tilde{b}_3 + \tilde{t}_{23} - \tilde{b}_3 \tilde{t}_{23}) \cdot (1 - \tilde{b}_4) \end{aligned} \quad (2.10)$$

The interpretation of the last line is that this is the probability that LNL I and IV do not become involved, while LNL III gets infected through lymphatic drainage from either the main tumor or LNL II. The probability of LNL II remaining involved is 1, of course, which is why we take the respective term to the power of 0.

2.3 Marginalization

To calculate the likelihood function, we have to calculate the probability of a given diagnostic observation. To that end, we first calculate the probability of observing

a given diagnosis $\mathbf{z} = \zeta_j$ at a fixed time-step t . As depicted in ??, we must consider every possible evolution of a patient's disease that leads to the observed diagnosis. Mathematically, this means that we need to marginalize over all such paths. And here is where the HMM-formalism comes in very useful, because this marginalization happens automatically when we multiply the transition matrix with itself and eventually with the observation matrix:

$$P(\mathbf{Z} = \zeta_j | t) = [\boldsymbol{\pi}^\top \cdot (\mathbf{A})^t \cdot \mathbf{B}]_j \quad (2.11)$$

where the $\boldsymbol{\pi}$ is the column vector for the healthy starting state. \mathbf{A} is multiplied with itself t times and thereby produces a matrix that describes the transition probability from the healthy state to all possible states $\mathbf{x}[t]$ in exactly t time-steps marginalized over the actual pathway of the patient's disease. The index $[\dots]_j$ here means that from the resulting (row-)vector of probabilities we take the component that corresponds to the diagnose $\mathbf{z} = \zeta_j$.

So, essentially, eq. (2.11) first computes the probability vector of all possible true hidden states, given a time step t

$$P(\mathbf{X} = \xi_i | t) = [\boldsymbol{\pi}^\top \cdot (\mathbf{A})^t]_i \quad (2.12)$$

and then multiplies it with the respective observation probability vector, which is a column of the \mathbf{B} matrix, to finally marginalize over all possible true hidden states – effectively a sum over i in eq. (2.12) – at the time t of diagnosis.

The problem that the number of time-steps until diagnosis is unknown cannot be solved in such an elegant fashion. Therefore, we must resort to brute force marginalization and introduce a prior $p(t)$, which is a discrete distribution over a finite number of time-steps. It describes the prior probability that a patient's cancer is diagnosed at a particular time-step t . To get the probability of a diagnosis \mathbf{z} we must compute

$$P(\mathbf{Z} = \zeta_j) = \sum_{t=0}^{t_{\max}} p(t) \cdot P(\mathbf{Z} = \zeta_j | t) = \left[\sum_{t=0}^{t_{\max}} p(t) \cdot \boldsymbol{\pi}^\top \cdot (\mathbf{A})^t \cdot \mathbf{B} \right]_j \quad (2.13)$$

While the choice of the time-prior may seem unclear at this point, its role for including T-stage into this model will be discussed in section 2.5.

2.4 Inference of model parameters

In the formalism of the last sections, the P_{HMM} depends implicitly through P_{BN} on parameters $\theta = \{\tilde{b}_v, \tilde{t}_{pv} \mid v \leq V, p \in \text{pa}(v)\}$, which – as mentioned – are now probability rates and have therefore a slightly different interpretation. Due to the marginalization over time-steps in eq. (2.13) the likelihood function additionally depends on the choice and parametrization of the prior $p(t)$. The parameters are to be inferred from a dataset of lymphatic progression patterns in a cohort of patients. We still assume that for each patient we record for every LNL v whether it is involved according to only one diagnostic modality. In other words, for each patient we observe one of the 2^V possible diagnoses. As mentioned before, we will expand this to multiple diagnostic modalities further down in ??.

Formally, we can then express the dataset \mathcal{Z} of N patients as vector \mathbf{f} of the number of patients f_i for which the diagnosis corresponds to the observational state ζ_i . The likelihood $P(\mathcal{Z} | \theta)$ of observing this dataset, given a particular choice of parameters, is then given by

$$P(\mathcal{Z} | \theta) = \prod_{i=1}^{2^V} P(\zeta_i | \theta)^{f_i} \quad (2.14)$$

with the probability $P(\zeta_i | \theta)$ specified by eq. (2.13). The product runs formally over all possible observational states. In reality, f_i will likely be zero for a number of rare or implausible states that are not in the dataset. Note that $\sum_i f_i = N$.

By Bayes' rule, the posterior distribution of those parameters is

$$P(\theta | \mathcal{Z}) = \frac{P(\mathcal{Z} | \theta) P(\theta)}{\int P(\mathcal{Z} | \theta') P(\theta') d\theta'} \quad (2.15)$$

where $P(\theta)$ is the prior over these parameters. Since they are exclusively probability rates, they must all come from the interval $[0, 1] \in \mathbb{R}$. In this work we will choose the most uninformative prior

$$p(\theta) = \begin{cases} 1 & \text{if } \theta_r \in [0, 1]; \forall r \leq E \\ 0 & \text{otherwise} \end{cases} \quad (2.16)$$

where E is the number of edges in the DAG we use to represent the lymphatic system. While it is easy to compute the likelihood, it is not feasible to efficiently calculate the normalization constant in the denominator of eq. (2.15). Hence, we will use Markov-chain Monte Carlo (MCMC) sampling methods to estimate the parameters θ and their uncertainty.

2.5 Incorporation of T-stage

We have introduced the HMM with the promise that it could handle the concept of T-stages through its explicit modeling of dynamic processes. To keep up with that, we will now explain how this is achieved using the time-prior $p(t)$.

