

Modified Gauss-Newton Algorithms under Noise

Krishna Pillutla*
University of Washington
Seattle, WA, USA
pillutla@cs.washington.edu

Vincent Roulet†
University of Washington
Seattle, WA, USA
vroulet@uw.edu

Sham M. Kakade
Harvard University
Cambridge, MA, USA
sham@seas.harvard.edu

Zaid Harchaoui
University of Washington
Seattle, WA, USA
zaid@uw.edu

Abstract—Gauss-Newton methods and their stochastic version have been widely used in machine learning and signal processing. Their nonsmooth counterparts, modified Gauss-Newton or prox-linear algorithms, can lead to contrasting outcomes when compared to gradient descent in large-scale statistical settings. We explore the contrasting performance of these two classes of algorithms in theory on a stylized statistical example, and experimentally on learning problems including structured prediction. In theory, we delineate the regime where the quadratic convergence of the modified Gauss-Newton method is active under statistical noise. In the experiments, we underline the versatility of stochastic (sub)-gradient descent to minimize such nonsmooth composite objectives.

Index Terms—Gauss-Newton, Nonsmooth, Composite problems

I. INTRODUCTION

Arising from the literature on non-linear least squares [1], [2], the Gauss-Newton method was proposed to tackle generic compositional problems of the form $\min_{w \in \mathbb{R}^d} f(\phi(w))$ by linearizing the inner function ϕ around the current iterate and solving the resulting subproblem [3].

The Gauss-Newton method and its variants such as the Levenberg-Marquardt method [4], [5] have been applied successfully in phase retrieval [6], [7], [8], nonlinear control [9], [10], and non-negative matrix factorization [11]. Modern machine learning problems such as deep learning possess a similar compositional structure, which makes Gauss-Newton-like algorithms potential good candidates [12], [13], [14]. However, in such problems, we are often interested in the generalization performance on unseen data. It is unclear whether the additional cost of solving the subproblems can be amortized by the superior efficiency of Gauss-Newton-like algorithms.

In this paper, we investigate whether modified Gauss-Newton methods or prox-linear algorithms with incremental gradient inner loops are superior to direct stochastic subgradient algorithms for nonsmooth problems with a compositional objective and a finite-sum structure. We present a statistical example and quantify when the quadratic convergence of the exact prox-linear method is not active before hitting the noise level of the problem. We present synthetic experiments that delineate the regimes where the stochastic subgradient methods outperform the prox-linear method. We also compare

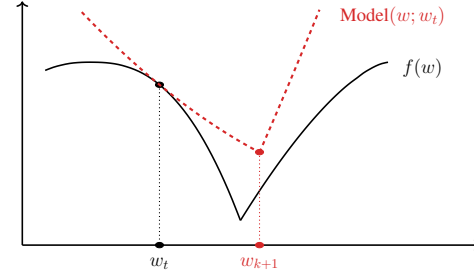


Fig. 1: The modified Gauss-Newton (a.k.a. the prox-linear) method builds a convex model of the objective $F(w)$ around w_t and finds w_{t+1} by minimizing this model.

these algorithms on a structured prediction problem with a convolutional neural network (end-to-end path planning). Experimental results suggest that modified Gauss-Newton methods or prox-linear algorithms offer marginal gains in some settings, and confirm the versatility of direct stochastic subgradient algorithms to tackle complex learning problems. All the proofs are given in the appendix.

II. PROBLEM SETTING AND OPTIMIZATION ALGORITHMS

Given $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}^k$ smooth, and $f : \mathbb{R}^k \rightarrow \mathbb{R}$ convex and Lipschitz, we consider finite-sum compositional minimization problems of the form

$$F(w) := \frac{1}{n} \sum_{i=1}^n f(\phi_i(w)) \quad (1)$$

For *multi-output regression* of input $x_i \in \mathbb{R}^p$ to output $y_i \in \mathbb{R}^k$, we take $\phi_i(w) = \varphi(x_i; w) - y_i$ as the residual of a predictor $\varphi(\cdot; w)$. We take a nonsmooth loss function such as $f(u) = \|u\|_2$ (ℓ_2 loss without the square), which is applicable in robust regression problems. A more sophisticated example is *structured prediction*, the prediction of a combinatorial object such as a sequence. Here, $\phi_i(w)$ is a score for each structured output, and f is the structural hinge loss [15], [16], [17], computed efficiently using dynamic programming [18]. For applications, see e.g. [19], [20], [21].

We compare two families of optimization algorithms, which are stochastic and nonsmooth versions of gradient descent and the Gauss-Newton algorithm. The **stochastic subgradient method** (abbreviated *SGD*) is the nonsmooth and stochastic analogue of gradient descent. In each iteration t , SGD samples an element i_t from the available n uniformly at random and

This work was supported by NSF DMS-2023166, NSF CCF-2019844, NSF DMS-2052239, NSF DMS-2134012, NSF DMS-2133244, NIH, CIFAR-LMB, and faculty research awards.

*† Now at Google Research.

