Investigating Variance Definitions for Stochastic Mirror Descent with Relative Smoothness

Anonymous Author(s) Affiliation Address email

Abstract

Mirror Descent is a popular algorithm, that extends Gradients Descent (GD) beyond 1 the Euclidean geometry. One of its benefits is to enable strong convergence 2 guarantees through smooth-like analyses, even for objectives with exploding or 3 vanishing curvature. This is achieved through the introduction of the notion of 4 relative smoothness, which holds in many of the common use-cases of Mirror 5 descent. While basic deterministic results extend well to the relative setting, most 6 existing stochastic analyses require additional assumptions on the mirror, such as 7 strong convexity (in the usual sense), to ensure bounded variance. In this work, we 8 revisit Stochastic Mirror Descent (SMD) proofs in the (relatively-strongly-) convex 9 and relatively-smooth setting, and introduce a new (less restrictive) definition 10 of variance which can generally be bounded (globally) under mild regularity 11 assumptions. We then investigate this notion in more details, and show that it 12 13 naturally leads to strong convergence guarantees for stochastic mirror descent. Finally, we leverage this new analysis to obtain convergence guarantees for the 14 Maximum Likelihood Estimator of a Gaussian with unknown mean and variance. 15

16 1 Introduction

¹⁷ The central problem of this paper is to solve optimization problems of the following form:

$$\min_{x \in C} f(x), \text{ where } f(x) = \mathbb{E}\left[f_{\xi}(x)\right], \tag{1}$$

where C is a closed convex subset of \mathbb{R}^d , and f_{ξ} are differentiable convex functions (stochasticity is on the variable ξ). The problems that we will consider typically arise from machine-learning use-cases, meaning that the dimension d can be very large. Therefore, first-order methods are popular for solving these problems, since they usually scale well with the dimension.

In standard machine learning setups, computing a gradient of f is very costly (or even impossible), since it requires computing gradients for all individual examples in the dataset. Yet, gradients of f_{ξ} are relatively cheap, and arbitrarily high precisions are generally not required. This makes Stochastic Gradient Descent (SGD) the method of choice [4]. Using a step-size $\eta > 0$, the SGD update from point $x \in \mathbb{R}^d$ can be written as $x_{\text{SGD}}^+ = \arg \min_{u \in C} \{\eta \nabla f_{\xi}(x)^\top u + \frac{1}{2} || u - x ||^2\}$. While the standard Euclidean geometry leading to Gradient Descent (GD) fits many use-cases quite

well, several applications are better solved with *Mirror Descent* (MD), a generalization of GD which
allows to better capture the geometry of the problem. For instance, the Kullback-Leibler divergence
might be better suited to discriminating between probability distributions than the (squared) Euclidean
norm, and this is something that one can leverage using MD with entropy as a mirror. As a matter
of fact, many standard algorithms can be interpreted as MD, *i.e.*, as generalized first-order methods.
This is for instance the case in statistics, where Expectation Minimization and Maximum A Posteriori

estimators can be interpreted as running MD with specific mirror and step-sizes [15, 17]. Mirror descent can also be used to solve Poisson inverse problems, which have many applications in

astronomy and medicine [3], to reduce the communication cost of distributed algorithms [24, 12],

³⁷ or to solve convex quartic problems [6]. In the online learning community as well, many standard

algorithms such as Exponential Weight Updates or Follow-The-Regularized-Leader can be interpreted

³⁹ as running mirror descent [21, 13]. There are still many open questions regarding the convergence

40 guarantees for most of the algorithms mentioned above. Therefore, progress on the understanding of

41 MD can lead to a plethora of results on these applications, and more generally to a more consistent

theory for Majorization-Minimization algorithms. This paper is a stepping stone in this direction.

43 Let us now introduce the *mirror map*, or *potential* function h, together with the *Bregman divergence*

44 with respect to h, which is defined for $x, y \in \text{dom } h$ as $D_h(x, y) = h(x) - h(y) - \nabla h(y)^\top (x-y)$. We

⁴⁵ now introduce the Stochastic Mirror Descent (SMD) update, which can be found in its deterministic

⁴⁶ form in, *e.g.*, Nemirovskij and Yudin [22]. SMD consists in replacing the squared Euclidean norm

47 from the SGD update by the Bregman divergence with respect to the mirror map h:

$$x^{+}(\eta,\xi) = \arg\min_{u\in C} \left\{ \eta \nabla f_{\xi}(x)^{\top} u + D_{h}(u,x) \right\}.$$
 (2)

⁴⁸ Note that since $D_{\|\cdot\|^2}(x, y) = \|x - y\|^2$, one can recover SGD by taking $h = \frac{1}{2} \|\cdot\|^2$. In this sense, ⁴⁹ SMD can be viewed as standard SGD, but changing the way distances are computed, and so the ⁵⁰ geometry of the problem. Yet, this change significantly complicates the convergence analysis of the ⁵¹ method, since the Bregman divergence, *in general:* (*i*) does not satisfy the triangular inequality, (*ii*) is

i not symmetric, (*iii*) is not translation-invariant, (*iv*) is not convex in its second argument.

This means that analyzing mirror descent methods requires quite some care, and that many standard (S)GD results do not extend to the mirror setting. For instance, one can prove that mirror descent cannot be *accelerated* in general [8]. Similarly, applying techniques such as variance-reduction requires additional assumptions [7]. To ensure that $x^+(\eta, \xi)$ exists and is unique, we first make the

⁵⁷ following blanket assumption throughout the paper:

Assumption 1. Function $h : \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is twice continuously differentiable and strictly convex on *C*. For every $y \in \mathbb{R}^d$, the problem $\min_{x \in C} h(x) - x^\top y$ has a unique solution, which lies in int *C*, and all f_{ξ} are convex.

Note that the regularity assumption on *h* could be relaxed, as discussed in Section 3, but we choose a rather strong one to make sure all the objects we will manipulate are well-defined. Interestingly, while mirror descent changes the way distances are computed to move away from the Euclidean geometry, standard analyses of mirror descent methods, and in particular in the online learning community, still require strong convexity and Lipschitz continuity with respect to norms [5, Chapter 4]. It is only recently that a *relative smoothness* assumption was introduced to study mirror descent [2, 20], together with the corresponding relative strong convexity.

Definition 1. The function f is said to be L-relatively smooth and μ -relatively strongly convex with respect to h if for all $x, y \in C$: $\mu D_h(x, y) \le D_f(x, y) \le LD_h(x, y)$. To lighten notation, we will omit the dependence on h and simply write that f is L-rel.-smooth unless clearly specified.

⁷¹ Definition 1 extends the standard smooth and strongly convex assumptions that correspond to the case ⁷² $h = \frac{1}{2} \| \cdot \|^2$, so that for all $x \in C$, $\nabla^2 h(x) = I$ the identity matrix. These assumptions allow MD ⁷³ analyses to generalize standard GD analyses, and in particular to obtain similar linear and sublinear ⁷⁴ rates, with constant step-size and conditions adapted to the *relative* assumptions.

While the basic deterministic setting is now well-understood under relative assumptions, a good 75 understanding of the stochastic setting remains elusive. In particular, as we will see in more details in 76 77 the related work section, all existing proofs somehow require the mirror h to be globally strongly convex with respect to a norm, or have non-vanishing variance. The only case that can be analyzed 78 tightly is under *interpolation* (there exists a point that minimizes all stochastic functions), or when 79 using Coordinate Descent instead of SMD [10, 11]. This is a major weakness, as the goal of relative 80 smoothness is precisely to avoid comparisons to norms. Indeed, even when these "absolute" regularity 81 assumptions hold, the smoothness and strong convexity constants are typically very loose, and the 82 theory is not representative of the observed behaviour of the algorithms. 83

However, as hinted at earlier, this was expected: acceleration is notoriously hard to achieve for mirror
 descent (and even impossible in general [8]), and variance reduction typically encounters the same

problems [7]. For stochastic updates, this comes from the fact that it is impossible to disentangle the 86 stochastic gradient from the effect of the curvature of h at the point at which it is applied. 87

Contribution and outline. The main contribution of this paper is to introduce a new analysis for 88 mirror descent, with a variance notion which is provably bounded under mild regularity assumptions: 89 typically, the same as those required for the deterministic case. We introduce our new variance 90 notion, and compare it with standard ones from the literature in Section 2. This new analysis is both 91 simpler and tighter than existing ones, as shown in Section 3. Finally, we use our results to analyse 92 the convergence of the Maximum Likelihood and Maximum A Posteriori estimators for a Gaussian 93 with unknown mean and variance in Section 4, and show that it is the first generic stochastic mirror 94 descent analysis that obtains meaningful finite-time convergence guarantees in this case. 95

2 Variance Assumptions 96

We now focus on the various variance assumptions under which Stochastic Mirror Descent is analyzed. 97 Some manipulations require technical lemmas, such as the duality property of the Bregman divergence 98 or the Bregman co-coercivity lemma, which can be found in Appendix A. 99

We start by introducing our variance definition, prove a few good properties for it, and then compare 100 it with the existing ones to highlight their shortcomings. The two key properties we would like to 101 ensure (and which are not satisfied by other definitions) are: (i) boundedness without strong convexity 102 of h or restricting the SMD iterates, and (ii) finiteness for $\eta \to 0$ (with the appropriate scaling). 103

2.1 New variance definition 104

Let $\eta > 0$, and recall that $x^+(\eta, \xi)$ is the result of a SMD step from x using function f_{ξ} with step-size 105 η (Equation (2)). From now on, when clear from the context, we will simply denote this point x^+ . 106 Yet, although the dependence is now implicit, do keep in mind that x^+ is a stochastic quantity that is 107 not independent from ξ nor η , as this is critical in most results. Under Assumption 1, x^+ writes: 108

$$\nabla h(x^{+}) = \nabla h(x) - \eta \nabla f_{\xi}(x).$$
(3)

Similarly, we denote by $\overline{x^+}$ the deterministic Mirror Descent update, which is such that $\nabla h(\overline{x^+}) =$ 109

 $\nabla h(x) - \eta \nabla f(x)$. We also introduce $h^*: y \mapsto \arg \max_{x \in C} x^\top y - h(x)$ the convex conjugate of h, 110 which verifies $\nabla h^*(\nabla h(x)) = x$. Let us now define the key function 111

$$f_{\eta}(x) = f(x) - \frac{1}{\eta} \mathbb{E}\left[D_h(x, x^+)\right].$$
 (4)

Definition 2. We define the variance of the stochastic mirror descent iterates given by (2) as 112 $\sigma_{\star,\eta}^2 = \frac{1}{\eta} \sup_{x \in C} \left(f(x_\star) - f_\eta(x) \right) = \frac{f^\star - f_\eta^\star}{\eta}, \text{ where } f^\star \text{ and } f_\eta^\star \text{ are respectively the inf. of } f \text{ and } f_\eta.$ 113

We now state various bounds on $\sigma_{\star,\eta}^2$, to help understand its behaviour. We start by positivity, which is an essential property that justifies the square in the definition. 114

- 115
- **Proposition 2.1** (Positivity). For all $\eta > 0$, $\sigma_{\star,\eta} \ge 0$. 116
- This result follows from $f_{\eta}(x) \leq f(x)$, since $D_h(x, x^+) \geq 0$ for all $x \in C$ by convexity of h. 117
- **Stochastic functions after a step.** We first upper bound $\sigma_{\star,\eta}^2$ directly in terms of f_{ξ} . 118
- **Proposition 2.2.** If f_{ξ} is L-rel.-smooth and $\eta \leq 1/L$, then $\sigma_{\star,\eta}^2 \leq \frac{1}{\eta} (f(x_{\star}) \min_{x \in C} \mathbb{E}[f_{\xi}(x^+)])$. 119

Proof. Since $D_h(x, x^+) = \langle \nabla h(x^+) - \nabla h(x), x^+ - x \rangle - D_h(x^+, x)$, then $D_h(x, x^+) = -\eta \nabla f_{\xi}(x)^{\top}(x^+ - x) - D_h(x^+, x) = \eta \left(D_{f_{\xi}}(x^+, x) - f_{\xi}(x^+) + f_{\xi}(x) \right) - D_h(x^+, x)$. The relative smoothness of f_{ξ} and the step-size condition imply that $\eta D_{f_{\xi}}(x^+, x) \leq D_h(x^+, x)$, leading to 120 121 122 $\frac{1}{n}D_h(x,x^+) \leq f_{\xi}(x) - f_{\xi}(x^+)$, and the result follows. 123

This bound offers a new point of view on the variance, which can be bounded as the difference 124 between the optimum of f, and the optimum of a related function, in which we make one mirror 125 descent step before evaluating each f_{ξ} . 126

Finiteness. Proposition 2.2 implies the following: 127

Corollary 2.3. If f_{ξ} is *L*-relatively-smooth w.r.t. *h* and admits a minimum $x_{\star}^{\xi} \in \text{int } C$ a.s., then for all $\eta \leq 1/L$, $\sigma_{\star,\eta}^2 \leq \frac{f(x_{\star}) - \mathbb{E}[f_{\xi}(x_{\star}^{\xi})]}{\eta}$. In particular, $\sigma_{\star,\eta}^2$ is finite.

This result directly comes from the fact that $\min_{x \in C} \mathbb{E}[f_{\xi}(x^+)] \geq \mathbb{E}[\min_{x \in C} f_{\xi}(x^+)] \geq \mathbb{E}[f_{\xi}(x^{\xi})]$. It shows that the standard regularity assumptions for the convergence of stochastic mirror descent guarantee that *the variance as introduced in Definition 2 remains bounded*. This is a strong result, that justifies the supremum in the variance definition. Indeed, **most other variance definitions require additional assumptions for the variance to remain bounded after the supremum.** Instead, we *globalize* the variance definition, by taking the supremum over the right quantity to ensure that it remains bounded over the whole domain without having to explicitly assume it.

Note that the bound from Corollary 2.3 has already been investigating in other settings for stochastic optimization [19], as discussed in Section 2.2. While useful to show boundedness, this bound has a major drawback, which is that it explodes when the step-size η vanishes. This does not reflect what happens in practice, which is why we investigate finer bounds on $\sigma_{\star,\eta}^2$.

Gradient norm at optimum. A usual way of formulating variance is to express it as the norm of the difference between stochastic gradients and the deterministic gradients. While the previous bounds highlight dependencies on the gradient steps (through evaluations at x^+), none of them really corresponds to "the size of the stochastic gradients at optimum". The key subtlety is that when using mirror descent, it is important to also specify the point at which these gradients are applied, and the following proposition gives a bound of this flavor on $\sigma_{\star,\eta}^2$. In this section, x_{η} denotes the minimizer of f_{η} when it exists and is in int C. Otherwise, unless explicitly stated, results involving x_{η} can be replaced by a limit for $x \to x_{\eta}$.