The core idea is to assume that early T-stage and late T-stage tumors share the same patterns of metastatic progression, except that late T-stage tumors are on average diagnosed at a later point in time, and thereby also show, on average, higher LNL involvement. Formally, this can be described by assuming a different time-prior $p_T(t)$ for every T-stage T . On the other hand, the transition matrix \mathbf{A} is assumed to be the same for all T-stages.

For the inference of model parameters, the training data is split into subgroups according to T-stage. We now define a column-vector \mathbf{f}_T separately for each T-stage, which counts the number of patients in the dataset that were diagnosed with one of the possible observational states and a given T-stage. The log-likelihood from which we want to sample is then simply a sum of the likelihoods as above, where the essential difference is that we equip each marginalization over time with a different time-prior $p_T(t)$, according to its T-stage:

$$\log P(\mathcal{Z} | \theta) = \sum_{T=1}^4 \log \left[\sum_{t=0}^{t_{\max}} p_T(t) \cdot \boldsymbol{\pi}^\top \cdot (\mathbf{A})^t \cdot \mathbf{B} \right] \cdot \mathbf{f}_T \quad (2.17)$$

The logarithm must be taken element-wise for the resulting row-vector inside the square brackets. The only data-dependent term here is the vector \mathbf{f}_T counting the occurrences of all possible observations. It is again important to note that the only difference between the part of the log-likelihood for the different T-stages is the exact shape or parametrization of the time-prior. The transition probabilities, and hence also the transition matrix \mathbf{A} , are the same for all T-stages. For this to work, we rely on the assumption that different typical patterns of nodal involvement for the same primary tumor location are caused mainly by different progression times

At this point, it makes sense to briefly introduce a notation of the above equation that is more suitable for the actual programmatic implementation of the inference and the extension we will discuss later. We can rewrite the term in the square brackets of eq. (2.17) by using the matrix

$$\mathbf{\Lambda} := P(\mathbf{X} | \mathbf{t}) = \begin{pmatrix} \boldsymbol{\pi}^\top \cdot (\mathbf{A})^0 \\ \boldsymbol{\pi}^\top \cdot (\mathbf{A})^1 \\ \vdots \\ \boldsymbol{\pi}^\top \cdot (\mathbf{A})^{t_{\max}} \end{pmatrix} \quad (2.18)$$

where row number t corresponds to the vector $\boldsymbol{\pi}^\top \cdot (\mathbf{A})^t$, i.e. the probabilities for all possible hidden states, given the diagnose time. So, the element Λ_{ti} corresponds to the probability $P(\xi_i | t)$ of a patient arriving in the i th state after t time steps. With this, we can rewrite the term in the square brackets of eq. (2.17) purely as a product of vectors and matrices:

$$\sum_{t=0}^{t_{\max}} p_T(t) \cdot \boldsymbol{\pi}^\top \cdot (\mathbf{A})^t = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \quad (2.19)$$

with $p_T(\mathbf{t}) = (p_T(0) \ p_T(1) \ \dots \ p_T(t_{\max}))$. The matrix $\mathbf{\Lambda}$ implicitly depends on the spread probabilities, while each of the $p_T(\mathbf{t})$ depends on the respective parametrization of the time prior. They are the only objects that depend on the parameters θ and they are independent of the data.

2.6 Sampling

With a parameter set $\theta = (\{\tilde{b}_v\}, \{\tilde{t}_{rv}\}_{r \in \text{pa}(v)}) \ \forall v \leq V$, we can assess the risk of nodal involvement, given a diagnosis \mathbf{z} , of a new patient. Using Bayes' law, the risk for a certain LNL v being involved is given by the conditional probability

$$\begin{aligned} R(X_v = 1 | \mathbf{z}, \theta) &= \frac{P(\mathbf{Z} = \mathbf{z} | X_v = 1, \theta) P(X_v = 1 | \theta)}{P(\mathbf{Z} = \mathbf{z} | \theta)} \\ &= \sum_{i: \xi_{iv}=1} \frac{P(\mathbf{Z} = \mathbf{z} | \xi_i, \theta) P(\xi_i | \theta)}{P(\mathbf{Z} = \mathbf{z} | \theta)} \end{aligned} \quad (2.20)$$

Note that in the second line, we have explicitly written out the marginalization over all hidden states ξ_i that have LNL v involved. We have written the state of LNL v in the state ξ_i as ξ_{iv} . The denominator can be computed using eq. (2.13), which already includes the marginalization over all hidden states ξ_i .

The process of sampling randomly generates L sets of parameters $\theta = (\theta_1 \ \theta_2 \ \dots \ \theta_L)$. They are therefore random variables and so is the risk $R(X_v | \mathbf{z}, \theta)$ since it is a function of θ . Using the Monte Carlo estimator, we can therefore compute the moments of the distribution over the risk, including e.g. the expectation value

$$\mathbb{E}_{\theta} [R(X_v = 1 | \mathbf{z})] = \frac{1}{L} \sum_{k=1}^L R(X_v = 1 | \mathbf{z}, \theta_k) \quad (2.21)$$

In the result sections below, we compute the individual risks for a large enough number L of sampled parameters. Thereby, we can compute histograms for the risk that will approach the real probability density of the respective risk for $L \rightarrow \infty$. This provides additional information on the uncertainty in the predicted risk resulting from uncertainty in the model parameters.

