

Dual Stochastic Natural Gradient Descent and convergence of interior half-space gradient approximations

Borja SÁNCHEZ-LÓPEZ
Jesús CERQUIDES

IIIA-CSIC, Campus UAB, Bellaterra, Spain

Abstract

The multinomial logistic regression (MLR) model is widely used in statistics and machine learning. Stochastic gradient descent (SGD) is the most common approach for determining the parameters of a MLR model in big data scenarios. However, SGD has slow sub-linear rates of convergence [1]. A way to improve these rates of convergence is to use manifold optimization [2]. Along this line, stochastic natural gradient descent (SNGD), proposed by Amari [3], was proven to be Fisher efficient when it converged. However, SNGD is not guaranteed to converge and it is computationally too expensive for MLR models with a large number of parameters.

Here, we propose a stochastic optimization method for MLR based on manifold optimization concepts which (i) has per-iteration computational complexity is linear in the number of parameters and (ii) can be proven to converge.

To achieve (i) we establish that the family of joint distributions for MLR is a dually flat manifold and we use that to speed up calculations. Sánchez-López and Cerquides [4] have recently introduced convergent stochastic natural gradient descent (CSNGD), a variant of SNGD whose convergence is guaranteed. To obtain (ii) our algorithm uses the fundamental idea from CSNGD, thus relying on an independent sequence to build a bounded approximation of the natural gradient. We call the resulting algorithm dual stochastic natural gradient descent (DNSGD). By generalizing a result from Suneag et al. [5], we prove that DNSGD converges. Furthermore, we prove that the computational complexity of DNSGD iterations are linear on the number of variables of the model.

Keywords: Multinomial logistic regression, Stochastic gradient descent, Natural gradient, Convergence, Riemannian manifold, Computational complexity

1 Introduction

Multinomial Logistic Regression (MLR) is a widely used tool for classification. Some relevant examples solving real-world tasks are [6] for the image classifi-

cation branch, [7] for video recommendation tasks, or numerous examples in health and life sciences for analyzing nominal qualitative response variables, to name some [8, 9, 10, 11].

The justification of MLR goes beyond practical. In statistical decision theory, it is well known that MLR for the choice probability can be derived assuming that (i) the random utilities are independent and identical distributed (i.i.d.) across alternatives and that (ii) their common distribution is a Gumbel function [12]. Recent results [13] show that the Gumbel distribution for the choice variables is not necessary and that any distribution which is asymptotically exponential in its tail is sufficient to obtain the MLR model.

Cross-entropy, also known as log-loss, is the loss function most used in MLR and it is also convenient from a statistical decision theory standpoint. Once the form of the loss function is elicited [14, 15] and the inverse link function is understood as a mapping from scores to class probabilities, the log-loss is proved to be a proper composite loss together with logistic regression models [15, 16, 17]. This provides theoretical support for the usage of the log-loss from the standpoint of statistical decision theory.

Classification algorithms predict the value of a discrete variable (class) given some other variables (features). We use \mathcal{Y} for the class variable and $\mathcal{X} \in \Omega$ for the features. We assume a finite set of classes $\mathcal{Y} \in \{1, \dots, s\}$. We are interested in computing the unknown conditional probability distributions $P(\mathcal{Y} | \mathcal{X})$. This is accomplished by optimizing the expected risk function [18]. Stochastic gradient descent (SGD), even though it was introduced in the mid-twentieth century, is the most common approach for the task because it is fast and simple. The strategy of SGD is very intuitive: Assume P is an unknown probability distribution on $\mathcal{X} \times \mathcal{Y}$, \mathbb{L} is a differentiable function from \mathbb{R}^k to \mathbb{R} , and l is a differentiable loss function such that $\mathbb{L}(\eta) = \mathbb{E}_{z \sim P} [l(\eta, z)]$. If $\eta_0 \in \mathbb{R}^k$ is an initial estimation, SGD algorithm is defined by the update equation

$$\eta_{t+1} \leftarrow \eta_t - \gamma_t \nabla l(\eta_t, s_t) \tag{1}$$

where $\nabla l(\eta_t, s_t) = \frac{1}{|s_t|} \sum_{z \in s_t} \nabla l(\eta_t, z)$ is an approximation to the gradient of $\mathbb{L}(\eta_t)$ according to a sample set $s_t \sim P$ and γ_t is a positive number encoding the learning rate. After certain regularities on the function and the learning rate, SGD converges to the minimum [19]. However, usually convergence speed suddenly drops after a moderate solution quality has been reached, making the minimum unreachable from a practical point of view. Furthermore, SGD highly depends on the learning rate parameter, which can be difficult to tune, and it is vulnerable to the plateau phenomenon and to ill-conditioning. To face this issue, many SGD variants have arisen that basically modify the direction $\nabla l(\eta_t, s_t)$ using a positive semi-definite matrix M_t . Specifically, such generalization of SGD can be described by equation

$$\eta_{t+1} \leftarrow \eta_t - \gamma_t M_t \nabla l(\eta_t, s_t) \tag{2}$$

where M_t is commonly a preconditioning matrix capturing the local curvature or related information such as the Hessian matrix in Newton's method [20] or

the inverse of the Fisher Information Matrix in Stochastic Natural Gradient Descent (SNGD) [3]. Due to the increment of the computational complexity (equation 2 defines a second order method) a trade-off between quality of curvature information and computational cost is assumed. Some widely used examples are preconditioned SGD, diagonal approximations of the Hessian [21]; Adagrad [22], Adadelata [23], RMSProp [24] or Adam [25] that use the diagonal of the covariance matrix of the gradients.

Among existing preconditioning matrix algorithms, we focus our attention to SNGD and its variants. This kind of algorithm runs over a smooth manifold \mathcal{M} of dimension n [26] equipped with a metric g defined at every $p \in \mathcal{M}$. Metric g_p is the positive-definite tensor that expresses the local metric information at p . The pair (\mathcal{M}, g) is a Riemannian manifold of dimension n [26]. It is possible to choose a system of coordinates $\eta \in U \subset \mathbb{R}^n$ – or parametrization – to refer to points in \mathcal{M} . In this case the metric information of g at η is given by an n dimensional square matrix G_η , symmetric and positive-definite, in the base derived by the parametrization.

Assume \mathcal{M} is not standard \mathbb{R}^n , for instance when \mathcal{M} is a sphere of dimension n or when it is a Statistical manifold [27], that is, manifolds whose points refer to probability distributions of the same family, where usually the Fisher information metric (FIM) is chosen. In such cases, it is interesting to work in a Riemannian manifold for two reasons. First, because it provides correct notions of angles and local lengths which yields better updates. This is important when the gradient is the key tool for an optimization algorithm, since the gradient possesses these local magnitudes. And second, it allows to correctly define a direction in the space, in opposition to what happens with SGD, where the descent direction is parametrization dependent, in the sense that the gradient depends on the selected parametrization η .

Roughly speaking, the gradient of a function f at $p \in \mathcal{M}$ is the steepest direction of f at p . In a Riemannian manifold this is called the natural gradient by [3], noted as $\tilde{\nabla}f(p)$, and it is well defined since it takes into account the metric. If a parametrization η is fixed, Amari [3] proved that

$$\tilde{\nabla}f(\eta) = (G_\eta)^{-1}\nabla f(\eta) \tag{3}$$

where $\tilde{\nabla}f(\eta)$ is the natural gradient at η in the base derived by the parametrization.

SNGD follows the natural gradient instead of the gradient. In particular, it sets the preconditioning matrix $M_t = (G_{\eta_t})^{-1}$ in update equation 2, to follow the natural gradient according to equation 3

$$\begin{aligned} \eta_{t+1} &\leftarrow \eta_t - \gamma_t \tilde{\nabla}l(\eta_t, z_t) \\ &\leftarrow \eta_t - \gamma_t (G_{\eta_t})^{-1} \nabla l(\eta_t, z_t) \end{aligned} \tag{4}$$

This algorithm, or an approximation, usually speeds up learning in many problems, avoids the plateau effect and it defines parametrization independent directions. However, it faces two main problems:

- i) its computational complexity is high due to the need of either inverting a matrix or solving a linear system, and
- ii) it does not converge in some scenarios where SGD does [28] or it needs stronger assumptions such as compactness to stabilize [29].

Issue i) warns about the higher computational complexity order, which really is a problem for nowadays large-scale high-dimensional problems, and issue ii) refers to the convergence property.

The objective of this work is to propose a natural gradient optimization method [2] for MLR, the Dual Stochastic Natural Gradient Descent (DSNGD), whose convergence is proven and whose computational complexity order equals that of SGD. Therefore this article reveals a strategy to approximate SNGD without suffering issues i) and ii).

To deal with issue i) we establish that the family of joint distributions for MLR is a dually flat manifold and we use that to speed up calculations. To overcome issue ii) our algorithm uses the fundamental idea from CSNGD [4], relying on an independent sequence to build a bounded approximation of the natural gradient.

Section 2 introduces some related work with essential concepts needed to solve issues i) and ii) for our natural gradient based algorithm. Section 3 defines DSNGD. Sections 4 and 5 face issues i) and ii) respectively for the discrete case, that is, when $\mathcal{X} = \{1, \dots, m\}$, and they prove discrete DSNGD is convergent and as fast as SGD, in terms of complexity order.

2 Related work

2.1 Dually flat manifolds

Issue i) is addressed in this paper by restricting to dually flat manifolds (DFM) [30, 31]. The computational cost of natural gradient can be significantly reduced if the ambient space is a dually flat manifold, as one can see for instance for mirror descent [32].

DFM are built after two dual connections – conjugate connections – that are flat, that is, where Riemann-Christoffel curvature vanishes. As it is proved in [30], in such manifolds, there exist two dual parametrizations η and η^* , related by the Legendre transform of a convex function $F(\eta)$, such that

$$\begin{aligned} \eta^* &= \nabla F(\eta) \\ \nabla^2 F(\eta) &= G_\eta \end{aligned} \tag{5}$$

considering that η and η^* refer to the same point. This leads to a key property

of DFM, starting by applying the chain rule to equation 3

$$\begin{aligned}
\tilde{\nabla} f(\eta) &= (G_\eta)^{-1} \nabla f(\eta) \\
&= (G_\eta)^{-1} \nabla \eta^*(\eta) \nabla f(\eta^*) \\
&= (G_\eta)^{-1} \nabla \nabla F(\eta) \nabla f(\eta^*) \\
&= (G_\eta)^{-1} G_\eta \nabla f(\eta^*) \\
&= \nabla f(\eta^*)
\end{aligned} \tag{6}$$

for any differentiable function f defined in \mathcal{M} . Equation 6 is proved for linear exponential families, a well known DFM, in [33] and [34]. Therefore, equation 6 states that the natural gradient equals the gradient in its dual parametrization. Here one deduces a strategy to compute the natural gradient, or an approximation, without paying the costs of matrix inversion or linear system solving. Summing up, the main idea that solves issue ii) is equation 6. For example, in the case of SNGD in a DFM, one equivalently writes SNGD as

$$\eta_{t+1} \leftarrow \eta_t - \gamma_t \nabla l(\eta_t^*, z_t) \tag{7}$$

Equations 4 and 7 define SNGD, but the latter avoids matrix inversions and linear system solving.

