
Stochastic model-based minimization of weakly convex functions

Fengjing Zhu

Institutes of Science and Development

Chinese Academy of Sciences

Beijing

zhufengjing24@mailsucas.ac.cn

Abstract

1 This report reviews the paper *Stochastic Model-Based Minimization of Weakly*
2 *Convex Functions*, which develops a general framework for minimizing weakly
3 convex and possibly nonsmooth functions via stochastic models. The proposed
4 algorithm iteratively samples a stochastic one-sided model of the objective and
5 performs a proximal minimization step. Under standard Lipschitz and weak con-
6 vexity assumptions, the authors prove that the expected norm of the gradient of
7 the Moreau envelope converges at a rate of $\mathcal{O}(k^{-1/4})$. The framework unifies
8 and establishes complexity guarantees for several key algorithms, including the
9 stochastic proximal point, proximal subgradient, and prox-linear methods. This
10 report summarizes the algorithmic structure, theoretical analysis, and practical
11 implications of the results. For code and reproducibility, see: [link](#).

1 Introduction

13 Optimization is fundamental in machine learning, signal processing, and operations research. Many
14 real-world problems, such as robust regression and neural network training, are nonconvex and
15 nonsmooth, making classical optimization techniques inadequate.

16 A key class of interest is *weakly convex functions*, which allow convexity after adding a quadratic
17 term. To address stochastic minimization under such settings, Davis and Drusvyatskiy propose a
18 unified stochastic model-based framework. This approach replaces exact gradients with one-sided
19 model approximations and introduces regularization via proximal steps.

20 Their algorithm covers stochastic proximal point, prox-linear, and subgradient methods, and evalu-
21 ates stationarity via the Moreau envelope. It achieves the first $\mathcal{O}(k^{-1/4})$ convergence guarantee for
22 this broad problem class.

23 The overall framework is summarized in Figure 1.

2 Related Works

25 Nonsmooth, nonconvex optimization problems arise in diverse domains such as phase retrieval,
26 covariance estimation, and neural network training. Classical methods often struggle due to non-
27 differentiability and non-convexity. Recent advances including subgradient methods (Davis et al.
28 2018), proximal algorithms (Duchi and Ruan 2019), and convex relaxations (Chen et al. 2015) offer
29 partial remedies, but typically depend on structural assumptions or careful initialization.

30 The *Moreau envelope* plays a central role in analyzing such problems, providing a smooth surrogate
31 whose gradient quantifies stationarity. Davis and Drusvyatskiy leveraged this in a stochastic model-
32 based framework, establishing the first $\mathcal{O}(k^{-1/4})$ convergence rate for weakly convex functions.

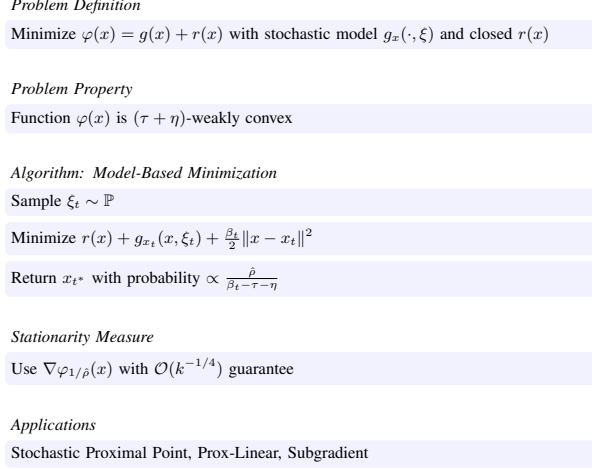


Figure 1: Overview: from problem setup to algorithm, analysis, and applications.

Extensions of the envelope to *Bregman distances* (Kan and Song 2012; Bauschke et al. 2022) have broadened its use in variational analysis and imaging (Durmus et al. 2018), though most results remain limited to convex settings.

In application, *Nonnegative Matrix Factorization (NMF)* (Gillis 2017) exemplifies constrained non-convex problems, but remains challenging under noise and incomplete observations.

