# Maxitive Donsker–Varadhan Formulation for Possibilistic Variational Inference

**Jasraj Singh**[*]
Nanyang Technological University

**Shelvia Wongso**[*]
Nanyang Technological University

**Jeremie Houssineau**
Nanyang Technological University

**Badr-Eddine Chérief-Abdellatif**
Sorbonne Université

## Abstract

Variational inference (VI) is a cornerstone of modern Bayesian learning, enabling approximate inference in complex models that would otherwise be intractable. However, its formulation depends on expectations and divergences defined through high-dimensional integrals, often rendering analytical treatment impossible and necessitating heavy reliance on approximate learning and inference techniques. Possibility theory, an imprecise probability framework, allows to directly model epistemic uncertainty instead of leveraging subjective probabilities. While this framework provides robustness and interpretability under sparse or imprecise information, adapting VI to the possibilistic setting requires rethinking core concepts such as entropy and divergence, which presuppose additivity. In this work, we develop a principled formulation of possibilistic variational inference and apply it to a special class of exponential-family functions, highlighting parallels with their probabilistic counterparts and revealing the distinctive mathematical structures of possibility theory.

## 1 Introduction

Variational inference (VI) has become a cornerstone of modern Bayesian learning, enabling approximate inference in complex models that were previously intractable. For many years, purely Bayesian methods were hampered by computational cost, but advances in Monte Carlo and especially variational approximations have brought them within reach (Blei et al., 2017; Salimans et al., 2015). In VI, one typically posits a simple (parametric) variational family $\mathcal{Q}$ and then finds the probability distribution $q \in \mathcal{Q}$ that best approximates the true posterior $q_{\text{add}}^\star$ by maximizing the evidence lower bound (ELBO). This framework underlies many tools today – e.g. variational autoencoders (Kingma & Welling, 2013), some Bayesian neural networks (Blundell et al., 2015; Kingma et al., 2015) – and blurs the line between tractable and intractable Bayesian analyses.

Introduced by Zadeh (1978), possibility theory, an imprecise probability framework, offers a complementary approach for modelling uncertainty. Unlike probability, possibility does not require additivity: events have a degree of plausibility ($\leq 1$) and a dual *necessity*, with logical rules based on min/max operators. This allows handling epistemic uncertainty naturally. For example, the analogue of probability distributions, referred to as possibility functions, can easily be made fully uninformative (Hieu et al., 2025), a difficult endeavour in probability theory. The flexibility and scalability of possibility theory makes it suitable for challenging applications such as control problems under uncertainty (Xue et al., 2025) and point process-based inference (Houssineau, 2021).

---

[*]Equal contribution.

A key challenge arises when we try to adapt variational methods to this possibilistic setting. Standard VI relies on notions like Shannon entropy, which do not directly carry over: for example, the ELBO includes an entropy term, but there is no obvious "possibilistic entropy" analogous to the Shannon entropy of a probability distribution (Shannon & Weaver, 1949). Similarly, divergences like Kullback–Leibler (KL) require additive measures (Kullback & Leibler, 1951). As a result, traditional VI objectives cannot be used out-of-the-box. Nonetheless, recent work has sought approximations in the possibilistic context. For instance, Cella & Martin (2025) developed a variational-like approximation for inferential models (Martin & Liu, 2013), and use a Monte Carlo-based strategy to search over the chosen family of possibility functions. In that spirit, our goal is to formulate a VI framework for general possibilistic models: we define an objective that aligns a tractable candidate possibility function with the (target) posterior possibility function, and which can be optimized in lieu of maximizing a probabilistic ELBO.

We present a novel maxitive analogue of the classical Donsker–Varadhan (DV) in Section 4, with Theorem 2 providing our main theoretical result. In Section 5, we develop a practical framework for possibilistic VI within a special class of exponential-family functions, highlighting their structural parallels with the probabilistic case and giving rise to special mathematical structures. These developments provide a foundation for extending variational reasoning to possibility theory, offering new avenues for inference under epistemic uncertainty.

## 2 Primer on Possibility Theory

Possibility theory provides a dedicated representation of epistemic uncertainty where, similarly to probabilities, each event is assigned a degree of possibility between $0$ and $1$ and, unlike probabilities, the main operation on possibilities is the maximum/supremum rather than the sum/integral.

To define epistemic uncertainty formally, we consider a sample space $\Omega$ whose elements characterise all the possible values of the relevant unknown quantities. Instead of equipping $\Omega$ with a probabilistic structure, we simply describe an unknown parameter $\theta_0$ in set $\Theta$ via a (deterministic) *uncertain variable* $\boldsymbol{\theta} : \Omega \to \Theta$. If an element $\omega \in \Omega$ were the correct one, then $\boldsymbol{\theta}(\omega)$ would be true value of the parameter $\theta_0$. To *describe* the available information about $\boldsymbol{\theta}$, we define a *possibility function* $f_{\boldsymbol{\theta}}$ (a.k.a. possibility distribution) on $\Theta$ as a non-negative function that verifies $\sup_{\theta \in \Theta} f(\theta) = 1$. The possibility of an event $\boldsymbol{\theta} \in A$ for some $A \subseteq \Theta$ is then $\sup_{\theta \in A} f_{\boldsymbol{\theta}}(\theta)$.

**Marginalization and conditioning:** Possibility functions behave similar to probability mass functions (p.m.f.), except that summation is replaced by a maximum. For instance, if $\boldsymbol{\theta}$ and another uncertain variable $\boldsymbol{\psi}$ on a set $\Psi$ are jointly described by a possibility function $f_{\boldsymbol{\theta},\boldsymbol{\psi}}$, then the marginal possibility function describing $\boldsymbol{\theta}$ is characterised by

$$f_{\boldsymbol{\theta}}(\theta) = \sup_{\psi \in \Psi} f_{\boldsymbol{\theta},\boldsymbol{\psi}}(\theta, \psi), \qquad \forall \theta \in \Theta.$$

Similarly, for a fixed $\psi \in \Psi$ satisfying $f_{\boldsymbol{\psi}}(\psi) > 0$, the conditional possibility function of $\boldsymbol{\theta}$ given $\boldsymbol{\psi} = \psi$ is characterised by

$$f_{\boldsymbol{\theta}}(\theta \mid \boldsymbol{\psi} = \psi) = \frac{f_{\boldsymbol{\theta},\boldsymbol{\psi}}(\theta, \psi)}{f_{\boldsymbol{\psi}}(\psi)}, \qquad \forall \theta \in \Theta.$$

As is standard with probability distributions, we will often omit which uncertain variable is being described by a possibility function and simply write, e.g. $f$ instead of $f_{\boldsymbol{\theta}}$.

**Possibilistic moments:** Based on Hieu et al. (2025), the expected value for a possibility function $f$ is given by

$$\mathbb{E}_f^\star[\boldsymbol{\theta}] \doteq \operatorname*{argmax}_{\theta \in \Theta} f(\theta),$$

which is a set in general, and satisfies $\mathbb{E}_f^\star[T(\boldsymbol{\theta})] = T(\mathbb{E}_f^\star[\boldsymbol{\theta}])$ for any mapping $T$ – a property shared with the maximum likelihood estimate. When $\mathbb{E}_f^\star[\boldsymbol{\theta}]$ is a singleton $\{\theta^\star\}$, we do not make a distinction between this singleton and the element $\theta^\star$. In this case, if $f$ is twice differentiable at $\theta^\star$, we define the precision as $\mathcal{I}_f(\boldsymbol{\theta}) \doteq \mathbb{E}_f^\star[-\nabla^2 \log f(\boldsymbol{\theta})] = -\nabla^2 \log f(\theta)|_{\theta=\theta^\star}$, where the operator $\nabla^2$ is the Hessian and *precision* is to be understood as inverse covariance matrix.

2

**Example 1.** The normal possibility function with expected value $\mu$ and positive definite covariance matrix $\Sigma$ is defined as

$$f(\theta) = \overline{\mathrm{N}}\left(\theta; \mu, \Sigma\right) \doteq \exp\left(-\frac{1}{2}(\theta - \mu)^{\top}\Sigma^{-1}(\theta - \mu)\right).$$

It verifies $\mathbb{E}_f^{\star}[\boldsymbol{\theta}] = \mu$ and $\mathcal{I}_f(\boldsymbol{\theta}) = \Sigma^{-1}$, as is usually leveraged in the Gaussian approximation (a.k.a. Laplace approximation).

