
Sequential Bayesian Learning: A Cognitive Neuroscience Framework

Robert T. Lange *

Einstein Center for Neurosciences Berlin
www.rob-lange.com
robert.lange17@imperial.ac.uk

Miro Grundei

Neurocomputation and Neuroimaging Unit
Free University Berlin
miro.grundei@fu-berlin.de

Sam Gijssen

Neurocomputation and Neuroimaging Unit
Free University Berlin
samgijssen@gmail.com

1 Introduction

The world is inherently sequential due to its temporal nature. It follows that our own survival relies on making sense of temporally as well as spatially ordered sensory data. Bayesian methods provide a probabilistic framework for integrating sensory information over time. Within cognitive neuroscience there has been a long-lasting tradition of theoretical frameworks which formalize this: The *Bayesian Brain Hypothesis* (BBH; Knill and Pouget [3]) postulates that the brain updates its posterior belief based on integrating its prior distribution (or previous posterior) with new likelihood evidence from the most recent sensory percept. Furthermore, the *Free Energy Principle* (FEP; Friston [2]) provides an extension to this more general hypothesis. It postulates that every self-organizing organism seeks to minimize surprise, i.e. free energy. Based on the log model evidence decomposition into the negative variational free energy and Kullback-Leibler (KL) divergence, this requires a Variational Inference (VI) scheme for approximate Bayesian model inversion. By iteratively minimizing the free energy term, one is able to obtain an ever better approximation to the log model evidence, on which we can subsequently perform Bayesian model comparison.

Still cognitive computational neuroscience lacks a general and unifying framework which allows experimentalists to easily generate sequential data, to build expressive regressors, and to analyze the quality of the neural correlates. In this piece of work we try to overcome such limitations. We introduce a general paradigm for generating feature- as well as temporally-dependent sequences of trials based on a simple graphical model. This general paradigm incorporates the roving paradigm a special case and allows for flexible design choices. Furthermore, we outline how one is able to model agents that learn the data-generating process in a sequential fashion. Based on the resulting posterior updating rules, we are able to compute different surprise measures which have recently become prominent in the cognitive computational neuroscience literature [1]. As a final contribution we outline a simple log model evidence-based model comparison scheme for the analysis of the resulting regressors.

The following document is structured as follows: First, we introduce the necessary background required in Bayesian modeling in order for us to introduce the special case of *Sequential Bayesian Learning* (SBL). Afterwards, we outline the general data-generating process and introduce agents which learn such sequences. Based on their posterior estimates and realizations one is able to calculate surprise measures. Given neural activity recordings the surprise measures can then be used

*This report got drafted during a lab rotation in the Winter of 2018/2019 and connects with an EEG study modeling mismatch-negativity in the somatosensory cortex. For further algorithmic details and replication please view <https://github.com/RobertTLange/SequentialBayesianLearning>.

as regressors in a Bayesian regression model comparison scheme. Finally, we provide an example EEG experiment which utilizes the general framework.

2 Background

Probabilistic methods allow one to elegantly formulate this form of belief updating in terms of a simple computational heuristic: Bayes' theorem. Compared to traditional frequentist statistical methods it accounts for the innate uncertainty associated with the statistical relationship between measured observations $o \in \mathbf{R}^d$ and the hidden/unobserved/latent state $s \in \mathbf{R}^d$.

At time t , the agent combines his prior over the distribution of the hidden state s_t with the likelihood of the observed state y_t :

$$p(s_t|y_t) = \frac{p(y_t, s_t)}{p(y_t)} = \frac{p(y_t|s_t)p(s_t)}{\int_{-\infty}^{\infty} p(y_t|s_t)p(s_t)dS} \propto p(y_t|s_t)p(s_t)$$

The computational procedure can often be formulated as a precision-weighted prediction error correction. The updated posterior then forms the prior distribution for time $t + 1$:

$$p(s_{t+1}) := p(s_t|y_t)$$

The beauty lies within the computational simplicity and its wide applicability.²

2.1 Bayesian Operations: Filtering, Smoothing, Decoding and Evaluation

- Filtering: $p(s_t|y_1, \dots, y_t)$
- Smoothing: $p(s_t|y_1, \dots, y_T)$
- Decoding: $\arg \max_{s_1, \dots, s_T} p(s_1, \dots, s_T|y_1, \dots, y_T)$
- Evaluation: $p(y_1, \dots, y_T)$