Proposition 2.4. If
$$f$$
 is L -rel.-smooth, $\eta \leq 1/L$ and $x_{\star} \in \text{int } C$, $\sigma_{\star,\eta}^2 \leq \frac{1}{\eta^2} \mathbb{E}\left[D_h\left(\overline{x_{\eta}^+}, x_{\eta}^+\right)\right]$

This can be considered as the Mirror Descent equivalent of $\mathbb{E}\left[\|\nabla f_{\xi}(x_{\star})\|^{2}\right]$. Yet, a key difference is that stochastic gradients are evaluated at point x_{η} instead of x_{\star} , and $\nabla f(x_{\eta}) \neq 0$ in general.

152 *Proof.* For all x, applying the duality property of the Bregman divergence leads to:

$$\mathbb{E}\left[D_{h}(x,x^{+})\right] = \mathbb{E}\left[D_{h^{*}}(\nabla h(x^{+}),\nabla h(x))\right] = \mathbb{E}\left[D_{h^{*}}(\nabla h(x) - \eta \nabla f_{\xi}(x),\nabla h(x))\right]$$
$$= \mathbb{E}\left[D_{h^{*}}(\nabla h(x) - \eta \nabla f(x),\nabla h(x))\right] + \mathbb{E}\left[D_{h^{*}}(\nabla h(x) - \eta \nabla f_{\xi}(x),\nabla h(x) - \eta \nabla f(x))\right]$$
$$= \mathbb{E}\left[D_{h^{*}}(\nabla h(x) - \eta \left[\nabla f(x) - \nabla f(x_{\star})\right],\nabla h(x))\right] + \mathbb{E}\left[D_{h}\left(\overline{x^{+}},x^{+}\right)\right],$$

where the last equality comes from the Bregman bias-variance decomposition Lemma [23]. We then use the Bregman cocoercivity Lemma [7] to obtain: $\mathbb{E}[D_h(x,x^+)] \leq \eta D_f(x,x_\star) + \mathbb{E}\left[D_h\left(\overline{x^+},x^+\right)\right]$. All these technical results can be found in Appendix A. In the end, $f_\eta(x) \geq f(x_\star) - \frac{1}{\eta}\mathbb{E}\left[D_h(\overline{x^+},x^+)\right]$, and this is in particular true for $x = x_\eta$.

Limit behaviour. A first observation is that both the $D_h(x, x^+)$ term in the definition of f_η and our variance definition are scaled by η^{-1} . Yet, they remain finite when $\eta \to 0$. While this is clear in the Euclidean setting, this property holds more generally, as shown in the two following results.

160 **Proposition 2.5.** Let
$$x \in C$$
 and $\eta_0 > 0$ s.t. $\mathbb{E}D_h(x, x^+(\eta_0, \xi)) < \infty$. Then, $f_\eta(x) \xrightarrow{\eta \to 0} f(x)$.

Note that uniform convergence of f_{η} to f would require that there exists $\eta > 0$ such that sup_{$x \in C$} $D_h(x, x^+)$ is finite, which we cannot guarantee in general (it does not hold for $f = g = \frac{1}{2} \|\cdot\|^2$ defined on \mathbb{R}^d for instance). Denote $\|x\|_A^2 = x^\top Ax$, then:

Proposition 2.6 (Small step-sizes limit). If f_{ξ} are *L*-rel.-smooth and *f* has a unique minimizer x_{\star} and for some $\eta_0 > 0$, $x_{\eta} = \arg \min f_{\eta}(x)$ exists and is in int *C* for $\eta \le \eta_0$,

$$\lim_{\eta \to 0} \sigma_{\star,\eta}^2 = \lim_{\eta \to 0} \frac{1}{\eta^2} \mathbb{E} \left[D_h(x_{\star}^+, x_{\star}) \right] = \frac{1}{2} \mathbb{E} \left[\| \nabla f_{\xi}(x_{\star}) \|_{\nabla^2 h(x_{\star})^{-1}}^2 \right].$$
(5)

This variance is actually the best we can hope for in the Bregman setting, which indicates the 166 relevance of Definition 2. Indeed, this term exactly correspond to the variance one would obtain 167 when making infinitesimal SMD steps from x_{\star} , *i.e.*, the norm of the stochastic gradients at optimum 168 in the geometry given by $\nabla^2 h(x_\star)^{-1}$. 169

2.2 Standard Assumptions 170

We now compare Definition 2 with several variance assumptions from the literature. Note that they 171 typically "only" require the bounds to hold for all iterates over the trajectory. However, in the absence 172 of proof that the iterates stay in certain regions of the space, suprema over the whole domain are 173 required for all variance definitions. 174

Euclidean case. Let us now take a step back and look at the Euclidean case, $h = \frac{1}{2} \|\cdot\|^2$, and assume that f is L-smooth. Writing Equation (3) with this specific h and replacing x_η by a supremum, we obtain $\sigma_{\star,\eta}^2 \leq \sup_{x \in C} \mathbb{E} \left[\frac{1}{2} \|\nabla f(x) - \nabla f_{\xi}(x)\|^2 \right]$, which is a common though debatable variance assumption. Indeed, it involves a maximum over the domain, and is in particular not bounded in 175 176 177 178 general even for simple examples like Linear Regression. Yet, we can recover another standard 179 variance assumption by assuming the smoothness of all f_{ξ} [9], which writes $\sigma_{\star,n}^2 \leq \mathbb{E} \left[\|\nabla f_{\xi}(x_{\star})\|^2 \right]$. 180

This result is obtained by writing that $\|\nabla f_{\xi}(x)\|^2 \leq 2\|\nabla f_{\xi}(x) - \nabla f_{\xi}(x_{\star})\|^2 + 2\|\nabla f_{\xi}(x_{\star})\|^2$, and bounding the first term using smoothness. In particular, we see that standard Euclidean variance definitions are natural bounds of $\sigma_{\star,\eta}^2$. Detailed derivations can be found in Appendix B. 181 182 183

Divergence between stochastic and deterministic gradients. An early variance definition for SMD 184 in the relative setting comes from Hanzely and Richtárik [10], who define σ_{sym}^2 as: 185

$$\sigma_{\text{sym}}^2 = \frac{1}{\eta} \sup_{x \in C} \mathbb{E}\left[\left\langle \nabla f(x) - \nabla f_{\xi}(x), x^+ - \overline{x^+}\right\rangle\right] = \frac{1}{\eta^2} \sup_{x \in C} \mathbb{E}\left[D_h\left(x^+, \overline{x^+}\right) + D_h\left(\overline{x^+}, x^+\right)\right],$$

186

187 188

where we recall that $\overline{x^+}$ is such that $\nabla h(\overline{x^+}) = \nabla h(x) - \eta \nabla f(x)$. We remark two main things when comparing σ_{sym}^2 with Proposition 2.4: (i) $\sigma_{\star,\eta}^2$ is not symmetrized, and contains only one of the two terms, and (ii) the bound only needs to hold at x_η instead of for all $x \in C$. As a result, we directly obtain that $\sigma_{\star,\eta}^2 \leq \sigma_{\text{sym}}^2$, and σ_{sym}^2 is actually infinite in most cases, whereas $\sigma_{\star,\eta}^2$ is usually finite, as 189 seen above. 190

Stochastic gradients at optimum. Dragomir et al. [7] define the variance as: 191

$$\sigma_{DEH}^{2} = \sup_{x \in C} \frac{1}{2\eta^{2}} \mathbb{E} \left[D_{h^{*}} (\nabla h(x) - 2\eta \nabla f_{\xi}(x_{\star}), \nabla h(x)) \right] = \sup_{x \in C} \mathbb{E} \left[\|\nabla f_{\xi}(x_{\star})\|_{\nabla^{2}h^{*}(z(x))}^{2} \right],$$

where $z(x) \in [\nabla h(x), \nabla h(x) - \eta \nabla f_{\xi}(x_{\star})]$ The main interest of this definition is that stochastic 192 gradients are only taken at x_{\star} . In particular, this variance is 0 in case there is interpolation (all 193 stochastic functions share a common minimum). However, this quantity can blow up if h is not 194 strongly convex, since in this case $\nabla^2 h^*$ is not upper bounded (indeed, smoothness of the conjugate 195 is ensured by strong convexity of the primal function [14]). Following similar derivations, but after 196 the supremum has been taken, we arrive at: 197

Proposition 2.7. If f is L-relatively-smooth w.r.t. h, then for $\eta < 1/(2L)$ and some $z_{\eta} \in$ 198 $[\nabla h(x_{\eta}), \nabla h(x_{\eta}) - \eta \nabla f_{\xi}(x_{\star})], \text{ the variance can be bounded as } \sigma^{2}_{\star,\eta} \leq \mathbb{E} \left[\|\nabla f_{\xi}(x_{\star})\|^{2}_{\nabla^{2}h^{*}(z_{\eta})} \right].$ 199

In particular, we obtain a finite bound without having to restrict the space. 200

Functions variance. Another variance definition that appears in the SGD literature is of the form 201 $f(x_{\star}) - \mathbb{E} \left| f_{\xi}(x_{\star}^{\xi}) \right|$, using the optima of the stochastic functions [19]. Unfortunately, the results 202 derived with this definition do not obtain a vanishing variance term when $\eta \to 0$, unlike most other 203 variance definitions, and contrary to what is observed in practice, that smaller step-sizes reduce 204 the variance. The vanishing variance term can be obtained by rescaling by $1/\eta$ (so considering 205 $\left(f(x_{\star}) - \mathbb{E}\left|f_{\xi}(x_{\star}^{\xi})\right|\right)/\eta$ instead), but this variance definition would explode for $\eta \to 0$. This is 206 because using such a definition would come down to performing the supremum step within the 207 expectation from Proposition 2.2, using that $f_{\xi}(x^+) \geq f_{\xi}(x^{\xi})$, which is a very crude bound. Instead, 208 Corolary 2.3 directly shows that our variance definition is tighter than this one, and in particular (i) it is 209 bounded for all $\eta > 0$, (ii) it remains finite as $\eta \to 0$ even with the proper rescaling (Proposition 2.6). 210

Relation to *c***-transform.** Mirror descent can be viewed as an alternate minimization method on 211 transforms of f [18]. This point of view subsumes many methods, including the Newton Method or 212 Mirror Descent. Central to their analysis is the notion of c-transform $f^c(y) = \sup_{x \in C} f(x) - c(x, y)$, a standard quantity from optimal transport [25]. It turns out that for $\eta \leq 1/L$, f_{η} is actually 213 214 linked to the c-transform as $f_{\eta}(x) = \mathbb{E} \left| f_{\xi}^{c}(x^{+}) \right|$, where we use the cost $c(x,y) = \frac{1}{\eta} D_{h}(x,y)$. 215 Since $f(x_{\star}) = f^c(x_{\star}) = \arg \min_{x \in C} f(\overline{x^+})$, denoting $\mathcal{T}_c(g) = g^c(\nabla h^*(\nabla h(x) - \eta \nabla g(x)))$, we have that $\sigma^2_{\star,\eta} = \frac{1}{\eta}(\min_{x \in C} \mathcal{T}_c(\mathbb{E}[f_{\xi}])(x) - \min_{x \in C} \mathbb{E}[\mathcal{T}_c(f_{\xi})](x))$. We recognize the structure of a 216 217 variance, as the difference between an operator applied to the expectation of a random variable, and 218 the expectation of the operator applied to the random variable. Yet, compared to standard (Euclidean) 219 analyses of SGD, it does not simply corresponds to the variance of the stochastic gradients (at 220 optimum), and bears a more complex form. 221

In this section, we have highlighted the connections with other definitions, and argued that f_{η} (and its minimum) is a relevant quantity. In particular, Definition 2 is the only definition that allows boundedness of the variance notion both after a supremum step over the iterates (and without strong convexity of h) and in the $\eta \to 0$ limit with the proper rescaling.

226 3 Convergence Analysis

Now that we have (extensively) investigated $\sigma_{\star,\eta}^2$, and the various interpretations that come from different bounds, we are ready to state the convergence results. Some proofs in this section are just sketched, but complete derivations can be found in Appendix C.

230 3.1 Relatively Strongly Convex setting.

Recall that $f_{\eta}^{\star} = \inf_{x \in C} f_{\eta}(x)$. Starting from an arbitrary $x^{(0)}$, the sequence $(x^{(k)})_{k \geq 0}$ is built as $x^{(k+1)} = (x^{(k)})^+$ for $k \in \{0, T\}$ for some T > 0

Theorem 3.1. If f is μ -relatively-strongly-convex with respect to h, under a constant step-size η , the iterates obtained by SMD (Equation (3)) verify

$$\eta \left[\mathbb{E} \left[f_{\eta}(x^{(T)}) \right] - f_{\eta}^{\star} \right] + \mathbb{E} \left[D_{h}(x_{\star}, x^{(T+1)}) \right] \le (1 - \eta \mu)^{T+1} D_{h}(x_{\star}, x^{(0)}) + \frac{\eta \sigma_{\star, \eta}^{2}}{\mu}.$$
(6)

Note that the (relatively) strongly-convex theorem has a standard form, and recovers usual MD results if we remove the variance, and standard SGD results if we take $h = \frac{1}{2} \| \cdot \|^2$.

237 *Proof of Theorem 3.1.* We start from a variation of Dragomir et al. [7, Lemma 4]:

$$\mathbb{E}\left[D_h(x_\star, x^+)\right] - D_h(x_\star, x) + \eta D_f(x_\star, x) = -\eta[f(x) - f(x_\star)] + \mathbb{E}\left[D_h(x, x^+)\right]$$
(7)

$$= \eta \left[f(x_{\star}) - \left(f(x) - \frac{1}{\eta} \mathbb{E} \left[D_h(x, x^+) \right] \right) \right] = \eta \left[f(x_{\star}) - f_\eta(x) \right]$$
(8)

$$= -\eta \left[f_{\eta}(x) - f_{\eta}^{\star} \right] + \eta \left[f(x_{\star}) - f_{\eta}^{\star} \right].$$
(9)

Using that $D_f(x_\star, x) \ge \mu D_h(x_\star, x)$, and remarking that $f(x_\star) - f_\eta^\star = \eta \sigma_{\star,\eta}^2$, we obtain:

$$\eta \left[f_{\eta}(x) - f_{\eta}^{\star} \right] + \mathbb{E} \left[D_{h}(x_{\star}, x^{+}) \right] \le (1 - \eta \mu) D_{h}(x_{\star}, x) + \eta^{2} \sigma_{\star, \eta}^{2}.$$
(10)

At this point, we can neglect the $\eta \left[f_{\eta}(x) - f_{\eta}^{\star} \right] \ge 0$ terms and chain the inequalities for $x = x^{(t)}$ for t from 0 to T to obtain the result.