**Notations:** We define the set $\mathcal{F}(\Theta) \doteq \{f : \Theta \to [0, 1] : \sup_{\theta \in \Theta} f(\theta) = 1\}$ of possibility functions over $\Theta$. As opposed to the set $\mathcal{P}(\Theta)$ of probability distributions, $\mathcal{F}(\Theta)$ is a pre-ordered set when equipped with the partial order $\preceq$ defined as $f \preceq g \iff f(\theta) \leq g(\theta), \forall \theta \in \Theta$. In particular, any two elements $f$ and $g$ of $\mathcal{F}(\Theta)$ have a least upper bound $f \vee g$ defined for any $\theta \in \Theta$ as $(f \vee g)(\theta) = \max\{f(\theta), g(\theta)\}$. There is a greatest element in $\mathcal{F}(\Theta)$, the function equal to $1$ everywhere, which we denote by $\mathbf{1}$. For any subset $\mathcal{G}$ of $\mathcal{F}(\Theta)$, $\max \mathcal{G}$ denotes the maximal element of $\mathcal{G}$. The minimal element $\min \mathcal{G}$ can also be considered when it exists. Table 1 summarizes other key notations.

## 3   Bayesian Inference and the Donsker-Varadhan Variational Formula

We first provide a review of probabilistic variational formulation of Bayesian inference and its variational approximations. We consider the problem of learning about an unknown parameter $\theta_0$ in a set $\Theta$, which can be thought of either as the true parameter in a statistical procedure or as the solution of an optimisation problem. In this context, it is usual to alternate between i) the optimisation viewpoint where the main objects are a loss $\ell$ and a regulariser $R$ on $\Theta$, and ii) the inference viewpoint where the main objects are a likelihood $L \propto \exp(-\ell)$, that is, $L(\theta) \propto \exp(-\ell(\theta))$ for all $\theta \in \Theta$, and a prior $\pi$ with density $\propto \exp(-R)$ with respect to Lebesgue's measure. This leads to the following formula for the (generalized) Bayesian posterior $q_{\mathrm{add}}^{\star}$ with respect to $\pi$:

$$q_{\mathrm{add}}^{\star}(\mathrm{d}\theta) \doteq \frac{\exp(-\ell(\theta))\,\pi(\mathrm{d}\theta)}{\int \exp(-\ell(\theta'))\pi(\mathrm{d}\theta')} = \frac{\exp\left(-\hat{\ell}(\theta)\right)\mathrm{d}\theta}{\int \exp\left(-\hat{\ell}(\theta')\right)\mathrm{d}\theta'},$$

where $\hat{\ell} \doteq \ell + R$ is the regularised loss.

We start with the following classical variational formula, known since at least the 1950s (see e.g. Kullback (1959, Exercise 8.28) for the finite case), generally attributed in the general setting to Donsker & Varadhan (1976), and later rediscovered in statistics by Zellner (1988). We refer the interested reader to Catoni (2004, Page 160) for a proof.

**Theorem 1** (Donsker and Varadhan's variational formula). Let $(\Theta, \mathcal{T})$ be a measurable space and let $\nu$ be a probability measure on $\Theta$. For any measurable function $h : \Theta \to \mathbb{R}$ such that $\int e^h \mathrm{d}\nu < +\infty$, we have

$$\log \int e^h \mathrm{d}\nu = \sup_{\rho \in \mathcal{P}(\Theta)} \left\{\int h \mathrm{d}\rho - \mathrm{KL}\left(\rho \| \nu\right)\right\},$$

where $\mathcal{P}(\Theta)$ denotes the set of probability measures on $\Theta$ and $\mathrm{KL}(\rho\|\nu)$ is the Kullback–Leibler divergence of $\rho$ with respect to $\nu$, with the convention $\infty - \infty = -\infty$.

Moreover, the supremum on the right-hand side is achieved exactly at the Gibbs measure $\nu_h$ whose density with respect to $\nu$ is

$$\frac{\mathrm{d}\nu_h}{\mathrm{d}\nu}(\theta) \doteq \frac{e^{h(\theta)}}{\int e^h \mathrm{d}\nu}.$$

**(Generalized) Bayesian inference via Donsker-Varadhan.** Theorem 1 provides an information-theoretic characterization of the Bayesian posterior. By setting $h = -\ell$ and $\nu = \pi$ (the prior), we obtain that the normalizing constant $Z_{\mathrm{add}} = \int \exp(-\ell)\,\mathrm{d}\pi$ (also called the *evidence* in the statistical setting) satisfies

$$\log Z_{\mathrm{add}} = \sup_{q \in \mathcal{P}(\Theta)} \mathrm{ELBO}(q), \qquad \mathrm{ELBO}(q) \doteq -\mathbb{E}_q[\ell(\boldsymbol{\theta})] + \mathrm{KL}(q\|\pi).$$

3

The ELBO (*Evidence Lower Bound*) provides a lower bound on the log-evidence. Furthermore, the exact posterior can then be expressed as the solution of the infinite-dimensional optimization problem

$$q_{\text{add}}^{\star} = \underset{q \in \mathcal{P}(\Theta)}{\operatorname{argmin}} \left\{ \mathbb{E}_{\theta \sim q}[\ell(\theta)] + \text{KL}(q \| \pi) \right\}. \tag{1}$$

Thus, Donsker and Varadhan's lemma simultaneously identifies the Bayesian posterior as the optimizer of a regularized expected loss (the negative ELBO) and expresses the log-evidence as the supremum of the ELBO over all probability measures.

**(Generalized) Variational inference via Donsker-Varadhan.** In practice, computing $q_{\text{add}}^{\star}$ exactly is generally intractable. *Generalized Variational Inference* (GVI) consists in restricting the optimization to a tractable family $\mathcal{Q} \subset \mathcal{P}(\Theta)$, turning the infinite-dimensional problem (1) into a simpler, typically parametric, optimization problem.

This restriction naturally induces a gap between the true log-evidence and the ELBO achievable within $\mathcal{Q}$, which can be quantified via

$$\log Z_{\text{add}} = \text{ELBO}(q) + \text{KL}(q \| q_{\text{add}}^{\star}), \qquad q \in \mathcal{Q}.$$

Maximizing the ELBO within $\mathcal{Q}$ therefore corresponds exactly to minimizing the KL divergence to the exact posterior restricted to the chosen family. In this sense, GVI can be interpreted both as the best tractable approximation (in the reverse KL sense) of the exact posterior within $\mathcal{Q}$, and as the optimal information-processing rule within $\mathcal{Q}$.

## 4 Possibilistic Inference and a Maxitive Donsker–Varadhan Analogue

In this section, we establish a maxitive analogue of the Donsker–Varadhan variational formula and use it to derive a possibilistic inference framework. We first observe that many standard regularizers naturally attain their minimum at 0, so that the prior $\pi = \exp(-R)$ is directly a possibility function whereas $\pi$ could be improper as a probability distribution if $\exp(-R)$ is not integrable. This leads to the following formula for the (generalized) posterior $g_{\text{max}}^{\star}$ in possibility theory:

$$g_{\text{max}}^{\star}(\theta) \doteq \frac{\exp(-\ell(\theta))\pi(\theta)}{\sup_{\theta'} \exp(-\ell(\theta'))\pi(\theta')} = \frac{\exp(-\hat{\ell}(\theta))}{\sup_{\theta'} \exp(-\hat{\ell}(\theta'))}.$$

The function $\pi(\theta)$ now denotes a prior possibility function over $\Theta$, typically $\pi(\theta) = \exp(-R(\theta))$. The normalizing constant in the maxitive posterior is defined as

$$Z_{\text{max}} \doteq \sup_{\theta' \in \Theta} \exp(-\ell(\theta'))\pi(\theta'),$$

and is referred to as the the maxitive marginal likelihood. Contrary to its probabilistic analogue $Z_{\text{add}}$, which measures the overall fit of the model to the data, the quantity $Z_{\text{max}}$ quantifies the *consistency* between the prior and the likelihood. This notion is central in robust inference and can be used, for instance, to detect outliers (Houssineau & Nott, 2022).

**Theorem 2** (Maxitive Donsker-Varadhan formula)**.** Let $\pi : \Theta \to [0, 1]$ be a possibility function. For any function $\ell : \Theta \to \mathbb{R}_+$, we have

$$\log \sup_{\theta \in \Theta} e^{-\ell(\theta)}\pi(\theta) = \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ -\ell(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}, \tag{2a}$$

$$= \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ -\ell(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}, \tag{2b}$$

where $\mathcal{F}(\Theta)$ denotes the set of all possibility functions on $\Theta$, with the convention $\infty \times 0 = 0$.

