
0.1 Comparing bilateral models

Up to this point we have largely argued that the mixing parameter makes intuitive sense because of the thought experiment, where we moved the primary tumor from a clearly lateralized position closer and closer to the mid-sagittal plane, until it was perfectly symmetric w.r.t. that plane. However, we now need to actually test whether our arguments hold. For that, we decided to compare three models:

- Model \mathcal{M}_{ag} , which is agnostic to the tumor's extension e over the mid-sagittal plane and treats the contralateral base spread in the same way for all patients.
- Model \mathcal{M}_{α} that uses the linear combination of the ipsilateral base probabilities and the contralateral ones for the patients without mid-plane extension to describe the spread for tumors which do extend over that plane.
- Model $\mathcal{M}_{\text{full}}$, going even further by defining a completely independent set of contralateral base probabilities for the patients whose tumor extends over the mid-sagittal plane.

Essentially, we now want to know which of these three models does the best job of describing the data. Intuitively, one would argue that it must be $\mathcal{M}_{\text{full}}$, but this model is also more complex than the other two. A natural choice for a metric that incorporates both the accuracy of the model and a penalty for model complexity – often also called *Occam's razor* – is the *model evidence* [1].

Model evidence and Bayes factor

In Bayesian terms, we would like to know which model \mathcal{M} has the highest probability $P(\mathcal{M} \mid \mathcal{D})$ given a dataset \mathcal{D} . This probability is given by

$$P(\mathcal{M} \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \mathcal{M}) P(\mathcal{M})}{P(\mathcal{D})} \quad (1)$$

If a priori all models we want to consider have the same probability $P(\mathcal{M})$ and we only make pairwise comparisons between models, then we can restrict ourselves to computing the *Bayes factor*:

$$K_{1v2} = \frac{P(\mathcal{M}_1 \mid \mathcal{D})}{P(\mathcal{M}_2 \mid \mathcal{D})} = \frac{P(\mathcal{D} \mid \mathcal{M}_1) P(\mathcal{M}_1)}{P(\mathcal{D} \mid \mathcal{M}_2) P(\mathcal{M}_2)} = \frac{P(\mathcal{D} \mid \mathcal{M}_1)}{P(\mathcal{D} \mid \mathcal{M}_2)} \quad (2)$$

On the right side in the above equation, we see the ratio of the two model's evidences, which are merely their respective likelihoods, marginalized over all parameters:

$$P(\mathcal{D} \mid \mathcal{M}) = \int p(\mathcal{D} \mid \theta, \mathcal{M}) p(\theta \mid \mathcal{M}) d\theta \quad (3)$$

So, if we can compute this model evidence – commonly also called *marginal likelihood* or *partition function* Z from physics – for our models \mathcal{M}_{ag} , \mathcal{M}_{α} and $\mathcal{M}_{\text{full}}$, the respective pairwise Bayes factors will indicate which of them is *most likely* to be the true one, given the observed data, in the probabilistic sense. Note that this does not mean it *is* the true data-generating model and not even that we should

believe it is the true one. But only that among the models investigated, this one is probably the best.

Harold Jeffreys gives a scale for interpreting values of the Bayes factor [7]:

| K_{1v2} | $\ln K_{1v2}$ | support for \mathcal{M}_1 |
|----------------------|---------------|---|
| $< 10^0$ | < 0 | negative evidence (supports \mathcal{M}_2) |
| 10^0 to $10^{1/2}$ | 0 to 1.15 | barely worth a mention |
| $10^{1/2}$ to 10^1 | 1.15 to 2.3 | substantial |
| 10^1 to $10^{3/2}$ | 2.3 to 3.45 | strong |
| $10^{3/2}$ to 10^2 | 3.45 to 4.6 | very strong |
| $> 10^2$ | > 4.6 | decisive |

We have also listed the natural logarithm $\ln K_{1v2}$ of the Bayes factor here, because what we will actually be doing is compute differences in the log-evidences.

Thermodynamic integration

Due to the integration over all model parameters, the quantity eq. (3) is usually impossible to calculate by brute force integration, even for models with only around a dozen parameters, as is the case for ours. Unless analytical solutions exist – which is rarely the case – it is often prohibitively expensive to compute the model evidence. For this reason, a large amount of approximation methods has been developed; [6] names only a few of those methods that can be used in the context of Markov-chain Monte Carlo (MCMC). Another method that is applicable in the context of MCMC is thermodynamic integration (TI), which is very well introduced in [1] and only roughly sketched out in this section.

