On algorithmic efficiency and statistical optimality in empirical risk minimization

Dmitrii M. Ostrovskii http://ostrodmit.github.io

TTI-Chicago April 8, 2019





Part I. Efficient Algorithms for Multiclass Classification

Joint with Dmitry Babichev (equal contribution) and Francis Bach



D. Babichev, D. M. Ostrovskii, F. Bach. Efficient Primal-Dual Algorithms for Large-Scale Multiclass Classification. *arXiv:1902.03755*.

Multiclass linear classification

- Dataset (x_i, y_i) , $i \in [n] = \{1, 2, ..., n\}$.
- $x_i \in \mathbb{R}^d$ features, $y_i \in \{e_1, ..., e_k\}$ class a vertex of simplex Δ_k .
- Large-scale: d, n, k up to $10^6 10^9$, as in NLP applications.



• Goal: $U^* \in \mathbb{R}^{d \times k}$ encoding k classifiers (empirical minimizer!)

• Loss $\ell : \mathbb{R}^k \times \Delta_k \to \mathbb{R}$ – misfit of predicting the (soft) label $y \in \Delta_k$ from the output of k classifiers $U^{\top} x \in \mathbb{R}^k$.

Multiclass linear classification

- Dataset (x_i, y_i) , $i \in [n] = \{1, 2, ..., n\}$.
- $x_i \in \mathbb{R}^d$ features, $y_i \in \{e_1, ..., e_k\}$ class a vertex of simplex Δ_k .
- Large-scale: d, n, k up to $10^6 10^9$, as in NLP applications.



• Goal: $U^* \in \mathbb{R}^{d \times k}$ encoding k classifiers (empirical minimizer!)

• Loss $\ell : \mathbb{R}^k \times \Delta_k \to \mathbb{R}$ – misfit of predicting the (soft) label $y \in \Delta_k$ from the output of k classifiers $U^{\top}x \in \mathbb{R}^k$.

Goal: find an optimum U^* with sublinear iteration cost O(d + n + k).

Challenge and previous work

Input: $X \in \mathbb{R}^{n \times d}$, $Y \in \Delta_k^{\otimes n}$; **output**: $U \in \mathbb{R}^{dk} \Rightarrow O(dn + dk + nk)$. **Biclass case**: k = O(1), linear time O(dn), sublinear time O(d + n).

- **Dual sampling**: sample $i \in [n]$, compute $\nabla_u \ell(u^\top x_i, y_i)$ in O(d).
 - SGD, Pegasos [Shalev-Shwartz et al., 2011], SVRG [Johnson and Zhang, 2013], SAGA [Defazio et al., 2014], SAG [Schmidt et al., 2017], SDCA [Shalev-Shwartz and Zhang, 2013], Frank-Wolfe [Lacoste-Julien et al., 2012],...

Challenge and previous work

Input: $X \in \mathbb{R}^{n \times d}$, $Y \in \Delta_k^{\otimes n}$; **output**: $U \in \mathbb{R}^{dk} \Rightarrow O(dn + dk + nk)$. **Biclass case**: k = O(1), linear time O(dn), sublinear time O(d + n).

- **Dual sampling**: sample $i \in [n]$, compute $\nabla_u \ell(u^\top x_i, y_i)$ in O(d).
 - SGD, Pegasos [Shalev-Shwartz et al., 2011], SVRG [Johnson and Zhang, 2013], SAGA [Defazio et al., 2014], SAG [Schmidt et al., 2017], SDCA [Shalev-Shwartz and Zhang, 2013], Frank-Wolfe [Lacoste-Julien et al., 2012],...
- Saddle-point approach: recast as a quasi-bilinear saddle-point problem with variables in ℝ^d and ℝⁿ, solve by mirror descent or mirror prox. Partial gradients in O(dn), accelerated to O(d + n).
 - [Grigoriadis and Khachiyan, 1995; Juditsky and Nemirovski, 2011; Xiao et al., 2017],...
 - Explicit bounds on sampling variance, certificate on duality gap.

Challenge and previous work

Input: $X \in \mathbb{R}^{n \times d}$, $Y \in \Delta_k^{\otimes n}$; **output**: $U \in \mathbb{R}^{dk} \Rightarrow O(dn + dk + nk)$.

Biclass case: k = O(1), linear time O(dn), sublinear time O(d + n).