2.7 Risk assessment of microscopic involvement

With a parameter set $\theta = (\{\tilde{b}_v\}, \{\tilde{t}_{rv}\}_{r \in \text{pa}(v)}) \ \forall v \leq V$, we can assess the risk of nodal involvement, given a diagnosis \mathbf{z} , of a new patient. Using Bayes' law, the risk for a certain LNL v being involved is given by the conditional probability

$$\begin{aligned} R(X_v = 1 | \mathbf{z}, \theta) &= \frac{P(\mathbf{Z} = \mathbf{z} | X_v = 1, \theta) P(X_v = 1 | \theta)}{P(\mathbf{Z} = \mathbf{z} | \theta)} \\ &= \sum_{i: \xi_{iv}=1} \frac{P(\mathbf{Z} = \mathbf{z} | \xi_i, \theta) P(\xi_i | \theta)}{P(\mathbf{Z} = \mathbf{z} | \theta)} \end{aligned} \quad (2.22)$$

Note that in the second line, we have explicitly written out the marginalization over all hidden states ξ_i that have LNL v involved. We have written the state of LNL v in the state ξ_i as ξ_{iv} . The denominator can be computed using eq. (2.13), which already includes the marginalization over all hidden states ξ_i .

The process of sampling randomly generates L sets of parameters $\theta = (\theta_1 \ \theta_2 \ \dots \ \theta_L)$. They are therefore random variables and so is the risk $R(X_v | \mathbf{z}, \theta)$ since it is a function of θ . Using the Monte Carlo estimator, we can therefore compute the moments of the distribution over the risk, including e.g. the expectation value

$$\mathbb{E}_{\theta} [R(X_v = 1 | \mathbf{z})] = \frac{1}{L} \sum_{k=1}^L R(X_v = 1 | \mathbf{z}, \theta_k) \quad (2.23)$$

In the result sections below, we compute the individual risks for a large enough number L of sampled parameters. Thereby, we can compute histograms for the risk that will approach the real probability density of the respective risk for $L \rightarrow \infty$. This provides additional information on the uncertainty in the predicted risk resulting from uncertainty in the model parameters.

2.8 Inference and risk assessment for incomplete diagnoses

A diagnosis is often not complete, meaning that not all LNLs might have been checked with a diagnostic modality. E.g., 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

$$d_v \in \{0, 1, \emptyset\} \quad (2.24)$$

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

$$\text{match}(\mathbf{d}, \mathbf{z}) := \begin{cases} \text{true} & \text{if } d_v = z_v \vee d_v = \emptyset; \forall v \\ \text{false} & \text{else} \end{cases} \quad (2.25)$$

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 eq. (2.22), we can compute the risk for an incomplete observation as

$$\begin{aligned} R(X_v = 1 \mid \mathbf{d}, \theta) &= \frac{P(\mathbf{d} \mid X_v = 1, \theta) P(X_v = 1 \mid \theta)}{P(\mathbf{d} \mid \theta)} \\ &= \sum_{i: \xi_{iv}=1} \frac{P(\mathbf{d} \mid \xi_i, \theta) P(\xi_i \mid \theta)}{P(\mathbf{d} \mid \theta)} \end{aligned} \quad (2.26)$$

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

$$\begin{aligned} P(\mathbf{d} \mid \xi_i, \theta) P(\xi_i \mid \theta) &= \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} P(\zeta_j \mid \xi_i, \theta) P(\xi_i \mid \theta) \\ &= \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} B_{ij} \left[p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \right]_i \end{aligned} \quad (2.27)$$

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. (2.26) 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. (2.26) becomes a marginalization over all possible diagnoses that are not available to us or that we deem unimportant

$$P(\mathbf{d} \mid \theta) = \sum_{\{j: \text{match}(\mathbf{d}, \zeta_j)\}} \left[p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \right]_j \quad (2.28)$$

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

$$c_i^{\mathbf{d}} = \text{match}(\mathbf{d}, \zeta_i) \quad (2.29)$$

where every *true* corresponds to a 1 and every *false* to a 0. This way we can rewrite eq. (2.28) in the following way:

$$P(\mathbf{d} \mid \theta) = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \cdot \mathbf{B} \cdot \mathbf{c}^{\mathbf{d}} \quad (2.30)$$

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

$$\mathbf{C} = (\mathbf{c}^{\mathbf{d}_1} \ \mathbf{c}^{\mathbf{d}_2} \ \dots \ \mathbf{c}^{\mathbf{d}_N}) \quad (2.31)$$

so that the row-vector of likelihoods reads as

$$P(\mathcal{D} \mid \theta) = (P(\mathbf{d}_n \mid \theta))_{n \in [1, N]} = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \cdot \mathbf{B} \cdot \mathbf{C} \quad (2.32)$$

2.9 Multiple diagnostic modalities

Throughout the last sections, we have only dealt with diagnoses from a single modality. In practice, however, most patients undergo screening for metastases using different modalities, like CT, MRI or FNA. The sensitivities and specificities of these might vary greatly and by combining them in a probabilistically rigorous way, we may gain a additional information.