2.2 Mirror descent

DSNGD described in this article is by no means the only algorithm taking profit from dual space properties. Mirror descent algorithm [32] makes use of dual parametrizations. As proved in [34], mirror descent in a dually flat manifold is nothing else than SNGD run in the dual space.

According to [35], mirror descent follows below update rule.

$$\begin{aligned}
\eta_t &\leftarrow \nabla F^*(\eta_t^*) \\
\eta_{t+1}^* &\leftarrow \nabla F(\eta_t) - \gamma_t \nabla l(\eta_t, s_t)
\end{aligned} \tag{8}$$

where F is a convex function and F^* is the Legendre transform of F . Even though both DSNGD and mirror descent rely on duality, there are clear differences between them: (i) mirror descent is normally defined for off-line learning, (ii) DSNGD has its convergence guaranteed, and (iii) mirror descent keeps a sequence of points in the manifold expressed both in primal and dual parametrizations while DSNGD has two sequences moving in the primal and dual space which are not necessarily connected. Difference (iii) is specially relevant since it requires mirror descent to rely on the computation of duals, while DSNGD can be run in spaces where we do not even know how to efficiently compute the dual coordinates of a point when given its primal coordinates.

2.3 Multinomial logistic regression

As described above, our strategy to reduce the per-iteration computational complexity of DSNGD relies on the fact that the family of joint distributions for MLR is a dually flat manifold.

The main assumption of MLR [36] is that the log-odds ratio of the class posteriors $P(\mathcal{Y} | \mathcal{X})$ is an affine function of the features \mathcal{X} . Banerjee [36] proved (Theorem 2) that a class of distributions fulfills that assumption if and only if for each value of \mathcal{Y} , the class of conditional distributions $P(\mathcal{X} | \mathcal{Y})$ belongs to the same linear exponential family (LEF)¹. In section 3.1 we use these results to prove that the class of joint distributions $P(\mathcal{Y}, \mathcal{X})$ is also a LEF. It is well known that a LEF is a DFM [30]. Usually, finding the minimum expected risk MLR parameters is formulated as an optimization problem in \mathbb{R}^k which is solved by means of SGD. Instead, we propose to formulate the problem as a manifold optimization problem [2], over the manifold of probability distributions $P(\mathcal{Y}, \mathcal{X})$ fulfilling the main assumption of MLR. Since we will prove that this manifold is dually flat, this formulation of the problem will allow us to capture the curvature information of the manifold efficiently.

2.4 Convergent Stochastic Natural Gradient Descent

In [4], a convergent variant for SNGD, namely CSNGD, is presented. CSNGD becomes stable in every toy scenario presented, unlike SNGD which fails in those same situations. Moreover, from a practical point of view it inherits the convergence speed of SNGD. CSNGD is defined with the update rule

$$\eta_{t+1} \leftarrow \eta_t - \gamma_t (G_{\zeta_t})^{-1} \nabla l(\eta_t, z_t) \quad (9)$$

where $\{\zeta_t\}_{t \in \mathbb{N}}$ can be any convergent sequence in \mathbb{R}^k . Both SNGD and CSNGD work by progressively building a sequence $\{\eta_t\}_{t \in \mathbb{N}}$. However, CSNGD additionally maintains an independent sequence ζ_t , which is only required to be convergent. This difference allows CSNGD to converge to the unique minimum of a convex function after some reasonable conditions on the learning rate parameter. Precisely, SNGD does not converge due to the inverse matrix $(G_{\eta_t})^{-1}$ in equation 4, since eigenvalues of that matrix are unbounded. CSNGD forces convergence and eigenvalue confinement of sequence $\{(G_{\zeta_t})^{-1}\}_{t \in \mathbb{N}}$ because of the convergence of sequence $\{\zeta_t\}_{t \in \mathbb{N}}$ and continuity property, stabilizing the algorithm. This idea is used in section 3 to define DSNGD and it allows us to prove its convergence in section 5.

2.5 Convergence of interior half-space gradient approximations

In [5] Sunehag et al. provide a variable metric stochastic approximation theory. One of the key results in that paper is Theorem 3.2 which proves convergence

¹For a definition of linear exponential family see [37]

given that we take a scaling matrix B_t at step t of the algorithm, provided that the spectrum of their (possibly non-convergent) scaling matrices is uniformly bounded from above by a finite constant and from below by a strictly positive constant. Moreover it assumes the step direction at iteration t is some Y_t modified by the scaling matrix. Vector Y_t is drawn from a family of random variables Y defined for all η , and $Y_t = Y(\eta_t)$. The result is stated here.

Theorem 1 (Theorem 3.2 in [5]). *Let $\mathbb{L} : \mathbb{R}^k \rightarrow \mathbb{R}$ be a twice differentiable function with a unique minimum $\bar{\eta}$ and $\eta_{t+1} = \eta_t - \gamma_t B_t Y_t$ where B_t is symmetric and only depends on information available at time t . Then η_t converges to $\bar{\eta}$ almost surely if the following conditions hold*

C.1 $(\forall t) \mathbb{E}_t Y_t = \nabla \mathbb{L}(\eta_t)$

C.2 $(\exists K)(\forall \eta) \|\nabla_{\eta}^2 \mathbb{L}(\eta)\| \leq 2K$

C.3 $(\forall \delta > 0) \inf_{\mathbb{L}(\eta) - \mathbb{L}(\bar{\eta}) > \delta} \|\nabla \mathbb{L}(\eta)\| > 0$

C.4 $(\exists A, B)(\forall t) \mathbb{E}_t \|Y_t\|^2 \leq A + B \cdot \mathbb{L}(\eta_t)$

C.5 $(\exists m, M : 0 < m < M < \infty) (\forall t) mI \prec B_t \prec MI$, where I is the identity matrix;

C.6 $\sum_t \gamma_t^2 < \infty, \sum_t \gamma_t = \infty$

\mathbb{E}_t in conditions **C.1** and **C.4** notes the conditional expectation given observations until time t . That is, $\mathbb{E}_t X = \mathbb{E}[X | \mathcal{F}_t]$, where $\mathcal{F}_t = \{\eta_1, \dots, \eta_t\}$ in this case. We recall Robbins-Siegmund Theorem [38] below, which is the key tool for proving both Theorem 1 and our generalization.

Theorem 2 (Robbins-Siegmund). *Let (Ω, \mathcal{F}, P) be a probability space and $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$ a sequence of sub- σ -fields of \mathcal{F} . Let U_t, β_t, ϵ_t and $\zeta_t, t = 1, 2, \dots$ be non-negative \mathcal{F}_t -measurable random variables such that*

$$\mathbb{E}(U_{t+1} | \mathcal{F}_t) \leq (1 + \beta_t)U_t + \epsilon_t - \zeta_t, \quad t = 1, 2, \dots \quad (10)$$

Then on the set $\{\sum_t \beta_t < \infty, \sum_t \epsilon_t < \infty\}$, U_t converges almost surely to a random variable, and $\sum_t \zeta_t < \infty$ almost surely.

3 Dual Stochastic Natural Gradient Descent

Recall from the introduction that the main idea to reduce the computational complexity of DSNGD is to define our optimization problem over a DFM. We start by establishing in section 3.1 that the family of joint distributions $P(\mathcal{Y}, \mathcal{X})$ satisfying the core MLR assumption is a LEF and hence a DFM. Then, we rely on duality to provide an efficient computation of the natural gradient of the log-loss function in section 3.2. Finally, we provide the DSNGD algorithm in section 3.3.

3.1 MLR generative model. The joint distribution

The next result proves that the the family of joint distributions $P(\mathcal{Y}, \mathcal{X})$ satisfying the core MLR assumption is a LEF and hence a DFM.

Proposition 1. *The log-odds ratio of the class posteriors $P(\mathcal{Y} | \mathcal{X})$ is an affine function of the features \mathcal{X} if and only if the joint distribution $P(\mathcal{Y}, \mathcal{X})$ belongs to LEF.*

Furthermore, there exist the LEF natural parametrization of the joint distribution

$$\begin{aligned} P_\eta(x, y) &= \frac{\exp S(y)^\top \alpha + T(x)^\top \beta_y}{\lambda(\eta)} \\ \lambda(\eta) &= \int_x \sum_y \exp S(y)^\top \alpha + T(x)^\top \beta_y \end{aligned} \quad (11)$$

where $\eta = (\alpha, \beta)$, $\alpha \in \mathbb{R}^{s-1}$, $\beta \in \mathbb{R}^{s \times t}$, β_y is the y -th row of β and

$$\begin{aligned} T : \Omega &\rightarrow \mathbb{R}^t \\ S : [1, \dots, s] &\rightarrow \mathbb{R}^{s-1} \end{aligned} \quad (12)$$

are sufficient and minimal statistics of \mathcal{X} and \mathcal{Y} respectively.

The proof relies strongly on theorem 2 in [36] and can be found in appendix A. This is convenient for our purpose, because, if we recall section 2.1, in a DFM the costs of natural gradient computations can be highly reduced, based on the property shown by equation 6. Next, we provide the dually flat parametrization of $P(\mathcal{Y}, \mathcal{X})$.

3.1.1 Dually flat parametrization of the joint distribution

We have seen that $P(\mathcal{Y}, \mathcal{X})$ is a LEF and that we can choose the natural parametrization of equation 11. With a linear transformation, S can become a canonical statistic, that means, $S(i)_j = \delta_{i=j}$ for $1 \leq i < s$ and $S(s) = 0$. For simplicity, we fix statistic S to be canonical from now on. The conditional probability distributions with η parametrization are

$$P_\eta(y | x) = \frac{\exp S(y)^\top \alpha + T(x)^\top \beta_y}{\sum_y \exp S(y)^\top \alpha + T(x)^\top \beta_y} \quad (13)$$

As [30] proves, the exponential family manifold is built after the convex function $F(\eta) = \log \lambda(\eta)$. The reference proves that this Riemannian manifold derived from $F(\eta)$, according to equation 5, has the Fisher information metric, as is usually considered for statistical manifolds, defined as

$$G_\eta = -\mathbb{E}_{x, y \sim P_\eta} [\nabla^2 \log P_\eta(y, x)] \quad (14)$$

Equation 5 also reveals the dual parametrization η^* . For LEF, it is called the expectation parametrization and it is shown below. For more properties of the

dual parametrization see [30]. To simplify the notation, if $x = (x_1 \ \cdots \ x_n)$, we note $\nabla_x = \left(\frac{\partial}{\partial x_1} \ \cdots \ \frac{\partial}{\partial x_n}\right)^\top$. So for every $i \in \{1, \dots, s\}$ write

$$\begin{aligned}\alpha^* &= \nabla_\alpha F(\eta) = \sum_y S(y)P_\eta(y) = \mathbb{E}_\mathcal{Y}[S(y)] = (P_\eta(\mathcal{Y} = 1), \dots, (P_\eta(\mathcal{Y} = s - 1)))^\top \\ \beta_i^* &= \nabla_{\beta_i} F(\eta) = P_\eta(\mathcal{Y} = i) \int_{\mathcal{X}} T(x)P_\eta(x | \mathcal{Y} = i) = P_\eta(\mathcal{Y} = i)\mathbb{E}_{\mathcal{X}|\mathcal{Y}=i}[T(x)]\end{aligned}\tag{15}$$

Define $\eta^* = (\alpha^*, \beta^*)$ with $\beta = (\beta_1^*, \dots, \beta_s^*)$ the dual parameterization, or equivalently, the expectation parameters.