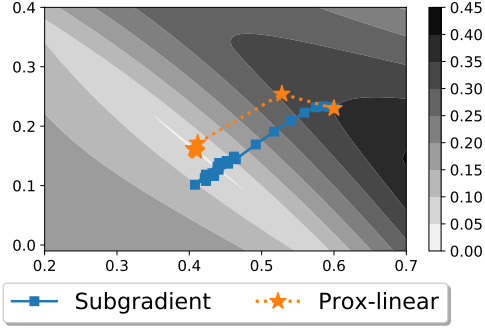


Fig. 2: A numerical comparison of gradient descent and the exact prox-linear method on a nonsmooth and nonconvex function $F: \mathbb{R}^2 \rightarrow \mathbb{R}$. The prox-linear method builds a more accurate model of the objective function, especially around points of nonsmoothness.

takes a step in the direction of its subgradient $v_t \in \partial(f \circ \phi_{i_t})(w_t)$:

$$w_{t+1} = w_t - \gamma v_t, \quad (2)$$

where γ is the learning rate and $\partial(f \circ \phi_i)$ denotes the regular (or Fréchet) subdifferential. For (1), it takes a simple form $\partial(f \circ \phi_i)(w) = \nabla \phi_i(w)^\top \partial f(\phi_i(w))$, where ∂f refers to the convex subdifferential of f and $\nabla \phi_i$ refers to the Jacobian of ϕ_i [22, Theorem 10.6]. In deep learning, the subgradient $v \in \partial(f \circ \phi_i(w))$ requires the computation of the vector-Jacobian product, readily given by reverse-mode automatic differentiation implemented in software such as PyTorch.

The **modified Gauss-Newton** method [5], also known as the **prox-linear** method [3], [23], [12], applied to (1) proceeds by finding approximate solutions of a partially linearized approximation of the objective with an additional regularization term. It computes iterates of the form

$$w_{t+1} = \arg \min_{w \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n f(\phi_i(w_t) + \nabla \phi_i(w_t)^\top (w - w_t)) + \frac{\kappa}{2} \|w - w_t\|_2^2 \right\} \quad (3)$$

As explained in Fig. 1, the prox-linear method creates a *convex model* of F around w_t by linearizing the inner function as $\phi_i(w) \approx \phi_i(w_t) + \nabla \phi_i(w_t)^\top (w - w_t)$. The next iterate (3) is the minimizer of the model plus a proximal term. Compare this with the subgradient method, where the model is $F(w_t) + v^\top (w - w_t)$, where $v \in \partial F(w_t)$; see also Fig. 2. In practice, it might not be possible to solve the subproblem (3) exactly. We consider using accelerated incremental algorithms such as Casimir-SVRG [24]. Computationally, each iteration of the inner loop requires having access to Jacobian-vector product $\nabla \phi_i(w)v$ for some vector v which are most efficiently computed via forward-mode automatic differentiation.

III. TRADEOFFS OF THE PROX-LINEAR METHOD IN STATISTICAL SETTINGS

Gauss-Newton methods and their variants are known to enjoy quadratic local convergence, provided the subproblems are solved exactly. In statistical learning problems, it is not meaningful to optimize beyond the noise level of the problem. If the noise level of the problem is large, the quadratic convergence may not be useful. We formalize this in a stylized example.

We start with a typical quadratic local convergence result of the prox-linear method in the overparameterized regime $d > nk$. In particular, we assume that the minimal singular value $\sigma_{\min}(\nabla \phi(w)^\top)$ of the transposed Jacobian of $\phi = (\phi_1; \dots; \phi_n)$ is strictly positive, or that the Jacobian $\nabla \phi$ is surjective.

Proposition 1. *Consider problem (1) where f is ℓ -Lipschitz, convex, and μ -sharp for some $\mu > 0$ (see long version for a precise definition). Suppose the function $\phi(w) = (\phi_1(w); \dots; \phi_n(w)) \in \mathbb{R}^{nk}$ is L -smooth and satisfies $\sigma_{\min}(\nabla \phi(w)^\top) \geq \nu > 0$ for any $w \in \mathbb{R}^d$. Then, the sequence $(w_t)_{t=0}^\infty$ produced by the exact prox-linear algorithm (3) with $\kappa = L\ell$ starting from arbitrary $w_0 \in \mathbb{R}^d$ converges globally to its minimum value $F(w_t) \rightarrow F^* := \min F$. Further, as soon as an iterate w_j satisfies $F(w_j) - F^* \leq (\mu\nu)^2 / (L\ell n^{3/2})$, the subsequence $(w_t)_{t=j}^\infty$ converges quadratically as*

$$F(w_{t+1}) - F^* \leq \frac{L\ell n^{3/2}}{2(\mu\nu)^2} (F(w_t) - F^*)^2. \quad (4)$$

Statistical Setting. Suppose we are given n input-output pairs (x_i, y_i) where $y_i = \varphi(x_i; \bar{w}) + \xi_i$ is given from a parameterized nonlinear function $\varphi(\cdot; \bar{w}): \mathbb{R}^p \rightarrow \mathbb{R}^k$ with parameters $\bar{w} \in \mathbb{R}^d$, and $\xi_i \sim \mathcal{N}(0, \sigma^2 I_k)$ is i.i.d. Gaussian noise. We wish to recover the true signal \bar{w} 's predictions, i.e. find \hat{w} such that $\varphi(x_i; \hat{w}) \approx \varphi(x_i; \bar{w})$ for each i . We instantiate (1) with $\phi_i(w) = \varphi(x_i; w) - y_i$ and $f = \|\cdot\|_2$ in order to solve this problem. This is different from the related choice of $f = \|\cdot\|_2^2$, which corresponds to non-linear least squares regression in the fixed design setting.