Luckily, the introduced formalism requires very little changes to be able to incorporate multiple diagnostic modalities. Let $\mathcal{O} = \{\text{CT}, \text{MRI}, \text{FNA}, \dots\}$ be the set of modalities. Then we can extend the collection of observed binary random variables (RVs) \mathbf{z} from a single modality

$$\mathbf{z} = (x_v)_{v \in [1, V]} = (x_1 \ \dots \ x_V) \quad (2.33)$$

to multiple diagnostic modalities

$$\mathbf{z} = (x_v^k)_{\substack{v \in [1, V] \\ k \in [1, |\mathcal{O}|]}} = \begin{pmatrix} x_1^1 & \dots & x_V^1 & x_1^2 & \dots & x_V^{|\mathcal{O}|} \end{pmatrix} \quad (2.34)$$

where k enumerates the elements in the set \mathcal{O} . We can use ζ_j again and this time the counting variable j goes from 1 to $2^{V \cdot |\mathcal{O}|}$. Notice that this means the observation matrix \mathbf{B} is not square anymore. Also, it now contains the sensitivities and specificities of all the modalities in \mathcal{O} . If we had separate square observation matrices \mathbf{B}^k for each diagnostic modality, the new total matrix' rows B_{i*} would be the outer products of the individual observation matrices:

$$B_{i*} = B_{i*}^1 \otimes B_{i*}^2 \otimes \dots \otimes B_{i*}^{|\mathcal{O}|} \quad (2.35)$$

Completely analogous to how we enlarged the vector of binary RVs \mathbf{z} , we can also extend the vectors \mathbf{c} and \mathbf{d} and then immediately use the entire formalism of the section before to model lymphatic progression with potentially incomplete diagnoses from multiple modalities. However, we will drop this way of continuously enumerating the observations in the next section again, because there is a slightly more efficient and elegant way to do it. This section only served to show that it is naturally possible to extend the formalism to combine findings from different diagnostic modalities.

2.10 Combining modalities and data

Note that the matrix \mathbf{B} – and also the matrix \mathbf{C} – can get very large very quickly: The former is of size $2^V \times 2^{V \cdot |\mathcal{O}|}$ and the latter has dimensions $2^{V \cdot |\mathcal{O}|} \times N$, meaning

both grow exponentially with the number of LNLs *and* diagnostic modalities. And although neither \mathbf{B} nor \mathbf{C} depend on the parameters θ , meaning their product can be precomputed, we can simply iterate over all patients, possible hidden states and available diagnostic modalities to compute $\mathbf{\Omega} := \mathbf{B} \cdot \mathbf{C}$ directly, which saves us building up and multiplying matrices with potentially millions of entries.

To compute this matrix $\mathbf{\Omega}$, we first abandon the just-introduced way of combining diagnoses for all modalities into one large vector and separate them again, so that we have complete and incomplete observations ζ_j^k and \mathbf{d}_n^k respectively for each modality, where $n \in [1, N]$ enumerates the patients in the data.

$$\begin{aligned} \Omega_{mn} &= P(\mathbf{d}_n \mid \xi_m) = \prod_{k=1}^{|\mathcal{O}|} P(\mathbf{d}_n^k \mid \xi_m) \\ &= \prod_{k=1}^{|\mathcal{O}|} \left[\sum_{j: \text{match}(\mathbf{d}_n^k, \zeta_j^k)} P(\zeta_j \mid \xi_m) \right] = \prod_{k=1}^{|\mathcal{O}|} \left[\sum_{j: \text{match}(\mathbf{d}_n^k, \zeta_j^k)} B_{mj}^k \right] \end{aligned} \quad (2.36)$$

Now, the elements Ω_{mn} encode the observation likelihood of patient n 's diagnose \mathbf{d}_n given their true state of involvement is ξ_m . Finally, with this the row-vector of likelihoods of a cohort of patients, given the model's spread parameters, becomes

$$P(\mathcal{D} \mid \theta) = p_T(\mathbf{t}) \cdot \mathbf{\Lambda} \cdot \mathbf{\Omega} \quad (2.37)$$

Again, the objects $p_T(\mathbf{t})$ and $\mathbf{\Lambda}$ depend on the parameters and hence need to be recalculated for every sample drawn during MCMC inference. $\mathbf{\Omega}$ depends only on the patient data \mathcal{D} and must therefore only be computed once at the beginning of the learning round.

Chapter 3

Bilateral hidden Markov model

In the previous chapter we have set up the formalism to deal with only one side of the neck. Implicitly, we have assumed that to be the ipsilateral side, i.e. the side of the sagittal plane where the primary tumor is located. This is because we assume lymphatic drainage to a process that is somewhat symmetric w.r.t. the sagittal plane, which means there can only be limited lymph flow across this plane. But depending on the tumor's location and lateralization, drainage and hence metastatic spread to the contralateral lymphatic system of the neck may also occur. In current clinical practice, a bilateral neck dissection or irradiation is often prescribed when the tumor is close to the mid-sagittal plane. So, ideally, we would like to model the risk for involvement in both sides of the neck at the same time.

The formalism of chapter 2 can easily be applied to the contralateral side and given respective training data for the sampling process, it would learn the appropriate spread probabilities to and among the contralateral LNLs just as it would learn the ones for the ipsilateral side. From clinical experience, the contralateral involvement is usually less severe than the ipsilateral one, and hence we would expect the contralateral spread to be less probable as well.

However, combining two such unilateral models naively would make the assumption that ipsi- and contralateral spread are independent, which seems unlikely: If we know a patient has advanced metastases in the contralateral neck nodes, the risk to find similarly or even more advanced disease in ipsilateral neck nodes should probably be higher than if the contralateral neck were healthy. In other words, we are now looking for the joint probability $P(\mathbf{X}^i, \mathbf{X}^c | \mathbf{Z}^i, \mathbf{Z}^c)$, where the superscripts i and c indicate the ipsi- and contralateral side respectively.

The following section will pick up the unilateral formalism, extend and modify it to come up with a less naive bilateral model.

3.1 Expanding the unilateral model

If we start by dissecting this joint conditional probability in the following way

$$P(\mathbf{X}^i, \mathbf{X}^c | \mathbf{Z}^i, \mathbf{Z}^c) = \frac{P(\mathbf{Z}^i, \mathbf{Z}^c | \mathbf{X}^i, \mathbf{X}^c) \cdot P(\mathbf{X}^i, \mathbf{X}^c)}{P(\mathbf{Z}^i, \mathbf{Z}^c)} \quad (3.1)$$

we notice right away that the likelihood on the right factorizes: Given the true states of involvement in the two sides of the neck, their respective diagnoses must

be independent. Furthermore, the two factors are already given by their corresponding observation matrices \mathbf{B}^i and \mathbf{B}^c .