2.2 Surprise

$$PS(y_t) := -\ln p(y_t|s_t) = -\ln p(y_t|y_1, \dots, y_{t-1})$$

The KL divergence is an asymmetric measure of how much two probability distributions differ. For a random variable x continuous distributions it is defined as

$$KL(P||Q) := \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)}$$

$$BS(y_t) := KL(p(s_{t-1}|y_{t-1}, \dots, y_1)||p(s_t|y_t, \dots, y_1))$$

$$CS(y_t) := KL(p(s_t)||\hat{p}(s_t|y_t))$$

Rob: Define the KL divergence, talk asymmetry, what about other more symmetric measures of divergence?

²The theorem follows directly from the product and sum rule of probability.

Rob: Talk about intuition and formal background of confidence-corrected surprise

2.3 Conjugacy

Conjugacy describes a simple mathematical relationship between the prior probability distribution and the likelihood: The updated posterior distribution which combines prior and likelihood information follows the same (differently parametrized) distribution as the prior distribution.

Popular conjugate pairs include the following:

Conjugate Prior	Likelihood	Posterior
Beta	Bernoulli	Beta
Dirichlet	Categorical	Dirichlet
Univ. Gaussian/Inverse Gamma	Univ. Gaussian	Univ. Gaussian/Inverse Gamma
Multiv. Gaussian/Inverse Wishart	Multiv. Gaussian	Multiv. Gaussian/Inverse Wishart

Please view the supplementary material for a proof of the general conjugate exponential family of conjugacy.

Rob: Provide
general proof
- See Jordan
notes

2.4 Approximate Bayesian Inference

As soon as we diverge from the simple conjugate-pair case, Bayesian inference faces severe computational challenges. There is no longer a closed-form expression for the posterior distribution, since the integration needed to compute the normalizing constant becomes computationally infeasible.

Instead, one has to revert to an approximate estimation of the posterior. The two most common frameworks for doing so are Markov Chain Monte Carlo (MCMC) and Variational Inference (VI) methods.

3 A General Data Generating Process

Throughout this text we will make use of a simple and general data-generation process. We differentiate between a simple first-order Markov paradigm in which a latent state s_t follows a Markov chain and emits observation o_t which also depend on its predecessor o_{t-1} .

3.1 First-Order Markov Dependencies

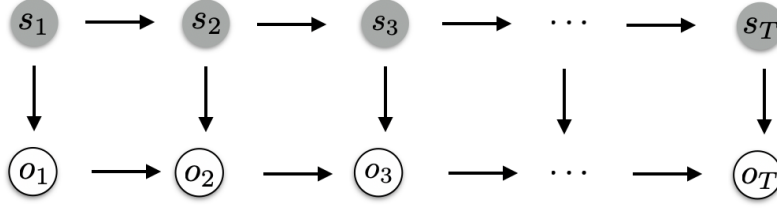


Figure 1: Graphical Model of Data-Generating Process with 1st order Markov Dependency.

$$p(s_{1:T}, o_{1:T}) = p(s_1)p(o_1|s_1) \prod_{t=2}^T p(s_t|s_{t-1})p(o_t|o_{t-1}, s_t)$$

- State space: $s \in \mathcal{S} = \{1, \dots, K\}$
- Observation space: $o \in \mathcal{O} = \{1, \dots, M\}$
- Initial state distribution: $p(s_1) = \{\frac{1}{K}, \dots, \frac{1}{K}\} \in [0, 1]^K, \sum_{j=1}^K p(s_1 = j) = 1$
- Initial obs. distribution: $p(o_1|s_1) = p(o_1) = \{\frac{1}{M}, \dots, \frac{1}{M}\} \in [0, 1]^M, \sum_{j=1}^M p(o_1 = j) = 1$
- Transitions: $p(s_t|s_{t-1}) = \mathbf{A} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$
 - $\mathbf{A}_{ij} = p(s_t = j | s_{t-1} = i); \sum_{j=1}^{|\mathcal{S}|} \mathbf{A}_{ij} = 1 \forall i = 1, \dots, |\mathcal{S}|$ and $\mathbf{A}_{ij} \geq 0 \forall i, j = 1, \dots, |\mathcal{S}|$
- Emissions: $p(o_t|o_{t-1}, s_t) = \mathbf{B} \in \mathbb{R}^{|\mathcal{O}| \times |\mathcal{O}| \times |\mathcal{S}|}$
 - $\mathbf{B}_{ijk} = p(o_t = j | o_{t-1} = i, s_t = k) := p_{ji}^{(k)}$
 - $\sum_{j=1}^{|\mathcal{O}|} \mathbf{B}_{ijk} = 1 \forall i = 1, \dots, |\mathcal{O}|, k = 1, \dots, |\mathcal{S}|$
 - $\mathbf{B}_{ijk} \geq 0 \forall i, j = 1, \dots, |\mathcal{O}|, k = 1, \dots, |\mathcal{S}|$