We consider *early stopping* of the optimization once we reach the noise level. That is, we stop the optimization once the objective value $F(w_t)$ in iteration t falls below $F(\bar{w})$. In this setting, we now show that the prox-linear method can enjoy quadratic local convergence only when the noise level σ of the problem is small enough. To this end, we make a general assumption on the radius R of local quadratic convergence; Prop. 1 provides a concrete lower bound on R .

Proposition 2. *Fix some $\delta \in (0, 1)$ and consider problem (1) with ϕ_i and f as defined above with the output dimension $k \geq 4 \log(2n/\delta)$. Suppose that (i) $w \mapsto \varphi(x_i; w)$ is L -smooth for each $i \in [n]$, and, (ii) φ can interpolate the data so that $\varphi(x_i; w^*) = y_i$ for each $i \in [n]$ for some $w^* \in \mathbb{R}^d$. Suppose*

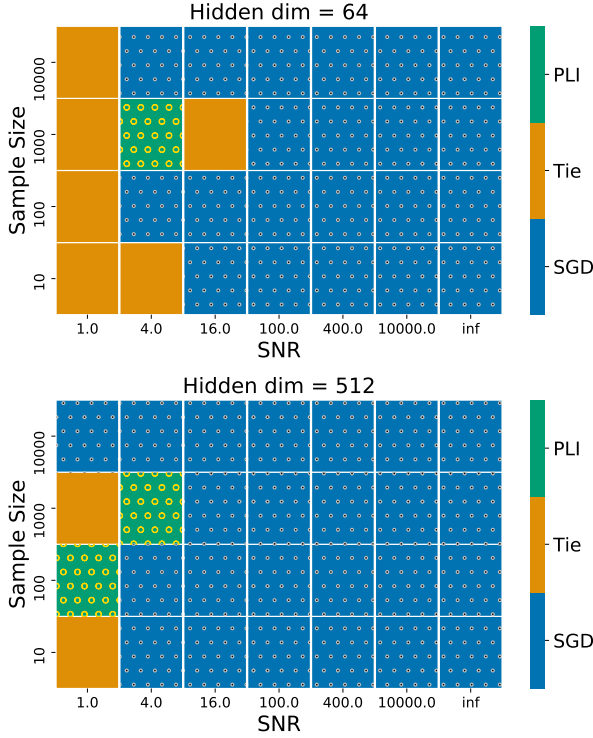


Fig. 3: Synthetic multi-output regression: stochastic subgradient method (SGD) vs. the prox-linear with incremental gradient inner-loop (PLI) while varying the number of samples n and the signal-to-noise ratio (SNR). We highlight the method which finds the smallest test ℓ_2 loss.

there exists a scalar $R > 0$ and an integer j such that for all integers $t \geq j$, we have

$$F(w_t) - \min F \leq R, \quad \text{and} \quad (5)$$

$$F(w_{t+1}) - \min F \leq \frac{1}{2R} (F(w_t) - \min F)^2.$$

Then, we have the following with probability at least $1 - \delta$. If the noise level satisfies $\sigma > \tilde{O}(R/(k^{1/2} - k^{1/4}))$, then the first iterate w_t enjoying quadratic convergence (5) satisfies $F(w_t) < F(\bar{w})$. Conversely, if the noise level satisfies $\sigma < \tilde{O}(R/(k^{1/2} + k^{1/4}))$, then the first iterate w_t enjoying quadratic convergence (5) satisfies $F(w_t) > F(\bar{w})$.

Prop. 2 shows that the potential advantages of the prox-linear method in terms of local quadratic convergence may not be relevant in some statistical problems with high noise.

IV. EXPERIMENTS

We consider 3 setups: multi-output regression, structured prediction, and solving non-linear equations. All hyper-parameters are tuned by grid search.

Synthetic Multi-output Regression. We consider a regression task of predicting output $y \in \mathbb{R}^k$ from input $x \in \mathbb{R}^p$, given a synthetic dataset $\{(x_i, y_i)\}_{i=1}^n$ of input-output pairs of varying size n where $p = 128$ and $k = 10$. We sample each input as $x_i \sim \mathcal{N}(0, \Sigma)$, where the covariance Σ exhibits a $1/j^2$ spectral

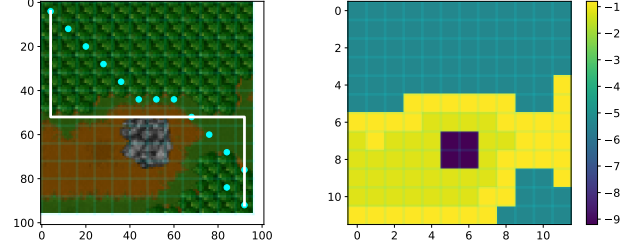


Fig. 4: Planning example. **Left:** a map and its best path in (solid) white. **Right:** corresponding rewards.

decay. The output is generated as $y_i = \varphi^*(x_i; w^*) + \sigma \xi_i$, where $\varphi^*(\cdot; w^*)$ is a multilayer perceptron (MLP) with one hidden layer of width 256, and standard normal weights w^* , while ξ_i is sampled from a standard Laplace distribution in \mathbb{R}^k and σ is the noise scale, which we vary. We define the signal-to-noise ratio (SNR) of a problem instance as $\text{SNR} = \|w^*\|^2 / \sigma^2$. Finally, the loss and evaluation measure we use is the nonsmooth ℓ_2 loss.