- **Dual sampling**: sample $i \in [n]$, compute $\nabla_u \ell(u^\top x_i, y_i)$ in O(d).
 - SGD, Pegasos [Shalev-Shwartz et al., 2011], SVRG [Johnson and Zhang, 2013], SAGA [Defazio et al., 2014], SAG [Schmidt et al., 2017], SDCA [Shalev-Shwartz and Zhang, 2013], Frank-Wolfe [Lacoste-Julien et al., 2012],...
- Saddle-point approach: recast as a quasi-bilinear saddle-point problem with variables in ℝ^d and ℝⁿ, solve by mirror descent or mirror prox. Partial gradients in O(dn), accelerated to O(d + n).
 - [Grigoriadis and Khachiyan, 1995; Juditsky and Nemirovski, 2011; Xiao et al., 2017],...
 - Explicit bounds on sampling variance, certificate on duality gap.

Multiclass case: all those become O(dk) or O(dk + nk), i.e. *linear*.

Reduction to saddle-point problem

We consider the class of Fenchel-Young losses [Blondel et al., 2018]:



• *k*-softmax loss (logistic):

$$f_y(v) = h(v) := -\sum_{i \in [n]} v_i \log(v_i)$$

• *k*-hinge loss (SVM) [Shalev-Shwartz and Ben-David, 2014]: $f_{v}(v) = 1 - v^{\top}v.$

Reduction to saddle-point problem

We consider the class of Fenchel-Young losses [Blondel et al., 2018]:



• *k*-softmax loss (logistic):

$$f_y(\mathbf{v}) = h(\mathbf{v}) := -\sum_{i \in [n]} v_i \log(v_i)$$

• k-hinge loss (SVM) [Shalev-Shwartz and Ben-David, 2014]:

$$f_y(v) = 1 - v^\top y.$$

Saddle-Point Formulation

$$\min_{\|U\|_{\mathscr{U}} \leq R^*} \max_{V \in \mathcal{V}} \underbrace{\frac{1}{n} \operatorname{tr} \left[{}_{k} (V - Y)^{\top} {}_{n} X_{d} U_{k} \right]}_{\text{bi-affine}} + \underbrace{\lambda \|U\|_{\mathscr{U}}}_{\text{simple}} + \underbrace{\mathcal{F}(V)}_{=\frac{1}{n} \sum_{i \in [n]} f_{y_{i}}(v_{i}), \text{ simple}} \\ \text{and the dual feasible set } \mathcal{V} = \Delta_{k}^{\otimes n} \text{ is the direct product of simplices.}$$

Mirror descent: primer



(Composite) gradient descent: given stepsizes $\{\gamma_t\}$, iterate

$$U^{t+1} = \operatorname*{argmin}_{U \in \mathcal{U}} \left\{ \langle \nabla f(U^t), U \rangle + \Psi(U) + \frac{1}{\gamma_t} \frac{\|U - U^t\|_2^2}{2} \right\}$$

Mirror descent: replace $||U - U^t||_2^2$ with **Bregman divergence**:

$$D_{\phi_{\mathcal{U}}}(U, U^{t}) := \phi_{\mathcal{U}}(U) - \phi_{\mathcal{U}}(U^{t}) - \langle \nabla \phi_{\mathcal{U}}(U^{t}), U - U^{t} \rangle$$

for some potential $\phi_{\mathcal{U}}(\cdot)$ generalizing $\frac{1}{2} \| \cdot \|_2^2$ and such that:

Mirror descent: primer



(Composite) gradient descent: given stepsizes $\{\gamma_t\}$, iterate

$$U^{t+1} = \operatorname*{argmin}_{U \in \mathcal{U}} \left\{ \langle \nabla f(U^t), U \rangle + \Psi(U) + \frac{1}{\gamma_t} \frac{\|U - U^t\|_2^2}{2} \right\}$$

Mirror descent: replace $||U - U^t||_2^2$ with **Bregman divergence**:

$$D_{\phi_{\mathcal{U}}}(U, U^{t}) := \phi_{\mathcal{U}}(U) - \phi_{\mathcal{U}}(U^{t}) - \langle \nabla \phi_{\mathcal{U}}(U^{t}), U - U^{t} \rangle$$

for some potential $\phi_{\mathcal{U}}(\cdot)$ generalizing $\frac{1}{2} \| \cdot \|_2^2$ and such that:

- $\phi_{\mathcal{U}}$ is 1-strongly convex w.r.t. the given norm $\|\cdot\|_{\mathscr{U}}$.
- "Fitting" geometry $\Omega_{\mathcal{U}} = \widetilde{O}(\operatorname{radius}^2_{\|\cdot\|_{\mathscr{U}}}(\mathcal{U})).$
- Step easily computable usually a quasi-separable problem.