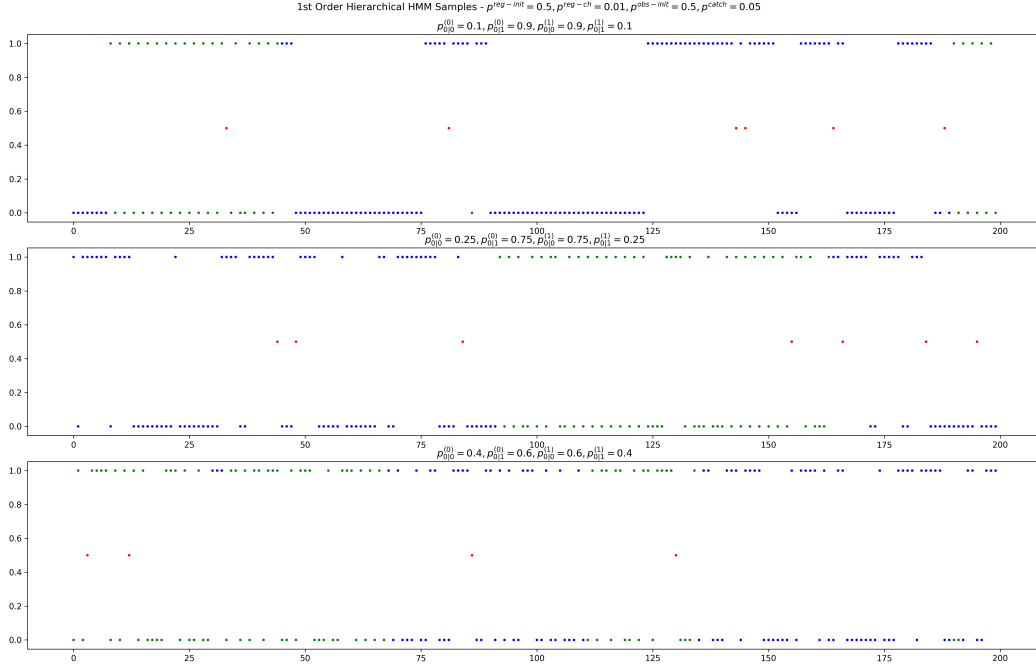


Figure 2: Example Sequences sampled from Graphical Model with first-order Markov dependency structure

3.2 Second-Order Markov Dependencies

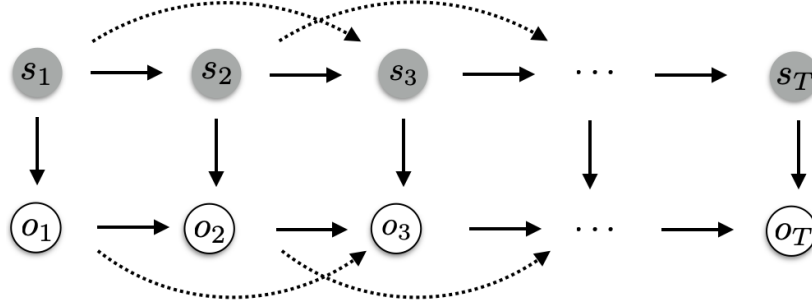


Figure 3: Graphical Model of Data-Generating Process with 2nd order Markov Dependency on the observation and hidden level.

