

SCORE: Approximating Curvature Information under Self-Concordant Regularization

Adeyemi D. Adeoye & Alberto Bemporad (IMT School for Advanced Studies, Lucca, Italy)

Problem

Data: $\{(x_n, y_n)\}_{n=1}^N$, $x_n \in \mathbb{R}^{n_p}$, $y_n \in \mathbb{R}^d$.

Objective:

$$\min_{\theta \in \mathbb{R}^{n_w}} \mathcal{L}(\theta) \triangleq \underbrace{\sum_{n=1}^N \ell(y_n, \hat{y}_n)}_{f(\theta)} + \lambda \underbrace{\sum_{j=1}^{n_w} r_j(\theta_j)}_{h(\theta)}. \quad (1)$$

- ℓ is the output-fit loss function.
- r_j , define a **regularization term** on θ , $\lambda > 0$.

Assumptions

- f, h are respectively γ_l, γ_a -strongly convex.
- f, h have γ_u, γ_b -Lipschitz continuous first derivatives g_f, g_h , respectively.
- f, h have γ_f, γ_h -Lipschitz continuous second derivatives H_f, H_h , respectively.
- λh is M_h -self-concordant:

$$|u^T (\partial^3 h(\theta)[u]) u| \leq 2M_h (u^T g_h u)^{3/2},$$

for any $u \in \mathbb{R}^{n_w}$, $M_h > 0$. Examples are the ℓ^2 -norm and pseudo-Huber function.

Main Contributions

- We present the **self-concordant regularization (SCORE)** framework for **(overparameterized) convex models** of problem (1).
- SCORE extends the idea of *Newton decrement* in [1] (or its extension in [2]) for self-concordant functions, and requires **only the regularization function** to be self-concordant, which is **less restrictive** in practice.
- We employ a more **general** (than popular Woodbury) **matrix identity** to scale the **generalized Gauss-Newton (GGN)** computations.

Curvature Approximation

- **Second-order (curvature) information in the data** is highly desirable and highly expensive.

Newton update:

$$\underbrace{\theta_{k+1} - \theta_k}_{\delta\theta} = -\rho \underbrace{(\mathbf{H}_f + \lambda \mathbf{H}_h)}_{n_w \times n_w}^{-1} (g_f + g_h).$$

- Its **GGN approximation** is **efficiently retrieved** in SCORE.
- We choose a mini-batch size $m < N$ so that $dm + 1 < n_w$ (possibly $dm \ll n_w$).

$$\delta\theta = -\rho \mathbf{H}_h^{-1} \mathbf{J}^T (\underbrace{\lambda \mathbf{I} + \mathbf{Q} \mathbf{J} \mathbf{H}_h^{-1} \mathbf{J}^T}_{(dm+1) \times (dm+1)})^{-1} e,$$

$$\mathbf{J} = \begin{bmatrix} \partial_{\theta} \hat{y} \\ \lambda \partial_{\theta} r \end{bmatrix}, \mathbf{Q} = \begin{bmatrix} \partial_{\hat{y}}^2 \ell & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, e = \begin{bmatrix} \partial_{\hat{y}} \ell \\ \mathbf{1} \end{bmatrix}.$$

- \mathbf{H}_h is always diagonal! ✓
- $\mathbf{Q}(1, 1) \equiv \mathbf{I}$ if ℓ is the squared loss ✓

GGN with SCORE

- At each mini-batch sampling time $k \in [1, t_m]$, $t_m \triangleq \lceil \frac{N}{m} \rceil$, compute the update θ_{t+1} using GGN-SCORE algorithm.

Algorithm 1 GGN-SCORE (mini-batch step)

- Input:** variables vector θ_k , data $\{(x_n, y_n)\}_{n=1}^m$, \mathbf{H}_h , \mathbf{Q} , \mathbf{J} , e , parameters $\alpha_k, M_h, \lambda > 0$
- Output:** variables vector θ_{k+1}
- Compute $g_h = \partial_{\theta_k} h(\theta_k)$
- Choose $\eta_k = (g_h^T \mathbf{H}_h^{-1} g_h)^{1/2}$
- Set $\rho_k = \frac{\alpha_k}{1 + M_h \eta_k}$
- Set $\mathbf{G} = \mathbf{H}_h^{-1} \mathbf{J}^T (\lambda \mathbf{I} + \mathbf{Q} \mathbf{J} \mathbf{H}_h^{-1} \mathbf{J}^T)^{-1} e$
- Compute $\theta_{k+1} = \theta_k - \rho_k \mathbf{G}$

Key Remark

There exist $\beta, \tilde{\beta}, K_1 > 0$, with $\mathbf{Q} \leq K_1 \mathbf{I}$, such that $\|e\| \leq \beta$, $\|\mathbf{J}\| \leq \tilde{\beta}$, and hence

$$\|\lambda \mathbf{I} + \mathbf{Q} \mathbf{J} \mathbf{H}_h^{-1} \mathbf{J}^T\| \leq \lambda + (K/\gamma_a), \quad \boxed{K = K_1 \tilde{\beta}^2}.$$