The concept of TI originates from the field of statistical mechanics and can be motivated from that standpoint. And although this path is certainly more educational and might convey a deeper understanding w.r.t. thermodynamics and information theory, we will take a more direct approach by starting with what we want to compute and subtracting a 0 from it:

$$\begin{aligned}
 \ln Z &:= \ln p(\mathcal{D} \mid \mathcal{M}) = \ln \int p(\mathcal{D} \mid \theta, \mathcal{M}) p(\theta \mid \mathcal{M}) d\theta - \ln 1 \\
 &= \ln \int p(\mathcal{D} \mid \theta, \mathcal{M}) p(\theta \mid \mathcal{M}) d\theta - \underbrace{\ln \int p(\theta \mid \mathcal{M}) d\theta}_{\ln Z_0}
 \end{aligned} \tag{4}$$

Writing it as this difference between two different log-evidences $\ln Z$ and $\ln Z_0$ itself does not get us far. But if we could somehow parametrize a differentiable path between the two, then maybe the integration

$$\ln Z - \ln Z_0 = \int_0^1 \frac{d}{d\beta} \ln Z_\beta d\beta \tag{5}$$

we end up with can actually be computed. Just by inspection of eq. (4) and eq. (5), one can see that on such differentiable path could be built using what we are going to call the *power posterior* $p_\beta(\theta \mid \mathcal{D}, \mathcal{M})$:

$$\begin{aligned}
 \ln Z_\beta &= \ln \int p_\beta(\theta \mid \mathcal{D}, \mathcal{M}) d\theta \\
 &= \ln \int p(\mathcal{D} \mid \theta, \mathcal{M})^\beta p(\theta \mid \mathcal{M}) d\theta
 \end{aligned} \tag{6}$$

with the derivative

$$\begin{aligned}
\frac{d}{d\beta} \ln Z_\beta &= \int \frac{p(\mathcal{D} \mid \theta, \mathcal{M})^\beta p(\theta \mid \mathcal{M})}{Z_\beta} \ln p(\mathcal{D} \mid \theta, \mathcal{M}) d\theta \\
&= \mathbb{E} [\ln p(\mathcal{D} \mid \theta, \mathcal{M})]_{p_\beta(\theta \mid \mathcal{D}, \mathcal{M})} \\
&\approx \frac{1}{S} \sum_{i=1}^S \ln p(\mathcal{D} \mid \hat{\theta}_{\beta_i}, \mathcal{M}) = A_{\text{MC}}(\beta)
\end{aligned} \tag{7}$$

The solution to computing the evidence now lies in sight: Using MCMC, we can draw samples from the power posterior p_β and use those samples to compute the expectation over the (unmodified) likelihood. Doing this for a sufficient number of steps in the interval $[0, 1]$ and integrating over the resulting $A_{\text{MC}}(\beta)$ will then yield an approximation to the log-evidence.

$$\ln Z \approx \frac{1}{2} \sum_{j=1}^{N-1} (\beta_{j+1} - \beta_j) (A_{\text{MC}}(\beta_{j+1}) + A_{\text{MC}}(\beta_j)) \tag{8}$$

This approximation gets better with larger values for S and N . But also how the β_j are chosen is crucial for computing a good estimate: Usually, the $A_{\text{MC}}(\beta)$ – which can be seen as accuracy terms – rise steeply for increasing β close to 0, while levelling off towards $\beta = 1$. It therefore makes sense to distribute the ladder of these values unevenly. A common choice, that we employed as well, was $\beta_j = x_j^5$ where the x_j are linearly spaced within the interval $[0, 1]$. This yields a very fine resolution for the first steps and gets successively coarser towards the end of the interval.

Implementation

To compare the introduced models \mathcal{M}_{ag} , \mathcal{M}_α and $\mathcal{M}_{\text{full}}$, we performed TI with a ladder of 64 β values distributed as a power-5 series. For each of the steps in the ladder, we performed an ensemble sampling round using the `emcee` [5] Python package. The size of the ensemble – consisting of so-called walkers that allow sampling in parallel and mutually influence each other’s proposals – was chosen to be 20 times the number of dimensions of the parameter space. We set the sampling algorithm to propose new samples according to a mixture of two procedures: with 80% probability it selected a differential evolution move [9] and with 20% probability a snooker move, also based on differential evolution [3]. The reason for this choice was that in previous experiments, this combination of proposals yielded the fastest convergence of the chain. Every one of the 64 sampling rounds consisted of a burn-in phase lasting 1000 steps, followed by 250 steps of which every fifth step was kept for later analysis. This might seem like a relatively short chain, but since the change of the posterior we sampled from only changed very slightly from β_j to β_{j+1} , fewer steps are required to reach convergence.

The models were trained on the combination of two datasets: One from our institution, the University Hospital Zurich, which has been published and described in great detail in a separate publication [8]. The other was kindly provided to us by researchers of the Centre Léon Bérard in Lyon, France and was the underlying data

for one of their papers [2]. Both datasets have been published in our repository `lydata`.

Different modalities were used to obtain the diagnoses for the patients in the two datasets. For the inference process, we combined all available diagnostic modalities using sensitivity and specificity values from the literature [4] using a maximum likelihood estimate. We treated this resulting "consensus diagnosis" as if it were the ground truth, i.e., we set its sensitivity and specificity to 1 respectively. Our motivation to do so was that this allows us to compare predictions of the model with data prevalences to see which of the model exhibits more flexibility in adapting to the data. If we had directly provided the models with all available diagnostic modalities and allowed it to combine them itself, as outlined in ??, this would not have been possible. Also, in this case we are not interested in learning the exact distribution over the posterior parameters of the model, i.e. the probability rates for spread along arcs of the lymphatic graph, but rather how well the different models are able to adapt to realistic data and make use of the additional information provided via the tumor's extension over the mid-sagittal plane.