$$p(s_{1:T}, o_{1:T}) = p(s_1)p(o_1|s_1)p(s_2|s_1)p(o_2|s_2, o_1) \prod_{t=3}^T p(s_t|s_{t-1}, s_{t-2})p(o_t|o_{t-1}, o_{t-2}, s_t)$$

- State space: $s \in \mathcal{S} = \{1, \dots, K\}$
- Observation space: $o \in \mathcal{O} = \{1, \dots, M\}$
- Initial state distribution: $p(s_1) = p(s_2|s_1) = \{\frac{1}{K}, \dots, \frac{1}{K}\} \in [0, 1]^K, \sum_{j=1}^K p(s_1 = j) = 1$
- Initial obs. distribution: $p(o_1|s_1) = (o_2|o_1, s_1) = \{\frac{1}{M}, \dots, \frac{1}{M}\} \in [0, 1]^M, \sum_{j=1}^M p(o_1 = j) = 1$
- Transitions: $p(s_t|s_{t-1}) = \mathbf{A} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{S}|}$

- $\mathbf{A}_{ijk} = p(s_t = j | s_{t-1} = i, s_{t-2} = k)$
- $\sum_{j=1}^{|\mathcal{S}|} \mathbf{A}_{ijk} = 1 \forall i, k = 1, \dots, |\mathcal{S}|$
- $\mathbf{A}_{ijk} \geq 0 \forall i, j, k = 1, \dots, |\mathcal{S}|$
- Emissions: $p(o_t | o_{t-1}, o_{t-2}, s_t) = \mathbf{B} \in \mathbb{R}^{|\mathcal{O}| \times |\mathcal{O}| \times |\mathcal{O}| \times |\mathcal{S}|}$
 - $\mathbf{B}_{ijkl} = p(o_t = j | o_{t-1} = i, o_{t-2} = k, s_t = l) := p_{j|il}^{(k)}$
 - $\sum_{j=1}^{|\mathcal{O}|} \mathbf{B}_{ijkl} = 1 \forall i = 1, \dots, |\mathcal{O}|, k = 1, \dots, |\mathcal{S}|$
 - $\mathbf{B}_{ijkl} \geq 0 \forall i, j, k = 1, \dots, |\mathcal{O}|, l = 1, \dots, |\mathcal{S}|$

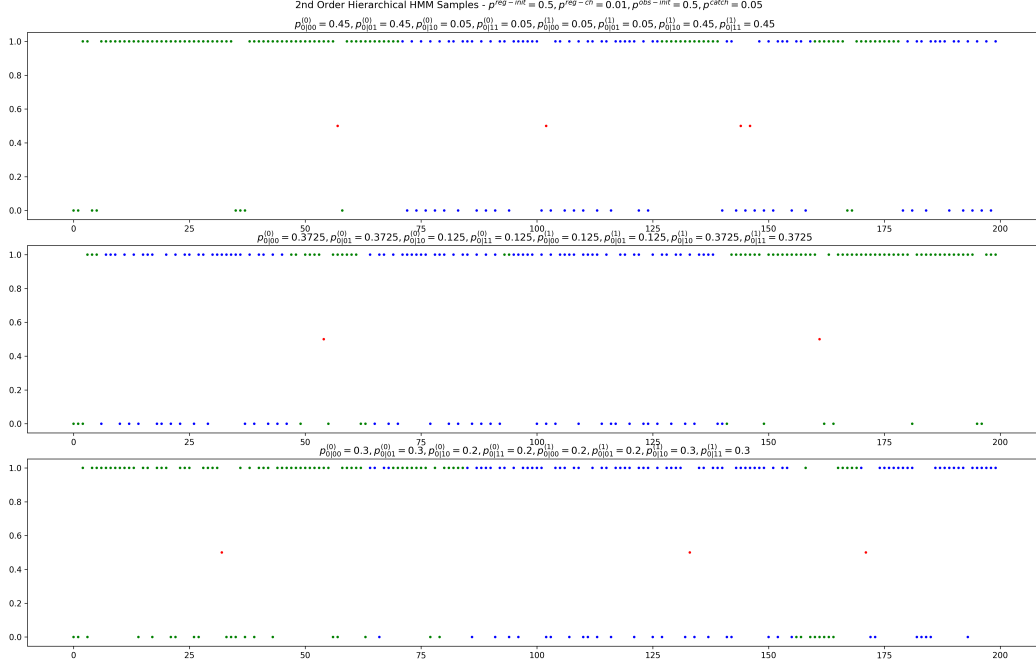


Figure 4: Example Sequences sampled from Graphical Model with second-order Markov dependency structure