We vary the number of samples n and the SNR (equivalently, σ) and compare the two methods introduced in Sec. II: the stochastic subgradient method (SGD) and prox-linear with incremental gradient inner loop (PLI). We tune hyperparameters to achieve the smallest loss on a held-out validation dataset in 100 epochs and report the test loss. We run the experiment in two regimes: (a) under-parameterized, where the model is a MLP with 64 hidden units, and, (b) over-parameterized, where the MLP has 512 hidden units, compared to the 256 hidden units of $\varphi^*(\cdot; w^*)$. We see in Fig. 3 that **SGD tends to outperform PLI overall**, especially in the high SNR regime. In the low SNR regime, PLI and SGD are mostly tied in their performance, exhibiting very similar test errors.

Path Planning as Structured Prediction. Among all monotonic paths from the top left corner to the bottom right corner of a grid, our task is to find the path that maximizes the rewards collected on each tile it passes through. Specifically, we consider images generated in the game Warcraft II [25]; see Fig. 4. Each tile corresponds to some terrain such as water, desert, grass, or rock with a fixed reward (grass > desert > water > rock). As long as the rewards can directly be observed, this task can be solved by dynamic programming. In this experiment, the rewards are not directly observed; they are computed as the transformation of the raw pixels of each tile by a convolutional neural network. Our goal is to learn the reward function from a dataset of random maps with their associated optimal path. Given a map x with associated best path y , denote by $\psi(x, y, y'; w)$ the score of a path y' parameterized by w . Our objective is

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max_{y' \in \mathcal{Y}} \psi(x_i, y_i, y'; w) + \frac{\mu}{2} \|w\|_2^2, \quad (6)$$

where $(x_i)_{i=1}^n$ are the maps, $(y_i)_{i=1}^n$ are their best paths, and $\mu \geq 0$ is a regularization parameter.

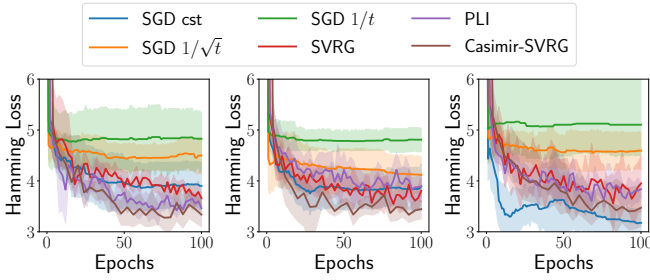


Fig. 5: Planning as a structured prediction problem. We plot the Hamming loss, which measures how good the predicted path is to the actual shortest path on unseen grids. From left to right: $\mu = 1/n, 10^{-2}/n, 10^{-4}/n$

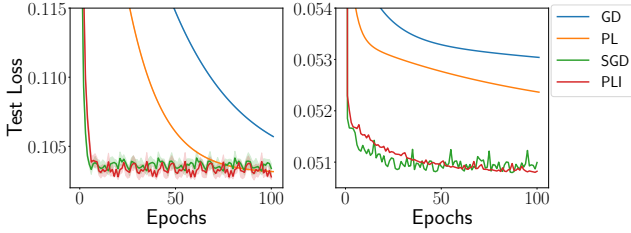


Fig. 6: Solving stochastic nonlinear equations. We plot the “test loss”, which is the objective value on a separate testing set. **Left:** ijccn1 dataset. **Right:** covtype dataset

The methods we consider are (i) stochastic subgradient methods [26], denoted SGD, with various learning rates strategies $\gamma_t = \gamma_0$, $\gamma_t = \gamma_0/\sqrt{t}$ and $\gamma_t = \gamma_0/t$, (ii) a variance-reduced stochastic sub-gradient method, denoted SVRG [27], (iii) an accelerated algorithm on the Moreau-envelope of the objective as described in [24], denoted Casimir-SVRG, (iv) a prox-linear algorithm with incremental inner loop as described in (3), denoted PLI. For SGD, subgradients of (6) can be computed by estimating the highest reward path y' associated with a given sample (x_i, y_i) for a feature map parameterized by the current parameters w . Both Casimir-SVRG and PL require smoothing the objective (6). We take inf-convolution of the max by a squared ℓ_2 norm, which can be approximated by returning the top- K shortest paths for the given score function [24].

In Fig. 5, we observe that SGD with constant step-size carefully tuned can perform as well as more sophisticated methods such as the modified Gauss-Newton method. Most importantly, for a small regularization parameter ($\mu = 10^{-4}/n$), **SGD yields the best test Hamming loss**, which in this task is the target metric.