This proof is quite simple, and naturally follows from Lemma C.1. One can also note that *relative smoothness of f is not required to obtain Theorem 3.1*, which has no condition on the step-size. This is not a typo, but reflects the fact that *step-size conditions are needed to obtain a bounded variance*. Indeed, the variance as defined here entangles aspects tied with the error due to discretization (which is usually dealt with using smoothness), and the error due to stochasticity. This is natural, as the stochastic noise vanishes in the continuous limit ($\eta \rightarrow 0$). Besides, the magnitude of the updates depends both on where the stochastic gradient is applied and on the step-size. Yet, the simplicity of

- the proof is partly due to this entanglement, meaning that we have deferred some of the complexityto the bounding of the variance term.
- Also note that Theorem 3.1 uses constant step-sizes, but Equation (10) can be used with time-varying step-sizes, as is done for instance in the proof of Theorem 4.3. A variant of Theorem 3.1 in which the discretization error is partly removed from the notion of variance writes:
- **Corollary 3.2.** Let f be μ -strongly-convex and L-relatively-smooth with respect to h, and $f_+^{\star} = \inf_{x \in C} \mathbb{E} [f_{\xi}(x^+)]$. If $\eta \leq 1/L$, the SMD iterates (Equation (3)) with constant step-size η verify

$$\eta \left[\mathbb{E} \left[f_{\xi}((x^{(T)})^{+}) \right] - f_{+}^{\star} \right] + \mathbb{E} \left[D_{h}(x_{\star}, x^{(T+1)}) \right] \leq (1 - \eta \mu)^{T+1} D_{h}(x_{\star}, x^{(0)}) + \frac{\eta}{\mu} \left[\frac{f(x_{\star}) - f_{+}^{\star}}{\eta} \right]$$

This alternate version is obtained using that $f_{\eta}(x) \geq \mathbb{E}[f_{\xi}(x^+)]$, a key step from the proof of Proposition 2.2 (see (8)). In the deterministic case, $f_{+}^{\star} = f(x_{\star})$, and we recover standard results.

257 3.2 Convex setting.

- Let us now consider the convex case, meaning that $\mu = 0$.
- **Theorem 3.3.** If f is convex, the iterates obtained by SMD using a constant step-size $\eta > 0$ verify

$$\frac{1}{T+1}\sum_{k=0}^{T} \mathbb{E}\left[f_{\eta}(x^{(k)}) - f_{\eta}^{\star} + D_{f}(x_{\star}, x^{(k)})\right] \le \frac{D_{h}(x_{\star}, x^{(0)})}{\eta(T+1)} + \eta \sigma_{\star, \eta}^{2}.$$
 (11)

This theorem is obtained by summing Equation (9) for $x = x^{(k)}$ for all $k \in \{1, ..., T\}$ and rearranging the terms. Note that varying step-size results can be obtained in the same way.

This case differs from standard convex analyses, in that we obtain a control on $f_{\eta}(x^{(k)}) - f_{\eta}^{\star} + D_f(x_{\star}, x^{(k)})$ instead of the usual $f(x^{(k)}) - f(x_{\star})$. One of the main consequences is that we cannot get a control on the average iterate since Bregman divergences are in general not convex in their second argument, and f_{η} is not necessarily convex. This non-standard result is a direct consequence of our choice of variance definition, but it is actually a quantity that naturally arises in the analysis. Note that a variant involving f_{\pm}^{\star} can be obtained in the same lines as Corollary 3.2.

Controlling f_{η} . The results in this section do not directly control the function gap $f(x) - f^*$, but rather the transformed one $f_{\eta}(x) - f^*_{\eta}$. Yet, the continuity result (in η) from Proposition 2.5 shows that the bounds we provide can still be interpreted as relevant function values for small η .

Controlling $D_f(x_\star, x^{(k)})$. An interesting property of $D_f(x_\star, x^{(k)})$ is that it can be linked with the size of the gradients of f, as shown by the following result.

Proposition 3.4. If $\nabla f(x_{\star}) = 0$ and f is L-relatively smooth with respect to h then for all $x \neq x_{\star}$, $D_f(x_{\star}, x) \ge LD_{h^*}\Big(\nabla h(x_{\star}) + \frac{\nabla f(x)}{L}, \nabla h(x_{\star})\Big) > 0.$

This is a Bregman equivalent of controlling the gradient squared norm, with the additional benefit that the reference point at which we apply the gradient is the optimum x_{\star} . Besides, Proposition 3.4 shows that $D_f(x_{\star}, x) > 0$ for $x \neq x_{\star}$ without requiring f to be strictly convex (only h).

Minimal assumptions on h. Note that the theorems in this section do not actually require h to satisfy Assumption 1, but only that iterations can be written in the form of Equation 3 (which is guaranteed by Assumption 1). While Assumption 1 allows for instance to use the Bregman cocoercivity lemma with any points, or ensures that $\nabla^2 h$ is well-defined, which we leverage extensively in Section 2, our theorems are much more general than this, and include applications such as proximal gradient mirror descent (next remark) or the MAP for Gaussian Parameters Estimation (next section).

Stochastic Mirror Descent with a Proximal term. Note that our results can be directly extended to handle a proximal term (similarly to the Euclidean proximal gradient algorithm), to handle composite objectives of the form f + g (and in particular projections, for cases in which g is the indicator of a convex set). More details can be found in Appendix E.

4 MAP For Gaussian Parameters Estimation.

So far, we have proposed new variance definitions for the analysis of stochastic mirror descent, and we have shown that they compare favorably to existing ones, while leading to simple convergence proofs. In this section, we investigate the open problem formulated by Le Priol et al. [17], which is to find non-asymptotic convergence guarantees for the KL-divergence of the Maximum A Posteriori (MAP) estimator. In particular, this example highlights the relevance of the infimum step on f_{η} , since it gives the first generic analysis that obtains meaningful finite time convergence rates.

295 4.1 MAP and MLE of exponential families.

We now rapidly review the formalism of exponential families. More details can be found in Le Priol 296 et al. [17], and Wainwright et al. [26, Chapter 3]. Let X be a random variable, and T a deterministic 297 298 function, then the density of an exponential family for a sample x writes $p_{\theta}(x) = p(x|\theta) =$ $\exp(\langle \theta, T(x) \rangle - A(\theta))$, where A is often referred to as the log-partition function. In this case, θ is called 299 the natural parameter, and T is the sufficient statistic. Function A is convex, and we can thus establish 300 a form of duality through convex conjugacy. The *entropy* writes $A^*(\mu) = \max_{\theta' \in \Theta} \langle \mu, \theta' \rangle - A(\theta')$. Parameter μ is called the *mean* parameter, and the standard MAP estimator can be derived for $n_0 \in \mathbb{N}$, $\mu_0 \in \mathbb{R}$ as $\mu_{\text{MAP}}^{(n)} = \frac{n_0 \mu^{(0)} + \sum_{i=1}^n T(X_i)}{n_0 + n}$. The Maximum Likelihood Estimator (MLE) corresponds to 301 302 303 taking $n_0 = 0$. An interesting observation is that $\mu_{MAP}^{(n)}$ can be obtained recursively for n > 0, as $\mu_{MAP}^{(0)} = \mu^{(0)}, \eta_n = (n+n_0)^{-1}, \mu_{MAP}^{(n+1)} = \mu_{MAP}^{(n)} - \eta_n \nabla g_{X_n}(\mu_{MAP}^{(n)})$, with $\nabla g_{X_n}(\mu) = \mu - T(X_n)$. In terms of primal variable $\theta^{(n)} = \nabla A^*(\mu_{MAP}^{(n)})$, the MAP writes: 304 305 306

$$\nabla A(\theta^{(n+1)}) = \nabla A(\theta^{(n)}) - \eta \nabla f_{X_n}(\theta^{(n)}), \tag{12}$$

where $f_{X_n}(\theta) = A(\theta) - \langle \theta, T(X_n) \rangle$, so that $f(\theta) = A(\theta) - \langle \theta, \mu_* \rangle$. We recognize stochastic mirror descent iterations, with mirror A and stochastic gradients ∇f_X . Similar results on the MLE can be obtained by taking $n_0 = 0$. This key observation implies that **convergence guarantees on the MAP** and the MLE can be deduced from stochastic mirror descent convergence guarantees.

While this appears as an appealing way to obtain convergence guarantees for the MAP, Le Priol et al. [17] observe that none of the existing SMD results obtain meaningful rates for the convergence of the MAP for general exponential families. In particular, none of them recover the O(1/n) asymptotic convergence rate for estimating a Gaussian with unknown mean and covariance.

This is due to the variance definitions used in the existing analyses, that all have issues (not uniformly bounded over the domain, not decreasing with the step-size...) as discussed in Section 2. Our analysis fixes this problem, and thus yields finite-time guarantees for the MAP estimator for the estimation of a Gaussian with unknown mean and covariance. This shows the relevance of Assumption 2.

319 4.2 Full Gaussian (unknown mean and covariance)

The main problem studied in Le Priol et al. [17] is that of the one-dimensional full-Gaussian 320 case, where the goal is to estimate the mean and covariance of a Gaussian from i.i.d. samples 321 $X_1, \ldots, X_n \sim \mathcal{N}(m_\star, \Sigma_\star)$, with $\Sigma_\star > 0$. Note that although notation Σ is usually reserved for 322 the covariance matrix of a multivariate Gaussian, we use it for a scalar value here to highlight the 323 distinction with $\sigma_{\star,\eta}^2$, the variance from stochastic mirror descent. In this case, the sufficient statistics 324 write $T(X) = (X, X^2)$, and the log-partition and entropy functions are, up to constants, $A(\theta) =$ 325 $\frac{\theta_1^2}{-4\theta_2} - \frac{1}{2}\log(-\theta_2), A^*(\mu) = -\frac{1}{2}\log(\mu_2 - \mu_1^2), \text{ for } \theta \in \Theta = \mathbb{R} \times \mathbb{R}^*_- \text{ and } \mu \in \{(u, v), u^2 < v\}.$ The goal is to estimate $D_A(\theta, \theta_*)$, for which Le Priol et al. [17] show that only partial solutions 326 327 exist: results are either asymptotic, or rely on the objective being (approximately) quadratic. Note 328 that there is a relationship between natural parameters, mean parameters, and (m, Σ^2) , the mean and 329 covariance of the Gaussian we would like to estimate. In the following, we will often abuse notations, 330 and write for instance $D_A(\tilde{\theta}, \theta)$ in terms of (m, Σ^2) and $(\tilde{m}, \tilde{\Sigma}^2)$ rather than θ and $\tilde{\theta}$. We now state a 331 few results, for which detailed derivations can be found in Appendix F. More specifically: 332

$$D_A(\tilde{\theta}, \theta) = -\frac{1}{2} \log \left(\frac{\Sigma^2}{\tilde{\Sigma}^2}\right) - \frac{\tilde{\Sigma}^2 - \Sigma^2}{2\tilde{\Sigma}^2} + \frac{(\tilde{m} - m)^2}{2\tilde{\Sigma}^2}.$$

The update formulas for the parameters are given by: 333

$$m^+ = (1 - \eta)m + \eta X,$$
 $(\Sigma^2)^+ = (1 - \eta) \left[\Sigma^2 + \eta (m - X)^2\right].$ (13)

Therefore, MAP iterations are well-defined although A does not verify Assumption 1. 334

Proposition 4.1. The iterations (12) are well-defined for $\eta < 1$ in the sense that if $\theta^{(n)} \in \Theta = \mathbb{R} \times \mathbb{R}^*_-$, 335

then $\nabla A(\theta^{(n)}) - \eta \nabla f_{X_n}(\theta^{(n)}) \in \text{Range}(\nabla A)$ almost surely, so that $\theta^{(n+1)} \in \Theta$ is well-defined 336

almost surely. Besides, f_{ξ} is 1-relatively-smooth and 1-relatively-strongly-convex with respect to A. 337

This result is a direct consequence of the fact that $D_{f_{\xi}} = D_f = D_A$ for all ξ , and the fact that 338

 $\nabla A(\theta) - \eta \nabla f_{X_n}(\theta) = (1 - \eta) \nabla A(\theta) + \eta T(X_n) \in \{(u, v), u^2 < v\} \text{ if } \nabla A(\theta) \in \{(u, v), u^2 < v\}.$ Proposition 4.1 means that we can apply Theorem 3.1, so the next step is to bound the variance $\sigma^2_{\star,\eta}$. 339

340

$$f_{\eta}(\theta) - f(\theta_{\star}) = \frac{1}{2\eta} \mathbb{E}\left[\log\left((1-\eta)\left(1+\eta\frac{(m-X)^2}{\Sigma^2}\right)\right)\right] - \frac{1}{2}\log\left(\frac{\Sigma_{\star}^2}{\Sigma^2}\right).$$
(14)

We now use this expression to to lower bound f_{η}^{\star} and so upper bound $\sigma_{\star,\eta}^2$. 341

Lemma 4.2. Let $(m_{\eta}, \Sigma_{\eta}^2)$ be the minimizer of f_{η} . Then, for $\eta < 1/3$, $m_{\eta} = m_{\star}$, $\Sigma_{\star}^2 \ge \Sigma_{\eta}^2 \ge (1-3\eta)\Sigma_{\star}^2$. In particular, the variance $\sigma_{\star,\eta}^2$ verifies $\sigma_{\star,\eta}^2 \le -\frac{1}{2\eta}\log(1-3\eta)$. For $1/3 < \eta \le 1-\varepsilon$, $\sigma_{\star,\eta}^2 \le c_{\varepsilon}$, where c_{ε} is a numerical constant that only depends on ε . 342

343

344

Note that we show in this example that Σ_{η}^2 is arbitrarily close to Σ_{\star}^2 as $\eta \to 0$, which is expected. 345

Theorem 4.3. Let $\Gamma \ge 0$ be a numerical constant and $\Gamma = 0$ if $n_0 > 3$. The MAP estimator satisfies: 346

$$\mathbb{E}\left[D_A(\theta_\star, \theta^{(n)})\right] \le \frac{n_0 D_A(\theta_\star, \theta^{(0)}) + \frac{3}{2}\log(1 + \frac{n+1}{n_0}) + \Gamma}{n+n_0}$$

Numerical constants are not optimized. Theorem 4.3 gives an anytime result on the convergence of 347 the MAP estimator for all $n \ge 0, n_0 \ge 1$ directly from the general SMD convergence theorem. Yet, 348 the open problem from Le Priol et al. [17] is not completely solved still, as discussed below. 349

Reverse KL bound. We obtain a bound on $D_A(\theta_\star, \theta^{(n)})$, instead of $D_A(\theta^{(n)}, \theta_\star) = f(\theta) - f(\theta_\star)$. 350 $D_A(\theta^{(n)}, \theta_{\star})$ can be controlled asymptotically thanks to the bound on $f_{\eta}(\theta^{(n)}) - f_{\eta}(\theta_{\eta})$, and $f_{\eta} \to f$ when $\eta = 1/n \to 0$, but we might also be able to exploit this control over the course of the iterations. 351 352

Asymptotic convergence. Theorem 4.3 leads to a $O(\log n/n)$ asymptotic convergence rate instead 353 of the expected O(1/n) [17]. This indicates that the $f_{\eta_n}(\theta^{(n)}) - f_{\eta_n}^{\star}$ terms should not be neglected. 354 Indeed, $\theta^{(n)}$ actually has a lot of structure in this example, since $\nabla A(\theta^{(n)}) = \frac{1}{n} \sum_{k=1}^{n} T(X_k)$. The SMD analysis is oblivious to this structure, hence the gap. Note that we can get rid of the $\log n$ factor 355 356 and recover the right O(1/n) rate from the same analysis by using a slightly different estimator than the MAP (or MLE). This is done by setting the step-size as $\eta_n = \frac{2}{n+1}$ for n > 1, and the analysis of this variant follows Lacoste-Julien et al. [16], as detailed in Appendix F.3. 357 358 359