4 Sequential Bayesian Agents

In the following section we will derive updating equations as well as analytically tractable expressions for the surprise measures of interest. More specifically, we introduce two different classes of agents: A conjugate Categorical-Dirichlet (CD) agent as well as a Hidden Markov Model (HMM) agent. They differ in their degree of complexity and allow for different types of hidden state tracking:

- **CD agent:** The hidden state \mathbf{s}_t is assumed to be shared across time and hence static:

$$p(\mathbf{s}_t | \mathbf{s}_{t-1}) = \delta_{\mathbf{s}_{t-1}}(\mathbf{s}_t) \Leftrightarrow \mathbf{s}_t = \mathbf{s}_{t-1} = \mathbf{s} \quad \forall t = 2, \dots, T.$$

- **HMM agent:** The hidden state s_t is a discrete variable and assumed to follow a first-order Markov chain. For a set of discrete hidden states K ($s_t \in 1, \dots, K$), the transition dynamics are given by the row-stochastic matrix $\mathbf{A} \in \mathbb{R}^{K \times K}$ with $a_{ij} \geq 0$ and $\sum_{j=1}^K a_{ij} = 1$:

$$p(s_t | s_{t-1}) = \mathbf{A} \Leftrightarrow p(s_t^j | s_{t-1}^j) = a_{ij} \quad \forall t = 1, \dots, T.$$

We differentiate between three model settings:

1. **Stimulus Probability (SP) Inference:** The model of the agent does not capture any Markov dependency. The current observation o_t only depends on the hidden state s .
2. **Alternation Probability (AP) Inference:** The model captures a limited form of first-order Markov dependency, where the probability the event of altering observations $d_t = \mathbf{1}_{o_t \neq o_{t-1}}$ is estimated given the hidden state \mathbf{s} and o_{t-1} .
3. **Transition Probability (TP) Inference:** The model accounts for full first-order Markov dependency and estimates separate alternation probabilities depending on the previous state o_{t-1} and s_t , i.e. $p(o_t | o_{t-1}, s_t)$.

Hence, the degrees of freedom and level of abstractions with which the agent is able to model the received sequence differs between the agents as well as the tracked statistic of interest.

4.1 Categorical-Dirichlet Agent

The Categorical-Dirichlet agent is part of the Bayesian conjugate pairs. It models the likelihood of the observations with the help of the Categorical distribution with $\{1, \dots, M\}$ different possible realizations per sample y_t . Given the probability vector $\mathbf{s} = \{s_1, \dots, s_M\}$ defined on the $M - 1$ dimensional simplex \mathcal{S}_{M-1} with $s_i > 0$ and $\sum_{j=1}^M s_j = 1$, the probability density function is given by

$$p(y_t = j | s_1, \dots, s_M) = s_j$$

Furthermore, the prior distribution over the hidden state \mathbf{s} is given by the Dirichlet distribution which is parametrized by the probability vector $\alpha = \{\alpha_1, \dots, \alpha_M\}$:

$$p(s_1, \dots, s_M | \alpha_1, \dots, \alpha_M) = \frac{\Gamma(\sum_{j=1}^M \alpha_j)}{\prod_{j=1}^M \Gamma(\alpha_j)} \prod_{j=1}^M s_j^{\alpha_j - 1}.$$

Hence, we have

$$\textbf{Prior: } s_1, \dots, s_M \sim \text{Dir}(\alpha_1, \dots, \alpha_M)$$

$$\textbf{Likelihood: } y \sim \text{Cat}(s_1, \dots, s_M)$$

Given a sequence of observations y_1, \dots, y_t the agent then combines the likelihood evidence with their prior beliefs in order to refine their posterior estimates over the latent variable space:

$$\begin{aligned}
p(s_1, \dots, s_M | y_1, \dots, y_t) &\propto p(s_1, \dots, s_M | \alpha_1, \dots, \alpha_M) \prod_{i=1}^t p(y_i | s_1, \dots, s_M) \\
&= \prod_{j=1}^M s_j^{\alpha_j - 1} \prod_{i=1}^t \prod_{j=1}^M s_j^{\mathbf{1}_{\{y_i=j\}}} \\
&= \prod_{j=1}^M s_j^{\alpha_j - 1 + \sum_{i=1}^t \mathbf{1}_{\{y_i=j\}}}
\end{aligned}$$