Solving Non-linear Equations. Gauss-Newton-type methods can be applied to stochastic non-linear equations of the form

$$\min_{w \in \mathbb{R}^d} f\left(\frac{1}{n} \sum_{i=1}^n \phi_i(w)\right) \quad (7)$$

where f is a convex, possibly nonsmooth, Lipschitz function such as $\|\cdot\|_1$ and the inner mappings are smooth and typically

of the form $\phi_i(w) = \phi(x_i, w) - y_i$ [13], [14]. Problem (7) can be interpreted as ensuring that, on average, the non-linear mapping $\phi(\cdot, w)$ maps the inputs x_i to the targets y_i .

Algorithms. Denoting $\phi(w) = \frac{1}{n} \sum_{i=1}^n \phi_i(w)$, a natural baseline algorithm is to compute iterates as

$$w_{t+1} = w_t - \gamma \widehat{\nabla} \phi(w_t)^\top \nabla f(\widehat{\phi}(w_t)), \quad (8)$$

where $\widehat{\nabla} \phi(w)$ and $\widehat{\phi}(w)$ are approximations of $\nabla \phi(w)$ and $\phi(w)$ respectively that can be approximated from a mini-batch [28]; we call this “SGD”. Note that the minibatch subgradient estimates can be biased since the outer function f can be non-linear. A modified Gauss-Newton or prox-linear method adapted to the inner finite-sum performs the iterations

$$w_{t+1} = \arg \min_{w \in \mathbb{R}^d} \left\{ f\left(\widehat{\phi}(w_t) + \widehat{\nabla} \phi(w_t)(w - w_t)\right) + \frac{M}{2} \|w - w_t\|_2^2 \right\} \quad (9)$$

where each sub-problem can be solved by incremental algorithms such as the accelerated dual proximal gradient ascent [13], [14].

Experiment. We consider the experimental setting of [13]. The objective is to solve (7) where f is the Huber loss, a smooth surrogate of the nonsmooth ℓ_1 norm, and inner mappings ϕ are the concatenation of four different losses, i.e., $\phi_i(w) = (\ell_1(x_i^\top w, y_i), \dots, \ell_4(x_i^\top w, y_i))$, where the formulations of the losses can be found in [13]. The samples (x_i, y_i) are drawn from the datasets `ijccn1` or `covtype` from the LIBSVM repository [29].

We consider (i) a gradient descent denoted GD, (ii) a modified Gauss-Newton or prox-linear method denoted PL, (iii) a baseline of the form (8), denoted SGD, (iv) an incremental Gauss-Newton or prox-linear method as described in (9), denoted PLI for consistency. In Fig. 6, we observe that PL outperforms GD as expected in the batch setting. However, this advantage is longer present in the incremental setting. Here, we find that the **SGD baseline (8) performs on par with the Gauss-Newton variant PLI**.

REFERENCES

- [1] J. Nocedal and S. Wright, *Numerical Optimization*. Springer Science & Business Media, 2006.
- [2] Å. Björck, *Numerical methods for least squares problems*. SIAM, 1996.
- [3] J. V. Burke, “Descent methods for composite nondifferentiable optimization problems,” *Mathematical Programming*, vol. 33, no. 3, pp. 260–279, 1985.
- [4] K. Levenberg, “A method for the solution of certain non-linear problems in least squares,” *Quarterly of applied mathematics*, vol. 2, no. 2, pp. 164–168, 1944.
- [5] Y. Nesterov, “Modified Gauss–Newton scheme with worst case guarantees for global performance,” *Optimisation methods and software*, vol. 22, no. 3, pp. 469–483, 2007.
- [6] A. Cichocki and S.-i. Amari, *Adaptive blind signal and image processing: learning algorithms and applications*. John Wiley & Sons, 2002.
- [7] J. L. Herring, J. Nagy, and L. Ruthotto, “Gauss–Newton Optimization for Phase Recovery from the Bispectrum,” *IEEE Transactions on Computational Imaging*, vol. 6, pp. 235–247, 2019.
- [8] A. Repetti, E. Chouzenoux, and J.-C. Pesquet, “A nonconvex regularized approach for phase retrieval,” in *2014 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2014, pp. 1753–1757.