The special case of the MLE. The MLE corresponds to $n_0 = 0$, which is not handled in our analysis 360 since the first step corresponds to $\eta = 1$, which necessarily results in $\theta_2^{(1)} = -\infty$ (which corresponds 361 to $\Sigma^2 = 0$, as can be seen from (13)). If we consider that mirror descent is run from $\theta^{(1)}$, then we 362 obtain $\mathbb{E}\left[D_A(\theta_\star, \theta^{(1)})\right] = \infty$ in general, where the expectation is over the value of the first sample 363 drawn. Therefore, we need to start the SMD analysis at $\theta^{(2)}$ to fit the MLE into this framework, and in 364 particular we need to be able to evaluate $\mathbb{E}\left[D_A(\theta_*, \theta^{(2)})\right]$. This is further discussed in Appendix F.4. 365

Conclusion 5 366

This paper introduces a new notion of variance for the analysis of stochastic mirror descent. This 367 notion, based on the fact that a certain function f_{η} admits a minimum, is less restrictive than existing 368 ones, has the right asymptotic scaling with the step-size and is bounded regardless of the trajectory of 369 the iterates without further assumptions. 370

We strongly believe that our analysis of SMD opens up new perspectives. As an example, we use our 371 SMD results to show convergence of the MAP for estimating a Gaussian with unknown mean and 372 covariance. As evidenced in Le Priol et al. [17], all existing generic analyses of stochastic mirror 373 descent failed to obtain such results. 374

375 **References**

- [1] H. H. Bauschke, J. M. Borwein, et al. Legendre functions and the method of random bregman projections. *Journal of convex analysis*, 4(1):27–67, 1997.
- [2] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond lipschitz gradient continuity: first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.
- [3] M. Bertero, P. Boccacci, G. Desiderà, and G. Vicidomini. Image deblurring with poisson data: from cells to galaxies. *Inverse Problems*, 25(12):123006, 2009.
- [4] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010: 19th International Conference on Computational StatisticsParis France, August 22-27, 2010 Keynote, Invited and Contributed Papers*, pages 177–186. Springer, 2010.
- [5] S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- [6] R.-A. Dragomir and Y. Nesterov. Convex quartic problems: homogenized gradient method and preconditioning. *arXiv preprint arXiv:2306.17683*, 2023.
- [7] R. A. Dragomir, M. Even, and H. Hendrikx. Fast stochastic bregman gradient methods: Sharp
 analysis and variance reduction. In *International Conference on Machine Learning*, pages
 2815–2825. PMLR, 2021.
- [8] R.-A. Dragomir, A. B. Taylor, A. d'Aspremont, and J. Bolte. Optimal complexity and certifica tion of bregman first-order methods. *Mathematical Programming*, pages 1–43, 2021.
- [9] R. M. Gower, N. Loizou, X. Qian, A. Sailanbayev, E. Shulgin, and P. Richtárik. Sgd: General analysis and improved rates. In *International conference on machine learning*, pages 5200–5209.
 PMLR, 2019.
- [10] F. Hanzely and P. Richtárik. Fastest rates for stochastic mirror descent methods. *Computational Optimization and Applications*, 79:717–766, 2021.
- [11] H. Hendrikx, F. Bach, and L. Massoulié. Dual-free stochastic decentralized optimization with
 variance reduction. *Advances in neural information processing systems*, 33:19455–19466, 2020.
- [12] H. Hendrikx, L. Xiao, S. Bubeck, F. Bach, and L. Massoulie. Statistically preconditioned
 accelerated gradient method for distributed optimization. In *International conference on machine learning*, pages 4203–4227. PMLR, 2020.
- [13] D. Hoeven, T. Erven, and W. Kotłowski. The many faces of exponential weights in online
 learning. In *Conference On Learning Theory*, pages 2067–2092. PMLR, 2018.
- [14] S. Kakade, S. Shalev-Shwartz, A. Tewari, et al. On the duality of strong convexity and
 strong smoothness: Learning applications and matrix regularization. Unpublished Manuscript,
 http://ttic. uchicago. edu/shai/papers/KakadeShalevTewari09. pdf, 2(1):35, 2009.
- [15] F. Kunstner, R. Kumar, and M. Schmidt. Homeomorphic-invariance of em: Non-asymptotic
 convergence in kl divergence for exponential families via mirror descent. In *International Conference on Artificial Intelligence and Statistics*, pages 3295–3303. PMLR, 2021.
- [16] S. Lacoste-Julien, M. Schmidt, and F. Bach. A simpler approach to obtaining an o (1/t) convergence rate for the projected stochastic subgradient method. *arXiv preprint arXiv:1212.2002*, 2012.
- [17] R. Le Priol, F. Kunstner, D. Scieur, and S. Lacoste-Julien. Convergence rates for the map
 of an exponential family and stochastic mirror descent–an open problem. *arXiv preprint arXiv:2111.06826*, 2021.
- [18] F. Léger and P.-C. Aubin-Frankowski. Gradient descent with a general cost. *arXiv preprint arXiv:2305.04917*, 2023.

- [19] N. Loizou, S. Vaswani, I. H. Laradji, and S. Lacoste-Julien. Stochastic polyak step-size for
 sgd: An adaptive learning rate for fast convergence. In *International Conference on Artificial Intelligence and Statistics*, pages 1306–1314. PMLR, 2021.
- 424 [20] H. Lu, R. M. Freund, and Y. Nesterov. Relatively smooth convex optimization by first-order 425 methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- 426 [21] B. McMahan. Follow-the-regularized-leader and mirror descent: Equivalence theorems and
- 11 regularization. In Proceedings of the Fourteenth International Conference on Artificial
- *Intelligence and Statistics*, pages 525–533. JMLR Workshop and Conference Proceedings, 2011.
- [22] A. S. Nemirovskij and D. B. Yudin. Problem complexity and method efficiency in optimization.
 1983.
- [23] D. Pfau. A generalized bias-variance decomposition for bregman divergences. Unpublished
 Manuscript, 2013.
- 434 [24] O. Shamir, N. Srebro, and T. Zhang. Communication-efficient distributed optimization using
 435 an approximate newton-type method. In *International conference on machine learning*, pages
 436 1000–1008. PMLR, 2014.
- 437 [25] C. Villani et al. Optimal transport: old and new, volume 338. Springer, 2009.
- [26] M. J. Wainwright, M. I. Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.

440 A Technical results on Bregman divergences

As for the rest of this paper, Assumption 1 is assumed throughout this section. However, some of these results hold even with less regularity, and in particular do not require second order continuous

443 differentiability.

444 **Lemma A.1** (Duality). For all $x, y \in C$, it holds that:

$$D_h(x,y) = D_{h^*}(\nabla h(y), \nabla h(x))$$
(15)

- 445 See, *e.g.* Bauschke et al. [1, Theorem 3.7] for the proof.
- 446 **Lemma A.2** (Symmetrized Bregman). For all $x, y \in C$, it holds that:

$$D_h(x,y) + D_h(y,x) = \langle \nabla h(x) - \nabla h(y), x - y \rangle$$
(16)

- The proof immediately follows from the definition of the Bregman divergence. The following result corresponds to Dragomir et al. [7, Lemma 3].
- Lemma A.3 (Bregman cocoercivity). If a convex function f is L-relatively-smooth with respect to h, then for all $\eta \leq 1/L$,

$$D_{h^*}(\nabla h(x) - \eta \left[\nabla f(x) - \nabla f(y)\right], \nabla h(x)) \le \eta D_f(x, y).$$
(17)

451 Denoting $x^{+y} = \nabla h^* (\nabla h(x) - \eta [\nabla f(x) - \nabla f(y)])$, a tighter result actually writes:

$$D_h(x, x^{+y}) + \eta D_f(x^{+y}, y) \le \eta D_f(x, y).$$
(18)

- The proof of the tighter version is simply obtained by not using that $D_f(x^{+y}, y) \ge 0$ in the original proof. While we don't directly use it in this paper, it is sometimes useful. We now introduce the
- generalized bias-variance decomposition Lemma [23, Theorem 0.1].
- 455 **Lemma A.4.** If X is a random variable, then for all $u \in C$,

$$\mathbb{E}\left[D_{h^{*}}(X, u)\right] = D_{h^{*}}(\mathbb{E}\left[X\right], u) + D_{h^{*}}(X, \mathbb{E}\left[X\right]).$$
(19)

456 B Missing results on the variances

- 457 We start this section by proving the following lemma, which in particular ensures that $D_h(x, x^+)/\eta$
- increases with η (and so decreases as $\eta \to 0$).

459 **Lemma B.1.** Let $\phi_{\xi} : \eta \mapsto \frac{1}{\eta} D_h(x, x^+(\eta, \xi))$. Then, $\nabla \phi_{\xi}(\eta) = \frac{1}{\eta^2} D_h(x^+(\eta, \xi), x) \ge 0$.

460 *Proof.* First remark that since $\nabla h(x^+) = \nabla h(x) - \eta \nabla f_{\xi}(x)$, we can write

$$\begin{aligned} \nabla_{\eta} \left[D_{h}(x,x^{+}) \right] &= \nabla_{\eta} \left[h(x) - h(x^{+}) - \nabla h(x^{+})^{\top} (x - x^{+}) \right] \\ &= -\nabla h(x^{+})^{\top} \nabla_{\eta} x^{+} + \nabla f_{\xi}(x)^{\top} (x - x^{+}) + \nabla h(x^{+})^{\top} \nabla_{\eta} x^{+} \\ &= \nabla f_{\xi}(x)^{\top} (x - x^{+}) \\ &= \frac{1}{\eta} \left(\nabla h(x) - \nabla h(x^{+}) \right)^{\top} (x - x^{+}) = \frac{D_{h}(x,x^{+}) + D_{h}(x^{+},x)}{\eta} \end{aligned}$$

461 Then, the expression follows from

$$\nabla \phi_{\xi}(\eta) = \nabla_{\eta} \left[\frac{1}{\eta} D_h(x, x^+) \right] = \frac{1}{\eta} \nabla_{\eta} \left[D_h(x, x^+) \right] - \frac{1}{\eta^2} D_h(x, x^+) = \frac{1}{\eta^2} D_h(x^+, x).$$
(20)

462

Proof of Proposition 2.5. We now prove that $f_{\eta} \to f$ when $\eta \to 0$. To show this, we note that for any fixed $x \in \text{int } C$:

• For any fixed
$$\xi$$
, $\frac{1}{\eta}D_h(x,x^+) = \frac{\eta}{2}||\nabla f_{\xi}(x)||_{\nabla^2 h^*(z)}^2$ for $z \in [\nabla h(x), \nabla h(x) - \eta \nabla f_{\xi}(x)]$.
• Therefore, $\frac{1}{\eta}D_h(x,x^+) \to 0$ for $\eta \to 0$ since $\nabla^2 h^*(\nabla h(x)) = (\nabla^2 h(x))^{-1} < \infty$ by strict convexity of h .

• Let
$$\eta \le \eta_0$$
. Then, for all ξ , $\frac{1}{\eta}D_h(x, x^+(\eta, \xi)) \le \frac{1}{\eta_0}D_h(x, x^+(\eta_0, \xi))$ since the function
 $\eta \mapsto \frac{1}{\eta}D_h(x, x^+(\eta, \xi))$ is an increasing function (positive gradient using Lemma B.1).

470 •
$$\frac{1}{n_0} \mathbb{E} \left[D_h(x, x^+(\eta_0, \xi)) \right]$$
 is finite.

Then, using the dominated convergence theorem, we obtain that we can invert the integral (expectation) and the limit, so that $\lim_{\eta\to 0} \mathbb{E}\frac{1}{\eta}D_h(x,x^+) = \mathbb{E}\lim_{\eta\to 0} \frac{1}{\eta}D_h(x,x^+) = 0.$

473 Proof of Proposition 2.6. We prove this result by successively upper bounding and lower bounding 474 $\sigma_{\star,\eta}^2$, and making $\eta \to 0$.

475 *1* - Upper bound on $\sigma^2_{\star,\eta}$. One side is direct, by writing that $f(x_\eta) \ge f(x_\star)$:

$$\sigma_{\star,\eta}^2 = \frac{1}{\eta} \left(f(x_\star) - f(x_\eta) + \frac{1}{\eta} \mathbb{E} \left[D_h(x_\eta, x_\eta^+) \right] \right) \le \frac{1}{\eta^2} \mathbb{E} \left[D_h(x_\eta, x_\eta^+) \right].$$
(21)

From the proof of Proposition 2.5 we have pointwise convergence of f_{η} to f. Since f is convex and has a unique minimizer x_{\star} , then $x_{\eta} \to x_{\star}$ for $\eta \to 0$, which leads to the result.