Lastly, we restricted ourselves to modelling the lymph node levels (LNLs) II, III and IV, because contralaterally we rarely observe involvement outside those levels and it drastically speeds up the inference process.

Reproducibility

Each of the three models investigated here, are available in `lynference`, where we have run the respective pipeline, pushed it as a tagged commit to GitHub and published it as a release alongside the produced data in the form of a data version control (DVC) remote.

The `README.md` file in this repository explains how one can reproduce an experiment and where to find documentation on the settings and configurations used.

- Model \mathcal{M}_{ag} : `bilateral-v1`
- Model \mathcal{M}_{α} : `midline-with-mixing-v1`
- Model $\mathcal{M}_{\text{full}}$: `midline-without-mixing-v1`

Results

First, we wanted to make sure that all three models are still able to describe the ipsilateral spread sufficiently well. We have plotted the prevalence our trained models predict in the forms of histograms against the Beta-posterior of the observed prevalence in the data (fig. 1). These plots were created by computing the respective prevalence for samples drawn during the final 250 steps at the end of the TI process of which every fifth step was discarded.

The shown differences between the model's predictions are miniscule. For late T-stages (bottom row of fig. 1) it seems as if the model that is agnostic to the tumor's extension over the mid-sagittal plane slightly overestimates the prevalence, while the other two models seem to match them better or underestimate them by a small amount. Overall the fit of all models ipsilaterally is very good and shows

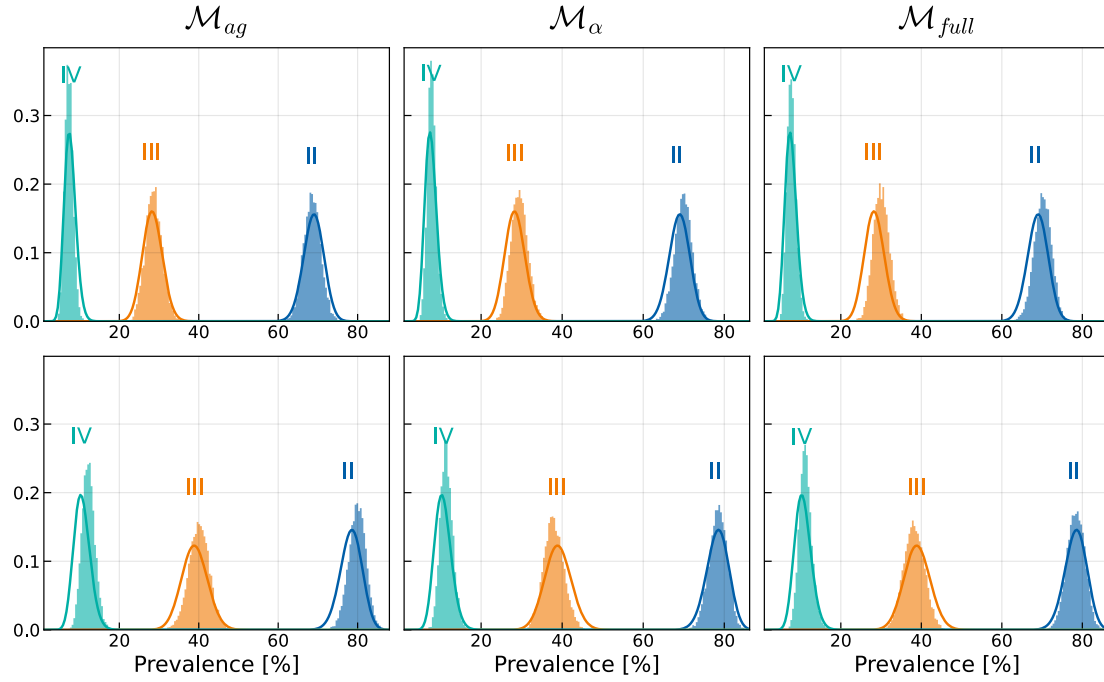


Figure 1: Predicted prevalences (shaded histograms) and posterior beta distributions of observed prevalences (solid lines) for the ipsilateral levels II (blue), III (orange) and IV (green). These prevalences have each been plotted for early T-stage patients (top row) and late T-stage (bottom row) and for the three models \mathcal{M}_{ag} (left column), \mathcal{M}_{α} (middle column) and for \mathcal{M}_{full} (right column). The differences between the models are negligible.

no indication that one model performs better than the other.

On the contralateral side, however, this does not hold anymore. Here, we do not only stratify the prevalence by T-stage, but also by midline extension. Naturally, this cannot be captured the agnostic model \mathcal{M}_{ag} since it has no method of modelling this. What is of interest to us here is how the mixing model \mathcal{M}_{α} and the full model \mathcal{M}_{full} fare against each other and whether their improvements in predicting contralateral spread are worth the additional complexity.