- [9] A. Sideris and J. E. Bobrow, "An efficient sequential linear quadratic algorithm for solving nonlinear optimal control problems," in *Proceedings of the 2005, American Control Conference, 2005*. IEEE, 2005, pp. 2275–2280.
- [10] V. Roulet, D. Drusvyatskiy, S. S. Srinivasa, and Z. Harchaoui, "Iterative Linearized Control: Stable Algorithms and Complexity Guarantees," in *ICML*, vol. 97, 2019, pp. 5518–5527.
- [11] K. Huang and X. Fu, "Low-Complexity Proximal Gauss-Newton Algorithm for Nonnegative Matrix Factorization," in *2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2019, pp. 1–5.
- [12] D. Drusvyatskiy and C. Paquette, "Efficiency of minimizing compositions of convex functions and smooth maps," *Mathematical Programming*, vol. 178, no. 1, pp. 503–558, 2019.
- [13] Q. Tran-Dinh, N. Pham, and L. Nguyen, "Stochastic Gauss-Newton Algorithms for Nonconvex Compositional Optimization," in *International Conference on Machine Learning*. PMLR, 2020, pp. 9572–9582.
- [14] J. Zhang and L. Xiao, "Stochastic Variance-Reduced Prox-Linear Algorithms for Nonconvex Composite Optimization," *arXiv preprint arXiv:2004.04357*, 2020.
- [15] K. Crammer and Y. Singer, "On the Algorithmic Implementation of Multiclass Kernel-based Vector Machines," *Journal of Machine Learning Research*, vol. 2, no. Dec, pp. 265–292, 2001.
- [16] B. Taskar, C. Guestrin, and D. Koller, "Max-margin Markov networks," in *Advances in Neural Information Processing Systems*, 2004, pp. 25–32.
- [17] I. Tschantzaris, T. Hofmann, T. Joachims, and Y. Altun, "Support vector machine learning for interdependent and structured output spaces," in *International Conference on Machine Learning*, 2004, p. 104.
- [18] A. J. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE Trans. Information Theory*, vol. 13, no. 2, pp. 260–269, 1967.
- [19] A. Rush, "Torch-struct: Deep structured prediction library," in *Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, July 2020, pp. 335–342.
- [20] M. J. F. Gales, S. Watanabe, and E. Fosler-Lussier, "Structured Discriminative Models for Speech Recognition: An Overview," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 70–81, 2012.
- [21] N. D. Ratliff, J. A. Bagnell, and M. Zinkevich, "Maximum Margin Planning," in *International Conference Machine Learning*, vol. 148, 2006, pp. 729–736.
- [22] R. T. Rockafellar and R. J.-B. Wets, *Variational Analysis*. Springer Science & Business Media, 2009, vol. 317.
- [23] J. V. Burke and M. C. Ferris, "A gauss—newton method for convex composite optimization," *Mathematical Programming*, vol. 71, no. 2, pp. 179–194, 1995.
- [24] K. Pillutla, V. Roulet, S. M. Kakade, and Z. Harchaoui, "A Smoother Way to Train Structured Prediction Models," in *NeurIPS*, 2018. [Online]. Available: <https://arxiv.org/pdf/1902.03228.pdf>
- [25] J. Guyomarch, "Warcraft II open-source map-editor," 2017. [Online]. Available: <http://github.com/war2/war2edit>
- [26] D. Davis and D. Drusvyatskiy, "Stochastic model-based minimization of weakly convex functions," *SIAM Journal on Optimization*, vol. 29, no. 1, pp. 207–239, 2019.
- [27] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," *Advances in neural information processing systems*, vol. 26, pp. 315–323, 2013.
- [28] M. Wang, E. X. Fang, and H. Liu, "Stochastic Compositional Gradient Descent: Algorithms for Minimizing Compositions of Expected-value Functions," *Mathematical Programming*, vol. 161, no. 1-2, pp. 419–449, 2017.
- [29] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Transactions on Intelligent Systems and Technology*, vol. 2, pp. 27:1–27:27, 2011, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [30] B. Laurent and P. Massart, "Adaptive Estimation of a Quadratic Functional by Model Selection," *Annals of Statistics*, pp. 1302–1338, 2000.

APPENDIX

A. Proof of Quadratic Convergence (Prop. 1)

Here, we give the full statement and a simple proof of Prop. 1, following standard techniques [5, Theorem 3]. We prove the proposition for the case $n = 1$ first by considering the problem $\min_{w \in \mathbb{R}^d} f \circ \phi(w)$. We then generalize to $n > 1$ for a proof of Prop. 1 in full generality.

We are interested primarily in the overparameterized case where $k \leq d$. Below, we denote $\nabla \phi(w) \in \mathbb{R}^{k \times d}$ as the Jacobian of ϕ at w . We impose the assumption that its minimal singular value is bounded away from 0 as $\sigma_{\min}(\nabla \phi(w)^\top) \geq \nu > 0$ for any $w \in \mathbb{R}^d$. This assumption implies the surjectivity of the Jacobian at each w . That is, for every $u \in \mathbb{R}^k$, there exists a $v \in \mathbb{R}^d$ such that $\nabla \phi(w)v = u$.

We also assume that the following minimum values are bounded from below:

$$f^* = \min_{u \in \mathbb{R}^k} f(u), \quad \text{and}, \quad (f \circ \phi)^* = \min_{w \in \mathbb{R}^d} f(\phi(w)).$$

We have the following statement.

Proposition 3. Consider the compositional problem $\min_{w \in \mathbb{R}^d} f \circ \phi(w)$ with following assumptions:

- (a) f is ℓ -Lipschitz continuous, convex and μ -sharp, i.e., $f(u) - f^* \geq \mu \text{dist}(u, U^*)$ for any $u \in \mathbb{R}^k$ with $\mu > 0$ and $\text{dist}(u, U^*)$ the Euclidean distance of u to $U^* = \arg \min_{u \in \mathbb{R}^k} f(u) \neq \emptyset$.
- (b) ϕ is L -smooth and satisfies $\sigma_{\min}(\nabla \phi(w)^\top) \geq \nu > 0$ for any $w \in \mathbb{R}^d$.