Mirror descent for saddle-point problems

Given convex sets ${\mathcal U}$ and ${\mathcal V}$, consider a saddle-point problem



- primal-dual variable: $W = (U, V) \in \mathcal{W} = \mathcal{U} \times \mathcal{V}$
- gradient field $G(W) = [\nabla_U f(U, V), -\nabla_V f(U, V)]$
- "balanced" joint potential $\phi_{\mathcal{W}}(W) = \frac{1}{\Omega_{\mathcal{U}}}\phi_{\mathcal{U}}(U) + \frac{1}{\Omega_{\mathcal{V}}}\phi_{\mathcal{V}}(V)$.

Mirror descent

$$\begin{split} \mathcal{W}^{t+1} &= \operatorname{prox}_{W^t}(\mathcal{G}(W^t)) \\ &:= \operatorname*{argmin}_{W \in \mathcal{W}} \left\{ \langle \mathcal{G}(W^t), W \rangle + \Psi(U) - \mathcal{F}(V) + \frac{D_{\phi_{\mathcal{W}}}(W, W^t)}{\gamma_t} \right\}, \end{split}$$

Mirror prox [Nemirovski and Yudin, 1983], faster convergence rate:

$$W^{t+1/2} = \operatorname{prox}_{W^t}(G(W^t)),$$

$$W^{t+1} = \operatorname{prox}_{W^t}(G(W^{t+1/2}))$$

Mirror descent for saddle-point problems

Given convex sets ${\mathcal U}$ and ${\mathcal V}$, consider a saddle-point problem



- primal-dual variable: $W = (U, V) \in \mathcal{W} = \mathcal{U} imes \mathcal{V}$
- gradient field $G(W) = [\nabla_U f(U, V), -\nabla_V f(U, V)]$
- "balanced" joint potential $\phi_{\mathcal{W}}(W) = \frac{1}{\Omega_{\mathcal{U}}}\phi_{\mathcal{U}}(U) + \frac{1}{\Omega_{\mathcal{V}}}\phi_{\mathcal{V}}(V)$.

Mirror descent

$$\begin{aligned} \mathcal{W}^{t+1} &= \operatorname{prox}_{\mathcal{W}^{t}}(\mathcal{G}(\mathcal{W}^{t})) \\ &:= \operatorname*{argmin}_{\mathcal{W} \in \mathcal{W}} \left\{ \langle \mathcal{G}(\mathcal{W}^{t}), \mathcal{W} \rangle + \Psi(\mathcal{U}) - \mathcal{F}(\mathcal{V}) + \frac{D_{\phi_{\mathcal{W}}}(\mathcal{W}, \mathcal{W}^{t})}{\gamma_{t}} \right\}, \end{aligned}$$

Stochastic mirror descent: replace $G(W^t)$ with cheaper **stochastic oracle** $\zeta(W^t)$ such that $\mathbf{E}[\zeta(W^t)] = G(W^t)$.

General convergence guarantee

Consider a saddle-point problem with quasi-bilinear objective:

$$f(U, V) = \frac{1}{n} \operatorname{tr} \left[(V - Y)^{\top} X U \right],$$

$$G(W) = \frac{1}{n} [X^{\top} (V - Y), -X U],$$

$$\zeta(W) = \frac{1}{n} [\eta_{V,Y}; -\xi_U]$$

Given norms $\|\cdot\|_{\mathscr{U}}, \|\cdot\|_{\mathscr{V}}$ and their dual norms $\|\cdot\|_{\mathscr{U}^*}, \|\cdot\|_{\mathscr{V}^*}$, define:

• Cross-Lipschitz constant $\mathcal{L}_{\mathscr{U},\mathscr{V}}$ of $G(\cdot)$:

$$\mathcal{L}_{\mathscr{U},\mathscr{V}} := \frac{1}{n} \sup_{\|U\|_{\mathscr{U}} \leqslant 1} \|XU\|_{\mathscr{V}^*}.$$

• Uniform bounds on "variances" of ξ_U and $\eta_{V,Y}$ over \mathcal{W} :

$$\bar{\sigma}_{\mathcal{U}}^2 \ge \frac{1}{n^2} \sup_{U \in \mathcal{U}} \mathbf{E}[\|XU - \xi_U\|_{\mathscr{V}^*}^2], \quad \bar{\sigma}_{\mathscr{V}}^2 \ge \dots$$