Moreover, the supremum in (2a) is achieved by any possibility function lower-bounding the posterior $g_{\text{max}}^{\star}$, that is

$$\underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmax}} \inf_{\theta \in \Theta} \left\{ -\ell(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\} = \left\{ g \in \mathcal{F}(\Theta) : g \leq g_{\text{max}}^{\star} \right\}. \tag{3}$$

Similarly, the infimum in (2b) is achieved by any possibility function upper-bounding $g_{\max}^\star$:

$$\underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmin}} \sup_{\theta \in \Theta} \left\{ -\ell(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\} = \left\{ g \in \mathcal{F}(\Theta) : g_{\max}^\star \leq g \right\}. \tag{4}$$

The proof is provided in Appendix B.1.

On the left-hand side of (2a) appears the log-consistency $\log Z_{\max}$, while the quantity optimized on the right-hand side defines a *consistency bound* (CBO). Two dual versions can be introduced:

$$\underline{\mathrm{CBO}}(g) \doteq \inf_{\theta} \left\{ -\ell(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}, \qquad \overline{\mathrm{CBO}}(g) \doteq \sup_{\theta} \left\{ -\ell(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.$$

These provide lower and upper bounds on the log-consistency, as shown in Theorem 2. For any $g \in \mathcal{F}(\Theta)$, we even have the decomposition

$$\log Z_{\max} = \underline{\mathrm{CBO}}(g) + D_{\max}(g \,\|\, g_{\max}^\star) = \overline{\mathrm{CBO}}(g) - D_{\max}(g_{\max}^\star \,\|\, g),$$

where the *max-relative entropy* between two possibility functions is defined as

$$D_{\max}(g \,\|\, f) \doteq \sup_{\theta \in \Theta} \log \frac{g(\theta)}{f(\theta)} \geq 0.$$

Maximizing (resp. minimizing) the lower (resp. upper) CBO with respect to $g$ is thus equivalent to minimizing the max-relative entropy to the Gibbs possibility function $g_{\max}^\star$. Any maximizer $g$ in (2a) satisfies $g \leq g_{\max}^\star$, whereas any minimizer $g$ in (2b) satisfies $g \geq g_{\max}^\star$. In particular, the standard posterior $g_{\max}^\star$ is always an optimizer, and pointwise bounds all other solutions (below or above, depending on the formulation).

**Canonical choice of maximizer:** Among all possibility functions $g$ that optimize the CBO, a natural canonical representative is the standard possibilistic posterior $g_{\max}^\star$. However, one must be careful with terminology: $g_{\max}^\star$ is an extremum of the CBO, which means it assigns the largest/smallest possibility degrees and is therefore the most/least *permissive* (i.e. informative) representative of the equivalence class of solutions. Choosing $g_{\max}^\star$ is natural because it is the sup/inf-envelope of all maximizers/minimizers, stable under maxitive combination, and i) *for the lower CBO*, does not introduce any extra, arbitrary information, ii) *for the upper CBO*, does not forgo any of the available information in the likelihood and/or in the prior.

**Recovering the posterior:** Possibilistic Bayesian inference can thus be understood as the optimisation of a possibilistic analogue of the ELBO: the CBO. Unlike the probabilistic case, this optimisation generally admits multiple solutions, but the standard maxitive posterior $g_{\max}^\star$ stands out as the natural choice – it is an extremum among the optimisers and serves as the canonical representative of the updated possibility distribution.

Recall that the standard Bayesian posterior is $\{q_{\mathrm{add}}^\star\} = \operatorname{argmax}_{q \in \mathcal{P}(\Theta)} \mathrm{ELBO}(q)$, while the two characterizations of the Bayesian posterior in possibility theory are:

$$g_{\max}^\star = \max\left( \underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmax}} \underline{\mathrm{CBO}}(g) \right), \tag{5a}$$

$$= \min\left( \underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmin}} \overline{\mathrm{CBO}}(g) \right). \tag{5b}$$

The additional $\max/\min$ operations in (5a) and (5b) add a layer of complexity to the recovery of the posterior possibility function. However, it holds that the posterior possibility function $g_{\max}^\star$ is the only element in

$$\underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmax}} \underline{\mathrm{CBO}}(g) \cap \underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmin}} \overline{\mathrm{CBO}}(g).$$

To exploit this, we can set a scalar $\alpha \in (0, 1)$, and express both optimisation problems over $\mathcal{F}(\Theta)$ as, e.g. an $\operatorname{argmin}$. It is then easy to see that

$$\{g_{\max}^\star\} = \underset{g \in \mathcal{F}(\Theta)}{\operatorname{argmin}} \, \alpha \overline{\mathrm{CBO}}(g) - (1 - \alpha) \underline{\mathrm{CBO}}(g), \tag{6}$$

where the first term penalises underestimation of the posterior and the second term penalises overestimation. $\alpha$ balances the degree of caution in the approximation – $\alpha \approx 1$ biases the optimization towards highly plausible but less informative estimates, and vice-versa when $\alpha \approx 0$. This achieves a similar effect as symmetrising the KL divergence.

5

**Relationship with VI:** The first two formulations, (5a) and (5b), to recover the posterior possibility function from an optimisation problem *on the set of possibility functions*, provide two different schemes for defining approximations of the posterior when we restrict the optimization to be on a subset $\mathcal{G}$ of $\mathcal{F}(\Theta)$. While (5a) is qualitatively related to standard VI and will typically underestimate the possibility when restricted to $\mathcal{G}$, (5b) provides an alternative that will instead overestimate possibilities. This latter behaviour is more in line with general statistical principles which tend to prefer pessimistic uncertainty estimates over optimistic ones. The third formulation, (6), strikes a balance between the two above extremes, with over/under-estimation being preferred as needed, depending on the value of $\alpha$.

**Relationship with Generalized VI:** In classical VI, the optimization is typically framed through the decomposition of the ELBO, which involves the reverse KL divergence, $\mathrm{KL}(q\|q_{\mathrm{add}}^\star)$. This yields a lower bound on the model evidence and has well-known properties such as the mode-seeking effect. Yet, in the probabilistic literature, many alternative divergences and bound constructions have been explored, leading to a variety of upper and lower bounds on the marginal likelihood. Notable examples include the $\chi$-divergence upper bounds (CUBO) (Dieng et al., 2017) and variational Rényi bounds (Li & Turner, 2016). In each of these cases, one can define dual optimization problems – minimisation to the left or to the right of the divergence – corresponding to lower or upper bounds providing a sandwiching of the true model evidence.

Remarkably, in the possibilistic framework, an analogous structure arises naturally, but with the max-relative entropy replacing the KL divergence. The lower and upper consistency bounds (CBOs) defined above correspond to two dual optimisation perspectives: maximizing the lower CBO or minimizing the upper CBO, which yield, respectively, under- and over-estimation of the posterior possibility degrees. This is conceptually analogous to probabilistic VI, where reverse or forward KL leads to mode-seeking versus mass-covering behaviours. In fact, one can interpret the upper CBO as a limiting case of a Rényi- or $\chi^n$-type bound when the divergence order tends to infinity, producing an extreme mass-covering (resp. "hyper mode-seeking") effect.

Finally, these dual bounds can be combined to define successive, balanced objectives, similar to approaches in probabilistic VI that mix ELBO and CUBO to control approximation error (Huggins et al., 2020). In the possibilistic setting, expression (6) provides a principled way to interpolate between the two extremes, penalising over- and under-estimation of possibility degrees in a single step. This construction thus generalises the notion of variational bounds to possibility theory, mirroring the rich family of bounds in probabilistic inference and offering a natural trade off between conservative and optimistic uncertainty estimates.