478 2 - Lower bound on $\sigma_{\star,\eta}^2$. By definition of x_η as the minimizer of f_η , we have $f_\eta(x_\eta) \le f_\eta(x_\star)$, and 479 so:

$$\sigma_{\star,\eta}^{2} = \frac{f(x_{\star}) - f_{\eta}(x_{\eta})}{\eta} \ge \frac{f(x_{\star}) - f_{\eta}(x_{\star})}{\eta} = \frac{1}{\eta^{2}} \mathbb{E} \left[D_{h}(x_{\star}, x_{\star}^{+}) \right].$$
(22)

480

481 Let us now prove the following proposition, which follows the proof from Dragomir et al. [7].

482 Proof of Proposition 2.7. Let us prove that $\sigma_{\star,\eta}^2 \leq \mathbb{E}\left[\|\nabla f_{\xi}(x_{\star})\|_{\nabla^2 h^*(z_{\eta})}^2\right]$. We start by

$$D_h(x,x^+) = D_{h^*}(\nabla h(x) - \eta \nabla f_{\xi}(x), \nabla h(x))$$
(23)

$$= D_{h^*}(\nabla h(x) - \eta \left[\nabla f_{\xi}(x) - \nabla f_{\xi}(x_{\star})\right] - \eta \nabla f_{\xi}(x_{\star}), \nabla h(x))$$
(24)

$$= D_{h^*}(\frac{(\nabla h(x) - 2\eta \left[\nabla f_{\xi}(x) - \nabla f_{\xi}(x_{\star})\right]) + (\nabla h(x) - 2\eta \nabla f_{\xi}(x_{\star}))}{2}, \nabla h(x)).$$
(25)

Using the convexity of D_{h^*} in its first argument and then the Bregman cocoercivity lemma, we obtain for $\eta \leq 1/2L$:

$$D_{h}(x,x^{+}) \leq \frac{1}{2} D_{h^{*}}(\nabla h(x) - 2\eta \left[\nabla f_{\xi}(x) - \nabla f_{\xi}(x_{\star})\right]), \nabla h(x))$$
(26)

$$+\frac{1}{2}D_{h^*}(\nabla h(x) - 2\eta \nabla f_{\xi}(x_*), \nabla h(x))$$
(27)

$$\leq \eta D_{f_{\xi}}(x, x_{\star}) + \frac{1}{2} D_{h^{\star}} (\nabla h(x) - 2\eta \nabla f_{\xi}(x_{\star}), \nabla h(x)).$$

$$(28)$$

485 Using that $\mathbb{E}\left[D_{f_{\xi}}(x, x_{\star})\right] = D_{f}(x, x_{\star})$ and applying this to $x = x_{\eta}$, we obtain

$$\begin{split} \sigma_{\star,\eta}^2 &= \frac{f(x_{\star}) - f_{\eta}(x_{\eta})}{\eta} \\ &= \frac{\mathbb{E}\left[D_{h}(x_{\eta}, x_{\eta}^{+})\right] + \eta f(x_{\star}) - \eta f(x_{\eta})}{\eta^{2}} \\ &\leq \frac{1}{2\eta^{2}} \mathbb{E}\left[D_{h^{*}}(\nabla h(x_{\eta}) - 2\eta \nabla f_{\xi}(x_{\star}), \nabla h(x_{\eta}))\right] + \frac{D_{f}(x_{\eta}, x_{\star}) + f(x_{\star}) - f(x_{\eta})}{\eta} \\ &= \frac{1}{2\eta^{2}} \mathbb{E}\left[D_{h^{*}}(\nabla h(x_{\eta}) - 2\eta \nabla f_{\xi}(x_{\star}), \nabla h(x_{\eta}))\right] = \frac{1}{2\eta^{2}} \times \mathbb{E}\left[\frac{1}{2} \|2\eta \nabla f_{\xi}(x_{\star})\|_{\nabla^{2}h^{*}(z_{\eta})}^{2}\right], \end{split}$$

and the result follows. The last inequality comes from the fact that if $x_{\star} = \arg \min_{x \in C} f(x)$, then - $\nabla f(x_{\star})$ is normal to C so $-\nabla f(x_{\star})^{\top}(x_{\eta} - x_{\star}) \leq 0$.

488 C Convergence results.

In this section, we detail the proofs of the various convergence theorems that were only sketched in the main text. We start by proving the first identity, which is a variation of *e.g.*, Dragomir et al. [7, Lemma 4], which we detail here for the sake of completeness.

Lemma C.1. Let $x^+ \in C$ be such that $\nabla h(x^+) = \nabla h(x) - \eta \nabla f_{\xi}(x)$, with f_{ξ} a random differentiable function such that $\mathbb{E}[f_{\xi}] = f$. Then, for all points $y \in C$,

$$\mathbb{E}\left[D_{h}(y,x^{+})\right] - D_{h}(y,x) + \eta D_{f}(y,x) = -\eta[f(x) - f(y)] + \mathbb{E}\left[D_{h}(x,x^{+})\right]$$
(29)

494 In particular, we can apply the result to $y = x_{\star}$.

Proof. We give a slightly different proof than Dragomir et al. [7], and in particular this version of the identity is slightly more direct (though maybe less insightful) and does not require $\nabla f(y) = 0$. We write:

$$\mathbb{E} \left[D_{h}(y,x^{+}) \right] = \mathbb{E} \left[h(y) - h(x^{+}) - \nabla h(x^{+})^{\top}(y-x^{+}) \right] \\ = \mathbb{E} \left[h(y) - h(x^{+}) - \nabla h(x^{+})^{\top}(y-x) - \nabla h(x^{+})^{\top}(x-x^{+}) \right] \\ = \mathbb{E} \left[h(y) - h(x) - \nabla h(x)^{\top}(y-x) + \eta \nabla f_{\xi}(x)^{\top}(y-x) \right] \\ \nabla h(x^{+})^{\top}(x-x^{+}) + h(x) - h(x^{+}) \\ = D_{h}(y,x) + \eta \nabla f(x)^{\top}(y-x) + \mathbb{E} \left[D_{h}(x,x^{+}) \right] \\ = D_{h}(y,x) - \eta D_{f}(y,x) + \eta \left[f(y) - f(x) \right] + \mathbb{E} \left[D_{h}(x,x^{+}) \right].$$

498

499 Proof of Corollary 3.2. We start back from Equation (8), and write, using that $f_{\eta}(x) \ge \mathbb{E}[f_{\xi}(x^+)]$ 500 (proof of Proposition 2.2):

$$\mathbb{E}\left[D_h(x_\star, x^+)\right] - D_h(x_\star, x) + \eta D_f(x_\star, x) = \eta \left[f(x_\star) - f_\eta(x)\right]$$
(30)

$$\leq \eta \left[f(x_{\star}) - \mathbb{E} \left[f_{\xi}(x^{+}) \right] \right]$$
(31)

$$\leq -\eta \left[\mathbb{E} \left[f_{\xi}(x^{+}) \right] - f_{\star}^{+} \right] + \eta^{2} \left(\frac{f(x_{\star}) - f_{\star}^{+}}{\eta} \right).$$
(32)

The result follows naturally from using the relative strong convexity of f, leading to:

$$\eta[\mathbb{E}\left[f_{\xi}(x^{+})\right] - f_{+}^{\star}] + \mathbb{E}\left[D_{h}(x_{\star}, x^{+})\right] \le (1 - \eta\mu)D_{h}(x_{\star}, x) + \eta^{2}\left[\frac{f(x_{\star}) - f_{+}^{\star}}{\eta}\right].$$
 (33)

⁵⁰² Then, we chain iterations as done for Theorem 3.1

Proof of Theorem 3.3. We also start from the same result as above, and write it for $x = x^{(k)}$, so that $x^+ = x^{(k+1)}$:

$$\mathbb{E}\left[D_h(x_{\star}, x^{(k+1)})\right] - D_h(x_{\star}, x^{(k)}) + \eta D_f(x_{\star}, x^{(k)}) = \eta \left[f(x_{\star}) - f_\eta(x^{(k)})\right]$$
(34)

$$\leq -\eta \left[f_{\eta}(x^{(k)}) - f_{\eta}(x_{\eta}) \right] + \eta^2 \sigma_{\star,\eta}^2. \tag{35}$$

Moving the f_{η} terms to the left, and summing this for k = 0 to T leads to:

$$\eta \sum_{k=0}^{T} \left[f_{\eta}(x^{(k)}) - f_{\eta}(x_{\eta}) + D_{f}(x_{\star}, x^{(k)}) \right] \le D_{h}(x_{\star}, x^{(0)}) - \mathbb{E} \left[D_{h}(x_{\star}, x^{(k+1)}) \right] + T \eta^{2} \sigma_{\star, \eta}^{2}.$$
(36)

- The final result is obtained by dividing by ηT , and the fact that $\mathbb{E}\left[D_h(x_\star, x^{(k+1)})\right] \ge 0.$
- ⁵⁰⁷ Proof of Proposition 3.4. We use Bregman cocoercivity (Lemma A.3) with $\eta = \frac{1}{L}$ between x_{\star} and ⁵⁰⁸ x (instead of x and x_{\star} as it had been done previously), which directly leads to:

$$D_{h^*}\left(\nabla h^*(x_\star) - \frac{1}{L}\left[\nabla f(x_\star) - \nabla f(x)\right]\right) \le \frac{1}{L}D_f(x_\star, x).$$
(37)

The first part of the proposition follows from the fact that $\nabla f(x_{\star}) = 0$. For the rest proof, we start with Inequality (32), which gives:

$$0 = D_{h^*}(\nabla h(x_\star) - \frac{1}{L}\nabla f(x), \nabla h(x_\star))$$
$$= D_h\left(x_\star, \nabla h^*\left(\nabla h(x_\star) - \frac{1}{L}\nabla f(x)\right)\right)$$

At this point, strict convexity of h leads to $\nabla h^* \left(\nabla h(x_\star) - \frac{1}{L} \nabla f(x) \right) = x_\star$, so that $\nabla f(x) = 0$ by applying ∇h on both sides.

513 **D** Variation on the convex case

In this section, we quickly illustrate that the result we obtain is tightly linked to the notion of variance that we define. As an example, a variation of Theorem 3.3 can be obtained with a control on $f(x) - f(x_*)$, but this requires a different notion of variance:

Theorem D.1. If f is convex, the iterates obtained by SMD using a constant step-sizes $\eta > 0$ verify

$$\mathbb{E}\left[f\left(\frac{1}{T}\sum_{k=0}^{T}x^{(k)}\right)\right] - f(x_{\star}) \le \frac{D_h(x_{\star}, x^{(0)})}{\eta T} + \eta \tilde{\sigma}_{\star, \eta}^2,\tag{38}$$

518 where

$$\tilde{\sigma}_{\star,\eta}^2 = \frac{1}{\eta} \max_{x \in C} \left\{ \frac{1}{\eta} \mathbb{E} \left[D_h(x, x^+) \right] - D_f(x_\star, x) \right\}.$$
(39)

Note that this alternative variance definition can be unbounded even when $\sigma_{\star,\eta}^2$ is bounded, as is the case for instance in the Gaussian MAP example. Besides, it does not inherit from most of the good properties of $\sigma_{\star,\eta}^2$ presented in Section 2, and cannot be compared to the other standard variance notions. The main case in which this alternative definition makes sense is the Euclidean case, in which $\tilde{\sigma}_{\star,\eta}^2$ can be bounded using cocoercivity.

524 *Proof of Theorem D.1.* This proof directly starts from Lemma C.1:

$$\mathbb{E}\left[D_h(x_\star, x^{(k+1)})\right] \tag{40}$$

$$= D_h(x_\star, x^{(k)}) - \eta D_f(x_\star, x^{(k)}) - \eta [f(x^{(k)}) - f(x_\star)] + \mathbb{E} \left[D_h(x^{(k)}, (x^{(k)})^+) \right]$$
(41)

$$= D_h(x_\star, x^{(k)}) - \eta[f(x^{(k)}) - f(x_\star)] + \eta \left[\frac{1}{\eta} \mathbb{E} \left[D_h(x^{(k)}, (x^{(k)})^+) \right] - D_f(x_\star, x^{(k)}) \right]$$
(42)

$$\leq D_h(x_\star, x^{(k)}) - \eta [f(x^{(k)}) - f(x_\star)] + \eta^2 \tilde{\sigma}_{\star,\eta}^2.$$
(43)

Summing this for k = 0 to T, and dividing by ηT we obtain:

$$\frac{1}{T}\sum_{k=0}^{T}f(x^{(k)}) - f(x_{\star}) \le \frac{D_h(x_{\star}, x^{(0)})}{\eta T} + \eta \tilde{\sigma}_{\star, \eta}^2$$
(44)

The result on the average iterate then follows from convexity of f and taking expectation on $x^{(k)}$.

527 E Stochastic Mirror Descent with a Proximal term

We are interested in this section in a variation of the original problem, where we would like to solve the following problem:

$$\min_{x \in C} f(x) + g(x),\tag{45}$$

where g is a convex proper lower semi-continuous function (but not necessarily differentiable). This

problem can be solved using the following stochastic proximal mirror descent algorithm:

$$x^{+} = \arg\min_{u \in C} g(u) + \nabla f_{\xi}(x)^{\top} u + \frac{1}{\eta} D_{h}(u, x).$$
(46)

This is a "proximal" version, which for instance corresponds to projected stochastic mirror descent if g is the indicator of a convex set. Under Assumption 1, the iterations write:

$$\nabla h(x^{+}) = \nabla h(x) - \eta \left[\nabla f_{\xi}(x) + \omega\right]$$
(47)

where $\omega \in \partial g(x^+)$, the subgradient of g at point x^+ . Equation (47) can be rewritten as

$$\nabla h(x^{+}) + \eta \omega = \nabla h(x) + \eta \omega_{x} - \eta \left[\nabla f_{\xi}(x) + \omega_{x}\right]$$
(48)

for any $\omega_x \in \partial g(x)$. In particular, (47) can be interpreted as a Stochastic mirror descent step with objective $f_{\xi} + g$ and mirror $h + \eta g$. While the mirror does not satisfy Assumption 1 (and in particular twice differentiability in case g is the indicator of a set), the iterations can still be written in the form of Equation (3). In particular, the theorems from Section 3 still apply, with the adapted variance definition involving function f + g and mirror $h + \eta g$. Similarly, f + g is $1/\eta$ relatively-smooth with respect to $h + \eta g$ as long as f is L-relatively-smooth with respect to h and $\eta \leq 1/L$.

541 We now prove an equivalent for Lemma C.1.

Lemma E.1. Let $x^+ \in C$ be such that $\nabla h(x^+) = \nabla h(x) - \eta [\nabla f_{\xi}(x) + \omega]$, with f_{ξ} a random differentiable function such that $\mathbb{E}[f_{\xi}] = f$ and $\omega \in \partial g(x^+)$ where g is a convex proper lower semi-continuous function. Then, for all $y \in C \cap \text{dom}g$,

$$\mathbb{E} \left[D_h(y, x^+) \right] = D_h(y, x) - \eta D_f(y, x) - \eta [f(x) - f(y)] \\ + \mathbb{E} \left[D_h(x, x^{+f}) - D_h(x^+, x^{+f}) \right] + \eta \omega^\top (y - x^+),$$

where x^{+f} is the point such that $\nabla h(x^{+f}) = \nabla h(x) - \eta \nabla f_{\xi}(x)$.

546 Proof. We write:

$$D_{h}(y,x^{+}) = h(y) - h(x^{+}) - \nabla h(x^{+})^{\top}(y - x^{+})$$

= $h(y) - h(x^{+f}) - \nabla h(x^{+f})^{\top}(y - x^{+}) + \eta \omega^{\top}(y - x^{+}) - h(x^{+}) + h(x^{+f})$
= $D_{h}(y,x^{+f}) - \nabla h(x^{+f})^{\top}(x^{+f} - x^{+}) + \eta \omega^{\top}(y - x^{+}) - h(x^{+}) + h(x^{+f})$
= $D_{h}(y,x^{+f}) - D_{h}(x^{+},x^{+f}) + \eta \omega^{\top}(y - x^{+})$

⁵⁴⁷ The result follows from applying Lemma C.1 to $D_h(y, x^{+f})$.