General convergence guarantee

Theorem (see [Juditsky and Nemirovski, 2011])

After T iterations of stochastic mirror descent with appropriate constant stepsize and uniform averaging, the expected duality gap is

$$\mathbf{\mathsf{E}}[\mathtt{Gap}^{\mathsf{T}}] \lesssim \frac{\mathcal{L}_{\mathscr{U},\mathscr{V}}\sqrt{\Omega_{\mathcal{U}}\Omega_{\mathcal{V}}}}{\sqrt{\mathcal{T}}} + \frac{\sqrt{\Omega_{\mathcal{U}}\bar{\sigma}_{\mathcal{V}}^2 + \Omega_{\mathcal{V}}\bar{\sigma}_{\mathcal{U}}^2}}{\sqrt{\mathcal{T}}}$$

• Cross-Lipschitz constant $\mathcal{L}_{\mathscr{U},\mathscr{V}}$ of $G(\cdot)$:

$$\mathcal{L}_{\mathscr{U},\mathscr{V}} := \frac{1}{n} \sup_{\|U\|_{\mathscr{U}} \leq 1} \|XU\|_{\mathscr{V}^*}.$$

• Uniform bounds on "variances" of ξ_U and $\eta_{V,Y}$ over \mathcal{W} :

$$\bar{\sigma}_{\mathcal{U}}^2 \ge \frac{1}{n^2} \sup_{\mathcal{U} \in \mathcal{U}} \mathbf{E}[\|X\mathcal{U} - \xi_{\mathcal{U}}\|_{\mathscr{V}^*}^2], \quad \bar{\sigma}_{\mathcal{V}}^2 \ge \dots$$

Saddle-Point Formulation

$$\min_{U \in \mathcal{U}} \max_{V \in \mathcal{V}} \frac{1}{n} \operatorname{tr} \left[(V - Y)^{\top} X U \right] + \underbrace{\lambda \| U \|_{\mathscr{U}}}_{\text{simple}} + \underbrace{\mathcal{F}_{Y}(V)}_{\text{simple}}$$
where \mathcal{U} is an $\| \cdot \|_{\mathscr{U}}$ -ball, and $\mathcal{V} = \Delta_{k}^{\otimes n}$ is the product of simplices.

Now we have to choose the setup $\|\cdot\|_{\mathscr{U}}, \|\cdot\|_{\mathscr{V}}, \phi_{\mathcal{U}}(\cdot), \phi_{\mathcal{V}}(\cdot)$ such that

$$\mathcal{L}_{\mathscr{U},\mathscr{V}}\sqrt{\Omega_{\mathcal{U}}\Omega_{\mathcal{V}}} pprox \sqrt{\Omega_{\mathcal{U}}ar{\sigma}_{\mathcal{V}}^2 + \Omega_{\mathcal{V}}ar{\sigma}_{\mathcal{U}}^2},$$

and both are small (hopefully not growing with d, n, k).

Choice of geometry: dual

Mixed norms with $p, q \in [1, \infty]$:

$$\|V\|_{\mathscr{V}} = \|V\|_{p \times q} = \left(\sum_{i \in [n]} \|V(i, :)\|_{q}^{p}\right)^{1/p},$$

We would like to choose p, q "compatible" with $\mathscr{V} = \Delta_k^{\otimes n}$. Intuition:

- When n = 1, we have $\mathscr{V} = \Delta_k$, and $\|\cdot\|_{p \times q} \equiv \|\cdot\|_q \Rightarrow |q = 1|$.
- When k = 2, we have 𝒴 = [-1,1]ⁿ cube, and || · ||_{p×q} ≡ || · ||_p. Correct choice known from [Nemirovski and Yudin, 1983]: p = 2.

$$\|\cdot\|_{\mathscr{V}} = \|\cdot\|_{2\times 1}, \quad \phi_{\mathcal{V}}(V) = -\sum_{i=1}^{n} h(V(i,:)),$$

where $h: \Delta_k \to \mathbb{R}$ is Shannon's entropy over rows.

D. M. Ostrovskii

Choice of geometry: primal

Apart from statistical considerations, we have algorithmic constraints:

- Need p = 1 to control variance when sampling features in the biclass case [Juditsky and Nemirovski, 2011].
- Need q = 1 to control variance when sampling classes (next).