Observe that $P(\mathcal{Y})$ is the categorical distribution (since \mathcal{Y} is discrete and finite) and therefore it is a LEF, where α^* are actually the expectation parameters. Moreover

$$\theta_i := \theta_i(\alpha^*, \beta_i^*) = \frac{\beta_i^*}{P_{\eta^*}(\mathcal{Y} = i)} = \mathbb{E}_{\mathcal{X}|\mathcal{Y}=i}[T(x)]\tag{16}$$

are the expectation parameters of the conditional distribution $P(\mathcal{X} | \mathcal{Y} = i)$.

3.2 Fast natural gradient of the log-loss

This section allows to compute the natural gradient of the log-loss function without having to use the metric matrix directly, but using both dual parametrizations instead.

Given $(y, x) \in \mathcal{Y} \times \mathcal{X}$ and $\eta \in \mathbb{R}^k$, the log-loss function is defined as

$$l(\eta, x, y) = -\log P_\eta(y | x)\tag{17}$$

Below result reveals $\tilde{\nabla}l(\eta, x, y)$ using both dual parametrizations η and η^* .

Proposition 2. *Let l be the log-loss function. Then, if $P(\mathcal{Y}, \mathcal{X})$ is a DFM*

$$\tilde{\nabla}l(\eta, x, y) = \nabla h(x, \eta^*) \cdot (q_\mathcal{Y}(x, P_\eta) - e_s(y))\tag{18}$$

where

$$q_\mathcal{Y}(x, P) = \begin{pmatrix} P(\mathcal{Y} = 1|x) \\ \vdots \\ P(\mathcal{Y} = s|x) \end{pmatrix},\tag{19}$$

$h(x, \eta^*) = (\log P_{\eta^*}(\mathcal{Y} = 1, x), \dots, \log P_{\eta^*}(\mathcal{Y} = s, x))$ and $e_s(k)$ is the k -th canonical s -dimensional vector.

The proof of proposition 2 is presented in Appendix B. To evaluate the computational complexity of using equation 18 we determine an expression of

$\nabla h(x, \zeta^*)$ with respect to the expectation parameters θ_y of \mathcal{X} given \mathcal{Y} already mentioned in equation 16. Below notation is used

$$K_i = \left(\begin{array}{c|c} & \begin{matrix} -1 \\ \vdots \\ -1 \end{matrix} \\ \hline Id^{i-1} & \end{array} \right) \quad d(x, y, \zeta^*) = \frac{1 - \theta_y^\top \nabla_{\theta_y} \log P_{\theta_y}(x | y)}{P_{\zeta^*}(y)} \quad (20)$$

and the proof is shown in Appendix C.

Proposition 3.

$$\nabla h(x, \zeta^*) = \begin{pmatrix} \nabla_{\alpha^*} h(x, \zeta^*) \\ \nabla_{\beta_1^*} h(x, \zeta^*) \\ \vdots \\ \nabla_{\beta_s^*} h(x, \zeta^*) \end{pmatrix} \quad (21)$$

where

$$\begin{aligned} \nabla_{\alpha^*} h(x, \zeta^*) &= K_s \cdot \text{diag}(d(x, 1, \zeta^*), \dots, d(x, s, \zeta^*)) \\ \nabla_{\beta_k^*} h(x, \zeta^*) &= \frac{\nabla_{\theta_k} \log P_{\theta_k}(x | \mathcal{Y} = k) \cdot e_s(k)^\top}{P_{\zeta^*}(\mathcal{Y} = k)} \end{aligned} \quad (22)$$

The complexity analysis of natural gradient is presented now, and the reader can find the proof in appendix E.

Proposition 4. *The computational complexity of the natural gradient $\tilde{\nabla} l(\eta, x, y)$ using proposition 2 is $O(s \cdot (A+t))$ where A is the cost of computing $\nabla_{\theta_y} \log P_{\theta_y}(x | y)$, s is the number of classes and t is the dimension of statistic T .*

Observe that the manifold dimension is $k = s - 1 + s \cdot t$ and therefore, a computation is linear on the number of the variables of the model if its complexity order is $O(k) = O(s(1+t)) = O(st)$. Therefore, the costs of computing the natural gradient can be reduced to linear if the cost A is low enough, precisely, if A is at most linear ($O(A) \leq O(k)$). This is the case when \mathcal{X} is discrete and finite (section 4).

3.3 DSNGD definition

DSNGD aims to solve the MLR optimization problem using the natural parametrization η of the LEF on $\mathcal{Y} \times \mathcal{X}$: If \bar{P} is an unknown probability distribution over $\mathcal{Y} \times \mathcal{X}$, optimize $\mathbb{L}(\eta) = \mathbb{E}_{x, y \sim \bar{P}} [l(\eta, x, y)]$ for $\eta \in \mathbb{R}^k$ where $l(\eta, x, y)$ is the log-loss function. The solution $\bar{\eta} \in \mathbb{R}^k$ to this problem refers to the conditional distributions $P_{\bar{\eta}}(\mathcal{Y} | \mathcal{X})$ that better fits the hidden conditional distributions $\bar{P}(\mathcal{Y} | \mathcal{X})$. To that end, we define a stochastic natural gradient based algorithm.

Using proposition 2, DSNGD moves by following the update equation

Definition 1 (DSNGD update).

$$\eta_{t+1} = \eta_t - \gamma_t \nabla h(x_t, \zeta_t^*) \cdot (q_{\mathcal{Y}}(x_t, P_{\eta_t}) - e_s(y_t)) \quad (23)$$

where $\{\zeta_t^*\}_{t \in \mathbb{N}}$ is a sequence in the expectation parametrization such that $\{\zeta_t\}_{t \in \mathbb{N}}$ converges.

Note that $q_{\mathcal{Y}}(x_t, P_{\eta_t})$ is a stable term (it only takes values between 0 and 1). Moreover, DSNGD forces the stability of the $\nabla h(x_t, \zeta_t^*)$ term, since ζ_t is a convergent sequence. This is the same strategy of CSNGD, and similarly, it is going to ensure the convergence of the algorithm in theorem 5. Observe that equation 23 is also well defined when the parameterization is not minimal (when T is not a minimal statistic), therefore DSNGD can be run in such general case, where S and T are not minimal. Steps taken by DSNGD are specified in Algorithm 1 below.

Algorithm 1: DSNGD

Result: η

- 1 $\eta \leftarrow \eta_0, \zeta^* \leftarrow \zeta_0^*, \gamma \leftarrow \gamma_0;$
- 2 **while** *observations x, y and stopping condition is false* **do**
- 3 $q \leftarrow q_{\mathcal{Y}}(x, P_{\eta});$
- 4 $grad_h \leftarrow \nabla h(x, \zeta^*);$
- 5 $d \leftarrow grad_h \cdot (q - e_s(y));$
- 6 $\eta \leftarrow \eta - \gamma \cdot d;$
- 7 update $\zeta^*;$
- 8 update γ
- 9 **end**

The sequence $\{\zeta_t^*\}_{t \in \mathbb{N}}$, or simply ζ_t^* as an abuse of notation, can be any sequence in the dual space whose dualized sequence $\{\zeta_t\}_{t \in \mathbb{N}}$ is convergent. For example, it can be constant. The resulting algorithm keeps track of two independent sequences; the main sequence η_t which estimates the solution $\bar{\eta}$ to the problem, and the sequence ζ_t^* selected with the convergence constraint and whose space is the dual. For example, assume the trivial case where $\mathcal{X} = \{0\}$ and $\mathcal{Y} = \{0, 1, 2\}$. The only conditional probability distribution of the problem is the Categorical distribution $P(\mathcal{Y} | \mathcal{X} = 0)$. This space is represented by \mathbb{R}^2 and its dual space is represented by the simplex S^2 . Then, the main sequence η_t moves in \mathbb{R}^2 while the independent sequence ζ_t^* traces its path in S^2 . Figure 1 illustrates iterations followed by η_t (instruction line 6 of the algorithm) and ζ_t^* (instruction line 7 of the algorithm) when running DSNGD for this simple example.

Recall that the sequence ζ_t^* can be chosen freely as long as its dual is convergent. However, recall that DSNGD is a natural gradient based algorithm. The algorithm effectively takes a natural gradient step only when η_t and ζ_t^* refer to the same probability distribution point, according to equation 23 and proposition 2. In section 5 there is our proof of DSNGD convergence to the solution $\bar{\eta}$, and if ζ_t^* is selected such that it also converges to the solution, then both sequences get closer along the optimization process, turning DSNGD steps into more accurate approximations of natural gradient steps. Therefore, in order to benefit from natural gradient speed up properties, it is recommended that sequence ζ_t^* converges to the solution $\bar{\eta}^* = \nabla F(\bar{\eta})$. For example, this can be accomplished by determining ζ_t^* using a maximum a posteriori estimator of the parameters of $P(\mathcal{Y}, \mathcal{X})$ obtained from data up to t .

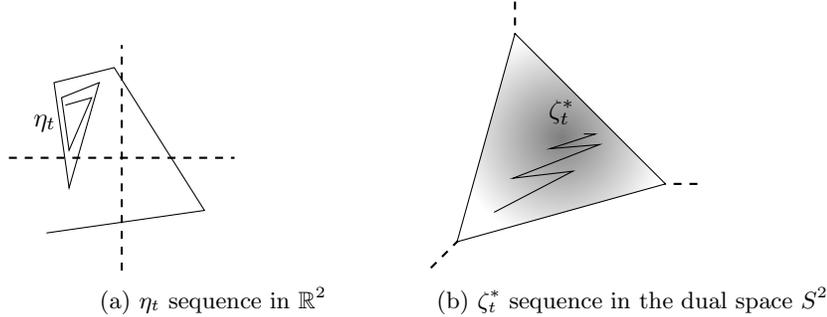


Figure 1: η_t and ζ_t^* sequences obtained in DSNGD where $\mathcal{X} = \{0\}$ and $\mathcal{Y} = \{0, 1, 2\}$

4 Discrete DSNGD and computational complexity

This section assumes that space \mathcal{X} is discrete, that is $\mathcal{X} = \{1, \dots, m\}$ for some $m \in \mathbb{N}$. For simplicity, we assume T to be also a canonical statistic, that is, $T(i)_j = \delta_{i=j} \in \mathbb{R}^{m-1}$ for $1 \leq i < m$ and $T(m) = 0$. A theorem deduces and proves that the complexity order for discrete DSNGD of one iteration is linear on the dimension of the parameter η . Let us show a simple example of discrete DSNGD to begin with.