3 Main Results

3.1 Problem Setting

We consider the stochastic composite optimization problem:

$$\min_{x \in \mathbb{R}^d} \varphi(x) := g(x) + r(x),$$

where $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is possibly nonconvex and nonsmooth, and $r : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is a proper closed function, often used for regularization or constraints.

The algorithm accesses g via a stochastic one-sided model $g_x(y, \xi)$, where $\xi \sim \mathbb{P}$. We assume:

- (A1) i.i.d. samples $\xi_t \sim \mathbb{P}$ are available;
- (A2) $\mathbb{E}_\xi[g_x(x, \xi)] = g(x)$, and $\mathbb{E}_\xi[g_x(y, \xi)] \geq g(y) - \frac{\tau}{2} \|y - x\|^2$;
- (A3) $y \mapsto g_x(y, \xi) + r(y)$ is η -weakly convex;
- (A4) g and $g_x(\cdot, \xi)$ are L -Lipschitz in expectation.

Then φ is $(\tau + \eta)$ -weakly convex, satisfying:

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) + \frac{\rho\lambda(1 - \lambda)}{2} \|x - y\|^2.$$

3.2 Algorithm Description

The method follows a stochastic model-based minimization framework, where at each step, a model of the objective is sampled and minimized with a proximal term. The procedure is shown in Algorithm 1.

Algorithm 1: Stochastic Model Based Minimization

Input: $x_0 \in \mathbb{R}^d$, parameters $\hat{\rho} > \tau + \eta$, sequence $\{\beta_t\}_{t=0}^T \subset (\hat{\rho}, \infty)$, iteration count T

1 **for** $t = 0$ **to** T **do**
2 Sample $\xi_t \sim \mathbb{P}$;
3 Set $x_{t+1} = \arg \min_x \{r(x) + g_{x_t}(x, \xi_t) + \frac{\beta_t}{2} \|x - x_t\|^2\}$;
4 Sample $t^* \in \{0, \dots, T\}$ with probability $\mathbb{P}(t^* = t) \propto \frac{\hat{\rho} - \tau - \eta}{\beta_t - \eta}$;
Output: x_{t^*}

3.3 Stationarity Measure via Moreau Envelope

To evaluate convergence, we use the *Moreau envelope*:

$$\varphi_\lambda(x) := \inf_y \left\{ \varphi(y) + \frac{1}{2\lambda} \|y - x\|^2 \right\},$$

which is differentiable if φ is ρ -weakly convex and $\lambda < 1/\rho$. Its gradient,

$$\nabla \varphi_\lambda(x) = \frac{1}{\lambda} (x - \text{prox}_{\lambda\varphi}(x)),$$

serves as a continuous stationarity measure.

Main Result. Under (A1)(A4), the output x_{t^*} satisfies:

$$\mathbb{E} [\|\nabla \varphi_{1/\hat{\rho}}(x_{t^*})\|^2] \leq \frac{\hat{\rho}(\varphi_{1/\hat{\rho}}(x_0) - \min \varphi) + 2\hat{\rho}^2 L^2 \sum_{t=0}^T \frac{1}{(\beta_t - \eta)(\beta_t - \hat{\rho})}}{\sum_{t=0}^T \frac{2(\hat{\rho} - \tau - \eta)}{\beta_t - \eta}}.$$

Convergence Rate. Choosing $\beta_t = \hat{\rho} + \frac{1}{\alpha\sqrt{T+1}}$ yields

$$\mathbb{E} [\|\nabla \varphi_{1/\hat{\rho}}(x_{t^*})\|^2] = \mathcal{O}(T^{-1/2}),$$

implying complexity $\mathcal{O}(\varepsilon^{-2})$ for ε -stationarity.