The sequence $(w_t)_{t=0}^\infty$ produced by the prox-linear algorithm (3) with $M = L\ell$ starting from arbitrary $w_0 \in \mathbb{R}^d$ converges globally as $(f \circ \phi)(w_t) \rightarrow (f \circ \phi)^* = f^*$. Furthermore, as soon as an iterate w_j satisfies $f(\phi(w_t)) - (f \circ \phi)^* \leq (\mu\nu)^2/(L\ell)$, the subsequence $(w_t)_{t=j}^\infty$ converges quadratically as

$$f(\phi(w_{t+1})) - (f \circ \phi)^* \leq \frac{L\ell}{2(\mu\nu)^2} (f(\phi(w_t)) - (f \circ \phi)^*)^2 \leq \frac{1}{2} (f(\phi(w_t)) - (f \circ \phi)^*).$$

Proof. For an iterate w_t of the prox-linear algorithm, denote $u_t^* = \text{Proj}_{U^*}(\phi(w_t))$ the Euclidean projection of $\phi(w_t)$ onto the set of minimizers of f such that $\text{dist}(\phi(w_t), U^*) = \|\phi(w_t) - u_t^*\|_2$.

Since the Jacobian $\nabla \phi(w_t)^\top$ is surjective, there exists v_t^* be such that $\nabla \phi(w_t)v_t^* = u_t^* - \phi(w_t)$. Furthermore, from the minimum singular value condition, there exists a choice of v_t^* such that $\|v_t^*\| \leq \|u_t^* - \phi(w_t)\|_2/\nu$ (see [5, Lemma 6] for a proof).

If $M \geq L\ell$, then the iterates of the prox-linear algorithm satisfy [12]

$$\begin{aligned} f(\phi(w_{t+1})) &\leq \min_{v \in \mathbb{R}^d} \left\{ f(\phi(w_t) + \nabla \phi(w_t)v) + \frac{M}{2} \|v\|_2^2 \right\} \\ &\stackrel{(i)}{\leq} \min_{s \in [0,1]} \left\{ f(\phi(w_t) + s\nabla \phi(w_t)v_t^*) + \frac{Ms^2}{2} \|v_t^*\|_2^2 \right\} \\ &\stackrel{(ii)}{\leq} \min_{s \in [0,1]} \left\{ f(\phi(w_t) + s(u_t^* - \phi(w_t))) + \frac{Ms^2}{2\nu^2} \|u_t^* - \phi(w_t)\|_2^2 \right\} \\ &\stackrel{(iii)}{\leq} \min_{s \in [0,1]} \left\{ sf^* + (1-s)f(\phi(w_t)) + \frac{Ms^2}{2(\nu\mu)^2} (f(\phi(w_t)) - f^*)^2 \right\}. \end{aligned} \quad (10)$$

Here, we (i) restricted the domain of the minimization to $v = sv_t^*$ with $s \in [0, 1]$, (ii) plugged in the definition of v_t^* and the bound on $\|v_t^*\|$, and, (iii) used the convexity and sharpness of f . Next, by subtracting f^* from both sides, we get

$$f(\phi(w_{t+1})) - f^* \leq \min_{s \in [0,1]} \left\{ (1-s)(f(\phi(w_t)) - f^*) + \frac{s^2 M}{2(\nu\mu)^2} (f(\phi(w_t)) - f^*)^2 \right\}.$$

If $f(\phi(w_t)) - f^* \leq (\mu\nu)^2/M$, the minimum in (10) is reached at $s = 1$ and we get

$$f(\phi(w_{k+1})) - f^* \leq \frac{M}{2(\nu\mu)^2} (f(\phi(w_t)) - f^*)^2 \leq \frac{1}{2} (f(\phi(w_t)) - f^*).$$

This is the quadratic convergence phase. On the other hand, if $f(\phi(w_t)) - f^* \geq (\mu\nu)^2/M$, then the minimum in (10) is reached at $s = (\nu\mu)^2/(M(f(\phi(w_t)) - f^*))$, and we have the bound

$$f(\phi(w_{t+1})) - f^* \leq f(\phi(w_t)) - f^* - \frac{(\mu\nu)^2}{2M}.$$

Since f is bounded from below, the sequence $f(\phi(w_t))$ converges to f^* . Hence, the minimum of the composite objective matches the minimum of the outer function, i.e., $f^* = (f \circ \phi)^*$. \square

We can now prove Prop. 1 as a corollary of Prop. 3.

Proof of Prop. 1. Consider the reduction $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{kn}$ and $\bar{f} : \mathbb{R}^{kn} \rightarrow \mathbb{R}$ given by

$$\phi(w) = (\phi_1(w); \dots; \phi_n(w)), \quad \text{and,} \quad \bar{f}(u_1; \dots; u_n) = \frac{1}{n} \sum_{i=1}^n f_i(u_i), \quad (11)$$

where we use semi-colons to denote the concatenation of vectors. The finite-sum problem (1) now reduces to $\min_w \bar{f} \circ \phi(w)$.