The joint probability of the hidden states $P(\mathbf{X}^i, \mathbf{X}^c)$ does not factorize in the same manner. But if we assume the lymphatic network to be symmetric and directed, there can be no direct connection between LNLs of the two sides of the neck, which means the probability for involvement of the ipsi- and contralateral side only correlate via the diagnose time t . Hence the joint probability is a sum of factorizing terms:

$$\begin{aligned} P(\mathbf{X}^i, \mathbf{X}^c) &= \sum_{t \in \mathbb{T}} p(t) \cdot P(\mathbf{X}^i, \mathbf{X}^c | t) \\ &= \sum_{t \in \mathbb{T}} p(t) \cdot P(\mathbf{X}^c | t) \cdot P^\top(\mathbf{X}^i | t) \end{aligned} \quad (3.2)$$

Note that the two row vectors of probabilities in the second line are multiplied using an outer product. Using the notation from the last section, We can write this in an algebraic way to effectively factorize this sum as follows

$$P(\mathbf{X}^c = \boldsymbol{\xi}_n, \mathbf{X}^i = \boldsymbol{\xi}_m) = [\boldsymbol{\Lambda}_c^\top \cdot \text{diag } p(\mathbf{t}) \cdot \boldsymbol{\Lambda}_i]_{n,m} \quad (3.3)$$

where the $\boldsymbol{\Lambda}$ are again matrices with rows of the conditional probabilities $P(\mathbf{X} | t)$ which can be computed as defined in eq. (2.18). Multiplying these two matrices – one for the contralateral side from the left and one for the ipsilateral side from the right – onto a diagonal matrix containing the time prior marginalizes over the diagnose time and results in a matrix where the value in row n and column m represents the probability to find the contralateral neck in state $\mathbf{X}^c = \boldsymbol{\xi}_n$ and the ipsilateral lymphatic system in state $\mathbf{X}^i = \boldsymbol{\xi}_m$.

Similarly, we can now multiply the observation matrices \mathbf{B} from the left and the right onto $P(\mathbf{X}^i, \mathbf{X}^c)$ to compute the bilateral equivalent of ??:

$$P(\mathbf{Z}^c = \boldsymbol{\zeta}_n, \mathbf{Z}^i = \boldsymbol{\zeta}_m) = [\mathbf{B}^\top \cdot \boldsymbol{\Lambda}_c^\top \cdot \text{diag } p(\mathbf{t}) \cdot \boldsymbol{\Lambda}_i \cdot \mathbf{B}]_{n,m} \quad (3.4)$$

Formally, all necessary terms can now be computed so that both inference and the subsequent risk prediction can be performed. However, in the next section we will go into more detail regarding how this was implemented.

3.2 Parameter symmetries and mid-line extension

Although it has been omitted, eqs. (3.1) to (3.4) are still functions of the same parameters as in the unilateral model, but each side now has their own set $\boldsymbol{\theta}^c$ and $\boldsymbol{\theta}^i$ of spread probabilities that are used to parameterize the transition matrices \mathbf{A}^c and \mathbf{A}^i respectively.

In principle, the spread probabilities of the two sides are entirely indepent, and a lateralized primary tumor certainly spreads to a different extend to the ipsi- versus the contralateral side. But the spread probabilities among the LNLs should be equal when assuming that the lymphatic network in the head and neck region is symmetric. This means

$$\begin{aligned} \tilde{b}_v^c &\neq \tilde{b}_v^i \\ \tilde{t}_{rv}^c &= \tilde{t}_{rv}^i \end{aligned} \quad \forall v \leq V, r \in \text{pa}(v) \quad (3.5)$$

Due to the reasonable assumption of a symmetric neck anatomy, we may avoid doubling the spread parameters when we model the bilateral lymphatic spread.

However, there are cases in which the primary tumor lies almost or exactly on the mid-sagittal plane of the patient. In such cases, we cannot reasonably distinguish between the ipsi- and contralateral side. Consequently, we must assume the base probability rates as well to be symmetric: $\tilde{b}_v^c = \tilde{b}_v^i$. This means there must be a continuous increase in the spread probabilities from the primary tumor to the contralateral LNLs if we were to move a patient's tumor from a clearly lateralized location closer and closer to that patient's mid-sagittal plane. Ideally, we would like to factor information about the tumor's "degree of asymmetry" into our model, e.g. by considering a normalized perpendicular distance from the mid-sagittal plane to the tumor's center of mass or by considering the tumor volume on either side of this plane. Data like this, however, is rarely available. What is frequently reported and also clinically considered as a risk factor for contralateral involvement is whether or not the tumor touches or extends over the mid-sagittal plane. With this binary variable (and the information on whether the tumor is central/symmetric w.r.t. to the sagittal symmetry plane) we can now distinguish three degrees of lateralizations:

1. \not{s}, \not{e} : The tumor does not cross or touch the mid-sagittal plane and is thus clearly lateralized. The base spread probabilities are $\{\tilde{b}_v^c\}$ and $\{\tilde{b}_v^{\not{c}, \not{e}}\}$.
2. \not{s}, e : The tumor is lateralized, but crosses or touches the mid-sagittal plane. We will discuss how to define the spread probabilities to the contralateral side below.
3. s, e : The tumor is symmetric w.r.t. to the sagittal plane, thus $\tilde{b}_v^{c,s} = \tilde{b}_v^i$

Note that $s(\not{s})$ and $e(\not{e})$ denote the two binary variables *symmetric* (or *not symmetric*) and *extending* (or *not extending*) over the mid-sagittal plane.