4.1 Example

Let $\mathcal{Y} = \{1, 2\}$ and $\mathcal{X} = \{1, 2\}$ and minimal and canonical statistics S and T . Let $\eta = (\alpha, \beta)$ be the natural parameter and $\zeta^* = (\alpha^*, \beta^*)$ be the independent dual parameter. Observe that in this case, α and α^* are 1-element vectors and β and β^* are 2-element not squared matrices. In this example we complete an iteration of discrete DSNGD algorithm, following the instructions listed in algorithm 1.

Let $(y, x) = (2, 1)$ be an observation. Statistics T and S are assumed to be canonical. Instruction line 3 consist on using equation 13 to compute

$$q_{\mathcal{Y}}(x, P_{\eta}) = \begin{pmatrix} P(\mathcal{Y} = 1 | x) \\ P(\mathcal{Y} = 2 | x) \end{pmatrix} = R \cdot \begin{pmatrix} \exp(\alpha_1 + \beta_1) \\ \exp \beta_2 \end{pmatrix} \quad (24)$$

where $R = \frac{1}{\exp(\alpha_1 + \beta_1) + \exp \beta_2}$. For instruction line 4, express function $h(x, \zeta^*)$ (use equation 15), then apply the gradient.

$$h(x = 1, \zeta^*) = (\log \beta_1^*, \log \beta_2^*) \rightarrow \nabla h(x = 1, \zeta^*) = \begin{pmatrix} 0 & 0 \\ \frac{1}{\beta_1^*} & 0 \\ 0 & \frac{1}{\beta_2^*} \end{pmatrix} \quad (25)$$

Proceed now with instruction line 5. It computes the approximation of the natural gradient and the direction that DSNGD uses for the η update.

$$\begin{aligned}
d &= \nabla h(x=1, \zeta^*) (q_{\mathcal{Y}}(x, P_\eta) - e_2(y=2)) \\
&= \begin{pmatrix} 0 & 0 \\ \frac{1}{\beta_1^*} & 0 \\ 0 & \frac{1}{\beta_2^*} \end{pmatrix} \cdot \begin{pmatrix} R \cdot \exp(\alpha_1 + \beta_1) \\ (R \cdot \exp \beta_2) - 1 \end{pmatrix} \\
&= \begin{pmatrix} 0 \\ \frac{R \cdot \exp(\alpha_1 + \beta_1)}{\beta_1^*} \\ \frac{(R \cdot \exp \beta_2) - 1}{\beta_2^*} \end{pmatrix}
\end{aligned} \tag{26}$$

Next instruction lines of the algorithm are standard to update the parameter vector $(\alpha, \beta_1, \beta_2)$ using direction d , so there is no need to go further. If for instance the observation is $(y, x) = (2, 2)$, then the approximation of the natural gradient is

$$d = \begin{pmatrix} \frac{R \exp \alpha_1}{\alpha_1^* - \beta_1^*} - \frac{R-1}{1-\alpha_1^* - \beta_2^*} \\ -\frac{R \cdot \exp \alpha_1}{\alpha_1^* - \beta_1^*} \\ \frac{R-1}{1-\alpha_1^* - \beta_2^*} \end{pmatrix} \tag{27}$$

where $R = \frac{1}{1 + \exp \alpha_1}$

Before analyzing the computational complexity of DSNGD, it is necessary to determine the generator of ζ_t^* sequence. Sequence ζ_t^* belongs to the dual space of the LEF distributions on $\mathcal{Y} \times \mathcal{X}$, and if S and T are canonical statistics then it implies that ζ_t^* are directly the probabilities $P(y, x)$ after equation 15. It is possible to select the well known maximum a posteriori (MAP) estimator with parameter $a \in \mathbb{R}$. This estimator is a simple counting of observations over the discrete space $\mathcal{Y} \times \mathcal{X}$ with an starting assumption of incidence of a for every event y, x . This estimator is linear and it clearly converges (to the solution).

First a similar result as proposition 3 is stated, taking into account the new assumption on \mathcal{X} . The proof of proposition below is found in appendix D.

Proposition 5. *Let $\mathcal{X} = \{1, \dots, m\}$ and let T be a minimal and canonical statistic. Then*

$$\begin{aligned}
\nabla_{\alpha^*} h(x, \zeta^*) &= \begin{cases} 0 & x \neq m \\ K_s \cdot \text{diag}\left(\frac{1}{P_{\zeta^*}(x, \mathcal{Y}=1)}, \dots, \frac{1}{P_{\zeta^*}(x, \mathcal{Y}=s)}\right) & x = m \end{cases} \\
\nabla_{\beta_y^*} h(x, \zeta^*) &= \frac{1}{P_{\zeta^*}(x, y)} \cdot \begin{cases} e_{m-1}(x) \cdot e_s(y)^\top & x \neq m \\ -\mathbf{1}_{m-1} \cdot e_s(y)^\top & x = m \end{cases}
\end{aligned} \tag{28}$$

where $\mathbf{1}_n \in \mathbb{R}^n$ is a vector filled with ones at every coordinate.

Now it is possible to analyze the computational complexity of discrete DSNGD. Below theorem proves that DSNGD, just as SGD, is a linear algorithm.

Theorem 3. *Let $\mathcal{X} = \{1, \dots, m\}$ and let T be a minimal and canonical statistic. Assume estimator ζ^* of DSNGD is linear. Then discrete DSNGD iterations have linear complexity order on the manifold dimension.*

Proof. Let $k = (s - 1) + s \cdot t$ be the dimension of η . Then $O(k) = O(st)$. Analyze the computational complexity of discrete DSNGD. That is, analyze the computational cost of instruction lines 3, 4, 5, 6 and 7 shown in Algorithm 1.

Complexity of instruction lines 3, 4 and 5 is given by proposition 4, which is $O(sA + st)$ where sA is the cost of computing $\nabla_{\theta_y} \log P_{\theta_y}(x | y)$ for all $y \in \mathcal{Y}$. Observe equations 15 and 16 assuming T canonical and write

$$\begin{aligned} \alpha^* &= (P_{\zeta^*}(\mathcal{Y} = 1), \dots, P_{\zeta^*}(\mathcal{Y} = s - 1))^\top \\ \theta_y &= (P_{\zeta^*}(\mathcal{X} = 1 | y), \dots, P_{\zeta^*}(\mathcal{X} = m - 1 | y))^\top \end{aligned} \tag{29}$$

Deduce then that $O(sA) = O(k)$.

Instruction line 6 adds k operations.

Finally, recall that a linear complexity order estimator is chosen for ζ_t^* sequence, implying that instruction line 7 is of linear order $O(k)$.

In conclusion, the computational complexity order of DSNGD is

$$O(sA + st) + O(k) + O(k) = O(k) \tag{30}$$

and therefore linear. \square

5 Discrete DSNGD convergence

In this section we prove the convergence of the discrete DSNGD. Discrete DSNGD refers to the case where $\mathcal{X} = \{1, \dots, m\}$ for some $m \in \mathbb{R}$. We start by generalizing Theorem 3.2 in Sunehag's et al. [5] (introduced above and referred to from now on as Theorem 1) in Section 5.1. This generalization provides enough flexibility so as to be used later to prove the convergence of DSNGD in Section 5.2.

5.1 Generalizing Sunehag et. al. variable metric stochastic approximation theory.

Theorem 1 is used to prove CSNGD convergence, however it can not be used to prove DSNGD convergence. First, because it demands the vector it follows to be factored as the product of a matrix B_t and a vector Y_t that approximates the gradient (condition **C.1**). But DSNGD is defined to directly approximate the natural gradient, without the gradient as reference. And second, even if DSNGD is written as the product of a matrix and a vector, matrix $\nabla h(x_t, \zeta_t^*)$ is not squared. So we need a more general convergence theorem.

Our result proves almost sure convergence of the sequence

$$\eta_{t+1} = \eta_t - \gamma_t Y(\eta_t, \mathcal{F}_t) \tag{31}$$

where $Y(\eta, \mathcal{F}_t)$ is a family of random vectors defined for every η and for every set

$$\mathcal{F}_t = \{(y_i, x_i) \mid i < t\} \quad (32)$$

As an abuse of notation write Y_t meaning the random variable $Y(\eta_t, \mathcal{F}_t) \in \mathbb{R}^n$. The main modification with respect to Theorem 1 is that we unify conditions C.1 and C.3

$$\begin{aligned} \mathbf{C.1} \quad & (\forall t) \quad \mathbb{E}_t Y_t = \nabla l(\eta_t) \\ \mathbf{C.3} \quad & (\forall \delta > 0) \quad \inf_{l(\eta) - l(\bar{\eta}) > \delta} \|\nabla l(\eta)\| > 0 \end{aligned} \quad (33)$$

to instead require

$$\mathbf{C.3} \quad (\forall \delta > 0) \quad \inf_{l(\eta_t) - l(\bar{\eta}) > \delta} \nabla l(\eta_t)^T \mathbb{E}_t [Y_t] > 0 \quad (34)$$

New condition **C.3** uses \mathbb{E}_t , referring to the conditional expectation given \mathcal{F}_t of equation 32, which is a generalization of the definition of \mathbb{E}_t in Sunehag [5].

Theorem 1 imposes that the expectation of the step taken must be the gradient and that the norm of the gradient must not approach to zero outside any environment of the minimum. Instead, we impose that the expectation of the step taken must not approach to the border of the half-space which has the gradient as its normal vector, unless we are approaching the minimum simultaneously. This is a more general condition. Furthermore, condition **C.5** on the maximum and minimum eigenvalues of the matrix B_t can also be removed. In fact, our result proves the convergence of algorithms with scaling matrices B_t whose spectrum is not bounded from below by a strictly positive number, as long as new version of condition **C.3** holds.

It is formally stated below. Proof is found in appendix F.

Theorem 4. *Let $l : \mathbb{R}^k \rightarrow \mathbb{R}$ be a twice differentiable function with a unique minimum $\bar{\eta}$ and $\eta_{t+1} = \eta_t - \gamma_t Y_t$. Then η_t converges to $\bar{\eta}$ almost surely if the following conditions hold*

$$\begin{aligned} \mathbf{C.2} \quad & (\exists K)(\forall \eta) \quad \|\nabla_{\eta}^2 l(\eta)\| \leq 2K \\ \mathbf{C.3} \quad & (\forall \delta > 0) \quad \inf_{l(\eta_t) - l(\bar{\eta}) > \delta} \nabla l(\eta_t)^T \mathbb{E} [Y_t] > 0 \\ \mathbf{C.4} \quad & (\exists A, B)(\forall t) \quad \mathbb{E} \|Y_t\|^2 \leq A + B l(\eta_t) \\ \mathbf{C.6} \quad & \sum_t (\gamma_t)^2 < \infty, \quad \sum_t \gamma_t = \infty \end{aligned}$$

5.2 Proving convergence

Next, we show how Theorem 4 can be used to prove DSNGD convergence in the discrete case. That is, we use it to prove the next result:

Theorem 5. *DSNGD in the canonical parametrization such that \mathcal{Y} and \mathcal{X} are discrete, converges almost surely to the optimum.*

The proof consists on showing that conditions **C.2**, **C.3**, **C.4** and **C.6** of theorem 4 hold. Condition **C.6** is assumed to hold, by just selecting an appropriate sequence of learning rates γ_t . Conditions **C.2** and **C.4** are proved in appendices G and H respectively. Proof of condition **C.3** is shown below.