## 5  Possibilistic VI with Exponential Families

We develop a general framework for possibilistic variational inference based on exponential families by considering their conjugate priors. We begin by defining the structure of exponential families in the possibilistic setting. Assuming that $\Theta$ is a subset of $\mathbb{R}^{d_\theta}$ and $\Lambda$ is a convex subset of $\mathbb{R}^{d_\lambda}$, we follow the probabilistic approach and define an exponential family as a set $\mathcal{G} = \{g_\lambda \in \mathcal{F}(\Theta) : \lambda \in \Lambda\}$ where $g_\lambda$ is of the form

$$g_\lambda(\theta) = \exp\left(\lambda^\top T(\theta) - A(\lambda) - B(\theta)\right), \qquad \forall \theta \in \Theta,$$

for a given base measure $B$ and sufficient statistics $T$ on $\Theta$ and the corresponding log-partition function $A$ on $\Lambda$, which ensures proper normalisation, that is

$$A(\lambda) = \log \sup_{\theta \in \Theta} \exp\left(\lambda^\top T(\theta) - B(\theta)\right) = \sup_{\theta \in \Theta} \lambda^\top T(\theta) - B(\theta).$$

We start by clarifying some properties of possibilistic exponential families:

**Fact 1:** There is no general correspondence between the components $T_+$ and $B_+$ of a probabilistic (canonical) exponential family and their analogue in possibility theory.

**Fact 2:** Although possibility functions are not densities and, therefore, do not require the definition of a reference measure, the base measure of a possibilistic exponentially family need not be an indicator function for the support of the family.

**Fact 3:** The log-partition function $A$ is simpler when compared to the expression of the log-partition $A_+$ for a probabilistic exponential family, that is

$$A_+(\lambda) = \log \int \exp\left(\lambda^\top T_+(\theta) - B_+(\theta)\right)\mathrm{d}\theta,$$

due to the fact that $\sup$ and $\log$ can be exchanged when defining $A$.

These properties are illustrated in the examples below.

**Example 2.** Consider a Bernoulli-style possibility function: set $\Theta = \{0,1\}$ and define the possibility of $\boldsymbol{\theta} = 0$ as $\alpha_0$ and the possibility of $\boldsymbol{\theta} = 1$ as $\alpha_1$, with $\max\{\alpha_0, \alpha_1\} = 1$ by construction. The corresponding possibility function is $f(\theta\,|\,\alpha_0, \alpha_1) = \alpha_0^{1-\theta}\alpha_1^{\theta}$. This can be expressed in an exponential family form as

$$g_\lambda(\theta) = \exp\left(\lambda^\top T(\theta)\right),$$

with $T(\theta) = (1 - \theta, \theta)$ and $\lambda = (\log \alpha_0, \log \alpha_1) \in \Lambda$, where $\Lambda = \{(\lambda_1, \lambda_2) \in \mathbb{R}^2 : \max\{\lambda_1, \lambda_2\} = 0\}$. These components differ significantly from the ones of a Bernoulli distribution. Although the parameters $\alpha_0$ and $\alpha_1$ are not directly related, the fact that their maximum is equal to $1$ allows to express this exponential family with a single variable: the odd ratio $r = \alpha_0/\alpha_1$. Indeed, it holds that

$$\alpha_0 = (r-1)^- + 1 \qquad \text{and} \qquad \alpha_1 = (r-1)^+ + 1,$$

with $(\cdot)^-$ and $(\cdot)^+$ denoting the negative part and positive part, respectively. Despite its simplicity, the Bernoulli possibility function is useful to express, for instance, possibilities of detection in tracking problems, see e.g. Ristic et al. (2020).

**Example 3.** Consider a Poisson-style possibility function: set $\Theta = \mathbb{N}_0$ and consider a parameter $\alpha \in \mathbb{R}^+$. The corresponding possibility function is

$$f(\theta\,|\,\alpha) = \frac{\alpha^{\theta - \lfloor\alpha\rfloor}\lfloor\alpha\rfloor!}{\theta!},$$

which can be expressed in an exponential family form as

$$g_\lambda(\theta) = \exp\left(\lambda^\top \theta - A(\lambda) - B(\theta)\right).$$

with $\lambda = \log \alpha$, $B(\theta) = \log(\theta!)$, and $A(\lambda) = \lfloor e^\lambda\rfloor\lambda - \log\left(\lfloor e^\lambda\rfloor!\right)$.

**Example 4.** Consider the univariate normal possibility function with known variance $\sigma^2$, $\overline{\mathrm{N}}(\theta; \mu, \sigma^2)$, for which we have $B(\theta) = \theta^2/(2\sigma^2)$, whereas the base measure of the corresponding univariate normal distribution is

$$\frac{\theta^2}{2\sigma^2} + \frac{1}{2}\log(2\pi\sigma^2).$$

## 5.1 Bayesian Inference with Conjugate Family

We consider a likelihood $p(\cdot\,|\,\theta)$, $\theta \in \Theta$, that is part of a regular and minimal exponential family with log-partition $A$, that is

$$p(x\,|\,\theta) = \exp(\theta^\top T(x) - A(\theta) - B(x)).$$

The following proposition rephrases known results on the relationship between exponential families and the Bregman divergence (Banerjee et al., 2005) in the context of possibilistic inference.

**Proposition 1.** The posterior possibility function under the uninformative prior on $\Theta$ is of the form

$$g_\lambda(\theta) = \exp\left(\lambda^\top \theta - A^\dagger(\lambda) - A(\theta)\right) \tag{7a}$$

$$= \exp\left(-D_A(\theta\|\theta^\star(\lambda))\right), \tag{7b}$$

with $\theta^\star(\lambda) \doteq \mathbb{E}_{g_\lambda}^\star[\boldsymbol{\theta}]$ the maximum likelihood estimator (MLE), $A^\dagger$ the Legendre transform of the log-partition $A$, and $D_A(\theta\|\theta') = A(\theta) - A(\theta') - \nabla_\theta A(\theta')^\top(\theta - \theta')$ the Bregman divergence.

Identifying the Bregman divergence in $g_\lambda$ allows us to leverage its properties. For instance, it is obvious from (7b) that $\pi(\cdot\,|\,\lambda)$ is a possibility function with expected value $\theta^\star(\lambda)$. We can also use the known relationship between $D_A$ and the KL divergence to see that $-\log g_\lambda(\theta)$ is the KL divergence between the likelihood at $\theta$ and the one at the MLE. We now highlight some important properties of possibility functions of the same form as $g_\lambda$ in the following proposition.

**Proposition 2.** Let $\mathcal{G}_A(\Theta)$ be the subset of $\mathcal{F}(\Theta)$ defined as

$$\mathcal{G}_A(\Theta) \doteq \{\theta \mapsto \exp(\lambda^\top \theta - A^\dagger(\lambda) - A(\theta)) : \lambda \in \nabla_\theta A(\Theta)\},$$

where $\Lambda \doteq \nabla_\theta A(\Theta)$ is the image of $\Theta$ by $\nabla_\theta A$, then, for any $g_\lambda \in \mathcal{G}_A(\Theta)$, it holds that

1. For any $\nu \geq 0$, $g_\lambda^\nu$ is a conjugate prior for $p(\cdot|\theta)$, and $g_\lambda^\nu \in \mathcal{G}_{\nu A}(\Theta)$.
2. The dual possibility function $f_\theta : \lambda \mapsto g_\lambda(\theta)$ is in $\mathcal{G}_{A^\dagger}(\Lambda)$.

The first point in Proposition 2 is a possibilistic analogue of the result of Diaconis & Ylvisaker (1979), with the additional features that the log-partition function is simply the Legendre transform of $A$ and that the $\nu$ parameter is a discount factor rather than an arbitrary parameter. This property is intuitive: $\nu$ can be interpreted as how many observations the prior is worth in terms of information content; in particular, when $\nu = 0$, it holds that $g_\lambda^\nu = \mathbf{1}$ and the prior is uninformative. The second point is a specific property of the considered class of conjugate prior possibility functions. Although it has additional properties, the possibilistic conjugate prior is often simply the renormalised version of its probabilistic counterpart.

We have the following properties for a given $g_\lambda \in \mathcal{G}_A(\Theta)$ as a possibility function describing the uncertain variable $\boldsymbol{\theta}$ in $\Theta$:

1. The expected value $\mathbb{E}_{g_\lambda}^\star[\boldsymbol{\theta}] = \nabla_\lambda A^\dagger(\lambda)$ is assumed to be the singleton $\{\theta^\star(\lambda)\}$. This assumption is weak as exponential families often have a unique mode except when they become fully uninformative, i.e. when $g_\lambda = \mathbf{1}$.
2. The precision $\mathcal{I}_\lambda \doteq \mathcal{I}_{g_\lambda}(\boldsymbol{\theta})$ verifies

$$\mathcal{I}_\lambda = \mathbb{E}_{g_\lambda}^\star\left[-\nabla_\theta^2 \log g_\lambda(\boldsymbol{\theta})\right] = -\nabla_\theta^2 \log g_\lambda(\theta^\star(\lambda)) = \nabla_\theta^2 A(\theta^\star(\lambda)) = (\nabla_\lambda^2 A^\dagger(\lambda))^{-1},$$

where $\nabla_\theta^2 \log g_\lambda(\theta^\star(\lambda))$ stands for $\nabla_\theta^2 \log g_\lambda(\theta)|_{\theta=\theta^\star(\lambda)}$.