We can infer that in case 2 the spread probabilities to the contralateral LNLs must be between the ones for the clearly lateralized (1) and the symmetric (3) case. Hence, we introduce a new "mixing" parameter α that defines the contralateral spread from tumor to the LNLs as a linear superposition between the two extremes:

$$\tilde{b}_v^{c,e} = \alpha \cdot \tilde{b}_v^i + (1 - \alpha) \cdot \tilde{b}_v^{c,\not{e}} \quad (3.6)$$

This new mixing parameter must be inferred from data just like the other spread probabilities and the parametrization of the time prior.

When using the learned parameters to predict the risk of a new patient g , the set of parameters for the risk computation $\hat{\theta}_g$ is compiled from the total set of inferred parameters $\hat{\theta} = \{\tilde{b}_v^i, \tilde{b}_v^{c,\not{e}}, \alpha, \tilde{t}_{rv}, p_T\}$, depending on the risk factors the patient presents with at the time of diagnosis. As always, for $\hat{\theta}$ we have $v \leq V$, $r \in \text{pa}(v)$ and the T-stage $T \in \{1, 2, 3, 4\}$. For example, if patient g has a T1 tumor that is clearly lateralized, their $\hat{\theta}_g$ may be computed in the following way:

$$\hat{\theta}_g = \left\{ \tilde{b}_v^i, \tilde{b}_v^c = \tilde{b}_v^{c,\not{e}}, \tilde{t}_{rv}, p_1 \right\} \quad (3.7)$$

while another patient m with a T3 tumor that clearly crosses the mid-sagittal plane would have the following set of parameters used for their risk prediction:

$$\hat{\theta}_m = \left\{ \tilde{b}_v^i, \tilde{b}_v^c = \alpha \cdot \tilde{b}_v^i + (1 - \alpha) \cdot \tilde{b}_v^{c,\not{e}}, \tilde{t}_{rv}, p_3 \right\} \quad (3.8)$$

In the actual computational implementation of this model, we essentially compute three different matrices $\mathbf{\Lambda}$ which are functions of different parameters:

$$\begin{aligned}\mathbf{\Lambda}_{\mathbf{i}} &= \mathbf{\Lambda} \left(\tilde{b}_v^{\mathbf{i}}, \tilde{t}_{rv} \right) \\ \mathbf{\Lambda}_{\mathbf{c}, \ell} &= \mathbf{\Lambda} \left(\tilde{b}_v^{\mathbf{c}, \ell}, \tilde{t}_{rv} \right) \\ \mathbf{\Lambda}_{\mathbf{c}, \mathbf{e}} &= \mathbf{\Lambda} \left(\alpha, \tilde{b}_v^{\mathbf{c}, \ell}, \tilde{b}_v^{\mathbf{i}}, \tilde{t}_{rv} \right)\end{aligned}\tag{3.9}$$

From those, the likelihoods of all patients in the training data can be computed when used with the respective p_T – that gives rise to the corresponding $\text{diag } p(\mathbf{t})$ – as in eq. (3.4).

Chapter 4

Possible further extension

In this section we will discuss possible extensions and improvements that we considered during the development. Ultimately, we decided against implementing them into the model we presented in the previous chapters due to reasons we will layout in the following subsections.

4.1 Trinary hidden random variables

As mentioned in the beginning, the reason for developing a probabilistic model of lymphatic tumor progression in the first place is that modern imaging modalities cannot detect tumor cells directly, but only the (macroscopic) changes they exert on their region of growth, e.g. when they cause lymph nodes to swell. In contrast, a histopathological examination of a resected malignancy or biopsy sample uses various staining techniques in combination with microscopes to detect carcinoma cells directly.

That raises the question whether clinical imaging and pathological examinations can both be modelled as observations of the same hidden random variable. The current resolution of the former is magnitudes away from being able to directly detect cancer cells and yet we consider the chance of detecting a microscopic involvement to be given by the sensitivity of MRI, CT, etc., the same one even as the chance for detecting macroscopic metastases. In other words: The sensitivity for detecting lesions far smaller than the voxel resolution of imaging is zero.

To account for this issue, one could categorize lymphatic involvement a little more finely: Instead of treating the true state of a LNL as a binary random variable representing a health node and a metastatic node respectively, we could consider micro- and macroscopic involvement separately. The hidden states were then modelled as *trinary* hidden variables (see eq. (1.1) for comparison):

$$\begin{array}{ll} \text{hidden} & \mathbf{X} = (X_v) \rightarrow \{0, \mu, M\}^V \\ \text{observed} & \mathbf{Z} = (Z_v) \rightarrow \{0, 1\}^V \end{array} \quad (4.1)$$

where μ and M now respectively represent *microscopic* and *macroscopic* involvement separately.

Diagnoses can still only be binary and in order to decide on a treatment, i.e. whether to irradiate the LNL in question, we also still only care about the distinction *cancerous* – meaning there are malign cells present – vs *healthy*. However,

the trinary hidden state represents the underlying reality much more precisely: To an imaging modality the hidden states 0 and μ are healthy states, since it cannot detect microscopic disease, and the respective probabilities for true and false negative observations is governed by the specificity of the modality. To a pathologist however, it is the other way around: The hidden states μ and M appear as truly involved and the probability to correctly identify these states as involved is given by the sensitivity (which we usually assume to be one for pathology).