Proof. Compute the gradient of $l(\eta)$ (see equation 60) and use proposition 2 to obtain $\mathbb{E}_t [Y_t]$ involved in condition **C.3**.

$$\begin{aligned}
\nabla l(\eta) &= \mathbb{E}_t [\nabla l(\eta, x, y)] \\
&= \sum_x \nabla h(x, \eta) \sum_y (q_Y(x) - e_s(y)) \bar{P}(x, y) \\
&= \sum_x \nabla h(x, \eta) \text{diff}_Y(x, \eta) \\
\mathbb{E}_t [Y_t] &= \sum_x \nabla h(x, \zeta^*) \text{diff}_Y(x, \eta)
\end{aligned} \tag{35}$$

where

$$\text{diff}_Y(x, \eta) = (q_Y(x, P_\eta) - q_Y(x, \bar{P})) \bar{P}(x) \tag{36}$$

Further evolve equation 35 to finally multiply $\nabla l(\eta)^\top \mathbb{E}_t [Y_t]$ and check condition **C.3**. Continue by developing $\nabla l(\eta)$ first, precisely compute $\nabla h(x, \eta)$. To simplify the notation, decompose $\nabla = (\nabla_\alpha, \nabla_{\beta_1}, \dots, \nabla_{\beta_s})$

$$\begin{aligned}
\nabla_\alpha h(x, \eta) &= S + u(P_\eta) \cdot (1, \dots, 1) & u(P) &= - \sum_y S(y) P(y) \\
\nabla_{\beta_y} h(x, \eta) &= T(x) e_s(y)^\top + v(y, P_\eta) \cdot (1, \dots, 1) & v(y, P) &= - \sum_x T(x) P(y, x)
\end{aligned} \tag{37}$$

Since $(1, \dots, 1) \cdot \text{diff}_Y(x, \eta) = 0$ then

$$\begin{aligned}
\nabla_\alpha l(\eta) &= \sum_x \nabla_\alpha h(x, \eta) \text{diff}_Y(x, \eta) \\
&= \sum_x S \cdot \text{diff}_Y(x, \eta) \\
&= S \cdot \text{diff}_Y(\eta) \\
\nabla_{\beta_y} l(\eta) &= \sum_x \nabla_{\beta_y} h(x, \eta) \text{diff}_Y(x, \eta) \\
&= \sum_x T(x) e_s(y)^\top \text{diff}_Y(x, \eta) \\
&= T \cdot \text{diff}_X(y, \eta)
\end{aligned} \tag{38}$$

where

$$\begin{aligned} \text{diff}_Y(\eta) &= \begin{pmatrix} P_\eta(\mathcal{Y} = 1) - \bar{P}(\mathcal{Y} = 1) \\ \vdots \\ P_\eta(\mathcal{Y} = s) - \bar{P}(\mathcal{Y} = s) \end{pmatrix} \\ \text{diff}_X(y, \eta) &= \begin{pmatrix} (P_\eta(\mathcal{Y} = y | \mathcal{X} = 1) - \bar{P}(\mathcal{Y} = y | \mathcal{X} = 1))\bar{P}(\mathcal{X} = 1) \\ \vdots \\ (P_\eta(\mathcal{Y} = y, \mathcal{X} = m) - \bar{P}(\mathcal{Y} = y | \mathcal{X} = m))\bar{P}(\mathcal{X} = m) \end{pmatrix} \end{aligned} \quad (39)$$

Now develop $\mathbb{E}_t Y_t$ further. Recall that the canonical parametrization is selected so plug in proposition 5 into equation 35. Decompose $\mathbb{E}_t = (\mathbb{E}_{t, \alpha^*}, \mathbb{E}_{t, \beta_1^*}, \dots, \mathbb{E}_{t, \beta_s^*})$

$$\begin{aligned} \mathbb{E}_{t, \alpha^*} [Y_t] &= \sum_x \nabla_{\alpha^*} h(x, \zeta^*) \text{diff}_Y(x, \eta) \\ &= K_s \cdot \text{diag}(d(m, 1, \zeta^*), \dots, d(m, s, \zeta^*)) \cdot \text{diff}_Y(m, \eta) \\ \mathbb{E}_{t, \beta_y^*} [Y_t] &= \sum_x \nabla_{\beta_y^*} h(x, \zeta^*) \text{diff}_Y(x, \eta) \\ &= K_m \cdot \text{diag}(d(1, y, \zeta^*), \dots, d(m, y, \zeta^*)) \cdot \text{diff}_X(y, \eta) \end{aligned} \quad (40)$$

Proceed now to check the condition. Develop the products until obtain

$$\begin{aligned} \nabla_\alpha l(\eta)^\top \mathbb{E}_{t, \alpha^*} [Y_t] &= \sum_y c(y) \\ \nabla_{\beta_y} l(\eta)^\top \mathbb{E}_{t, \beta_y^*} [Y_t] &= -c(y) + \sum_x d(x, y, \zeta^*) (P_\eta(y|x) - \bar{P}(y|x))^2 \bar{P}(x)^2 \end{aligned} \quad (41)$$

where $c(y) = d(m, y, \zeta^*) (P_\eta(y) - \bar{P}(y)) (P_\eta(y|x=m) - \bar{P}(y|x=m)) \bar{P}(x=m)$

Finally,

$$\begin{aligned} \nabla l(\eta)^\top \mathbb{E}_t [Y_t] &= \nabla_\alpha l(\eta)^\top \mathbb{E}_{t, \alpha^*} [Y_t] + \sum_y \nabla_{\beta_y} l(\eta)^\top \mathbb{E}_{t, \beta_y^*} [Y_t] \\ &= \sum_y c(y) + \sum_y -c(y) + \sum_x d(x, y, \zeta^*) (P_\eta(y|x) - \bar{P}(y|x))^2 \bar{P}(x)^2 \\ &= \sum_{y,x} d(x, y, \zeta^*) (P_\eta(y|x) - \bar{P}(y|x))^2 \bar{P}(x)^2 \end{aligned} \quad (42)$$

Notice in equation 42 that $\nabla l(\eta)^\top \mathbb{E}_t [Y_t]$ is a sum of positive numbers, and it vanishes only if $\eta = \bar{\eta}$. Also, since $d(x, y, \zeta^*) > 1$, observe that

$$\begin{aligned} \nabla l(\eta)^\top \mathbb{E}_t [Y_t] &> \sum_{y,x} (P_\eta(y|x) - \bar{P}(y|x))^2 \bar{P}(x)^2 \\ &= \sum_y \|\text{diff}_X(y, \eta)\|^2 \end{aligned} \quad (43)$$

To finish proving the result, let $\{\eta_i\}_{i \in \mathbb{N}}$ be a sequence such that

$$\sum_y \|\text{diff}_{\mathcal{X}}(y, \eta_i)\|^2 \xrightarrow{i \rightarrow \infty} 0 \quad (44)$$

since every term is positive, then for every $y \in \mathcal{Y}$

$$\|\text{diff}_{\mathcal{X}}(y, \eta_i)\|^2 \xrightarrow{i \rightarrow \infty} 0 \quad (45)$$

implying that $P_{\eta_i}(y|x) - \bar{P}(y|x) \xrightarrow{i \rightarrow \infty} 0$ for all x, y and that

$$l(\eta_i) - l(\bar{\eta}) \xrightarrow{i \rightarrow \infty} 0 \quad (46)$$

Hence it's proven

$$(\forall \delta > 0) \quad \inf_{l(\eta) - l(\bar{\eta}) > \delta} \sum_y \|\text{diff}_{\mathcal{X}}(y, \eta_i)\|^2 > 0 \quad (47)$$

and therefore, after equation 43, condition **C.3** holds. \square

6 Conclusion and future work

Natural gradient based algorithms behave erratically when tested in practical problems. However, as CSNGD shows, these kind of algorithms may stabilize once convergence is guaranteed. With this in mind, we defined DSNGD, which approximates the natural gradient at each step and whose convergence in the discrete case can be proved. To that end, we stated and proved a general result showing the convergence of interior half-space gradient approximations. Furthermore, we point out that this convergence result may prove the convergence of more general algorithms, since it doesn't require the expectation of the update's direction to factor as a symmetric positive-definite matrix and the gradient.

This paper concentrates on the theoretical aspects of DSNGD. We are currently working on a flexible implementation of the algorithm that can be easily set up for different LEF linked to the conditional distributions $P_{\eta}(\mathcal{X} | \mathcal{Y})$, including several commonly used discrete, continuous and multivariate distributions such as the normal, Poisson and exponential. Moreover, DSNGD can potentially be used in high dimensional scenarios due to its low computational complexity. The benefits of approximating the natural gradient are specially promising in this case, since the parameter space is potentially twisted and using metric information can be crucial for an algorithm's good performance. In preliminary empirical studies we are observing how it increasingly outperforms SGD as the manifold dimension grows larger. We plan to compare DSNGD against the most effective algorithms nowadays, in order to expose its weaknesses and reveal its strengths.

For the more theoretical part, in the future we plan to study the convergence of continuous and mixed DSNGD, that is when \mathcal{X} is continuous, and in cases where $\mathcal{X} = (\mathcal{X}_d, \mathcal{X}_c)$ is divided into a discrete and a continuous part.

7 Declarations

7.1 Funding

The work has been funded by EU Horizon 2020 under grant agreements 872944 (Crowd4SDG) and 825619 (Humane-AI-net), and by the Spanish Ministry of Science and Innovation through the CI-SUSTAIN project (PID2019-104156GB-I00).

7.2 Conflicts of interest/Competing interests

On behalf of all authors, the corresponding author states that there is no conflict of interest.

7.3 Availability of data and material

Not applicable

7.4 Code availability

Not applicable

7.5 Authors' contributions

The authors have contributed equally to this work.

Appendices

A Proof of Proposition 1

Proof. Prove first that if the logg-odds ratio of $P(\mathcal{Y} | \mathcal{X})$ is an affine function of \mathcal{X} then the joint distribution $P(\mathcal{Y}, \mathcal{X})$ belongs to LEF.