Entrywise
$$\ell_1$$
-norm: $\|U\|_{\mathscr{U}} = \sum_{j \in [d]} \sum_{\kappa \in [k]} |U_{j\kappa}|.$

• Convert $\mathcal{U} = \{U : \|U\|_1 \leq R^*\}$ into *"solid"* simplex in $\mathbb{R}^{2d \times k}$:

$$\widehat{\mathcal{U}} = \{ \widehat{\mathcal{U}} \in \mathbb{R}^{2d \times k} : \mathbb{1}_{2d \times k}^{\top} \widehat{\mathcal{U}} \leq R^* \}, \qquad u \in \mathbb{R}^1$$

$$\Leftrightarrow \quad \widehat{\mathcal{U}} = [\mathcal{U}_+; \mathcal{U}_-].$$

Potential: renormalized nonnegative entropy (1-strongly cvx on $\widehat{\mathcal{U}}$),

Partial and full sampling

$$G(U^t, V^t) = \frac{1}{n} [_d X^\top {}_n (V^t - Y)_k, -{}_n X_d U_k^t].$$

Computing $G(U^t, V^t)$ is hard: O(dnk) arithmetic operations (a.o.).

• Sample the rows of U^t and $V^t - Y$, i.e. features and examples:

$$\xi_U(\mathsf{p}) = X rac{e_j e_j^ op}{\mathsf{p}_j} U, \; \; ext{where} \; \; j \sim \mathsf{p} \in \Delta_d,$$

and similarly for $\eta_{V,Y}(q)$ with $q \in \Delta_n$. Computed in O(dk + nk).

Partial and full sampling

$$G(U^t, V^t) = \frac{1}{n} [_d X^\top {}_n (V^t - Y)_k, -{}_n X_d U_k^t].$$

Computing $G(U^t, V^t)$ is hard: O(dnk) arithmetic operations (a.o.).

• Sample the rows of U^t and $V^t - Y$, i.e. features and examples:

$$\xi_U(\mathsf{p}) = X rac{e_j e_j^ op}{\mathsf{p}_j} U, \; \; ext{where} \; \; j \sim \mathsf{p} \in \Delta_d,$$

and similarly for $\eta_{V,Y}(q)$ with $q \in \Delta_n$. Computed in O(dk + nk). • Further sample classes (conditionally):

$$\xi_{U}(\mathbf{p}, \mathbf{P}) = X \frac{\mathbf{e}_{j} \mathbf{e}_{j}^{\top}}{\mathbf{p}_{i}} U \frac{\mathbf{e}_{\kappa} \mathbf{e}_{\kappa}^{\top}}{\mathbf{P}_{j\kappa}}, \text{ where } \begin{cases} j \sim \mathbf{p}, \\ \kappa \sim \mathbf{P} \in \Delta_{k}^{\otimes d}, \end{cases}$$

and similarly for $\eta_{V,Y}(q,Q)$ with $Q \in \Delta_k^{\otimes n}$.

• Computed in O(d + n + k) a.o. as we only need one row of P,Q.

D. M. Ostrovskii

Variance control

Ideally, we are interested in solving

$$\inf_{\mathbf{p}\in\Delta_d, \mathbf{P}\in\Delta_k^{\otimes d}} \left\{ \sup_{U\in\mathcal{U}} \mathbf{E}[\|XU-\xi_U(\mathbf{p})\|_{\mathscr{V}^*}^2] \right\}, \quad \inf_{\mathbf{q}\in\Delta_n, \mathbf{Q}\in\Delta_k^{\otimes n}} \{...\}.$$

Lemma. Second moments $\sup_{U \in \mathcal{U}} \mathbf{E}[\|\xi_U(\mathbf{p}, \mathbf{P})\|_{\mathscr{V}^*}^2]$, ... minimized by: $\mathbf{p}_j^* \propto \|X(:,j)\|_2 \cdot \|U(j,:)\|_1, \qquad \mathbf{P}_{j\kappa}^* \propto |U_{j\kappa}|,$ $\mathbf{q}_i^* \propto \|X(i,:)\|_\infty \cdot \|V(i,:) - Y(i,:)\|_1, \qquad \mathbf{Q}_{i\kappa}^* \propto |V_{i\kappa} - Y_{i\kappa}|.$

Final bound

With this choice of geometry and sampling distributions:

$$\mathbf{E}[\operatorname{Gap}^{T}] = \frac{R^{*}}{\sqrt{T}} \cdot \widetilde{O}\left(\max_{j \in [d]} \mathbf{E}_{n}^{1/2}[\phi_{j}^{2}] + \mathbf{E}_{n}\left[\max_{j \in [d]} |\phi_{j}|\right]\right)$$

For light-tailed data distribution, the two terms are of the same order.