Hence, we see that the Dirichlet prior and Categorical likelihood pair follow the concept of conjugacy. More specifically, given an initial $\alpha^0 = \{\alpha_1^0, \dots, \alpha_M^0\}$ (set as a hyperparameter or optimized in an Empirical Bayes or Bayesian non-parametrics setting) we can compute the **filtering** step:

$$p(\mathbf{s}_t | y_1, \dots, y_t) = p(s_1, \dots, s_M | y_1, \dots, y_t) = \text{Dir}(\alpha^t) \text{ with } \alpha_j^t = \alpha_j^0 + \sum_{i=1}^t \mathbf{1}_{\{y_i=j\}}.$$

Furthermore, we can easily obtain the **posterior predictive** distribution needed to compute the predictive surprise measure by integrating over the space of latent states:

$$\begin{aligned}
p(y_t = x | y_1, \dots, y_{t-1}) &= \int p(y_t = x | s_1, \dots, s_M) p(s_1, \dots, s_M | y_1, \dots, y_{t-1}) d\mathcal{S}_M \\
&= \int_{s_x} \frac{\Gamma(\sum_{j=1}^M \alpha_j^t)}{\prod_{j=1}^M \Gamma(\alpha_j^t)} \prod_{j=1}^M s_j^{\alpha_j^t - 1} d\mathcal{S}_M \\
&= \frac{\Gamma(\sum_{j=1}^M \alpha_j^t)}{\prod_{j=1}^M \Gamma(\alpha_j^t)} \int \prod_{j=1}^M s_j^{\mathbf{1}_{\{x=j\}} + \alpha_j^t - 1} d\mathcal{S}_M \\
&= \frac{\Gamma(\sum_{j=1}^M \alpha_j^t)}{\prod_{j=1}^M \Gamma(\alpha_j^t)} \frac{\prod_{j=1}^M \Gamma(\mathbf{1}_{\{x=j\}} + \alpha_j^t)}{\Gamma(1 + \sum_{j=1}^M \alpha_j^t)} \\
&= \frac{\alpha_x^t}{\sum_{j=1}^M \alpha_j^t}
\end{aligned}$$

Finally, we can **evaluate** the likelihood of a specific sequence of events:

$$p(o_1, \dots, o_t) = \dots$$

We are now able to model different agents which track different statistics of the incoming sequence. We differ between three cases:

Rob: Add marginal distribution derivation

1. The stimulus probability model: $y_t = o_t \forall t = 1, \dots, T$.
2. The alternation probability model: $y_t = d_t \forall t = 2, \dots, T$
3. The transition probability model: $y_t = o_t \forall t = 1, \dots, T$ with a set of hidden parameters $\mathbf{s}^{(i)}$ for each transition from $o_{t-1} = i$.

The most simple agent models the identity of the stimulus, i.e. o_1, \dots, o_t .

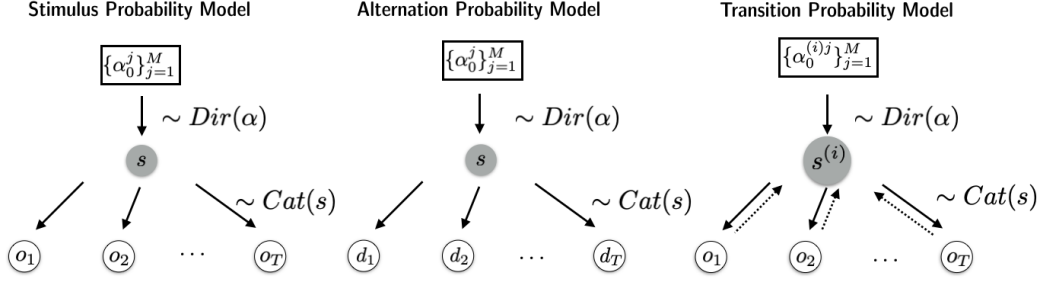


Figure 5: Categorical-Dirichlet Agent as Graphical Model. **Left.** The stimulus probability model in which the agent tracks the hidden state vector which determines the sampling process of the raw observations. **Middle.** The transition probability model in which the agent infers the hidden state distribution based on alternations of the observations. **Right.** The transition probability in which the agent assumes a different data-generating process based on the previous observations. Hence, he infers M sets of probability vectors α^i .