We have by definition that \bar{f} is convex. Next, \bar{f} is $\bar{\ell}$ -Lipschitz with $\bar{\ell} = \ell/\sqrt{n}$ since

$$\begin{aligned} |\bar{f}(u) - \bar{f}(u')| &\leq \frac{1}{n} \sum_{i=1}^n |f(u_i) - f(u'_i)| \leq \frac{\ell}{n} \sum_{i=1}^n \|u_i - u'_i\|_2 \\ &\leq \frac{\ell}{\sqrt{n}} \left(\sum_{i=1}^n \|u_i - u'_i\|_2^2 \right)^{1/2} = \frac{\ell}{\sqrt{n}} \|u - u'\|_2. \end{aligned}$$

Further, we argue that \bar{f} is $\bar{\mu}$ -sharp with $\bar{\mu} = \mu/n$. Note that $(U^*)^n$ is the argmin set of \bar{f} where U^* is the argmin set of f . Further, their minimum values satisfy $\bar{f}^* := \min \bar{f} = \min f = f^*$. Therefore, we have,

$$\bar{f}(u) - \bar{f}^* = \frac{1}{n} \sum_{i=1}^n (f(u_i) - f^*) \geq \frac{\mu}{n} \sum_{i=1}^n \text{dist}(u_i, U^*) \geq \frac{\mu}{n} \text{dist}(u, (U^*)^n),$$

where we used

$$\text{dist}(u, (U^*)^n) = \sqrt{\sum_{i=1}^n \text{dist}(u_i, U^*)^2} \leq \sum_{i=1}^n \text{dist}(u_i, U^*).$$

All the assumptions of Prop. 3 are met, and invoking it now completes the proof. \square

B. Proof of the Statistical Setting (Prop. 2)

We give the full statement of Prop. 2 and its proof.

Proposition 4. Fix some $\delta \in (0, 1)$ and consider problem (1) with g_i and f as defined above with the output dimension $k \geq 4 \log(2n/\delta)$. Suppose that (i) $w \mapsto \varphi(x_i; w)$ is L -smooth for each $i \in [n]$, and, (ii) the function φ can interpolate the data so that $\varphi(x_i; w^*) = y_i$ for each $i \in [n]$ for some $w^* \in \mathbb{R}^d$. Suppose there exists a scalar $R > 0$ and an integer j such that for all integers $t \geq j$, we have

$$F(w_t) - \min F \leq R, \quad \text{and} \quad F(w_{t+1}) - \min F \leq \frac{1}{2R} (F(w_t) - \min F)^2. \quad (12)$$

Then, we have the following with probability at least $1 - \delta$:

(a) If the noise level satisfies

$$\sigma > \frac{R}{\sqrt{k}} \left(1 - \left(\frac{4}{k} \log(2n/\delta) \right)^{1/4} \right)^{-1},$$

then the first iterate w_t enjoying quadratic convergence (12) satisfies $F(w_t) < F(\bar{w})$.

(b) Conversely, if the noise level satisfies

$$\sigma < \frac{R}{\sqrt{k}} \left(1 + \left(\frac{16}{k} \log(2n/\delta) \right)^{1/4} \right)^{-1},$$

then the first iterate w_t enjoying quadratic convergence (12) satisfies $F(w_t) > F(\bar{w})$.

Proof. First, under the interpolation assumption, we have that $0 \leq \min F \leq F(w^*) = 0$, so $\min F = 0$. Therefore, in this setting, quadratic convergence holds before the noise level if and only if $F(\bar{w}) \geq R$. To complete the proof, we show below that with probability at least $1 - \delta$ that

$$\sigma \sqrt{k} \left(1 - \left(\frac{4}{k} \log(2n/\delta) \right)^{1/4} \right) \leq F(\bar{w}) \leq \sigma \sqrt{k} \left(1 + \left(\frac{16}{k} \log(2n/\delta) \right)^{1/4} \right). \quad (13)$$

To this end, we simplify

$$F(\bar{w}) = \frac{1}{n} \sum_{i=1}^n \|\varphi(x_i; \bar{w}) - y_i\|_2 = \frac{1}{n} \sum_{i=1}^n \|\xi_i\|_2. \quad (14)$$

Noting that $\|\xi_i\|_2^2$ follows a χ^2 distribution with k degrees of freedom, a standard concentration argument shows that (see example [30, Lemma 1])

$$\mathbb{P} \left(\sigma^2 k \left(1 - 2\sqrt{\frac{\lambda}{k}} \right) \leq \|\xi_i\|_2^2 \leq \sigma^2 k \left(1 + 2\sqrt{\frac{\lambda}{k}} + \frac{2\lambda}{k} \right) \right) \geq 1 - 2\exp(-\lambda)$$

for any $\lambda > 0$. Next, we plug in $\lambda = \log(2n/\delta)$. Noting that $\lambda/k \leq 1/4$, we use the bound $\lambda/k \leq \sqrt{\lambda/k}$ and $\sqrt{1-x} \geq 1 - \sqrt{x}$ to get that

$$\sigma\sqrt{k} \left(1 - (4\lambda/k)^{1/4} \right) \leq \|\xi_i\|_2 \leq \sigma\sqrt{k} \sqrt{1 + 4\sqrt{\lambda/k}} \leq \sigma\sqrt{k} \left(1 + (16\lambda/k)^{1/4} \right)$$

with probability at least $1 - \delta/n$. Invoking the union bound over $i = 1, \dots, n$, we have with probability at least $1 - \delta$ that

$$\sigma\sqrt{k} \left(1 - \left(\frac{4}{k} \log(2n/\delta) \right)^{1/4} \right) \leq \|\xi_i\|_2 \leq \sigma\sqrt{k} \left(1 + \left(\frac{16}{k} \log(2n/\delta) \right)^{1/4} \right)$$

holds simultaneously for each $i \in [n]$. Plugging this into (14) completes the proof. \square