## 5.2 Variational Inference with Conjugate Family

We now show how to perform VI for possibilistic conjugate families, which is exact when the likelihood is in an exponential family distribution, as is commonly the case, e.g. Gaussian distribution for regression, and categorical distribution for classification. The secondary objective is to recover known optimisation techniques as pioneered by Khan & Rue (2023). We focus on the lower CBO since it is the one that is the closest to standard VI.

**Proposition 3.** Given a variational family $\mathcal{G}_A(\Theta)$, let $g_{\lambda_t}$ be the possibility function at the $t$-th step of the maximisation of CBO. Then a valid update rule based on sub-gradients is

$$\lambda_{t+1} = \lambda_t - \rho_t\big(\theta^\star(\lambda_t + \nabla_\theta \hat{\ell}(\underline{\theta}_t)) - \theta^\star(\lambda_t)\big),$$

where $\underline{\theta}_t \in \arg\max_{\theta \in \Theta} \hat{\ell}(\theta) + \lambda_t^\top \theta - A(\theta)$. An approximate explicit expression

$$\lambda_{t+1} \approx \lambda_t - \rho_t \mathcal{I}_{\lambda_t}^{-1} \nabla_\theta \hat{\ell}(\theta^\star(\lambda_t)), \tag{8}$$

follows from the approximations $\underline{\theta}_t \approx \theta^\star(\lambda_t) + \mathcal{I}_{\lambda_t}^{-1} \nabla_\theta \hat{\ell}(\underline{\theta}_t)$ and $\underline{\theta}_t \approx \theta^\star(\lambda_t)$ in the argument of $\nabla_\theta \hat{\ell}$.

**Corollary 1.** *Let $A(\theta) = \frac{1}{2}\theta^\top \Sigma \theta$ be the log-partition function of the multivariate normal distribution with known variance $\Sigma$, and let $\hat{\ell}_s(\mu)$ be the regularised loss parametrised by the standard parameter $\mu = \Sigma\theta$, then*

$$\lambda_{t+1} \approx \lambda_t - \rho_t \nabla_\mu \hat{\ell}_s(\lambda_t),$$

*which is a standard gradient descent update.*

Although similar to Khan & Rue (2023), our result leverages a different approximation: instead of using the delta-method (Dorfman, 1938), we simply assume that the gradients of the regularised loss are small.

**Example 5.** Suppose that $\exp(-\hat{\ell}_s)$ is the normal distribution $N(x; \mu, \Sigma + \Sigma_0)$ for some positive definite covariance $\Sigma_0$, then $\nabla_\lambda \hat{\ell}_s(\lambda_t) = (\Sigma + \Sigma_0)^{-1}(\lambda_t - x)$, so that the update rule becomes

$$\lambda_{t+1} \approx \lambda_t - \rho_t(\Sigma + \Sigma_0)^{-1}(\lambda_t - x).$$

This expression can be interpreted as follows: $\lambda_t$ is the current estimate for the sufficient statistics $x$, so that $\lambda_t - x$ is an estimation error term, which is scaled by the precision $(\Sigma + \Sigma_0)^{-1}$; directions in which the precision is low should receive smaller updates, and this is naturally achieved here.

**Corollary 2.** *Let* $A(\theta) = n\log(1+\exp(\theta))$ *be the log-partition function of the binomial distribution with known number of trials* $n$, *then the approximate update rule* ([8](#)) *becomes*

$$\lambda_{t+1} \approx \lambda_t - \rho_t \frac{1}{\lambda_t(1-\lambda_t/n)} \nabla_\theta \hat{\ell}\left(\log \frac{\lambda_t}{n-\lambda_t}\right).$$

*Reparameterizing the loss in terms of the standard parameter* $p = (1+\exp(-\theta))^{-1}$, *the approximate update rule is given by*

$$\lambda_{t+1} \approx \lambda_t - \frac{\rho_t}{n} \nabla_p \hat{\ell}_s(\lambda_t/n).$$

**Example 6.** Suppose that $\exp(-\hat{\ell}_s)$ is a binomial distribution with standard parameter $p \in (0,1)$ and $n$ trials, for which $x$ successes have been observed. It is useful in this case to use the parameter $\hat{p}_t = \lambda_t/n$ instead of $\lambda_t$, so we have $\hat{p}_{t+1} \approx \hat{p}_t - \rho_t/n^2 \nabla_p \hat{\ell}_s(\hat{p}_t)$. Computing $\nabla_p \hat{\ell}_s(\hat{p}_t)$ yields the update rule

$$\hat{p}_{t+1} \approx \hat{p}_t - \rho_t \frac{\hat{p}_t - x/n}{\mathbb{V}(X \mid \hat{p}_t)},$$

where $\mathbb{V}(X \mid p) = np(1-p)$ is the variance of the observation when the probability of success is $p$. This expression makes sense: the current estimate $\hat{p}_t$ is compared to the probability of success $x/n$ induced by the observation $x$, and $\mathbb{V}(X \mid \hat{p}_t)$ scales the update accordingly. Indeed, if $\hat{p}_t$ is close to $0$ (or close to $1$) then the estimation error $\hat{p}_t - x/n$ should be smaller since $x$ has lower variance.

## 6  Discussion

Possibilistic Bayesian inference is already an optimisation problem, which seems to limit the interest in the development of possibilistic VI at first sight. On the contrary, we have shown that there is a surprising amount of subtleties to be understood and of insights to be gained when exploring the properties of possibilistic VI in general as well as specifically through the lens of a possibilistic version of exponential families and conjugate priors. And indeed, even when closely following existing treatments of VI and its applications, many differences emerge, such as regularities in conjugate families and a new route leading to gradient descent.

## References

Arindam Banerjee, Srujana Merugu, Inderjit S Dhillon, and Joydeep Ghosh. Clustering with bregman divergences. *Journal of machine learning research*, 6(Oct):1705–1749, 2005.

David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.

Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1613–1622, Lille, France, 07–09 Jul 2015. PMLR.

Olivier Catoni. Statistical learning theory and stochastic optimization. saint-flour summer school on probability theory 2001 (jean picard ed.). *Lecture Notes in Mathematics. Springer*, 2:10, 2004.

Leonardo Cella and Ryan Martin. Computationally efficient variational-like approximations of possibilistic inferential models. *International Journal of Approximate Reasoning*, 186:109506, 2025.

Persi Diaconis and Donald Ylvisaker. Conjugate priors for exponential families. *The Annals of statistics*, pp. 269–281, 1979.

Adji Bousso Dieng, Dustin Tran, Rajesh Ranganath, John Paisley, and David Blei. Variational inference via $\chi$ upper bound minimization. *Advances in Neural Information Processing Systems*, 30, 2017.

Monroe David Donsker and S. R. Srinivasa Varadhan. Asymptotic evaluation of certain markov process expectations for large time. iii. *Communications on Pure and Applied Mathematics*, 28: 389–461, 1976.

Robert A Dorfman. A note on the delta-method for finding variance formulae. *Biometric Bulletin*, 1938.

Nong Minh Hieu, Jeremie Houssineau, Neil K Chada, and Emmanuel Delande. Decoupling epistemic and aleatoric uncertainties with possibility theory. In *The 28th International Conference on Artificial Intelligence and Statistics*, pp. 2899–2907. ML Research Press, 2025.

Jeremie Houssineau. A linear algorithm for multi-target tracking in the context of possibility theory. *IEEE Transactions on Signal Processing*, 69:2740–2751, 2021.

Jeremie Houssineau and David J Nott. Robust bayesian inference in complex models with possibility theory. *arXiv preprint arXiv:2204.06911*, 2022.

Jonathan Huggins, Mikolaj Kasprzak, Trevor Campbell, and Tamara Broderick. Validated variational inference via practical posterior error bounds. In *International Conference on Artificial Intelligence and Statistics*, pp. 1792–1802. PMLR, 2020.

Mohammad Emtiyaz Khan and Håvard Rue. The bayesian learning rule. *Journal of Machine Learning Research*, 24(281):1–46, 2023.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2013.