Note that by abuse of notation, if we denote
$$D_g(y, x^+) = g(y) - g(x^+) - \omega^\top (y - x^+)$$
, and
 $D_g(y, x) = g(y) - g(x) - \omega_x^\top (y - x)$ for any $\omega_x \in \partial g(x)$, then with a few lines of computations, and
noting in particular that $D_h(x, x^{+f}) - D_h(x^+, x^{+f}) = D_h(x, x^+) - [\nabla h(x^{+f}) - \nabla h(x^+)]^\top (x - x^+)$ we obtain:

$$\mathbb{E} \left[D_{h+\eta g}(y, x^{+}) \right] = D_{h}(y, x) - \eta D_{f}(y, x) - \eta [f(x) - f(y)] \\ + \mathbb{E} \left[D_{h}(x, x^{+}) \right] - \eta \mathbb{E} \left[\omega^{\top} (x - x^{+}) \right] + \eta \mathbb{E} \left[g(y) - g(x^{+}) \right] \\ = D_{h+\eta g}(y, x) - \eta D_{g}(y, x) - \eta D_{f}(y, x) - \eta [f(x) - f(y)] \\ + \mathbb{E} \left[D_{h+\eta g}(x, x^{+}) \right] + \eta \left[g(y) - g(x) \right]$$

In particular, we exactly recover the result of Lemma C.1 applied to the iterations in which we take (sub)-gradients of f + g with mirror $h + \eta g$, *i.e.*,

$$\mathbb{E}\left[D_{h+\eta g}(y,x^{+})\right] = D_{h+\eta g}(y,x) - \eta D_{f+g}(y,x) - \eta [(f+g)(x) - (f+g)(y)] + \mathbb{E}D_{h+g}(x,x^{+}).$$

Therefore, using the same sequence of derivations, Theorems 3.1 and 3.3 can be transposed directly

to the composite (f + g) setting by simply defining generalized Bregman divergences where the

⁵⁵⁶ gradient parts are replaced by the subgradients picked in the actual SMD steps.

⁵⁵⁷ While $h + \eta g$ does not necessarily satisfy Assumption 1, the key point is that iterations can be written ⁵⁵⁸ in the form of Equation (47), which is the case for instance if g is the indicator of a convex set.

Note that Corollary 3.2 also holds in the same way, since relative smoothness is only needed to obtain that $\eta D_{f_{\xi}+g}(x,x^+) \leq D_{h+\eta g}(x,x^+)$, which is equivalent to $\eta D_{f_{\xi}}(x,x^+) \leq D_h(x,x^+)$, which holds by *L*-relative smoothness of *f* with respect to *h* for $\eta \leq 1/L$.

⁵⁶² F Gaussian case with unknown covariance.

In this section, we prove the various results for Gaussian estimation with unknown mean and covariance. For the sake of brevity, we only prove the propositions, and refer the interested reader to, *e.g.*, Le Priol et al. [17] for standard results about the setting.

566 F.1 Instanciation in the Stochastic mirror descent setting

⁵⁶⁷ We first write what the various divergences are in our setting, together with the mirror updates and ⁵⁶⁸ finally the form of f_{η} . Following Le Priol et al. [17, Section 4.2], we write that:

$$\theta_1 = \frac{m}{\Sigma^2}, \qquad \theta_2 = -\frac{1}{2\Sigma^2}.$$
(49)

This allows us to express $A(\theta)$ in terms of (m, Σ^2) :

$$A(\theta) = -\frac{1}{2}\log(-\theta_2) - \frac{\theta_1^2}{4\theta_2} = \frac{1}{2}\log(2\Sigma^2) + \frac{1}{2}\frac{m^2}{\Sigma^2}$$
(50)

Proposition F.1. The Bregman divergence with respect to $\tilde{\theta}, \theta$ writes:

$$D_A(\tilde{\theta}, \theta) = -\frac{1}{2} \log \left(\frac{\Sigma^2}{\tilde{\Sigma}^2}\right) - \frac{\tilde{\Sigma}^2 - \Sigma^2}{2\tilde{\Sigma}^2} + \frac{(\tilde{m} - m)^2}{2\tilde{\Sigma}^2}.$$
(51)

571 *Proof.* We know that $\nabla A(\theta) = \mu = (m, m^2 + \Sigma^2)$. Therefore,

$$\nabla A(\theta)^{\top}(\tilde{\theta} - \theta) = m\left(\frac{\tilde{m}}{\tilde{\Sigma}^2} - \frac{m}{\Sigma^2}\right) - \frac{1}{2}(m^2 + \Sigma^2)\left(\frac{1}{\tilde{\Sigma}^2} - \frac{1}{\Sigma^2}\right)$$
(52)

$$= \frac{m\tilde{m}}{\tilde{\Sigma}^{2}} - \frac{m^{2}}{2\Sigma^{2}} - \frac{m^{2}}{2\tilde{\Sigma}^{2}} - \frac{1}{2} \left(\frac{\Sigma^{2}}{\tilde{\Sigma}^{2}} - 1\right)$$
(53)

$$= -\frac{(m-\tilde{m})^2}{2\tilde{\Sigma}^2} + \frac{\tilde{m}^2}{2\tilde{\Sigma}^2} - \frac{m^2}{2\Sigma^2} - \frac{\Sigma^2 - \tilde{\Sigma}^2}{2\tilde{\Sigma}^2}.$$
 (54)

572 Using Equation (50), we obtain:

$$D_A(\tilde{\theta}, \theta) = A(\tilde{\theta}) - A(\theta) - \nabla A(\theta)^\top (\tilde{\theta} - \theta)$$

= $\frac{1}{2} \log(2\tilde{\Sigma}^2) - \frac{1}{2} \log(2\Sigma^2) + \frac{\Sigma^2 - \tilde{\Sigma}^2}{2\tilde{\Sigma}^2} + \frac{(m - \tilde{m})^2}{2\tilde{\Sigma}^2},$

- 573 which finishes the proof.
- 574 In the Gaussian with unknown covariance, the sufficient statistics are:

$$T(X) = (X, X^2),$$
 (55)

- where $x \in \mathbb{R}$ is an observation drawn from $\mathcal{N}(m_{\star}, \Sigma_{\star})$.
- ⁵⁷⁶ Let us now prove the form on the updates, which corresponds to (13):
- **Proposition F.2.** In (m, Σ^2) parameters, the updates write:

$$m^{+} = (1 - \eta)m + \eta X, \tag{56}$$

$$(\Sigma^2)^+ = (1 - \eta) \left[\Sigma^2 + \eta (m - X)^2 \right].$$
(57)

Proof. Since the (stochastic) gradients write $g(\mu) = \mu - T(X)$, the iterations are defined by:

$$\mu_1^+ = (1 - \eta)\mu_1 + \eta X \tag{58}$$

$$\mu_2^+ = (1 - \eta)\mu_2 + \eta X^2.$$
(59)

Since $(\mu_1, \mu_2) = (m, m^2 + \Sigma^2)$, the update on m is immediate. For the update on Σ^2 , we write: $(\Sigma^2)^+ = \mu_2^+ - (m^+)^2$

$$= (1 - \eta)\mu_2 + \eta X^2 - ((1 - \eta)m + \eta X)^2$$

= $(1 - \eta)\Sigma^2 + (1 - \eta)m^2 + \eta X^2 - (1 - \eta)^2 m^2 - 2\eta (1 - \eta)Xm - \eta^2 X^2$
= $(1 - \eta)\Sigma^2 + \eta (1 - \eta)(m - X)^2$.

580

⁵⁸¹ We now use this to show that updates are well-defined.

Proof of Proposition 4.1. If $\theta_2 < 0$ then $\Sigma^2 > 0$ so for $\eta < 1$, $(\Sigma^2)^+ > 0$ almost surely so that $\theta_2^+ < 0$ and $|\theta_1^+| < \infty$. In particular, $\theta^+ \in \mathbb{R} \times \mathbb{R}_-^*$ so the update is well-defined. The rest of the proposition comes from the fact that $\nabla^2 f_{\xi} = \nabla^2 f = \nabla^2 A$.

We can now proceed to proving the form of f_{η} . We first start by writing that:

$$f(\theta) = A(\theta) - \theta^{\top}(m_{\star}, m_{\star}^2 + \Sigma_{\star}^2)$$
(60)

$$= \frac{1}{2}\log(2\Sigma^2) + \frac{1}{2}\frac{m^2}{\Sigma^2} - \frac{mm_{\star}}{\Sigma^2} + \frac{m_{\star}^2 + \Sigma_{\star}^2}{2\Sigma^2}.$$
 (61)

586 Therefore,

$$f(\theta) = \frac{1}{2}\log(2\Sigma^2) + \frac{\Sigma_{\star}^2}{2\Sigma^2} + \frac{(m - m_{\star})^2}{2\Sigma^2}$$
(62)

587 In particular,

$$f(\theta) - f(\theta_{\star}) = \frac{1}{2} \log \left(\frac{\Sigma^2}{\Sigma_{\star}^2}\right) + \frac{\Sigma_{\star}^2 - \Sigma^2}{2\Sigma^2} + \frac{(m - m_{\star})^2}{2\Sigma^2}$$
(63)

Note that, as expected, this corresponds to $D_A(\theta, \theta_*)$, that we can also compute through Proposition F.1. We now write:

$$D_{A}(\theta,\theta^{+}) = -\frac{1}{2}\log\left(\frac{(\Sigma^{2})^{+}}{\Sigma^{2}}\right) - \frac{\Sigma^{2} - (\Sigma^{2})^{+}}{2\Sigma^{2}} + \frac{(m-m^{+})^{2}}{2\Sigma^{2}}$$
(64)
$$= -\frac{1}{2}\log\left((1-\eta)\left[1+\eta\frac{(m-X)^{2}}{\Sigma^{2}}\right]\right) + \frac{(1-\eta)(\Sigma^{2}+\eta(m-X)^{2}) - \Sigma^{2}}{2\Sigma^{2}} + \frac{\eta^{2}(m-X)^{2}}{2\Sigma^{2}}$$
(65)
$$= -\frac{1}{2}\log\left((1-\eta)\left[1+\eta\frac{(m-X)^{2}}{\Sigma^{2}}\right]\right) - \frac{\eta}{2} + \eta(1-\eta)\frac{(m-X)^{2}}{2\Sigma^{2}} + \frac{\eta^{2}(m-X)^{2}}{2\Sigma^{2}}$$
(66)
$$= -\frac{1}{2}\log\left((1-\eta)\left[1+\eta\frac{(m-X)^{2}}{\Sigma^{2}}\right]\right) - \frac{\eta}{2} + \eta\frac{(m-X)^{2}}{2\Sigma^{2}}.$$
(67)

590 Therefore,

$$f(\theta) - \frac{D_A(\theta, \theta^+)}{\eta} - f(\theta_\star)$$
(68)
= $\frac{1}{2} \log\left(\frac{\Sigma^2}{\Sigma_\star^2}\right) + \frac{\Sigma_\star^2 - \Sigma^2}{2\Sigma^2} + \frac{(m - m_\star)^2}{2\Sigma^2} + \frac{1}{2\eta} \log\left((1 - \eta) \left[1 + \eta \frac{(m - X)^2}{\Sigma^2}\right]\right) + \frac{1}{2} - \frac{(m - X)^2}{2\Sigma^2}$ (69)
= $\frac{1}{2} \log\left(\frac{\Sigma^2}{\Sigma_\star^2}\right) + \frac{1}{2\eta} \log\left((1 - \eta) \left[1 + \eta \frac{(m - X)^2}{\Sigma^2}\right]\right) + \frac{\Sigma_\star^2}{2\Sigma^2} + \frac{(m - m_\star)^2}{2\Sigma^2} - \frac{(m - X)^2}{2\Sigma^2}.$ (70)

Finally,
$$\mathbb{E}\left[(m-X)^2\right] = (m-m_\star)^2 + \Sigma_\star^2$$
, and so:

$$f_\eta(\theta) - f(\theta_\star) = \frac{1}{2}\log\left(\frac{\Sigma^2}{\Sigma_\star^2}\right) + \frac{1}{2\eta}\mathbb{E}\left[\log\left((1-\eta)\left[1+\eta\frac{(m-X)^2}{\Sigma^2}\right]\right)\right],$$
(71)

which precisely corresponds to Equation (14). We now proceed to proving bounds on θ_{η} for $\eta < 1$.

593 F.2 Bounding the stochastic mirror descent variance $\sigma_{\star,\eta}^2$.

Now that we have an explicit form for f_{η} , we can characterize its minimizer θ_{η} , and use this to prove results on $f_{\eta}(\theta_{\eta})$, which will in turn lead to bounds on $\sigma_{\star,\eta}^2$. This is the core of Lemma 4.2.

596 *Proof.* Proof of Lemma 4.2. The proof will proceed in three different stages:

• Differentiating f_{η} with respect to m and Σ^2 .

• Using these expressions to obtain bounds on the (m_η, Σ_η^2) for which ∇f_η is 0.

• Plugging these bounds into the expression of f_{η} to bound Σ_{η}^2

600 **1** - **Differentiating** f_{η} . Before differentiating, we rewrite:

$$f_{\eta}(\theta) - f(\theta_{\star}) = \frac{1}{2} \log \left(\Sigma^{2}\right) + \frac{1}{2\eta} \mathbb{E} \left[\log \left(1 + \eta \frac{(m-X)^{2}}{\Sigma^{2}}\right) \right] - \frac{1}{2} \log \left(\Sigma^{2}_{\star}\right) + \frac{1}{2\eta} \log(1-\eta)$$
(72)
$$= -\frac{1-\eta}{2\eta} \log \left(\Sigma^{2}\right) + \frac{1}{2\eta} \mathbb{E} \left[\log \left(\Sigma^{2} + \eta (m-X)^{2}\right) \right] - \frac{1}{2} \log \left(\Sigma^{2}_{\star}\right) + \frac{1}{2\eta} \log(1-\eta)$$
(73)

Indeed, the two terms on the right are constant and so do not matter. If we differentiate in m, we obtain:

$$\nabla_m f_\eta(\theta) = \mathbb{E}\left[\frac{1}{2\eta} 2\eta \frac{m-X}{\Sigma^2} \frac{1}{\Sigma^2 + \eta(m-X)^2}\right] = \mathbb{E}\left[\frac{m-X}{\Sigma^2 + \eta(m-X)^2}\right].$$
 (74)

Now, differentiating in Σ^2 yields:

1

$$\nabla_{\Sigma^2} f_{\eta}(\theta) = -\frac{1-\eta}{2\eta\Sigma^2} + \frac{1}{2\eta} \mathbb{E}\left[\frac{1}{\Sigma^2 + \eta(m-X)^2}\right] = \frac{1}{2\Sigma^2} - \mathbb{E}\left[\frac{(m-X)^2}{2\Sigma^2(\Sigma^2 + \eta(m-X)^2)}\right].$$
 (75)

2 • **Obtaining bounds on** $(m_{\eta}, \Sigma_{\eta}^2)$. The solution to $\nabla_m f_{\eta}(\theta) = 0$ is $m = m_{\star}$. Indeed, it is direct to verify that in this case, $\mathbb{E}\left[\frac{\tilde{X}}{\Sigma^2 + \eta \tilde{X}^2}\right] = 0$ since $\tilde{X} = m_{\star} - X$ is symmetric (with respect to 0). For $m > m_{\star}, \mathbb{E}\left[\frac{\tilde{X}}{\Sigma^2 + \eta \tilde{X}^2}\right] > 0$ since we integrate the same values as the previous case, but now more mass is put on the positive values (and similarly for $m < m_{\star}$). Note that this is the case regardless of Σ_{η}^2 .