Main Theorem (Local Convergence Rate)

If $\alpha_k = \alpha = \frac{\sqrt{\gamma_a}}{\beta_1} (K + \lambda \gamma_a)$, $\beta_1 = \beta \tilde{\beta}$, then the iterates of GGN-SCORE satisfy the following descent properties:

$$\mathbb{E}[\mathcal{L}(\theta_{k+1})] \leq \mathcal{L}(\theta_k) - \left(\frac{\lambda \omega(\zeta_k)}{M_h^2} + \frac{\gamma_l \omega''(\tilde{\zeta}_k)}{2\gamma_a} - \xi \right),$$

$$\mathbb{E}\|\theta_{k+1} - \theta^*\|_{\theta_{k+1}} \leq \vartheta \|\theta_k - \theta^*\|_{\theta_k} + \frac{\gamma_u}{\beta_1} \|\theta_k - \theta^*\| + \frac{\gamma_g}{2} \|\theta_k - \theta^*\|^2,$$

$$\omega(t) \triangleq t - \ln(1+t), \quad \zeta_k \triangleq \frac{M_h}{1 + M_h \eta_k},$$

$$\tilde{\zeta}_k \triangleq M_h \eta_k, \quad \xi \triangleq \frac{2(\gamma_u + \lambda \gamma_b)}{\sqrt{\gamma_a}},$$

$$\vartheta \triangleq 1 + \frac{\lambda}{\sqrt{\gamma_a} \beta_1 (1 - M_h \|\theta_k - \theta^*\|_{\theta_k})}.$$

Experimental Results

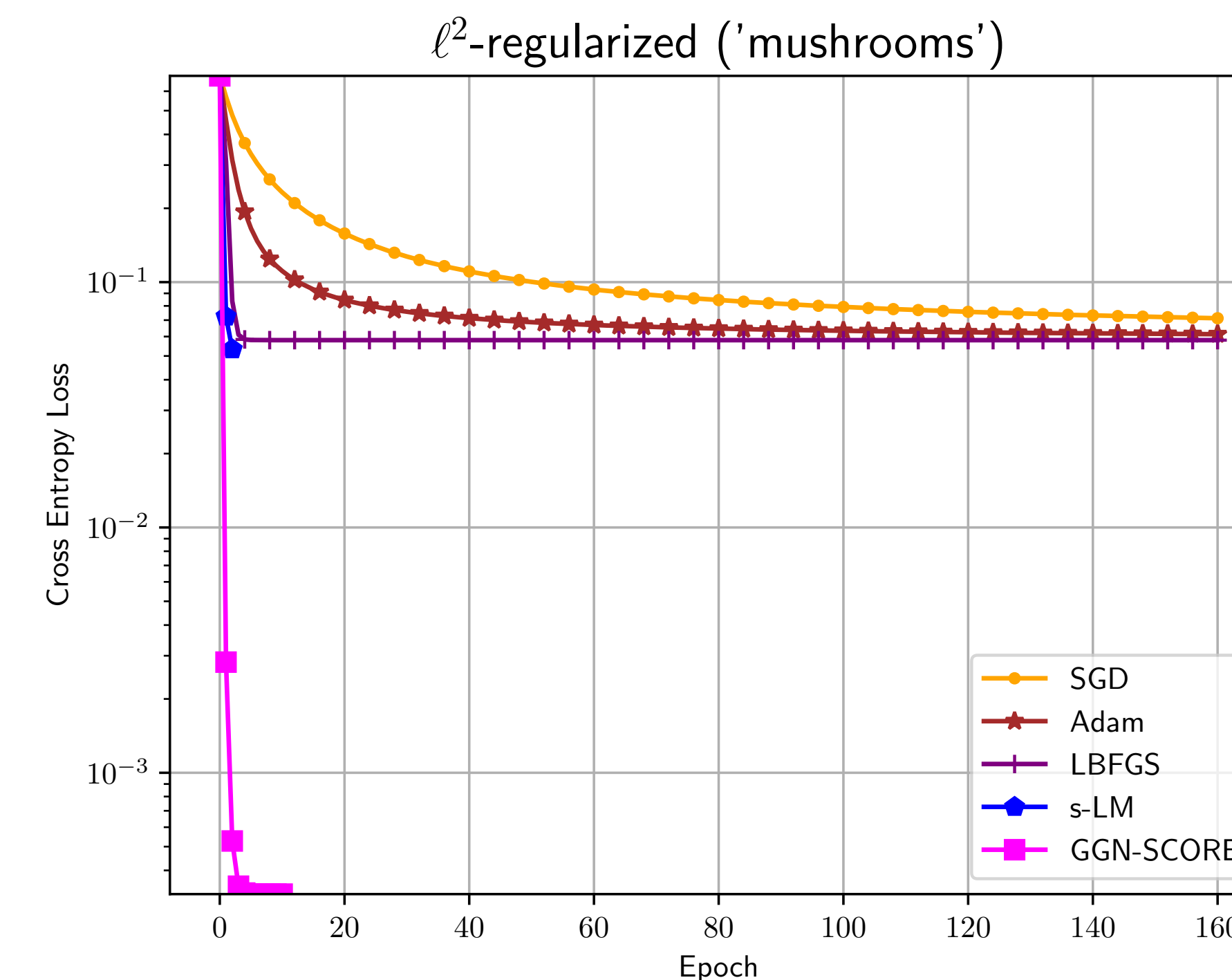


Figure 1: Convex, "overparameterized"