Now, it seems that with binary observations of an underlying trinary state we cannot expect to correctly infer the true state of the LNL. But imagine that we have data on diagnostic CTs and/or MRIs taken before a neck dissection. After the surgery, the resected LNLs are examined by a pathologist. The outcome is as follows:

$$\begin{aligned} Z_v^{\text{MRI}} &= 0 \\ Z_v^{\text{path}} &= 1 \end{aligned} \tag{4.2}$$

From the second observation we can immediately infer that the hidden state must be either μ or M , because it means tumor cells have been found in the LNL. The MRI diagnosis gives us the additional information that the probability for state μ is $P(X_v = \mu) = s_P$, i.e. the specificity. For state M it is $p(X_v = M) = 1 - s_N$, where s_N is the sensitivity. Obviously, we can infer the likely true hidden state by combining observations of different kinds of modalities: Imaging on the one hand and FNA and pathology on the other hand.

The inference above might seem unnecessary at first: Could the pathologist not simply distinguish the three hidden states? But the distinction between micro- and macroscopic metastases is actually done by the imaging modality. Any occult disease that cannot be detected by them is considered a microscopic disease. This ability depends on many factors, like presence of necrotic tissue and also characteristics of the machine used to obtain the scans, that make the distinction fuzzy. Clearly, pathologists cannot routinely think about whether the disease they see would already have been visible on a scan or not and consequently, the distinction between micro and macro is not done in pathology.

Chapter 5

Data on patterns of lymphatic progression

One critical aspect of our effort to model and predict the lymphatic tumor progression is the data we use to train the model. As previously explained, our model essentially consumes tables with rows of patients and columns involvement by LNL. Data in this relatively simple format has been extracted in the past to create studies like [4] or [21]. However, the authors then used the data to compute statistics of it – e.g. the prevalence of involvement – but stopped short of publishing that data in its raw format. From these statistics it is – with one exception [20] – usually not possible to reconstruct the correlations between the involvement of LNLs.

With almost no usable data, of course, our methodology for modelling lymphatic progression cannot be tested or applied. So, we decided to start at the University Hospital Zurich (USZ) to extract all patterns of lymphatic progression in patients with newly diagnosed oropharyngeal squamous cell carcinoma (OP-SCC) between 2013 and 2019. We then not only used that data for inference on it, but also published it freely, hoping that other researchers might find it useful and that it may even motivate them to share their data in a similar fashion in the future.

In the following sections, I will include large parts of the publication [15], in which we detailed the extraction of the dataset, its characteristics and how we made it available. It is important to note that the first authorship is shared in this publication: Jean-Marc Hoffmann, a radiation oncologist at the USZ, extracted most of the data from digital patient and imaging records. Bertrand Pouymayou, a medical physicist and postdoctoral researcher at the USZ built up a complex template for easier extraction and storage of the patient information. He also created the initial interface for viewing the data. My contribution to this work was the processing of the data, creating figures and tables for the publication, host the cohort in the form of a comma separated values (CSV) table in online repositories and, lastly, develop and deploy an online interface akin to what Bertrand Pouymayou had implemented earlier.

5.1 Detailed patient-individual reporting of lymph node involvement in oropharyngeal squamous cell carcinoma with an online interface

5.1.1 Abstract

Purpose/Objective

Whereas the prevalence of LNL involvement in HNSCC has been reported, the details of lymphatic progression patterns are insufficiently quantified. In this study, we investigate how the risk of metastases in each LNL depends on the involvement of upstream LNLs, T-category, Human Papillomavirus (HPV) status and other risk factors.

Results

We retrospectively analyzed patients with newly diagnosed OPSCC treated at a single institution, resulting in a dataset of 287 patients. For all patients, involvement of LNLs I-VII was recorded individually based on available diagnostic modalities (PET, MRI, CT, FNA) together with clinicopathological factors. To analyze the dataset, a web-based graphical user interface (GUI) was developed, which allows querying the number of patients with a certain combination of co-involved LNLs and tumor characteristics.

Results

The full dataset and GUI is part of the publication. Selected findings are: Ipsilateral level IV was involved in 27% of patients with level II and III involvement, but only in 2% of patients with level II but not III involvement. Prevalence of involvement of ipsilateral levels II, III, IV, V was 79%, 34%, 7%, 3% for early T-category patients (T1/T2) and 85%, 50%, 17%, 9% for late T-category (T3/T4), quantifying increasing involvement with T-category. Contralateral levels II, III, IV were involved in 41%, 19%, 4% and 12%, 3%, 2% for tumors with and without midline extension, respectively. T-stage dependence of LNL involvement was more pronounced in HPV negative than positive tumors, but overall involvement was similar. Ipsilateral level VII was involved in 14% and 6% of patients with primary tumors in the tonsil and the base of tongue, respectively.

Conclusions

Detailed quantification of LNL involvement in HNSCC depending on involvement of upstream LNLs and clinicopathological factors may allow for further personalization of elective clinical target volume (CTV-N) definition in the future.

5.1.2 Introduction

HNSCC spread through the lymphatic system of the neck and form metastases in regional lymph nodes. Therefore, the target volume in radiotherapy of HNSCC patients includes, in addition to the primary tumor, parts of the lymph drainage

volume [2], [10]. The nodal gross tumor volume nodal gross tumor volume (GTV-N) contains detectable macroscopic lymph node metastases, while the elective clinical target volume CTV-N contains parts of lymph drainage volume that is at risk of harboring microscopic tumor, i.e. occult metastases that are not yet visible with current imaging techniques.