According to theorem 2 in [36], assume that $P(\mathcal{X} | \mathcal{Y} = i)$ belongs to the same LEF for all $i \in \mathcal{Y}$. Also, since \mathcal{Y} is discrete and finite, $P(\mathcal{Y})$ is a categorical distribution and hence, it belongs to LEF. This means that there exist parameters $\bar{\alpha} \in \mathbb{R}^{s-1}$ and $\bar{\theta}_i \in \mathbb{R}^t$ for all $i \in \mathcal{Y}$ such that

$$\begin{aligned} P_{\bar{\alpha}}(\mathcal{Y} = i) &= \frac{\exp S(i)^\top \bar{\alpha}}{\sum_y \exp S(y)^\top \bar{\alpha}} \\ P_{\bar{\theta}_i}(x | \mathcal{Y} = i) &= \frac{\exp T(x)^\top \bar{\theta}_i}{\int_x \exp T(x)^\top \bar{\theta}_i} \end{aligned} \tag{48}$$

where S and T are sufficient statistics of \mathcal{Y} and \mathcal{X} respectively. If $\bar{\theta}$ is the matrix having $\bar{\theta}_i$ as i -th row, name $\bar{\eta} = (\bar{\alpha}, \bar{\theta})$ and write

$$\begin{aligned} P_{\bar{\eta}}(\mathcal{Y} = i, x) &= P_{\bar{\alpha}}(\mathcal{Y} = i) P_{\bar{\theta}_i}(x \mid \mathcal{Y} = i) \\ &= \frac{\exp S(i)^\top \bar{\alpha}}{\sum_y \exp S(y)^\top \bar{\alpha}} \frac{\exp T(x)^\top \bar{\theta}_i}{\int_x \exp T(x)^\top \bar{\theta}_i} \\ &= \frac{\exp S(i)^\top \bar{\alpha} + T(x)^\top \bar{\theta}_i}{\sum_y \exp S(y)^\top \bar{\alpha} \int_x \exp T(x)^\top \bar{\theta}_i} \end{aligned} \quad (49)$$

To prove the result, it is enough to find a change of variables from $\bar{\eta} = (\bar{\alpha}, \bar{\theta})$ to $\eta = (\alpha, \beta)$ satisfying $P_{\bar{\eta}}(y, x) = P_\eta(y, x)$ where

$$P_\eta(\mathcal{Y} = i, x) = \frac{\exp S(i)^\top \alpha + T(x)^\top \beta_i}{\int_x \sum_y \exp S(y)^\top \alpha + T(x)^\top \beta_y} \quad (50)$$

since η is the natural parametrization of a LEF.

In particular, the change of variables has to satisfy that $P_{\bar{\eta}}(x \mid \mathcal{Y} = i) = P_\eta(x \mid \mathcal{Y} = i)$ and $P_{\bar{\eta}}(y) = P_\eta(y)$. Start with the conditional probability and observe that

$$\begin{aligned} P_\eta(x \mid \mathcal{Y} = i) &= \frac{P_\eta(\mathcal{Y} = i, x)}{\int_x P_\eta(\mathcal{Y} = i, x)} = \frac{\exp S(i)^\top \alpha + T(x)^\top \beta_i}{\int_x \exp S(i)^\top \alpha + T(x)^\top \beta_i} \\ &= \frac{\exp T(x)^\top \beta_i}{\int_x \exp T(x)^\top \beta_i} \end{aligned} \quad (51)$$

Last equation matches exactly with equation 48 by just setting $\beta = \bar{\theta}$. To complete the change of variables continue by matching $P_{\bar{\eta}}(y) = P_\eta(y)$.

$$\begin{aligned} P_\eta(\mathcal{Y} = i) &= \frac{\int_x \exp S(i)^\top \alpha + T(x)^\top \beta_i}{\sum_j \int_x \exp S(j)^\top \alpha + T(x)^\top \beta_j} \\ &= \frac{\exp(S(i)^\top \alpha) \int_x \exp T(x)^\top \beta_i}{\sum_j \exp(S(j)^\top \alpha) \int_x \exp T(x)^\top \beta_j} \\ &= \frac{\exp S(i)^\top \alpha + \log A_i}{\sum_j \exp S(j)^\top \alpha + \log A_j} \end{aligned} \quad (52)$$

where $A_i = \int_x \exp T(x)^\top \beta_i$. Last equation must coincide with equation 48. That is

$$P_\eta(\mathcal{Y} = i) = P_{\bar{\eta}}(\mathcal{Y} = i) \iff \frac{\exp S(i)^\top \alpha + \log A_i}{\sum_j \exp S(j)^\top \alpha + \log A_j} = \frac{\exp S(i)^\top \bar{\alpha}}{\sum_y \exp S(y)^\top \bar{\alpha}} \quad (53)$$

To simplify, assume S is canonical. That is $S(i) = e_i$ is the i -th canonical vector for all $i \neq s$ and $S(s) = 0 \in \mathbb{R}^{s-1}$. Note that it is enough to prove that there exists a $\mu \in \mathbb{R}$ such that

$$S(i)^\top \alpha + \log A_i - \mu = S(i)^\top \bar{\alpha}, \quad \forall i \in \mathcal{Y} \quad (54)$$

because as a consequence, equation 53 clearly holds. In our case, it is $S(i)^\top \alpha = \alpha_i$ when $i \neq s$ and $S(i)^\top \alpha = 0$, and therefore the solution is

$$\alpha + \begin{pmatrix} \log A_1 \\ \vdots \\ \log A_{s-1} \end{pmatrix} - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \mu = \bar{\alpha} \quad (55)$$

$$\mu = \log A_s$$

and the proof is completed when S is canonical.

Prove now the result for a general sufficient statistic S . Equation 54 describes the below linear equations system

$$\mathbf{S}\alpha + \begin{pmatrix} \log A_1 \\ \vdots \\ \log A_{s-1} \end{pmatrix} - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \mu = \mathbf{S}\bar{\alpha} \quad (56)$$

$$S(s)^\top \alpha + \log A_s - \mu = S(s)^\top \bar{\alpha}$$

where \mathbf{S} is the matrix having $S(1), \dots, S(s-1)$ as rows. Since S is a sufficient statistic, assume without loss of generality that $S(1), \dots, S(s-1)$ are linearly independent vectors, and then \mathbf{S} is invertible. Finally, it is easy to check that the change of variables is

$$\alpha + \mathbf{S}^{-1} \left(\begin{pmatrix} \log A_1 \\ \vdots \\ \log A_{s-1} \end{pmatrix} - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \mu \right) = \bar{\alpha}$$

$$\mu = \frac{S(s)^\top \mathbf{S}^{-1} \begin{pmatrix} \log A_1 \\ \vdots \\ \log A_{s-1} \end{pmatrix} - \log A_s}{S(s)^\top \mathbf{S}^{-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} - 1} \quad (57)$$

The converse implication is straightforward. Assuming that $P(\mathcal{X}, \mathcal{Y})$ belongs to LEF, and therefore assuming equation 50, start by expressing the conditional probability distributions of \mathcal{Y} given \mathcal{X} in η .

$$P_\eta(y | x) = \frac{P_\eta(y, x)}{\sum_y P_\eta(y, x)} \quad (58)$$

$$= \frac{\exp S(y)^\top \alpha + T(x)^\top \beta_y}{\sum_y \exp S(y)^\top \alpha + T(x)^\top \beta_y}$$

and compute the log-odds ratio

$$\log \frac{P_\eta(x | \mathcal{Y} = k)}{P_\eta(x | \mathcal{Y} = h)} = S(k)^\top \alpha + T(x)^\top \beta_k - (S(h)^\top \alpha + T(x)^\top \beta_h) \quad (59)$$

$$= (S(k) - S(h))^\top \alpha + T(x)^\top (\beta_k - \beta_h)$$

which is clearly an affine function of features \mathcal{X} . \square

B Proof of Proposition 2

Proof. First, claim that

$$\nabla l(\eta, x, y) = \nabla h(x, \eta) \cdot (q_{\mathcal{Y}}(x, P_{\eta}) - e_s(y)) \quad (60)$$

Indeed,

$$\begin{aligned} \nabla \log P_{\eta}(y | x) &= \nabla \log P_{\eta}(y, x) - \nabla \log \sum_y P_{\eta}(y, x) \\ &= \nabla \log P_{\eta}(y, x) - \frac{\sum_y \nabla P_{\eta}(y, x)}{\sum_y P_{\eta}(y, x)} \\ &= \nabla \log P_{\eta}(y, x) - \frac{\sum_y P_{\eta}(y, x) \nabla \log P_{\eta}(y, x)}{\sum_y P_{\eta}(y, x)} \\ &= \nabla \log P_{\eta}(y, x) - \sum_y P_{\eta}(y | x) \nabla \log P_{\eta}(y, x) \\ &= \nabla h(y, x, \eta) - \mathbb{E}_{\mathcal{Y}|x}[\nabla h(y, x, \eta)] \end{aligned} \quad (61)$$

where $h(y, x, \eta) = \log P_{\eta}(y, x)$. Observe we can rewrite equation 61 as;

$$\nabla \log P_{\eta}(y | x) = -\nabla h(x, \eta) \cdot (q(x, \eta) - e_s(i)) \quad (62)$$

where $h(x, \eta) = (h(1, x, \eta), \dots, h(s, x, \eta))$ implying the claim. From equation 60 observe that

$$\tilde{\nabla} l(\eta, x, y) = \tilde{\nabla} h(x, \eta) \cdot (q_{\mathcal{Y}}(x, P_{\eta}) - e_s(y)) \quad (63)$$

Finally, since the log-loss is defined in a DFM, then use previous equation and equation 6 to finish the proof \square

C Proof of Proposition 3

Proof. To simplify, break $\nabla_{\eta^*} = (\nabla_{\alpha^*}, \nabla_{\beta_1^*}, \dots, \nabla_{\beta_s^*})$ and then it's clear that

$$\nabla h(x, \eta^*) = \begin{pmatrix} \nabla_{\alpha^*} h(x, \eta^*) \\ \nabla_{\beta_1^*} h(x, \eta^*) \\ \vdots \\ \nabla_{\beta_s^*} h(x, \eta^*) \end{pmatrix} \quad (64)$$

Start with $\nabla_{\alpha^*} h(x, \eta^*)$ expression. Observe that i -th column of $\nabla_{\alpha^*} h(x, \eta^*)$ is

$$\begin{aligned} \nabla_{\alpha^*} \log P_{\eta^*}(\mathcal{Y} = i, x) &= \nabla_{\alpha^*} \log P_{\alpha^*}(\mathcal{Y} = i) + \nabla_{\alpha^*} \log P_{\alpha^*}(x | \mathcal{Y} = i) \\ &= \nabla_{\alpha^*} \log P_{\alpha^*}(\mathcal{Y} = i) + d_{\alpha^*} \theta_i \nabla_{\theta_i} \log P_{\theta_i}(x | \mathcal{Y} = i) \end{aligned} \quad (65)$$

where in last step the chain rule is applied and $d_{\alpha^*}\theta_i$ stands for the Jacobian of θ_i with respect to α^* .