Exponentially weighted forgetting

$$p(\mathbf{s}_t | y_1, \dots, y_t) = p(s_1, \dots, s_M | y_1, \dots, y_t) = \text{Dir}(\alpha^t) \text{ with } \alpha_j^t = \alpha_j^0 + \sum_{i=1}^t e^{-\tau(t-i)} \mathbf{1}\{y_i = j\}.$$

- SP: $p(o_t | s_t)$, AP: $p(d_t | s_t)$, TP: $p(o_t | o_{t-1}, s_t)$
- Limited memory via exponential decay in parameter updates
- Closed-form posterior/suprise via Beta-Bernoulli conjugacy

In summary, the CD provides a simple conjugate model which assumes that the latent dynamics are static. By incorporating the exponential memory-decay parameter $\tau \in [0, 1]$ we are able to model the subject-specific cognitive capabilities. Let us now turn to the computation of surprise regressors for the Categorical-Dirichlet model.

Predictive Surprise

$$\begin{aligned} PS(y_t) &= -\ln p(y_t | y_1, \dots, y_{t-1}) \\ &= -\ln \left(\prod_{j=1}^M \left(\frac{\alpha_j^t}{\sum_{j=1}^M \alpha_j^t} \right)^{\mathbf{1}\{y_t=j\}} \right) \end{aligned}$$

Bayesian Surprise

$$BS(o_t) := KL(p(s_{t-1} | o_{t-1}, \dots, o_1) || p(s_t | o_t, \dots, o_1))$$

Confidence-Corrected Surprise

$$CS(o_t) := KL(p(s_t) || \hat{p}(s_t | o_t))$$

4.2 Hidden Markov Model Agent

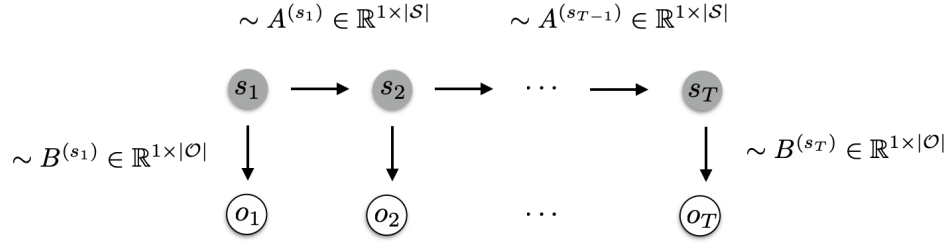


Figure 6: Hidden Markov Model Agent as Graphical Model.

Predictive Surprise

$$\begin{aligned}
 PS(y_t) &= -\ln p(y_t | y_1, \dots, y_{t-1}) \\
 &= -\ln \left(\prod_{j=1}^M \left(\frac{\alpha_x^t}{\sum_{j=1}^M \alpha_j^t} \right)^{\mathbf{1}_{\{y_t=j\}}} \right)
 \end{aligned}$$

Bayesian Surprise

$$BS(o_t) := KL(p(s_{t-1} | o_{t-1}, \dots, o_1) || p(s_t | o_t, \dots, o_1))$$

Confidence-Corrected Surprise

$$CS(o_t) := KL(p(s_t) || \hat{p}(s_t | o_t))$$

5 Bayesian Model Comparison

VI³

5.1 Variational Inference for Conjugate Exponential Family Models

5.2 Auto-Differentiation Variational Inference

³Original work on VI originates from classical spin glass models in statistical physics. In the machine learning (ML) community the negative free energy term has been better known as the evidence lower bound (ELBO). Hence, the FEP community refers to the minimization of the free energy, while the ML community, on the other hand, speaks about the maximization of the ELBO. In order to avoid confusion: The two are operationally equivalent.

6 An Example of an EEG Study

Now that we have established the theoretical framework, we are putting the above frameworks to the test.