Durk P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.

Solomon Kullback. *Information Theory and Statistics*. John Wiley & Sons, 1959.

Solomon Kullback and Richard Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79 – 86, 1951.

Yingzhen Li and Richard E Turner. Rényi divergence variational inference. *Advances in neural information processing systems*, 29, 2016.

Ryan Martin and Chuanhai Liu. Inferential models: A framework for prior-free posterior probabilistic inference. *Journal of the American Statistical Association*, 108(501):301–313, 2013.

Branko Ristic, Jeremie Houssineau, and Sanjeev Arulampalam. Target tracking in the framework of possibility theory: The possibilistic bernoulli filter. *Information Fusion*, 62:81–88, 2020.

Tim Salimans, Diederik Kingma, and Max Welling. Markov chain monte carlo and variational inference: Bridging the gap. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1218–1226, Lille, France, 07–09 Jul 2015. PMLR.

Claude E. Shannon and Warren Weaver. *The Mathematical Theory of Communication*. University of Illinois Press, Urbana, IL, 1949.

Zhirun Xue, Han Cai, Jeremie Houssineau, and Jingrui Zhang. Orbit-attitude coupled control for multi-target tracking based on partition pattern search. *IEEE Transactions on Aerospace and Electronic Systems*, 2025.

Lotfi Aliasger Zadeh. Fuzzy sets as a basis for a theory of possibility. *Fuzzy Sets and Systems*, 1(1): 3–28, 1978.

Arnold Zellner. Optimal information processing and bayes's theorem. *The American Statistician*, 42(4):278–280, 1988.

# A  Notations

Table 1: Summary of key notations

| Notation | Description |
| --- | --- |
| $\mathbb{E}_f^\star$ | Possibilistic expectation (mode) of $f$ |
| $\mathcal{I}_f$ | Possibilistic precision of $f$ |
| $A^\dagger$ | Legendre transform of $A$ |
| $\hat{\ell}$ | Regularized loss: $-\log L(\theta) - \log \pi(\theta)$ |
| $\overline{\text{CBO}}$ | Upper consistency bound |
| $\underline{\text{CBO}}$ | Lower consistency bound |
| $\theta^\star(\lambda)$ | Mode / expected value of $g_\lambda$, that is $\mathbb{E}_{g_\lambda}^\star[\boldsymbol{\theta}]$ |
| $\underline{\theta}(\lambda)$ | Set of minimizers of $\underline{\text{CBO}}(g_\lambda)$ |

# B  Proofs

In this section, we present the proofs for the main theoretical results in the main text.

## B.1  Theorem 2

We prove each equation of the theorem one after the other.

*Proof of Equation* (2a). Consider some fixed possibility function $g$.

Let us first assume that there exists some $\theta_g^\star \in \operatorname{argsup}_{\theta \in \Theta} g(\theta)$. We then have $g(\theta_g^\star) = 1$, so that:

$$
\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \geq \log e^{h(\theta_g^\star)} \pi(\theta_g^\star)
$$

$$
= \log \frac{e^{h(\theta_g^\star)} \pi(\theta_g^\star)}{g(\theta_g^\star)}
$$

$$
\geq \inf_{\theta \in \Theta} \log \frac{e^{h(\theta)} \pi(\theta)}{g(\theta)}
$$

$$
= \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.
$$

Now, if $\operatorname{argsup}_{\theta \in \Theta} g(\theta) = \varnothing$, then one can still define a sequence $(\theta_{g,n}^\star)_{n=1}^\infty \in \Theta^{\mathbb{N}}$ such that for any integer $n > 0$, $g(\theta_{g,n}^\star) \geq 1 - 1/n$. We can then write:

$$
\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \geq \log e^{h(\theta_{g,n}^\star)} \pi(\theta_{g,n}^\star)
$$

$$
= \log \frac{e^{h(\theta_{g,n}^\star)} \pi(\theta_{g,n}^\star)}{1 - \frac{1}{n}} + \log \left( 1 - \frac{1}{n} \right)
$$

$$
\geq \log \frac{e^{h(\theta_{g,n}^\star)} \pi(\theta_{g,n}^\star)}{g(\theta_{g,n}^\star)} + \log \left( 1 - \frac{1}{n} \right)
$$

$$
\geq \inf_{\theta \in \Theta} \log \frac{e^{h(\theta)} \pi(\theta)}{g(\theta)} + \log \left( 1 - \frac{1}{n} \right)
$$

$$
= \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\} + \log \left( 1 - \frac{1}{n} \right),
$$

so by letting $n \to +\infty$, we have

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \geq \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.$$

Consequently, the inequality above holds for any possibility function $g \in \mathcal{F}(\Theta)$, and by taking the supremum over $g \in \mathcal{F}(\Theta)$ in the right-hand side leads to:

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \geq \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.$$

Furthermore, the choice of possibility function $g(\theta) = g_{\max}^\star(\theta) = e^{h(\theta)} \pi(\theta) / \sup_{\theta'} e^{h(\theta')} \pi(\theta')$ transforms the inequality into an equality, which finally gives Formula (2a):

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) = \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.$$

$\square$

*Proof of Equation* (2b). Any possibility function $g$ takes its values in $[0, 1]$, so that for any possibility function $g$ and any parameter $\theta \in \Theta$, we have $\log g(\theta) \leq 0$. Thus, for any possibility function $g$ and any parameter $\theta$,

$$\log e^{h(\theta)} \pi(\theta) \leq h(\theta) - \log \frac{g(\theta)}{\pi(\theta)},$$

so taking the supremum over $\theta$ in both sides leads for any $g$ to:

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \leq \sup_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\},$$

and taking the infimum over possibility functions $g$ in the right-hand side leads to

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) \leq \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\},$$

Once again, $g = g_{\max}^\star$ transforms the inequality into an equality, which finally gives Formula (2b)

$$\log \sup_{\theta \in \Theta} e^{h(\theta)} \pi(\theta) = \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log \frac{g(\theta)}{\pi(\theta)} \right\}.$$

$\square$

*Proof of Equation* (3). We have already mentioned in the proof of Equation (2a) that

$$g_{\max}^\star(\theta) \doteq \frac{e^{h(\theta)} \pi(\theta)}{\sup_{\theta' \in \Theta} e^{h(\theta')} \pi(\theta')} \in \operatorname*{argmax}_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g(\theta)}{\pi(\theta)} \right) \right\},$$

so that

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g_{\max}^\star(\theta)}{\pi(\theta)} \right) \right\} = \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g(\theta)}{\pi(\theta)} \right) \right\}.$$

Furthermore, any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g' \leq g_{\max}^\star$ satisfies by monotonicity of the logarithmic function:

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} \geq \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g_{\max}^\star(\theta)}{\pi(\theta)} \right) \right\}.$$

Hence, combining the two lines above, we get for any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g' \leq g_{\max}^\star$:

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} \geq \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log \left( \frac{g(\theta)}{\pi(\theta)} \right) \right\}.$$

Since by definition of the supremum, we have the reverse inequality:

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} \leq \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\},$$

we finally have for any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g' \leq g^\star_{\max}$:

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} = \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\}.$$

This provides half of the proof, namely

$$\{ g \in \mathcal{F}(\Theta) : g \leq g^\star_{\max} \} \subset \operatorname*{argmax}_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\}.$$

To get the equality, we still have to show that any possibility function $g$ not satisfying $g \leq g^\star_{\max}$ does not belong to the argmax. To show this, let us consider some possibility function $g' \notin \{ g \in \mathcal{F}(\Theta) : g \leq g^\star_{\max} \}$, and show that

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} < \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\}.$$

To see this, recall that since $g' \not\leq g^\star_{\max}$, there exists at least one $\theta' \in \Theta$ such that $g'(\theta') > g^\star_{\max}(\theta')$. Therefore,

$$h(\theta') - \log\left( \frac{g'(\theta')}{\pi(\theta')} \right) < h(\theta') - \log\left( \frac{g^\star_{\max}(\theta')}{\pi(\theta')} \right),$$

Using the definition of the infimum, we have

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} \leq h(\theta') - \log\left( \frac{g'(\theta')}{\pi(\theta')} \right) < h(\theta') - \log\left( \frac{g^\star_{\max}(\theta')}{\pi(\theta')} \right).$$