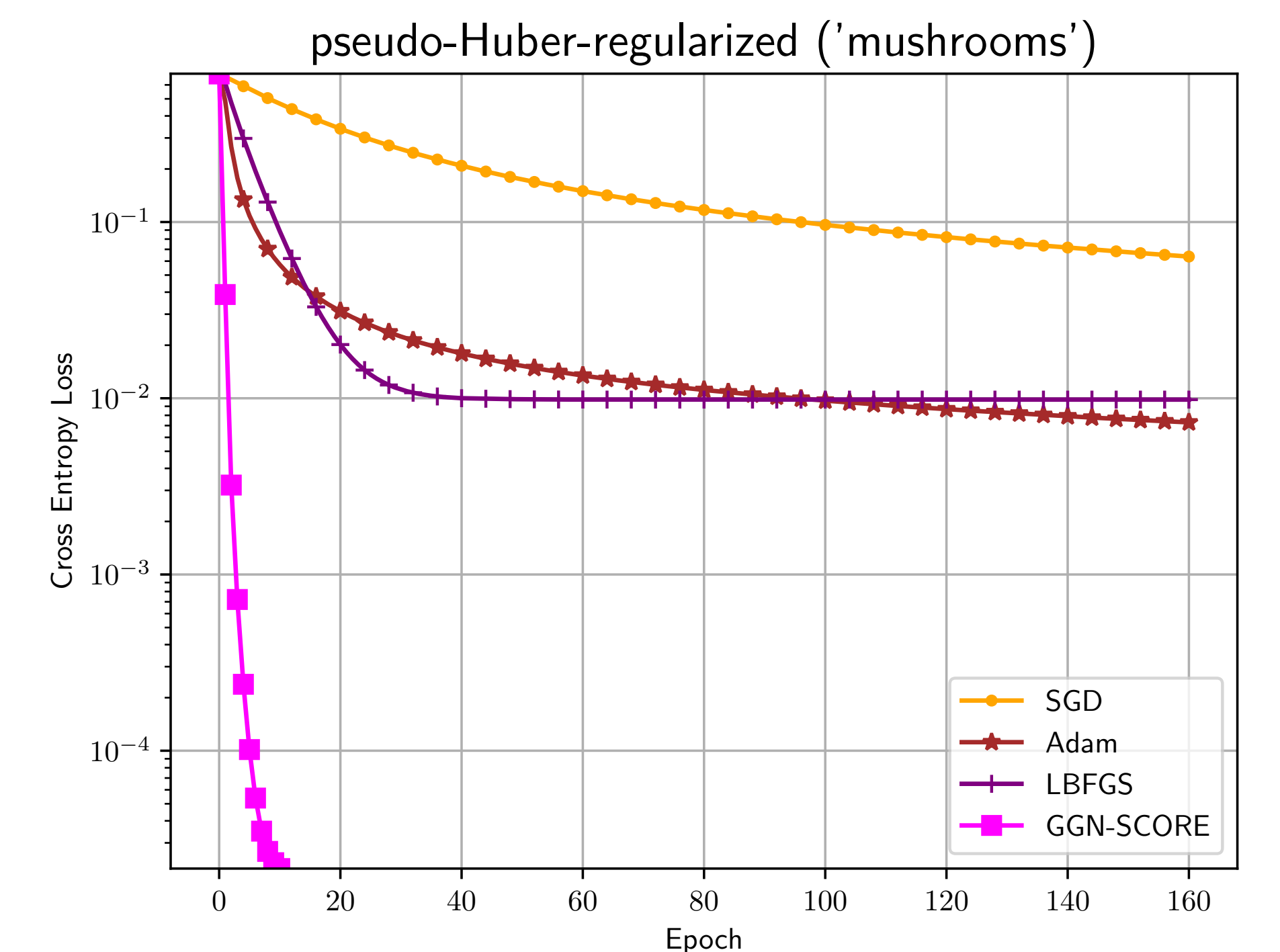


Figure 2: Nonconvex, "overparameterized"

References

- [1] Yurii Nesterov and Arkadii Nemirovskii. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- [2] Wenbo Gao and Donald Goldfarb. Quasi-newton methods: superlinear convergence without line searches for self-concordant functions. *Optimization Methods and Software*, 34(1):194–217, 2019.

Contact Information

Web adeyemiadeoye.github.io

Email adeyemi.adeoye@imtlucca.it