- Catch trial (intermediate) if $o_t = 2$ - Probability: p^{catch}
- Regime change if $o_t = 3$ - Probability: $p^{reg-change}$
- Let the 2-nd order dependencies be denoted by $p_{01}^0 = p(o_t = 0 | o_{t-1} = 0, o_{t-2} = 1)$.
- Matrix is fully specified by the four probabilities p_{00}^0, p
- E.g. given p_{00}^0 , we have $p_{00}^1 = 1 - p_{00}^0 - p^{catch}/4 - p^{reg-change}/4$.

$$\mathbb{P}_{s_t} = p(o_t | o_{t-1}, o_{t-2}, s_t) = \begin{matrix} & \begin{matrix} o_t=0 & o_t=1 & o_t=2 & o_t=3 \end{matrix} \\ \begin{matrix} o_{t-1}=0, o_{t-2}=0 \\ o_{t-1}=0, o_{t-2}=1 \\ o_{t-1}=1, o_{t-2}=0 \\ o_{t-1}=1, o_{t-2}=1 \\ o_{t-1}=2, o_{t-2}=0 \\ o_{t-1}=2, o_{t-2}=1 \\ o_{t-1}=3, o_{t-2}=0 \\ o_{t-1}=3, o_{t-2}=1 \\ o_{t-1}=0, o_{t-2}=2 \\ o_{t-1}=1, o_{t-2}=2 \\ o_{t-1}=2, o_{t-2}=2 \\ o_{t-1}=3, o_{t-2}=2 \\ o_{t-1}=2, o_{t-2}=3 \\ o_{t-1}=0, o_{t-2}=3 \\ o_{t-1}=1, o_{t-2}=3 \\ o_{t-1}=3, o_{t-2}=3 \end{matrix} & \begin{pmatrix} p_{00}^0 & p_{00}^1 & p^{catch}/4 & p^{reg-change}/4 \\ p_{01}^0 & p_{01}^1 & p^{catch}/4 & p^{reg-change}/4 \\ p_{10}^0 & p_{10}^1 & p^{catch}/4 & p^{reg-change}/4 \\ p_{11}^0 & p_{11}^1 & p^{catch}/4 & p^{reg-change}/4 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

- A^{s_t} is row-stochastic.
- Marginalizing yields the following relationships 0-th order dependency:

$$p(o_t = i) = \sum_{j,k} p(o_t = i | o_{t-1} = j, o_{t-2} = k) = \sum_{j,k} p_{jk}^i$$

- Furthermore, 1-st order dependencies are the following:

$$p(o_t = i | o_{t-1}) = \begin{cases} p_{10}^i + p_{11}^i, & \text{for } o_{t-1} = 1 \\ p_{00}^i + p_{01}^i, & \text{for } o_{t-1} = 0 \end{cases}$$

- What does slow or fast changing mean in this framework (1-st or 2-nd order)?
- How to go about catch trial? Go to third order? $p(o_t | o_{t-1}, o_{t-2} = 2) = p(o_t | o_{t-1}, o_{t-3})$

7 Conclusions

Add a table which compares the different model specifications!

References

- [1] FARAJI, M., K. PREUSCHOFF, AND W. GERSTNER (2018): “Balancing new against old information: the role of puzzlement surprise in learning,” *Neural computation*, 30, 34–83.
- [2] FRISTON, K. (2010): “The free-energy principle: a unified brain theory?” *Nature reviews neuroscience*, 11, 127.
- [3] KNILL, D. C. AND A. POUGET (2004): “The Bayesian brain: the role of uncertainty in neural coding and computation,” *TRENDS in Neurosciences*, 27, 712–719.

Todo list

Rob: Define the KL divergence, talk asymmetry, what about other more symmetric measures of divergence?	2
Rob: Talk about intuition and formal background of confidence-corrected surprise	2
Rob: Provide general proof - See Jordan notes	3
Rob: Add marginal distribution derivation	8
Rob: Add log model evidence decomposition	14

Supplementary Material

Mathematical Derivation

Notes on Reproduction

Please clone the repository <https://github.com/RobertTLange/SequentialBayesianLearning> and follow the instructions outlined below:

Rob: Add log model evidence decomposition