However, notice that by definition of $g^\star_{\max}$, the quantity in the right-hand side above does not depend on $\theta'$ since:

$$h(\theta') - \log\left( \frac{g^\star_{\max}(\theta')}{\pi(\theta')} \right) = h(\theta') - \log\left( \frac{\frac{e^{h(\theta')}\pi(\theta')}{\sup e^{h(\cdot)}\pi(\cdot)}}{\pi(\theta')} \right) = \log \sup_{\theta \in \Theta} e^{h(\theta)}\pi(\theta),$$

so combining the two lines above, we get

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} < \log \sup_{\theta \in \Theta} e^{h(\theta)}\pi(\theta),$$

By using Formula (2a), we can rewrite the quantity in the right-hand side:

$$\inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g'(\theta)}{\pi(\theta)} \right) \right\} < \sup_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\},$$

which is exactly what we wanted to show. Hence, $g'$ cannot belong to the set of maximisers. This proves that

$$\operatorname*{argmax}_{g \in \mathcal{F}(\Theta)} \inf_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\} = \{ g \in \mathcal{F}(\Theta) : g \leq g^\star_{\max} \},$$

which concludes the proof. $\qquad\square$

*Proof of Equation* (4). The proof of Equation (4) is very similar to the proof of Equation (3), and is only provided for the sake of completeness. Once again, we start from the following fact mentioned in the proof (2b):

$$g^\star_{\max}(\theta) \doteq \frac{e^{h(\theta)}\pi(\theta)}{\sup_{\theta' \in \Theta} e^{h(\theta')}\pi(\theta')} \in \operatorname*{argmin}_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left( \frac{g(\theta)}{\pi(\theta)} \right) \right\},$$

so that

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g_{\max}^{\star}(\theta)}{\pi(\theta)}\right) \right\} = \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}.$$

Furthermore, any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g_{\max}^{\star} \preceq g'$ satisfies by monotonicity of the logarithmic function:

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} \leq \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g_{\max}^{\star}(\theta)}{\pi(\theta)}\right) \right\}.$$

Hence, combining the two lines above, we get for any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g_{\max}^{\star} \preceq g'$:

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} \leq \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}.$$

Since by definition of the infimum, we have the reverse inequality:

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} \geq \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\},$$

we finally have for any possibility function $g' \in \mathcal{F}(\Theta)$ such that $g_{\max}^{\star} \preceq g'$:

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} = \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}.$$

This provides half of the proof, namely

$$\{g \in \mathcal{F}(\Theta) : g_{\max}^{\star} \preceq g\} \subset \operatorname*{argmin}_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}.$$

To get the equality, we now have to show that any possibility function $g$ not satisfying $g_{\max}^{\star} \preceq g$ does not belong to the argmin. To show this, let us consider some possibility function $g' \notin \{g \in \mathcal{F}(\Theta) : g_{\max}^{\star} \preceq g\}$, and show that

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} > \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\}.$$

To see this, recall that since $g_{\max}^{\star} \npreceq g'$, there exists at least one $\theta' \in \Theta$ such that $g'(\theta') < g_{\max}^{\star}(\theta')$. Therefore,

$$h(\theta') - \log\left(\frac{g'(\theta')}{\pi(\theta')}\right) > h(\theta') - \log\left(\frac{g_{\max}^{\star}(\theta')}{\pi(\theta')}\right),$$

Using the definition of the supremum, we now have

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} \geq h(\theta') - \log\left(\frac{g'(\theta')}{\pi(\theta')}\right) > h(\theta') - \log\left(\frac{g_{\max}^{\star}(\theta')}{\pi(\theta')}\right).$$

Again, the quantity in the right-hand side above does not depend on $\theta'$:

$$h(\theta') - \log\left(\frac{g_{\max}^{\star}(\theta')}{\pi(\theta')}\right) = h(\theta') - \log\left(\frac{\dfrac{e^{h(\theta')}\pi(\theta')}{\sup e^{h(\cdot)}\pi(\cdot)}}{\pi(\theta')}\right) = \log \sup_{\theta \in \Theta} e^{h(\theta)}\pi(\theta),$$

so combining the two lines above leads to

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} > \log \sup_{\theta \in \Theta} e^{h(\theta)}\pi(\theta),$$

Now using Formula (2b), we can rewrite the quantity in the right-hand side:

$$\sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g'(\theta)}{\pi(\theta)}\right) \right\} > \inf_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\},$$

which is precisely what we wanted to show. Hence, $g'$ cannot belong to the set of minimisers. This proves that

$$\operatorname*{argmin}_{g \in \mathcal{F}(\Theta)} \sup_{\theta \in \Theta} \left\{ h(\theta) - \log\left(\frac{g(\theta)}{\pi(\theta)}\right) \right\} = \{g \in \mathcal{F}(\Theta) : g_{\max}^{\star} \preceq g\},$$

which concludes the proof. $\qquad\square$

## B.2  Proofs of Results in Section 5

*Proof of Proposition 1.*  Given the probabilistic likelihood $p(x\,|\,\theta) = \exp(\theta^\top T(x) - A(\theta) - B(x))$, under the uninformative prior on $\Theta$, $\pi(\theta) = 1$, the posterior possibility is given by

$$\pi(\theta\,|\,x) = \frac{p(x\,|\,\theta)\pi(\theta)}{\max_{\theta'} p(x\,|\,\theta')\pi(\theta')} = \frac{p(x\,|\,\theta)}{\max_{\theta'} p(x\,|\,\theta')}$$
$$= \exp\left((\theta - \theta^\star)^\top T(x) - (A(\theta) - A(\theta^\star))\right),$$

where we denote the MLE by $\theta^\star \doteq \operatorname{argmax}_\theta p(x\,|\,\theta) = \operatorname{argmax}_\theta \theta^\top T(x) - A(\theta)$. The first-order condition for the maximum tells us that $T(x) = \nabla_\theta A(\theta^\star)$. Plugging this in the equation above, we get the result in (7b):

$$\pi(\theta\,|\,x) = \exp\left((\theta - \theta^\star)^\top \nabla_\theta A(\theta^\star) - (A(\theta) - A(\theta^\star))\right)$$
$$= \exp\left(-D_A(\theta\|\theta^\star)\right).$$

Next, we rewrite the posterior as

$$\pi(\theta\,|\,x) = \exp\left(\theta^\top T(x) - A(\theta) - \left(\theta^{\star\top} T(x) - A(\theta^\star)\right)\right)$$
$$= \exp\left(\theta^\top \lambda - A(\theta) - A'(\lambda)\right),$$

which is indeed of the same form as $g_\lambda$, with $\lambda \doteq T(x)$. From the definition of $\theta^\star$, we get the result in Equation (7a), that $A' = A^\dagger$.  □

*Proof of Proposition 2.*  Consider the likelihood $p(x\,|\,\theta) = \exp(\theta^\top T(x) - A(\theta) - B(x))$, and a choice of prior $g_{\lambda,\nu} \doteq g_\lambda^\nu$, with $g_\lambda \in \mathcal{G}_A(\Theta)$ defined as $g_\lambda(\theta) = \exp\left(\lambda^\top \theta - A^\dagger(\lambda) - A(\theta)\right)$, for some $\lambda \in \Lambda$. Denoting the maximum *a posteriori* (MAP) estimate by $\theta^\star \doteq \operatorname{argmax}_\theta p(x\,|\,\theta) g_{\lambda,\nu}(\theta)$, the posterior possibility is given by

$$g_{\lambda,\nu}(\theta\,|\,x) = \frac{p(x\,|\,\theta) g_\lambda^\nu(\theta)}{\max_{\theta'} p(x\,|\,\theta') g_\lambda^\nu(\theta')}$$
$$= \frac{\exp\left(\theta^\top (T(x) + \nu\lambda) - (\nu + 1) A(\theta)\right)}{\max_{\theta'} \exp\left(\theta'^\top (T(x) + \nu\lambda) - (\nu + 1) A(\theta')\right)}$$
$$= \exp\left(\theta^\top (T(x) + \nu\lambda) - (\nu + 1) A(\theta) - (\nu + 1) A^\dagger\left(\frac{T(x) + \nu\lambda}{\nu + 1}\right)\right)$$
$$= g_{\frac{T(x) + \nu\lambda}{\nu + 1}, \nu + 1}(\theta),$$

As a sanity check, note that with $\nu = 0$, $g$ is the uninformative prior, and we recover the posterior in Proposition 1.[2] Therefore, $g_{\lambda,\nu}$ is a valid conjugate prior for the likelihood $p(\cdot\,|\,\theta)$.