We are now interested in Σ_{η}^2 . Note that we will not get such a clean expression as for m_{η} , but only bounds. From its expression, we deduce that $\nabla_{\Sigma^2} f_{\eta}(\theta_{\eta}) = 0$ can be reformulated as:

$$\mathbb{E}\left[\frac{(m_{\eta}-X)^2}{\Sigma_{\eta}^2+\eta_{\eta}(m-X)^2}\right] = 1$$
(76)

⁶¹¹ For the upper bound, we simply write that:

$$1 = \mathbb{E}\left[\frac{(m_{\eta} - X)^2}{\Sigma_{\eta}^2 + \eta_{\eta}(m - X)^2}\right] \le \mathbb{E}\left[\frac{(m_{\eta} - X)^2}{\Sigma_{\eta}^2}\right] = \frac{\Sigma_{\star}^2}{\Sigma_{\eta}^2},\tag{77}$$

from which we deduce that $\Sigma_{\eta}^2 \leq \Sigma_{\star}^2$. Let us now introduce some $\alpha > 0$. We have that:

$$\mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}\right] = \mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\alpha-\alpha+\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}\right] = \mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\alpha}\frac{1}{1-1+\frac{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}{\alpha}}\right]$$

613 We now use that for $u \ge -1$, $\frac{1}{1+u} \ge 1-u$, and so:

$$\mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}\right] \geq \mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\alpha}\left(1-\left[-1+\frac{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}{\alpha}\right]\right)\right]$$
(79)
$$=\mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\alpha}\left(2-\frac{\Sigma_{\eta}^{2}}{\alpha}\right)-\eta\frac{(m_{\eta}-X)^{4}}{\alpha^{2}}\right].$$
(80)

Now, recall that $m_{\eta} = m_{\star}$, so $X - m_{\eta} \sim \mathcal{N}(0, \Sigma_{\star})$, leading to

$$I = \mathbb{E}\left[\frac{(m_{\eta} - X)^2}{\Sigma_{\eta}^2 + \eta(m_{\eta} - X)^2}\right] \ge \frac{\Sigma_{\star}^2}{\alpha} \left(2 - \frac{\Sigma_{\eta}^2}{\alpha}\right) - \eta \frac{3\Sigma_{\star}^4}{\alpha^2}.$$
(81)

615 Rearranging terms, we obtain:

$$\frac{\alpha^2}{\Sigma_\star^2} - 2\alpha \ge -\Sigma_\eta^2 - 3\Sigma_\star^2, \text{ so } \Sigma_\eta^2 \ge \frac{2\alpha\Sigma_\star^2 - \alpha^2}{\Sigma_\star^2} - 3\eta\Sigma_\star^2.$$
(82)

⁶¹⁶ We see that $\alpha = \Sigma_{\star}^2$ maximizes the right term, and we obtain the desired result, *i.e.*:

$$\Sigma_{\eta}^2 \ge (1 - 3\eta)\Sigma_{\star}^2. \tag{83}$$

⁶¹⁷ Unfortunately, we see that this bound is only informative for $3\eta < 1$. For the rest of the cases, we ⁶¹⁸ will use the Markov inequality instead, which writes for all a > 0:

$$\mathbb{P}\left(\frac{(m_{\eta} - X)^{2}}{\Sigma_{\eta}^{2} + \eta(m_{\eta} - X)^{2}} \ge a\right) \le \frac{1}{a} \mathbb{E}\left[\frac{(m_{\eta} - X)^{2}}{\Sigma_{\eta}^{2} + \eta(m_{\eta} - X)^{2}}\right] = \frac{1}{a}.$$
(84)

619 Yet,

$$\mathbb{P}\left(\frac{(m_{\eta}-X)^2}{\Sigma_{\eta}^2+\eta(m_{\eta}-X)^2} \ge a\right) = \mathbb{P}\left(\frac{(m_{\eta}-X)^2}{\Sigma_{\star}^2} \ge \frac{a}{1-\eta a}\frac{\Sigma_{\eta}^2}{\Sigma_{\star}^2}\right) = 2\mathbb{P}\left(\frac{X-m_{\star}}{\Sigma_{\star}} \ge \sqrt{\frac{a}{1-\eta a}\frac{\Sigma_{\eta}}{\Sigma_{\star}}}\right)$$
(85)

Therefore, denoting Φ the cumulative distribution function of the standard Gaussian, we have:

$$2\left(1 - \Phi\left(\sqrt{\frac{a}{1 - \eta a}}\frac{\Sigma_{\eta}}{\Sigma_{\star}}\right)\right) \le \frac{1}{a},\tag{86}$$

and since Φ^{-1} is an increasing function, this leads to:

$$\sqrt{\frac{a}{1-\eta a}} \frac{\Sigma_{\eta}}{\Sigma_{\star}} \ge \Phi^{-1} \left(1 - \frac{1}{2a} \right), \tag{87}$$

622 so that:

$$\Sigma_{\eta} \ge \sqrt{1 - \eta a} \frac{\Phi^{-1} \left(1 - \frac{1}{2a}\right)}{\sqrt{a}} \Sigma_{\star}$$
(88)

623 One can check that $\Phi^{-1}\left(1-\frac{1}{2a}\right)/\sqrt{a} < 1$ for all a, which is consistent with the fact that $\Sigma_{\eta}^{2} \leq \Sigma_{\star}^{2}$. 624 Also note that for $\eta = 1$, a non-trivial bound would require a < 1, but then $\Phi^{-1}\left(1-\frac{1}{2a}\right) \leq 0$ so 625 (as expected), we cannot get better than $\Sigma_{\eta}^{2} \geq 0$. However, the previous bounding (Equation (83)) is 626 more precise for small η since $\Phi^{-1}\left(1-\frac{1}{2a}\right)/\sqrt{a} < 1-c$ with c > 0 a constant regardless of a. In 627 particular, for any ε , by using any $1 < a < 1/(1-\varepsilon)$, we obtain that $\Sigma_{\eta}^{2} \geq \alpha_{\varepsilon} \Sigma_{\star}^{2}$ for some constant 628 α_{ε} that only depends on the a that we choose. In particular, we can handle the cases $\eta = 1/2$ and 629 $\eta = 1/3$ that gave trivial results $\Sigma_{\eta}^{2} \geq 0$ with the previous bounds.

The last part consists in proving that $f_{\eta}(\theta_{\eta}) - f(\theta_{\star}) \ge \frac{1}{2} \log \left(\frac{\Sigma_{\eta}^2}{\Sigma_{\star}^2} \right)$. To do so, we start back from $f_{\eta}(\theta) = \frac{1}{2} \log \left(\frac{\Sigma^2}{\Sigma_{\star}^2} \right) + \frac{1}{2} \mathbb{E} \left[\log \left((1 - m) \left[1 + m (m - X)^2 \right] \right) \right]$

$$f_{\eta}(\theta) - f(\theta_{\star}) = \frac{1}{2} \log \left(\frac{\Sigma^2}{\Sigma_{\star}^2} \right) + \frac{1}{2\eta} \mathbb{E} \left[\log \left((1-\eta) \left[1 + \eta \frac{(m-X)^2}{\Sigma^2} \right] \right) \right]$$

$$u \text{ that } \mathbb{E} \left[\log \left((1-\eta) \left[1 + \eta \frac{(m_{\eta}-X)^2}{\Sigma^2} \right] \right) \right] > 0. \text{ We start by the inequality } \log (1-\eta) \left[1 + \eta \frac{(m_{\eta}-X)^2}{\Sigma^2} \right] \right]$$

and show that $\mathbb{E}\left[\log\left((1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^2}{\Sigma_{\eta}^2}\right]\right)\right] \ge 0$. We start by the inequality $\log(1+x) \ge \frac{x}{1+x}$, leading to:

$$\mathbb{E}\left[\log\left((1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}}\right]\right)\right] \ge \mathbb{E}\left[\frac{(1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}}\right]-1}{(1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}}\right]}\right]$$
(89)

$$= \mathbb{E}\left[\frac{\eta(1-\eta)\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}} - \eta}{(1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}}\right]}\right]$$
(90)

$$= \eta \mathbb{E}\left[\frac{(1-\eta)(m_{\eta} - X)^2 - \Sigma_{\eta}^2}{(1-\eta)\left[\Sigma_{\eta}^2 + \eta(m_{\eta} - X)^2\right]}\right]$$
(91)

Recall that the optimality conditions for $(m_{\eta}, \Sigma_{\eta}^2)$ write:

$$1 = \mathbb{E}\left[\frac{(m_{\eta} - X)^2}{\Sigma_{\eta}^2 + \eta(m_{\eta} - X)^2}\right] = \frac{1}{\eta} \mathbb{E}\left[1 - \frac{\Sigma_{\eta}^2}{\Sigma_{\eta}^2 + \eta(m_{\eta} - X)^2}\right],$$
(92)

634 so that

$$\mathbb{E}\left[\frac{\Sigma_{\eta}^2}{\Sigma_{\eta}^2 + \eta(m_{\eta} - X)^2}\right] = 1 - \eta.$$
(93)

635 Combining these, we obtain that

$$\mathbb{E}\left[\log\left((1-\eta)\left[1+\eta\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}}\right]\right)\right] \ge \eta\mathbb{E}\left[\frac{(1-\eta)(m_{\eta}-X)^{2}-\Sigma_{\eta}^{2}}{(1-\eta)\left[\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}\right]}\right] = \eta\left(\mathbb{E}\left[\frac{(m_{\eta}-X)^{2}}{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}\right] - \frac{1}{1-\eta}\mathbb{E}\left[\frac{\Sigma_{\eta}^{2}}{\Sigma_{\eta}^{2}+\eta(m_{\eta}-X)^{2}}\right]\right) = 0,$$

- 636 which is the desired result.
- The final result is obtained by plugging the lower bounds for Σ_{η}^2 into this bound, leading to either $\sigma_{\star,\eta}^2 \leq -\frac{1}{2\eta} \log(1-3\eta)$ for $\eta < 1/3$ or $\sigma_{\star,\eta}^2 \leq -\frac{1}{2\eta} \log \alpha_{\varepsilon}$ for $\eta < 1-\varepsilon$.

640 F.3 Unrolling the recursions to derive actual convergence results.

641 F.3.1 Proof of Theorem 4.3

Now that we have bounded the stochastic mirror descent variance $\sigma_{\star,\eta}^2$ in this setting, we can plug it into Theorem 3.1 to obtain finite-time convergence guarantees on the MAP and MLE estimators.

644 *Proof of Theorem 4.3.* Starting from Theorem 3.1, we obtain:

$$D_A(\theta_\star, \theta^{(k+1)}) \le (1-\eta) D_A(\theta_\star, \theta^{(k)}) - \frac{\eta}{2} \log (1-3\eta) \le (1-\eta) D_A(\theta_\star, \theta^{(k)}) + \frac{3\eta^2}{2}, \quad (94)$$

where the right term is replaced by c_{ε} (where $c_{\epsilon} = -\frac{1}{2} \log \alpha_{\varepsilon}$) for $k \leq 3$. Taking $\eta = 1/k$ for k > 1and multiplying by k leads for k > 3 to:

$$kD_A(\theta_\star, \theta^{(k+1)}) \le (k-1)D_A(\theta_\star, \theta^{(k)}) + \frac{3}{2k}.$$
 (95)

⁶⁴⁷ Therefore, a telescopic sum leads to, for $n_0 > 0$:

$$(n+n_0)D_A(\theta_\star,\theta^{(n)}) \le n_0 D_A(\theta_\star,\theta^{(0)}) + \frac{3}{2} \sum_{k=n_0}^{n+n_0} \frac{1}{k} + 2c_{1/2},$$
(96)

and so, since $\sum_{k=n_0}^{n} \frac{1}{k} \le \log(n+n_0+1) - \log(n_0)$:

$$D_A(\theta_\star, \theta^{(n)}) \le \frac{n_0 D_A(\theta_\star, \theta^{(0)}) + (3/2)\log(1 + (n+1)/n_0) + \Gamma}{n+n_0},\tag{97}$$

- 649 where $\Gamma = 2c_{1/2}$ and we actually have $\Gamma = 0$ for $n_0 > 3$.
- 650 **F.3.2** O(1/n) convergence result.
- ⁶⁵¹ We now consider a different estimator (from the MAP and the MLE), which we construct in the ⁶⁵² following way:
- Choose $n_0 \ge 6$ and initial parameter $\tilde{\theta}^{(n_0)}$.
- Obtain $\tilde{\theta}^{(n)}$ by performing $n n_0$ stochastic mirror descent steps from $\tilde{\theta}^{(n_0)}$ with step-sizes $\eta_k = 2/(k+1)$ for $k \in \{n_0, ..., n\}$.

- This estimator is a modified version of the MAP, where n_0 controls how much weight we would like to put on the prior, and $\tilde{\theta}^{(n_0)}$ would typically be the same starting parameter as for the MAP estimator.
- This estimator is built so that we can use the convergence analysis from Lacoste-Julien et al. [16] and
- obtain a O(1/n) convergence rate. Note that we make the $n_0 \ge 6$ restriction for simplicity to ensure
- that $\sigma_{\star,\eta}^2 \leq 3/2$, but the result can be easily adapted to $n_0 \geq 2$.
- **Proposition F.3.** After $n n_0$ steps, this modified estimator $\tilde{\theta}^{(n)}$ verifies:

$$\mathbb{E}D_{h}(\theta_{\star},\tilde{\theta}^{(n)}) \leq \frac{2n_{0}(n_{0}-1)}{n(n-1)}D_{h}(\theta_{\star},\tilde{\theta}^{(n_{0})}) + \frac{6}{n}.$$
(98)

Proof. Let us note $D_k = \mathbb{E}\left[D_h(\theta_\star, \tilde{\theta}^{(k)})\right]$. In this case, using that $\sigma_{\star,\eta}^2 \leq 3/2$, Theorem 3.1 writes (since $\mu = 1$):

$$D_{k+1} \le (1 - \eta_k) D_k + \frac{3\eta_k^2}{2}.$$
(99)

At this point, we can multiply by k(k+1) on both sides, and take $\eta_k = \frac{2}{k+1}$ for $k \ge n_0$. Remarking that $1 - \eta_k = 1 - \frac{2}{k+1} = \frac{k-1}{k+1}$, we obtain that:

$$(k+1)kD_{k+1} \le k(k-1)D_k + \frac{6k}{k+1} \le k(k-1)D_k + 6.$$
 (100)

Unrolling this recursion from $k = n_0$ to k = n - 1 (since $(k + 1)kD_{k+1} = L_{k+1}$, where $L_k = k(k-1)D_k$), we obtain:

$$n(n-1)D_n \le n_0(n_0-1)D_{n_0} + \sum_{k=n_0}^{n-1} 6,$$
(101)

and the result follows by dividing by n(n-1), and using that $(n-n_0)/(n-1) \le 1$.