Proof Sketch. The proof applies Lemma 3.1 to bound the envelope value at x_{t+1} :

$$\mathbb{E}_t[\varphi_{1/\hat{\rho}}(x_{t+1})] \leq \varphi_{1/\hat{\rho}}(x_t) - \frac{\hat{\rho}(\hat{\rho} - \tau - \eta)}{2(\beta_t - \eta)} \|x_t - \hat{x}_t\|^2 + \frac{2\hat{\rho}^2 L^2}{(\beta_t - \eta)(\beta_t - \hat{\rho})}.$$

Summing over t and using $\varphi_{1/\hat{\rho}}(x_{t+1}) \geq \min \varphi$ leads to:

$$\mathbb{E} \left[\sum_{t=0}^T \frac{\hat{\rho} - \tau - \eta}{\beta_t - \eta} \|x_t - \hat{x}_t\|^2 \right] \leq 2 \cdot \frac{\varphi_{1/\hat{\rho}}(x_0) - \min \varphi}{\hat{\rho}} + 4L^2 \sum_{t=0}^T \frac{1}{(\beta_t - \eta)(\beta_t - \hat{\rho})}.$$

Since $\nabla \varphi_{1/\hat{\rho}}(x_t) = \hat{\rho}(x_t - \hat{x}_t)$, we obtain the final bound.

3.4 Special Cases

The framework recovers several classical methods by varying model type and assumptions, summarized in Table 1.

Table 1: Special Cases under the General Framework

Method	Model Type	Key Assumptions
<i>Stochastic Proximal Point</i>	Exact	$\tau = 0, \eta = \rho$
<i>Proximal Subgradient</i>	Linear	Bounded variance
<i>Prox-Linear</i>	Composite	Smooth c , convex h

67 3.5 Numerical Experiment

68 We consider the stochastic logistic regression problem:

$$\min_{x \in \mathbb{R}^d} \varphi(x) := g(x) + r(x),$$

69 where $g(x) = \mathbb{E}_\xi [\log(1 + \exp(-yz^\top x))]$ with $\xi = (z, y)$, $y \in \{-1, 1\}$, and $r(x) = \frac{\lambda}{2} \|x\|^2$ with
70 $\lambda = 0.1$.

71 We initialize at $x_0 = 0$ and run for $T = 100$ iterations, setting $\tau + \eta = 2$, $\hat{\rho} = 3$, $\beta_t > \hat{\rho}$.

72 The experiment is repeated 5 times to evaluate robustness. Results are shown in Figure 2 and 3.

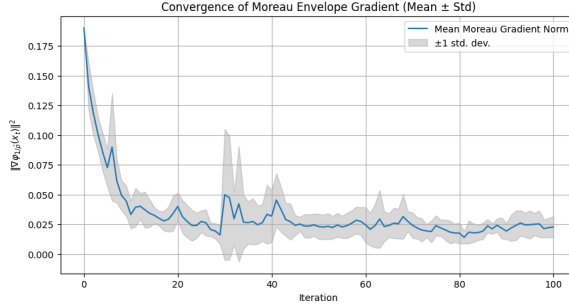


Figure 2: Convergence of Moreau Envelope Gradient

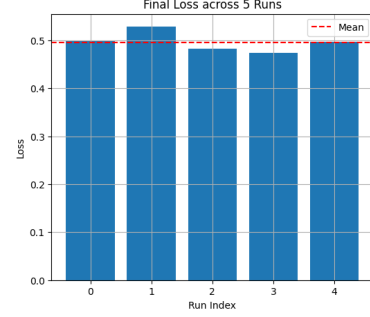


Figure 3: Final Loss across 5 Runs

73 4 Conclusion

74 This report examined the stochastic model-based minimization framework for weakly convex func-
75 tions proposed by Davis and Drusvyatskiy. Its main contribution is a unified algorithmic scheme that
76 generalizes several classical methods and provides the first complexity guarantees for convergence
77 to approximate stationary points. A central analytical tool is the Moreau envelope, whose gradient
78 yields a smooth measure of near-stationarity.