GTV-N definition is primarily performed through imaging techniques (PET-CT/MRI, MRI or CT) as well as FNA. Imaging criteria for lymph node metastases include size, round rather than oval shape, central necrosis, and FDG uptake as summarized by Biau et al [2]. Goel et al. gives an overview over clinical practice in PET/CT for the management of HNSCC [6]. However, all imaging techniques have finite sensitivity and specificity [17], [12], [19], i.e. they fail to detect small metastases or may incorrectly identify suspicious lymph nodes as tumor.

For standardized reporting of the location of lymph node metastases as well as delineation of the CTV-N, the lymph drainage system of the neck is divided into anatomically defined regions called LNL [8]. CTV-N definition amounts to the decision which LNLs to include into the elective CTV-N and is based on international consensus guidelines. Such guidelines were first published by Grégoire et al in 2000 and have been updated in 2006, 2014 and 2019 [2], [8], [7], [9]. Current recommendations for the selection of lymph node levels in OPSCC can be found in Table 2 of the guidelines published in 2019 by Biau et al. [2]. Current guidelines are primarily based on the prevalence of LNL involvement for a given primary tumor location, i.e. the percentage of patients diagnosed with metastases in each level. It is recommended that the elective CTV-N includes all LNLs that are involved in 10–15% of patients or more. Patients are primarily stratified by primary tumor location. For example, tumors of the soft palate, the posterior pharyngeal wall and the base of tongue show lymph node metastases on both sides via crossing lymph vessels. For this reason, even for lateralized tumors of these localizations, bilateral neck treatment is recommended. However, the lymphatic drainage of the tonsil is mainly unilateral, therefore an ipsilateral irradiation is recommended for lateralized low T-category (T1/T2) tumors (at least up to lymph node stage N2a). Volume-reduced elective nodal irradiation has been or is being investigated in several trials [3], [16].

While the general patterns of lymph drainage in the neck is understood and prevalence of LNL involvement has been reported in the literature [7], [4], [11], [1], the details of progression patterns in OPSCC are poorly quantified. How much does the risk of level IV involvement increase depending on whether levels II and III harbors macroscopic metastases? How much does the risk of involvement increase for late versus early T-category? Are progression patterns different for HPV positive versus HPV negative tumors? Answering these questions quantitatively may allow for further personalizing CTV-N definition based on an individual patient’s clinical presentation at the time of diagnosis.

The basis for better quantification of LNL involvement are detailed datasets of HNSCC patients for whom involvement is reported per individual LNL together with tumor and patient characteristics. For example, to answer the question of how much the risk in level IV increases depending on the involvement of upstream levels II and III, it is insufficient to only report the prevalence of LNL involvement in levels II, III, and IV. Instead, the observed frequency of certain involvement combinations must be known, e.g. how often levels II, III and IV are involved

simultaneously, versus how often only the levels II and III are involved without level IV. The contributions of this work are:

- We provide a dataset of lymphatic progression patterns in 287 OPSCC patients treated at our institution in whom involvement of LNLs together with tumor characteristics are reported on a patient-individual basis.
- To visualize and explore the complex dataset, a graphical user interface is provided that allows the user to query the number of patients who were diagnosed with a specific combination of simultaneously involved LNLs and tumor characteristics.

We hope that this work provides the basis for collecting large multicenter datasets of lymphatic progression patterns, which can then inform future guidelines on further personalized CTV-N definition.

5.1.3 Material & methods

Data curation

We included patients diagnosed with OPSCC (primary diagnosis) between 2013 and 2019 and treated at the department of radiation oncology and/or head and neck surgery of the USZ. Patients with prior radiotherapy or surgery to the neck were excluded, resulting in a dataset of 287 patients. Specific subsites of oropharyngeal cancer included the base of tongue, the tonsils as well as the oropharyngeal side of the vallecula and the posterior or lateral wall of the oropharynx. Patient information consisted of the date of birth, gender, the date of the 1st histological confirmation of the tumor, the performed treatment (surgery with neck dissection prior to RT/RCHT vs. surgery only vs. definitive radio(chemo)therapy), risk factors such as nicotine abuse and HPV-status (p16 pos/neg), the TNM-classification (UICC 7th edition until 2017, 8th edition since 2017), the position of the primary tumor (left/right neck) as well as positive vs. negative mid-sagittal plane extension. Further details are described in the accompanying data-in-brief article [13].

The analysis of the lymphatic spread included levels Ia, Ib, IIa, IIb, III, IV, V, VII and was performed separately for the diagnostic imaging modalities available for a patient (FDG PET-CT, FDG PET-MRI, MRI, CT) as well as FNA and radiotherapy planning CT if available. This was performed by 2 experienced radiation oncologists by reviewing radiology and pathology reports together with the diagnostic images. Criteria for considering a lymph node as malignant followed the description in Biau et al [2] and are described in detail in the data-in-brief article [13].

5.1.4 Data base

The full dataset is available as a CSV-file via the data-in-brief article linked to this publication [13] and on GitHub at <https://github.com/rmnldwg/lydata> in a folder named 2021-usz-oropharynx.

In addition, the dataset has been archived and given a persistent identifier: <https://doi.org/10.5281/zenodo.6024778>.

5.1.5 Graphical user interface

We developed an online GUI based on the Django framework [5] and provide it to explore the dataset. It allows the user to conveniently determine the number of patients that show a particular combination of co-involved LNLs and tumor characteristics. The GUI is available at <https://2021-oropharynx.lyprox.org>; its source code under MIT license is available on GitHub at <https://github.com/rmnlwdg/lyprox>. Documentation is provided within the GUI; a video demonstrating the use of the GUI is available in the supplementary materials.

Colophon

This thesis was made in L^AT_EX 2_ε using my blood, sweat and tears.

List of Figures

List of Tables