669 F.4 The case of the MLE

For the MLE estimator, directly applying the mirror descent approach would require using $\eta_0 = 1$, starting from an arbitrary $\theta^{(0)}$ (that would not affect the results anyway). The problem in this case is that $D_h(\theta_\star, \theta^{(1)})$ is infinite since $\Sigma^{(2)} = 0$. This also means that we cannot start the stochastic mirror descent algorithm from $\theta^{(1)}$, since the recursion would still involve the infinite $D_h(\theta_\star, \theta^{(1)})$. Therefore, in the case of the MLE, considering that the first two samples are $X^{(1)}$ and $X^{(2)}$, then the first two points are:

$$m^{(1)} = X^{(1)}, \Sigma^{(1)} = 0 \text{ and } m^{(2)} = \frac{X^{(1)} + X^{(2)}}{2}, (\Sigma^{(2)})^2 = \frac{(X^{(1)} - X^{(2)})^2}{4}.$$
 (102)

676 More generally, a direct recursion for the MLE leads to:

$$m^{(n)} = \frac{1}{n} \sum_{k=1}^{n} X^{(k)}, \ (\Sigma^{(n)})^2 = \frac{1}{n} \sum_{k=1}^{n} (X^{(k)} - m^{(n)})^2.$$
(103)

677 From this, we derive that:

$$\mathbb{E}\left[(\Sigma^{(n)})^2\right] = \mathbb{E}\left[(X^{(n)} - m^{(n)})^2\right]$$
(104)

$$= \mathbb{E}\left[\left(\left(1-\frac{1}{n}\right)X^{(n)}-\frac{n-1}{n}m^{(n-1)}\right)^2\right]$$
(105)

$$= \left(\frac{n-1}{n}\right)^{2} \mathbb{E}\left[\left(X^{(n)} - m_{\star} - (m^{(n-1)} - m_{\star})\right)^{2}\right]$$
(106)

$$= \left(\frac{n-1}{n}\right)^{2} \mathbb{E}\left[\left(X^{(n)} - m_{\star}\right)^{2} + \left(m^{(n-1)} - m_{\star}\right)^{2}\right]$$
(107)

$$= \left(\frac{n-1}{n}\right)^2 \left(\Sigma_\star^2 + \frac{1}{n-1}\Sigma_\star^2\right) = \left(1 - \frac{1}{n}\right)\Sigma_\star^2,\tag{108}$$

where (107) comes from the fact that $X^{(n)}$ and $m^{(n-1)}$ are independent with mean m_{\star} . Plugging this into the expression of $D_h(\theta_{\star}, \theta)$ for the MLE after *n* steps, we obtain:

$$D_h(\theta_\star, \theta^{(n)}) = -\frac{1}{2} \mathbb{E}\left[\log\frac{(\Sigma^{(n)})^2}{\Sigma_\star^2}\right].$$
(109)

Unfortunately, there is no closed-form for this expression for arbitrary n, hence the need for a more involved analysis, for instance through the mirror descent framework. For the case n = 2 however

682 (which is the one we are interested in), we obtain that

$$D_h(\theta_\star, \theta^{(2)}) = -\frac{1}{2} \mathbb{E}\left[\log\left(\frac{X^{(1)} - X^{(2)}}{2\Sigma_\star}\right)^2\right] = -\frac{1}{2} \mathbb{E}\left[\log\frac{Y^2}{2}\right],$$
(110)

where $Y = \frac{X^{(1)} - X^{(2)}}{\sqrt{2\Sigma_{\star}}} \sim \mathcal{N}(0, 1)$. Therefore, this can simply be treated as a constant that we can precisely evaluate numerically (for instance remarking that Y^2 is gamma distributed and using results on logarithmic expectations of gamma distributions).

For n > 2, it is tempting to use the convexity of $-\log$ to use a similar reasoning, but this only leads to a constant bound on $D_h(\theta_\star, \theta^{(n)})$.

688 NeurIPS Paper Checklist

	1	
689	1.	Claims
690		Question: Do the main claims made in the abstract and introduction accurately reflect the
691		paper's contributions and scope?
692		Answer: [Yes]
693		Justification: We have theorems for all the results we claim to have in the abstract.
694		Guidelines:
695		• The answer NA means that the abstract and introduction do not include the claims
696		made in the paper.
697		• The abstract and/or introduction should clearly state the claims made, including the
698		contributions made in the paper and important assumptions and limitations. A No or
699		NA answer to this question will not be perceived well by the reviewers.
700		• The claims made should match theoretical and experimental results, and reflect how
701		much the results can be expected to generalize to other settings.
702		• It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.
703	\mathbf{r}	Limitations
/04	۷.	
705		Question: Does the paper discuss the limitations of the work performed by the authors?
706		Answer: [Yes]
707		Justification: Limitations for each result are discussed after they are introduced, in particular
708		the fact that Theorem 4.3 does not completely solve the problem from Le Priol et. al. (2021).
709		Guidelines:
710		• The answer NA means that the paper has no limitation while the answer No means that
711		the paper has limitations, but those are not discussed in the paper.
712		• The authors are encouraged to create a separate "Limitations" section in their paper.
713		• The paper should point out any strong assumptions and how robust the results are to
714		violations of these assumptions (e.g., independence assumptions, noiseless settings,
715		should reflect on how these assumptions might be violated in practice and what the
716		implications would be
718		• The authors should reflect on the scope of the claims made e.g. if the approach was
719		only tested on a few datasets or with a few runs. In general, empirical results often
720		depend on implicit assumptions, which should be articulated.
721		• The authors should reflect on the factors that influence the performance of the approach.
722		For example, a facial recognition algorithm may perform poorly when image resolution
723		is low or images are taken in low lighting. Or a speech-to-text system might not be
724		used reliably to provide closed captions for online lectures because it fails to handle
725		technical jargon.
726		• The authors should discuss the computational efficiency of the proposed algorithms
727		and now they scale with dataset size.
728		• If applicable, the authors should discuss possible limitations of their approach to
729		While the outhout might fear that complete honesty shout limitations might be used by
730		• While the authors high rear that complete holiesty about miniations high be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover
731		limitations that aren't acknowledged in the paper. The authors should use their best
733		judgment and recognize that individual actions in favor of transparency play an impor-
734		tant role in developing norms that preserve the integrity of the community. Reviewers
735		will be specifically instructed to not penalize honesty concerning limitations.
736	3.	Theory Assumptions and Proofs
737		Question: For each theoretical result, does the paper provide the full set of assumptions and
738		a complete (and correct) proof?
739		Answer: [Yes]

740	Justification: Assumptions are clearly introduced, and all proofs can be found in the ap-
741	pendix.
742	Guidelines:
7/3	• The answer NA means that the naper does not include theoretical results
743	• All the theorems, formulas, and proofs in the paper should be numbered and cross
744	referenced
745	• All accumptions should be clearly stated or referenced in the statement of any theorems
746	• All assumptions should be clearly stated of referenced in the statement of any theorems.
747	• The proofs can either appear in the main paper or the supplemental material, but if
748	they appear in the supplemental material, the authors are encouraged to provide a short
749	
750	• Inversely, any informal proof provided in the core of the paper should be complemented
751	by formal proofs provided in appendix or supplemental material.
752	• Theorems and Lemmas that the proof relies upon should be properly referenced.
753	4. Experimental Result Reproducibility
754	Ouestion: Does the paper fully disclose all the information needed to reproduce the main ex-
755	perimental results of the paper to the extent that it affects the main claims and/or conclusions
756	of the paper (regardless of whether the code and data are provided or not)?
757	Answer: [NA]
	Institution No apportmental regulta
758	
759	Guidelines:
760	 The answer NA means that the paper does not include experiments.
761	• If the paper includes experiments, a No answer to this question will not be perceived
762	well by the reviewers: Making the paper reproducible is important, regardless of
763	whether the code and data are provided or not.
764	• If the contribution is a dataset and/or model, the authors should describe the steps taken
765	to make their results reproducible or verifiable.
766	• Depending on the contribution, reproducibility can be accomplished in various ways.
767	For example, if the contribution is a novel architecture, describing the architecture fully
768	might suffice, or if the contribution is a specific model and empirical evaluation, it may
769	be necessary to either make it possible for others to replicate the model with the same
770	dataset, or provide access to the model. In general, releasing code and data is often
771	one good way to accomplish this, but reproducibility can also be provided via detailed
772	instructions for now to replicate the results, access to a nosted model (e.g., in the case
773	of a large language model), releasing of a model checkpoint, of other means that are
//4	appropriate to the research performed.
775	• while NeurIPS does not require releasing code, the conference does require all submis-
776	sions to provide some reasonable avenue for reproducionity, which may depend on the
	(a) If the contribution is gringerille a general contribution the general could make it clear have
//8	(a) In the contribution is primarily a new algorithm, the paper should make it clear now
700	(b) If the contribution is primarily a new model architecture, the paper should describe
780	(b) in the contribution is primarily a new model architecture, the paper should describe
701	(a) If the contribution is a new model (e.g., a large language model), then there should
182	either be a way to access this model for reproducing the results or a way to reproduce
103	the model (e.g. with an open-source dataset or instructions for how to construct
785	the dataset).
786	(d) We recognize that reproducibility may be tricky in some cases in which case
787	authors are welcome to describe the particular way they provide for reproducibility
788	In the case of closed-source models, it may be that access to the model is limited in
789	some way (e.g., to registered users), but it should be possible for other researchers
790	to have some path to reproducing or verifying the results.
791	5. Open access to data and code
792	Ouestion: Does the paper provide open access to the data and code, with sufficient instruc-
793	tions to faithfully reproduce the main experimental results, as described in supplemental
79/	material?

795	Answer: [NA]
796	Justification: No experimental results.
797	Guidelines:
798	• The answer NA means that paper does not include experiments requiring code.
799	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
800	public/guides/CodeSubmissionPolicy) for more details.
801	• While we encourage the release of code and data, we understand that this might not be
802	possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
803	including code, unless this is central to the contribution (e.g., for a new open-source
804	benchmark).
805	• The instructions should contain the exact command and environment needed to run to
806	reproduce the results. See the NeurIPS code and data submission guidelines (https: //nips_cc/public/guides/CodeSubmissionPolicy) for more details
807	• The authors should provide instructions on data access and preparation including how
808 809	to access the raw data, preprocessed data, intermediate data, and generated data, etc.
810	• The authors should provide scripts to reproduce all experimental results for the new
811	proposed method and baselines. If only a subset of experiments are reproducible, they
812	should state which ones are omitted from the script and why.
813 814	• At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
815	• Providing as much information as possible in supplemental material (appended to the
816	paper) is recommended, but including URLs to data and code is permitted.
817	6. Experimental Setting/Details
818	Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
819	parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
820	results?
821	Answer: [NA]
822	
823	Guidelines:
824	 The answer NA means that the paper does not include experiments.
825	• The experimental setting should be presented in the core of the paper to a level of detail
826	that is necessary to appreciate the results and make sense of them.
827	• The full details can be provided either with the code, in appendix, or as supplemental
828	
829	7. Experiment Statistical Significance
830 831	Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
832	Answer: [NA]
833	Justification: No experimental results.
834	Guidelines:
835	• The answer NA means that the paper does not include experiments.
836	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
837	dence intervals, or statistical significance tests, at least for the experiments that support
838	the main claims of the paper.
839	• The factors of variability that the error bars are capturing should be clearly stated (for
840	example, train/test split, initialization, random drawing of some parameter, or overall
841	run with given experimental conditions).
842	• The method for calculating the error bars should be explained (closed form formula,
843	• The assumptions made should be given (a g. Normally, distributed arrays)
044	 The assumptions made should be given (c.g., Normally distributed efforts). It should be clear whether the error her is the standard deviation or the standard error.
846	of the mean.

847 848 849		• It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
850 851 852		• For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
853 854		• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
855	8.	Experiments Compute Resources
856 857 858		Question: For each experiment, does the paper provide sufficient information on the com- puter resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?
859		Answer: [NA]
860		Justification: No experimental results.
861		Guidelines:
862		• The answer NA means that the paper does not include experiments.
863 864		• The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
865 866		• The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
867		• The paper should disclose whether the full research project required more compute
868 869		than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
870	9.	Code Of Ethics
871 872		Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
873		Answer: [Yes]
874 875		Justification: Only theoretical results for the convergence of an optimization algorithm, with no foreseeable societal impact.
876		Guidelines:
877		• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
878 879		• If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
880 881		• The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
882	10.	Broader Impacts
883 884		Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?
885		Answer: [NA]
886 887		Justification: This is a theoretical work on an optimization algorithm, it has no foreseeable societal impact bey
888		Guidelines:
889		• The answer NA means that there is no societal impact of the work performed.
890 891		• If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
892		• Examples of negative societal impacts include potential malicious or unintended uses
893 894 895		(e.g., disinformation, generating take profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

896 897 898 899 900 901 902 903 904 905 906 907 908		 The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster. The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology. If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks).
909 910		mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
911	11.	Safeguards
912 913 914		Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?
915		Answer: [NA]
916		Justification: No model release.
917		Guidelines:
918		• The answer NA means that the paper poses no such risks.
919		• Released models that have a high risk for misuse or dual-use should be released with
920		necessary safeguards to allow for controlled use of the model, for example by requiring
921 922		that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
923		• Datasets that have been scraped from the Internet could pose safety risks. The authors
924		should describe how they avoided releasing unsafe images.
925		• We recognize that providing effective safeguards is challenging, and many papers do
926		not require this, but we encourage authors to take this into account and make a best faith effort
928	12.	Licenses for existing assets
929		Ouestion: Are the creators or original owners of assets (e.g. code data models) used in
930		the paper, properly credited and are the license and terms of use explicitly mentioned and
931		properly respected?
932		Answer: [NA]
933		Justification: Not using existing assets.
934		Guidelines:
935		• The answer NA means that the paper does not use existing assets.
936		• The authors should cite the original paper that produced the code package or dataset.
937		• The authors should state which version of the asset is used and, if possible, include a
938		URL.
939		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
940		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided
941		• If assets are released, the license, convright information, and terms of use in the
943		package should be provided. For popular datasets, paperswithcode.com/datasets
944		has curated licenses for some datasets. Their licensing guide can help determine the
945		license of a dataset.
946 947		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

948 949		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
950	13.	New Assets
951 952		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
953		Answer: [NA]
954		Justification: No new assets.
955		Guidelines:
056		• The answer NA means that the paper does not release new assets
950		 Researchers should communicate the details of the dataset/code/model as part of their
958		submissions via structured templates. This includes details about training, license,
959		limitations, etc.
960 961		• The paper should discuss whether and how consent was obtained from people whose asset is used.
962 963		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
964	14.	Crowdsourcing and Research with Human Subjects
965 966 967		Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?
968		Answer: [NA]
969		Justification: No human subjects or crowdsourcing.
970		Guidelines:
971		• The answer NA means that the paper does not involve crowdsourcing nor research with
972		human subjects.
973 974 975		• Including this information in the supplemental material is fine, but if the main contribu- tion of the paper involves human subjects, then as much detail as possible should be included in the main paper.
976		 According to the NeurIPS Code of Ethics, workers involved in data collection, curation.
977 978		or other labor should be paid at least the minimum wage in the country of the data collector.
979 980	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects
981 982 983 984		Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
985		Answer: [NA]
986		Justification:
987		Guidelines:
988		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects
990 991 992		 Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
993 994 995		• We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
996 997		• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.