79 The framework stands out for its generality and minimal assumptions: it requires only stochastic
80 one-sided models accurate in expectation, making it widely applicable. It achieves an $\mathcal{O}(k^{-1/4})$
81 convergence rate in the Moreau envelope gradient norm, improving known results in certain settings.

82 However, limitations remain. The theory focuses on first-order stationarity, with second-order or
83 global analysis still unexplored. Performance may be sensitive to hyperparameters such as β_t and
84 $\hat{\rho}$, suggesting the value of adaptive tuning. Future enhancements could address heavy-tailed noise,
85 variance reduction, or structured modeling (e.g., Hessian surrogates) to improve robustness.

86 Promising directions include applications to large-scale nonconvex learning problems such as deep
87 networks and stochastic reinforcement learning where this framework's combination of flexibility
88 and theoretical guarantees offers a compelling foundation for stochastic nonsmooth optimization.

89 5 Some possible readings

- 90 Attouch, H (1984). “Variational convergence for functions and operators”. *Applicable Mathematics*
91 *Series*.
- 92 Bauschke, H. H., M. D., and S. Lindstrom (2022). “Using BregmanMoreau envelopes for regular-
93 ization”. *Set-Valued and Variational Analysis*.
- 94 Chen, Y., Y. Chi, and A. J. Goldsmith (2015). “Exact and stable covariance estimation from quadratic
95 sampling via convex programming”. *IEEE Transactions on Information Theory*, vol. 61, no. 7,
96 pp. 4034–4059.
- 97 Davis, D., D. Drusvyatskiy, K. J. MacPhee, and C. Paquette (2018). “Subgradient methods for sharp
98 weakly convex functions”. *Journal of Optimization Theory and Applications*, vol. 179, no. 3,
99 pp. 962–982.
- 100 Davis, D., D. Drusvyatskiy, and C. Paquette (2020). “The nonsmooth landscape of phase retrieval”.
101 *IMA Journal of Numerical Analysis*, vol. 40, no. 3, pp. 2652–2695.
- 102 Duchi, J. C. and F. Ruan (2019). “Solving (most) of a set of quadratic equalities: composite op-
103 timization for robust phase retrieval”. *Information and Inference: A Journal of the IMA*, vol. 8,
104 no. 3, pp. 471–529.
- 105 Durmus, A., É. Moulines, and M. Pereyra (2018). “Efficient Bayesian computation by proximal
106 Markov chain Monte Carlo: when Langevin meets Moreau”. *SIAM Journal on Imaging Sciences*,
107 vol. 11, no. 1, pp. 473–506.
- 108 Gillis, N. (2017). “Introduction to Nonnegative Matrix Factorization”. *SIAG/OPT Views and News*,
109 vol. 25, no. 1. arXiv:1703.00663, pp. 7–16.
- 110 Kan, I and Y Song (2012). “BregmanMoreau envelopes and Bregman proximity operators”. *Journal*
111 *of Optimization Theory and Applications*, vol. 154, no. 3, pp. 904–930.
- 112 Kecis, I. and L. Thibault (2015). “Moreau envelopes of s-lower regular functions”. *Nonlinear Anal-*
113 *ysis: Theory, Methods & Applications*, vol. 127, pp. 157–181.
- 114 Moreau, J.-J. (1965). “Proximité et dualité dans un espace hilbertien”. *Bulletin de la Société Math-*
115 *ématique de France*, vol. 93, pp. 273–299.
- 116 Pérez-Aros, P. and E. Vilches (2021). “Moreau envelope of supremum functions with applications
117 to infinite and stochastic programming”. *SIAM Journal on Optimization*, vol. 31, no. 3, pp. 1635–
118 1657.

119 **Acknowledgement**

120 I would like to express my sincere gratitude to Professor Yaoliang Yu for his inspiring lectures
121 and invaluable guidance throughout the course. His clear and accessible teaching style, along with
122 the concise and well-structured lecture notes, made complex concepts in optimization theory and
123 stochastic algorithms much more approachable. Thanks to his instruction, I have taken my first
124 solid step into the world of optimization algorithms in data science.