Assume the canonical parametrization is used, then according to equation 15 write

$$\alpha^* = \begin{pmatrix} P_{\eta^*}(\mathcal{Y} = 1) \\ \vdots \\ P_{\eta^*}(\mathcal{Y} = s - 1) \end{pmatrix} \quad (66)$$

From equations 66 and 16 obtain

$$\begin{aligned} \nabla_{\alpha^*} \log P_{\alpha^*}(\mathcal{Y} = i) &= \frac{1}{P_{\alpha^*}(\mathcal{Y} = i)} \begin{cases} e_{s-1}(i) & i \neq s \\ (-\mathbf{1}) & i = s \end{cases} \\ d_{\alpha^*}\theta_i &= \frac{-1}{P_{\alpha^*}(\mathcal{Y} = i)} \begin{cases} e_{s-1}(i) \cdot \theta_i^\top & i \neq s \\ (-\mathbf{1}) \cdot \theta_i^\top & i = s \end{cases} \end{aligned} \quad (67)$$

where $e_{s-1}(i)$ is the i -th canonical $s-1$ dimensional vector and $\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$. From here deduce,

$$\begin{aligned} \nabla_{\alpha^*} h(x, \eta^*) &= K_{(s-1) \times s} \cdot \text{diag}(d(x, 1, \zeta^*), \dots, d(x, s, \zeta^*)) \\ d(x, y, \zeta^*) &= \frac{1 - \theta_y^\top \nabla_{\theta_y} \log P(x | y)}{P_{\zeta^*}(y)} \end{aligned} \quad (68)$$

The part $\nabla_{\beta_k^*} h(x, \eta^*)$ follows the same steps. Observe that i -th column of $\nabla_{\beta_y^*} h(x, \eta^*)$ is

$$\begin{aligned} \nabla_{\beta_y^*} \log P_{\eta^*}(\mathcal{Y} = i, x) &= \nabla_{\beta_y^*} \log P_{\beta_y^*}(x | \mathcal{Y}) \\ &= d_{\beta_y^*}\theta_i \nabla_{\theta_i} \log P_{\theta_i}(x | \mathcal{Y} = i) \\ &= \begin{cases} 0 & y \neq i \\ \frac{\nabla_{\theta_i} \log P_{\theta_i}(x | \mathcal{Y} = i)}{P_{\zeta^*}(y)} & y = i \end{cases} \end{aligned} \quad (69)$$

and therefore the claim is proved \square

D Proof of Proposition 5

Proof. Compute $\nabla_{\theta_y} \log P_{\theta_y}(x|y)$ and proposition 3 finishes the proof. Parameters $\theta_y = (\theta_{y,1}, \dots, \theta_{y,m-1})$ are the expectation parameters of the probability distribution $P_{\theta_y}(x|y)$ which belongs to LEF.

Recall that the canonical parametrization is taken and by equation 16 deduce

$$P_{\theta_y}(x|y) = \begin{cases} \theta_{y,x} & x \neq m \\ 1 - \sum_j \theta_{y,j} & x = m \end{cases} \quad (70)$$

which clearly implies

$$\nabla_{\theta_y} \log P_{\theta_y}(x|y) = \frac{1}{P_{\theta_y}(x|y)} \begin{cases} e_{m-1}(x) & x \neq m \\ -\mathbf{1}_{m-1} & x = m \end{cases} \quad (71)$$

Finally observe

$$d(x, y, \zeta^*) = \frac{1 - \theta_y^\top \nabla_{\theta_y} \log P_{\theta_y}(x|y)}{P_{\zeta^*}(y)} = \begin{cases} 0 & x \neq m \\ \frac{1}{P_{\theta_y}(x|y)P_{\zeta^*}(y)} & x = m \end{cases} \quad (72)$$

Substitute the computations in proposition 3 to finish the proof. \square

E Proof of Proposition 4

Proof. Let A be the cost of computing $\nabla_{\theta_k} \log P_{\theta_k}(x|y)$. Prove first the next claim: the number of operations required to compute $\nabla_{\zeta^*} h(x, \zeta^*)$ is

$$s \cdot (A + 3t + 2) - 1 \quad (73)$$

and hence $O(s \cdot (A + t))$.

Indeed, terms $\nabla_{\theta_k} \log P_{\theta_k}(x|y)$ for every $y \in \mathcal{Y}$ need $s \cdot A$ operations. The cost of computing $P_{\zeta^*}(y)$ for every $y \in \mathcal{Y}$ is $s - 1$ according to equation 15 (only the term $P_{\zeta^*}(s)$ requires operations). Obtain term $d(x, y, \zeta^*)$ after $2t + 1$ operations ($2t - 1$ for the scalar product of vectors, 1 for the subtraction in the numerator and 1 last operation for the division). Since this needs to be done for every $y \in \mathcal{Y}$ then $d(x, 1, \zeta^*), \dots, d(x, s, \zeta^*)$ is known with $s \cdot (2t + 1)$ operations. Now, $\nabla_{\alpha^*} h(x, \zeta^*)$ is obtained with the product of matrices M (which is almost the identity matrix) and a diagonal matrix, which does not require any operation (it is just a transformation). Finally $\nabla_{\beta_k^*} h(x, \zeta^*)$ demands for t divisions for every $y \in \mathcal{Y}$, and therefore for $s \cdot t$ operations. Hence, the claim is proved.

To previous analysis, add the costs represented by equation 18. That is, analyze the costs of computing $q_{\mathcal{Y}}(x, P_{\eta})$ and then the products shown in that equation.

The vector $q_{\mathcal{Y}}(x, P_{\eta})$ consist on computing $P_{\eta}(y|x)$ for every $y \in \mathcal{Y}$. Using equation 13, $q_{\mathcal{Y}}(x, P_{\eta})$ needs $2t + 1$ for scalar products $T(x)^\top \beta_y$, 1 subtraction in $S(y)^\top \alpha - T(x)^\top \beta_y$ (recall that S statistic is canonical), then 1 exponentiation and finally 1 division. This is done for every $y \in \mathcal{Y}$. The denominator is the same for every y so it can be computed just once with $s - 1$ sums. The total is

$$2ts + 5s - 1 \quad (74)$$

operations.

Finally, the operations described in equation 18 are 1 for $(q_{\mathcal{Y}}(x, P_{\eta}) - e_s(y))$, $s - 1$ for $\nabla_{\alpha^*} h(x, \zeta^*) \cdot \nabla(q_{\mathcal{Y}}(x, P_{\eta}) - e_s(y))$ product and t operations for $\nabla_{\beta_y^*} h(x, \zeta^*) \cdot \nabla(q_{\mathcal{Y}}(x, P_{\eta}) - e_s(y))$, this last one needs to be done for every $y \in \mathcal{Y}$. The total operations for this block it is then

$$s + s \cdot t \quad (75)$$

To conclude the proof, the total operations needed is

$$s \cdot (A + 6t + 8) - 2 \quad (76)$$

and the complexity order is $O(s \cdot (A + t))$

□

F Proof of Theorem 4

Proof. The proof uses Robbins-Siegmund theorem as key tool. Steps taken are closely inspired by those taken in the proof of Theorem 3.2 in [5].

Compute Taylor' second order approximation of $l(\eta_{t+1})$, and after condition **C.2** apply Taylor's inequality

$$l(\eta_{t+1}) \leq l(\eta_t) - \gamma_t \nabla l(\eta_t)^T Y_t + \gamma_t^2 K \|Y_t\|^2 \quad (77)$$

Therefore, applying the expectation conditioned to information at time t obtain

$$\mathbb{E}_t[l(\eta_{t+1})] \leq l(\eta_t) - \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] + \gamma_t^2 K \mathbb{E}_t\|Y_t\|^2 \quad (78)$$

Use bound of **C.4** to third term of right hand side

$$\mathbb{E}_t[l(\eta_{t+1})] \leq l(\eta_t) - \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] + \gamma_t^2 K(A + B l(\eta_t)) \quad (79)$$

Finally, substitute $U_t = l(\eta_t)$ and arrange terms to match with equation 10

$$\mathbb{E}_t[U_{t+1}] \leq (1 + B\gamma_t^2 K)U_t - \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] + \gamma_t^2 K A \quad (80)$$

Note that theorem 2 conditions are satisfied, since condition **C.6** implies $\sum_t \beta_t = \sum_t BK\gamma^2 = BK \sum_t \gamma^2 < \infty$ and $\sum_t \epsilon_t = \sum_t KA\gamma^2 < \infty$. Hence, Robbins-Siegmund theorem ensures that $U_t = l(\eta_t)$ converges almost surely to a random variable and

$$\sum_t \zeta_t = \sum_t \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] < \infty \quad (81)$$

Now prove that $\lim_t l(\eta_t) = l(\bar{\eta})$. If $l(\eta_t)$ converges to some different random variable, condition **C.3**, second condition of **C.6** and equation 81 lead to a contradiction. Indeed, if $\lim_t l(\eta_t) = v \neq l(\bar{\eta})$, use condition **C.3** and deduce that for a fixed $0 < \delta < v - l(\bar{\eta})$ there exists an N large enough and $\epsilon > 0$ such that

$$\nabla l(\eta_t)^T \mathbb{E}_t[Y_t] \geq \epsilon \quad (82)$$

for all $t > N$. Therefore, equation 81 becomes

$$\begin{aligned} \sum_t \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] &= \sum_t^N \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] + \sum_{t>N} \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] \\ &\geq \sum_t^N \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] + \sum_{t>N} \epsilon \gamma_t \\ &\geq \epsilon \sum_{t>N} \gamma_t \end{aligned} \quad (83)$$

Second condition in **C.6** applied to right hand side of above equation assures that

$$\sum_t \gamma_t \nabla l(\eta_t)^T \mathbb{E}_t[Y_t] = \infty \quad (84)$$

which contradicts equation 81.

Finally, it is only possible that $\lim_t l(\eta_t) = l(\bar{\eta})$ almost surely as we wanted to prove. \square

G Proof of condition C.2 in Theorem 5

Proof. Compute the hessian of

$$l(\eta) = \sum_{x,y} l(\eta, y, x) \bar{P}(y, x) \quad (85)$$

The gradient of $l(\eta, y, x)$ is

$$\begin{aligned} \nabla_\alpha l(\eta, y, x) &= S \cdot (q_{\mathcal{Y}}(x) - e_s(y)) \\ \nabla_{\beta_{y'}} l(\eta, y, x) &= (q_{\mathcal{Y}}(x)_{y'} - \delta_{y=y'}) \cdot T(x) \end{aligned} \quad (86)$$

where S is the matrix having $S(i)$ as i -th column for $i \in \mathcal{Y}$. Therefore, the hessian is

$$\begin{aligned} \nabla_\alpha^2 l(\eta, y, x) &= S \cdot (\text{diag}(q_{\mathcal{Y}}(x)) - q_{\mathcal{Y}}(x) \cdot q_{\mathcal{Y}}(x)^\top) \cdot S^\top \\ \nabla_{\beta_{y_2}} \nabla_{\beta_{y_1}} l(\eta, y, x) &= \nabla_{\beta_{y_2}} q_{\mathcal{Y}}(x)_{y_1} \cdot T(x) \\ &= -T(x) \cdot T(x)^\top q_{\mathcal{Y}}(x)_{y_1} (q_{\mathcal{Y}}(x)_{y_2} - \delta_{y_1=y_2}) \\ \nabla_\alpha \nabla_{\beta_{y'}} l(\eta, y, x) &= \nabla_\alpha q_{\mathcal{Y}}(x)_{y'} \cdot T(x) \\ &= T(x) \cdot (q_{\mathcal{Y}}(x) - e_s(y'))^\top \cdot S^\top \end{aligned} \quad (87)$$