Next, we show that $g_\lambda^\nu \in \mathcal{G}_{\nu A}(\Theta)$. It holds that

$$g_\lambda^\nu(\theta) = \exp\left(\nu\lambda^\top \theta - \nu A^\dagger(\lambda) - \nu A(\theta)\right)$$
$$= \exp\left(\nu\lambda^\top \theta - (\nu A)^\dagger(\nu\lambda) - \nu A(\theta)\right),$$

where the second equality follows from the following property of the convex conjugate:

$$\nu A^\dagger(\lambda) = \nu \sup_\theta \lambda^\top \theta - A(\theta) = \sup_\theta \nu\lambda^\top \theta - \nu A(\theta) = (\nu A)^\dagger(\nu\lambda).$$

We conclude that $g_\lambda^\nu \in \mathcal{G}_{\nu A}(\Theta)$ as required. Finally, we have

$$g_\lambda(\theta) = \exp\left(\lambda^\top \theta - A^\dagger(\lambda) - A(\theta)\right)$$
$$= \exp\left(\theta^\top \lambda - A^{\dagger\dagger}(\theta) - A^\dagger(\lambda)\right),$$

so the mapping $f_\theta : \lambda \to g_\lambda(\theta)$, which can be written as

$$f_\theta : \lambda \to \exp\left(\theta^\top \lambda - A^{\dagger\dagger}(\theta) - A^\dagger(\lambda)\right)$$

is in $\mathcal{G}_{A^\dagger}(\Lambda)$.  □

---

[2] We note here that when $\nu$ is an integer, it can be understood as encoding the number of *pseudo-observations*, as also in probabilistic conjugate priors.

*Proof of Proposition 3.* Let $\underline{\theta}(\lambda) \doteq \mathrm{argmax}_{\theta \in \Theta} \, \hat{\ell}(\theta) + \lambda^\top \theta - A(\theta)$ denote the solution set of the optimisation problem within $\underline{\mathrm{CBO}}$. A sub-gradient of $\underline{\mathrm{CBO}}$ is

$$\nabla_\lambda \underline{\mathrm{CBO}}(g_{\lambda_t}) = -\nabla_\lambda \log g_{\lambda_t}(\theta)|_{\theta = \underline{\theta}_t} = \nabla_\lambda A^\dagger(\lambda_t) - \underline{\theta}_t,$$

where $\underline{\theta}_t$ is an optimiser in $\underline{\theta}(\lambda_t)$, which satisfies $\nabla_\theta \hat{\ell}(\underline{\theta}_t) + \lambda_t - \nabla_\theta A(\underline{\theta}_t) = 0$, so that

$$\underline{\theta}_t = \nabla_\lambda A^\dagger(\lambda_t + \nabla_\theta \hat{\ell}(\underline{\theta}_t)) = \theta^\star(\lambda_t + \nabla_\theta \hat{\ell}(\underline{\theta}_t)).$$

Therefore, recalling that $\underline{\mathrm{CBO}}(g_{\lambda_t})$ is maximised, a valid update rule for the parameter $\lambda$ is

$$\lambda_{t+1} = \lambda_t + \rho_t \nabla_\lambda \underline{\mathrm{CBO}}(g_{\lambda_t}) = \lambda_t - \rho_t \big[\theta^\star(\lambda_t + \nabla_\theta \hat{\ell}(\underline{\theta}_t)) - \theta^\star(\lambda_t)\big].$$

We consider the first-order Taylor approximation of $\theta^\star$:

$$\theta^\star(\lambda_t + \nabla_\theta \hat{\ell}(\underline{\theta}_t)) \approx \theta^\star(\lambda_t) + \nabla_\lambda^2 A^\dagger(\lambda_t) \nabla_\theta \hat{\ell}(\underline{\theta}_t)$$
$$\approx \theta^\star(\lambda_t) + \nabla_\lambda^2 A^\dagger(\lambda_t) \nabla_\theta \hat{\ell}(\theta^\star(\lambda_t)),$$

where a zeroth-order approximation has been made in the argument of $\nabla_\theta \hat{\ell}$. Identifying the term $\nabla^2 A^\dagger(\lambda_t)$ with the inverse of the information $\mathcal{I}_{\lambda_t}$ completes the proof of the proposition. $\qquad \square$

*Proof of Corollary 1.* We consider the log-partition function $A(\theta) = \frac{1}{2}\theta^\top \Sigma \theta$ of the normal distribution with unknown mean and known variance $\Sigma$. The convex conjugate of $A(\theta)$ is

$$A^\dagger(\lambda) = \sup_{\theta \in \Theta} \left\{\theta^\top \lambda - \tfrac{1}{2}\theta^\top \Sigma \theta\right\} = \tfrac{1}{2}\lambda^\top \Sigma^{-1}\lambda.$$

Possibility functions in $\mathcal{G}_A(\Theta)$ are of the form

$$g_\lambda(\theta) = \exp\left(\lambda^\top \theta - \tfrac{1}{2}\lambda^\top \Sigma^{-1}\lambda - \tfrac{1}{2}\theta^\top \Sigma \theta\right)$$
$$= \exp\left(-\tfrac{1}{2}(\theta - \Sigma^{-1}\lambda)^\top \Sigma(\theta - \Sigma^{-1}\lambda)\right)$$
$$= \overline{\mathrm{N}}(\theta; \Sigma^{-1}\lambda, \Sigma^{-1}).$$

A convenient property of these exponential families, as in the probabilistic case, is that $\mathcal{I}_\lambda = \Sigma$ which does not depend on $\lambda$. We also obtain from the properties of the normal possibility function that $\theta^\star(\lambda) = \Sigma^{-1}\lambda$. Therefore, the approximate update rule can be expressed as

$$\lambda_{t+1} \approx \lambda_t - \rho_t \Sigma^{-1} \nabla_\theta \hat{\ell}(\Sigma^{-1}\lambda_t).$$

Finally, noticing that

$$\nabla_\theta \hat{\ell}(\theta) = \nabla_\theta \hat{\ell}_{\mathrm{s}}(\mu(\theta)) = \nabla_\theta \mu(\theta) \nabla_\mu \hat{\ell}_{\mathrm{s}}(\mu(\theta)) = \Sigma \nabla_\mu \hat{\ell}_{\mathrm{s}}(\Sigma\theta),$$

the update rule could also be expressed as

$$\lambda_{t+1} \approx \lambda_t - \rho_t \nabla_\mu \hat{\ell}_{\mathrm{s}}(\lambda_t).$$

$\qquad \square$

*Proof of Corollary 2.* We consider the log-partition function $A(\theta) = n \log(1 + \exp(\theta))$ of the binomial distribution with known number of trials $n$. The convex conjugate of $A(\theta)$ is

$$A^\dagger(\lambda) = \sup_{\theta \in \Theta} \left\{\theta^\top \lambda - n \log(1 + \exp(\theta))\right\}$$
$$= \lambda \log \tfrac{\lambda}{n} + (n - \lambda) \log\left(1 - \tfrac{\lambda}{n}\right).$$

Subsequently, we can also obtain the following:

$$\theta^\star = \nabla_\lambda A^\dagger(\lambda) = \log \tfrac{\lambda}{n - \lambda}$$
$$\mathcal{I}_\lambda^{-1} = \nabla_\lambda^2 A^\dagger(\lambda) = \tfrac{1}{\lambda(1 - \lambda/n)}$$

Therefore, the approximate update rule can be expressed as:

$$\lambda_{t+1} \approx \lambda_t - \rho_t \frac{1}{\lambda_t(1 - \lambda_t/n)} \nabla_\theta \hat{\ell}\left( \log \frac{\lambda_t}{n - \lambda_t} \right).$$

Recall that $p(\theta) = (1 + \exp(-\theta))^{-1} = \sigma(\theta)$, where $\sigma$ is the *sigmoid* function. Chain rule gives us

$$\nabla_\theta \hat{\ell}(\theta) = \nabla_\theta p(\theta) \nabla_p \hat{\ell}_s(\sigma(\theta))$$
$$= \sigma(\theta)(1 - \sigma(\theta)) \nabla_p \hat{\ell}_s(\sigma(\theta))$$

At $\theta = \log(\lambda/(n - \lambda))$, it holds that $\sigma(\theta) = \lambda/n$. Hence, the update is given by

$$\lambda_{t+1} \approx \lambda_t - \frac{\rho_t}{n} \nabla_p \hat{\ell}_s(\lambda_t/n).$$

$\square$