Observe how all matrices in equation 87 have their elements bounded once S and T statistics are fixed, since $\|q_{\mathcal{Y}}(x)\| \leq 1$. Therefore

$$\|\nabla^2 l(\eta, y, x)\| \leq 2K_{x,y} \quad (88)$$

for some positive numbers $K_{x,y}$. Define $K = \max_{x,y} K_{x,y}$. then finally

$$\begin{aligned} \|\nabla_\eta^2 l(\eta)\| &= \|\nabla^2 \sum_{y,x} l(\eta, y, x) \cdot \bar{P}(x, y)\| \\ &= \|\sum_{y,x} \nabla^2 l(\eta, y, x) \cdot \bar{P}(x, y)\| \\ &\leq \sum_{y,x} \|\nabla^2 l(\eta, y, x)\| \cdot \bar{P}(x, y) \\ &\leq \sum_{y,x} 2 \cdot K_{x,y} \cdot \bar{P}(x, y) \\ &\leq 2 \cdot K \sum_{y,x} \bar{P}(x, y) \\ &= 2 \cdot K \end{aligned} \quad (89)$$

□

H Proof of condition C.4 in Theorem 5

Proof. Observe that for any ϵ and t large enough there exists A_{x_t} such that

$$\begin{aligned} \|Y_t\|^2 &= (q_Y(x_t) - e_s(y_t))^\top \cdot h(x_t, \zeta_t^*)^\top h(x_t, \zeta_t^*) \cdot (q_Y(x_t) - e_s(y_t)) \\ &\leq A_{x_t} \|q_Y(x_t) - e_s(y_t)\|^2 \end{aligned} \quad (90)$$

where

$$A_{x_t} \geq \|h(x_t, \zeta_t^*)^\top h(x_t, \zeta_t^*)\| + \epsilon \quad (91)$$

This is because ζ_t converges and because of theorem 5. Now

$$\begin{aligned} \|q_Y(x_t) - e_s(y_t)\|^2 &= 1 - 2P_{\eta_t}(y_t|x_t) + \sum_y P_{\eta_t}(y|x_t)^2 \\ &\leq s + 1 \end{aligned} \quad (92)$$

therefore $\|Y_t\|^2 \leq A_{x_t}(s + 1)$ and

$$\begin{aligned} \mathbb{E}\|Y_t\|^2 &\leq \mathbb{E}[A_{x_t}(s + 1)] \\ &\leq A'(s + 1) = A \end{aligned} \quad (93)$$

where $A' = \max_x A_x$ and then condition C.4 holds. □

References

- [1] Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Robust Stochastic Approximation Approach to Stochastic Programming. *SIAM Journal on Optimization* **19**(4), 1574–1609 (2009). DOI 10.1137/070704277. URL <https://epubs.siam.org/doi/abs/10.1137/070704277>. Publisher: Society for Industrial and Applied Mathematics
- [2] Hu, J., Liu, X., Wen, Z.W., Yuan, Y.X.: A Brief Introduction to Manifold Optimization. *Journal of the Operations Research Society of China* **8**(2), 199–248 (2020). DOI 10.1007/s40305-020-00295-9. URL <https://doi.org/10.1007/s40305-020-00295-9>
- [3] Amari, S.i.: Natural Gradient Works Efficiently in Learning. *Neural Computation* **276**, 251–276 (1998)
- [4] Sánchez-López, B., Cerquides, J.: Convergent stochastic almost natural gradient descent. *Artificial Intelligence Research and Development- Proceedings of the 22nd International Conference of the Catalan Association for Artificial Intelligence* **319**, 54–63 (2019)

- [5] Sunehag, P., Trunpf, J., Vishwanathan, S.V.N., Schraudolph, N.: Variable Metric Stochastic Approximation Theory. In: Artificial Intelligence and Statistics, pp. 560–566 (2009). URL <http://proceedings.mlr.press/v5/sunehag09a.html>
- [6] Li, J., Bioucas-Dias, J.M., Plaza, A.: Spectral–spatial hyperspectral image segmentation using subspace multinomial logistic regression and markov random fields. *IEEE Transactions on Geoscience and Remote Sensing* **50**(3), 809–823 (2012). DOI 10.1109/TGRS.2011.2162649
- [7] Covington, P., Adams, J., Sargin, E.: Deep neural networks for youtube recommendations. In: Proceedings of the 10th ACM Conference on Recommender Systems, RecSys ’16, p. 191–198. Association for Computing Machinery, New York, NY, USA (2016). DOI 10.1145/2959100.2959190. URL <https://doi.org/10.1145/2959100.2959190>
- [8] Daniels, M.J., Gatsonis, C.: Hierarchical polytomous regression models with applications to health services research. *Statistics in Medicine* **16**(20), 2311–2325 (1997)
- [9] Bull, S.B., Lewinger, J.P., Lee, S.S.: Confidence intervals for multinomial logistic regression in sparse data. *Statistics in Medicine* **26**(4), 903–918 (2007)
- [10] Biesheuvel, C., Vergouwe, Y., Steyerberg, E., Grobbee, D., Moons, K.: Polytomous logistic regression analysis could be applied more often in diagnostic research. *Journal of clinical epidemiology* **61**(2), 125–134 (2008)
- [11] Leppink, J.: Multicategory nominal choices. In: *The Art of Modelling the Learning Process*, pp. 103–110. Springer (2020)
- [12] Ben-Akiva, M.E., Lerman, S.R., Lerman, S.R.: *Discrete Choice Analysis: Theory and Application to Travel Demand*. MIT Press (1985). Google-Books-ID: oLC6ZYPs9UoC
- [13] Tadei, R., Perboli, G., Manerba, D.: A Recent Approach to Derive the Multinomial Logit Model for Choice Probability. In: P. Daniele, L. Scrimali (eds.) *New Trends in Emerging Complex Real Life Problems: ODS, Taormina, Italy, September 10–13, 2018, AIRO Springer Series*, pp. 473–481. Springer International Publishing, Cham (2018)
- [14] Nock, R., Nielsen, F.: On the efficient minimization of classification calibrated surrogates. In: D. Koller, D. Schuurmans, Y. Bengio, L. Bottou (eds.) *Advances in Neural Information Processing Systems*, vol. 21, pp. 1201–1208. Curran Associates, Inc. (2009)
- [15] Reid, M.D., Williamson, R.C.: Composite Binary Losses. *Journal of Machine Learning Research* **11**(83), 2387–2422 (2010)

- [16] Vernet, E., Reid, M.D., Williamson, R.C.: Composite Multiclass Losses. In: J. Shawe-Taylor, R.S. Zemel, P.L. Bartlett, F. Pereira, K.Q. Weinberger (eds.) *Advances in Neural Information Processing Systems 24*, pp. 1224–1232. Curran Associates, Inc. (2011)
- [17] Nock, R., Menon, A.K.: Supervised Learning: No Loss No Cry. *arXiv:2002.03555 [cs, stat]* (2020)
- [18] Vapnik, V.: Principles of risk minimization for learning theory. In: *Proceedings of the 4th International Conference on Neural Information Processing Systems, NIPS'91*, pp. 831–838. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1991)
- [19] Bottou, L.: Online algorithms and stochastic approximations. In: D. Saad (ed.) *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK (1998). URL <http://leon.bottou.org/papers/bottou-98x>. Revised, oct 2012
- [20] Dennis, J.E., Schnabel, R.B.: Numerical methods for unconstrained optimization and nonlinear equations. No. 16 in *Classics in applied mathematics*. SIAM, Philadelphia, Pa (1996). OCLC: 845110213
- [21] Becker, S., Lecun, Y.: Improving the convergence of back-propagation learning with second-order methods. In: D. Touretzky, G. Hinton, T. Sejnowski (eds.) *Proceedings of the 1988 Connectionist Models Summer School, San Mateo*, pp. 29–37. Morgan Kaufmann (1989)
- [22] Duchi, J., Hazan, E., Singer, Y.: Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.* **12**(null), 2121–2159 (2011)
- [23] Zeiler, M.D.: ADADELTA: An Adaptive Learning Rate Method. *arXiv:1212.5701 [cs]* (2012). URL <http://arxiv.org/abs/1212.5701>. ArXiv: 1212.5701
- [24] Tieleman, T., Hinton, G.: Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: *Neural networks for machine learning 4*(2), 26–31 (2012)
- [25] Kingma, D.P., Ba, L.J.: Adam: A Method for Stochastic Optimization. In: *Proceedings of the 3rd International Conference on Learning Representations (ICLR)* (2015). URL <https://dare.uva.nl/search?identifier=a20791d3-1aff-464a-8544-268383c33a75>
- [26] Carmo, M.P.d.: *Riemannian geometry, corrected at 14th printing* 2013 edn. *Mathematics: theory & applications*. Birkhäuser, Boston Basel Berlin (2013)
- [27] Murray, M.K., Rice, J.W.: *Differential geometry and statistics*, vol. 48. CRC Press (1993)

- [28] Thomas, P.S.: Genga: A generalization of natural gradient ascent with positive and negative convergence results. 31st International Conference on Machine Learning, ICML 2014 **5**, 3533–3541 (2014)
- [29] Bonnabel, S.: Stochastic gradient descent on Riemannian manifolds. IEEE Transactions on Automatic Control **58**(9), 2217–2229 (2013). DOI 10.1109/TAC.2013.2254619. URL <http://arxiv.org/abs/1111.5280>. ArXiv: 1111.5280
- [30] Amari, S.i.: Information geometry and its applications, vol. 5416. Springer (2016)
- [31] Nielsen, F.: An elementary introduction to information geometry. arXiv:1808.08271 [cs, math, stat] (2018). URL <http://arxiv.org/abs/1808.08271>. ArXiv: 1808.08271
- [32] Nemirovskii A Yudin, D.: Problem Complexity and Method Efficiency in Optimization. A Wiley-Interscience publication. Wiley (1983). URL <https://books.google.es/books?id=6ULvAAAAMAAJ>
- [33] Masegosa, A.R.: Stochastic Discriminative EM. In: Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, UAI'14, pp. 573–582. AUAI Press, Arlington, Virginia, United States (2014). URL <http://dl.acm.org/citation.cfm?id=3020751.3020811>. Event-place: Quebec City, Quebec, Canada
- [34] Raskutti, G., Mukherjee, S.: The Information Geometry of Mirror Descent. IEEE Transactions on Information Theory **61**(3), 1451–1457 (2015). DOI 10.1109/TIT.2015.2388583
- [35] Beck, A., Teboulle, M.: Mirror descent and nonlinear projected subgradient methods for convex optimization. Operations Research Letters **31**(3), 167 – 175 (2003). DOI [https://doi.org/10.1016/S0167-6377\(02\)00231-6](https://doi.org/10.1016/S0167-6377(02)00231-6). URL <http://www.sciencedirect.com/science/article/pii/S0167637702002316>
- [36] Banerjee, A.: An Analysis of Logistic Models: Exponential Family Connections and Online Performance. In: Proceedings of the 2007 SIAM International Conference on Data Mining, Proceedings, pp. 204–215. Society for Industrial and Applied Mathematics (2007). DOI 10.1137/1.9781611972771.19. URL <https://epubs.siam.org/doi/abs/10.1137/1.9781611972771.19>
- [37] Wani, J.K.: On the linear exponential family. Mathematical Proceedings of the Cambridge Philosophical Society **64**(2), 481–483 (1968). DOI 10.1017/S0305004100043097
- [38] Robbins, H., Siegmund, D.: A convergence theorem for non negative almost supermartingales and some applications. In: J.S. Rustagi (ed.) Optimizing Methods in Statistics, pp. 233 – 257. Academic Press (1971)