Sublinear algorithm for multiclass SVM

Fully sampled oracle $\zeta_W = \frac{1}{n} [\eta_{V,Y}, -\xi_U]$ sparse, cost O(d + n + k). Challenges:

- Iterate updates $W^{t+1} = \operatorname{prox}_{W^t}(\zeta_{W^t})$ dense.
- Maintain distributions p,q that depend on the row norms:

 $p_j^* \propto \|X(:,j)\|_2 \cdot \|U(j,:)\|_1, \ \ q_i^* \propto \|X(i,:)\|_\infty \cdot \|V(i,:) - Y(i,:)\|_1.$

Key idea: Entropy-type potentials + affine composite terms = special **multiplicative** updates: each row **rescaled** (modulo one element).

- Hence our trick with "solid" simplex / renormalized entropy in U.
- Hinge loss has affine dual composite term $\mathcal{F}_{Y}(V) = 1 tr[V^{\top}Y]$.

Sublinear algorithm for multiclass SVM, cont.

Algorithm sketch:

- Initialize $U^0 \in \Delta_{2dk}, V^0 \in \Delta_k^{\otimes n}$ by uniform distributions.
- Compute initial sampling distributions p⁽⁰⁾, q⁽⁰⁾ from X.
- For *t* = 0, 1, ..., *T*
 - 1. Sample $i_t \sim q^{(t)}, j_t \sim p^{(t)}, \kappa_t \sim Q^{(t)}(i_t, :), \kappa'_t \sim P^{(t)}(j^t, :).$
 - 2. Compute $\zeta_{W^t} = \frac{1}{n} [\eta_{V^t, Y}, -\xi_{U^t}]$ (two columns κ_t and κ'_t).
 - 3. Lazy updates: update scaling factors for the rows of U and V.
 - 4. Update $||U(j_t,:)||_1, ||V(i_t,:) Y(i_t,:)||_1$ and thus $p^{(t)}$ and $q^{(t)}$.
 - 5. Explicitly update the sampled element $U(j_t, \kappa_t)$, $V(i_t, \kappa'_t)$.

Behind the scenes: lazy averaging of iterates in O(d + n + k) per iteration + O(dk + nk) postprocessing.

- Sublinear iterations O(d + n + k).
- Linear pre/postprocessing and memory: O(dn + nk + dk).
- O(dk) can be improved to $O(d\min(k, T))$.

Experiments and perspectives



Convergence of sublinear stochastic mirror descent (Full-SMD) vs. mirror prox (MP) and SGD, in natural (left) and logarithmic scale (right).

Perspectives:

(?) Faster (at least linear!) certificate for the duality gap.

- (?) Logistic and other Fenchel-Young losses.
- (?) Heavy-tailed data via Hadamard/DFT randomization technique.

Part II: Fast Rates for *M*-Estimators using Self-Concordance

D. M. Ostrovskii, F. Bach. Finite-Sample Analysis of *M*-Estimators using Self-Concordance. *arXiv:1810.06838*.

NOT IN THIS TALK:

U. Marteau-Ferey, D. M. Ostrovskii, A. Rudi, F. Bach. Beyond Least-Squares: Fast Rates for Regularized Empirical Risk Minimization through Self-Concordance. *arXiv:1902.03046*.



Statistical learning problem

Given some loss $\ell : \mathbb{R} \times \mathcal{Y} \to \mathbb{R}$, find a minimizer $\theta_* \in \mathbb{R}^d$ of expected risk:

$$\theta_* \in \operatorname{Argmin}_{\theta \in \mathbb{R}^d} L(\theta) := \mathbf{E}[\ell(x^{\top}\theta, y)],$$

where expectation $\mathbf{E}[\cdot]$ is w.r.t. the unknown distribution \mathcal{P} of $(x, y) \in \mathbb{R}^d \times \mathcal{Y}$. Since \mathcal{P} is unknown, θ_* can't be found; instead, it is estimated from **i.i.d. sample**:

$$(x_1, y_1), ..., (x_n, y_n) \sim \mathcal{P}$$
 (i.i.d.)

- Random-design classification, $\mathcal{Y} = \{0, 1\}$, and regression, $\mathcal{Y} = \mathbb{R}$.
- Performance of a candidate $\hat{\theta}$ measured by excess risk $L(\hat{\theta}) L(\theta_*)$.

Goal

• Empirical risk minimization: replace $L(\theta)$ with empirical risk:

$$\widehat{\theta}_n \in \operatorname*{Argmin}_{\theta \in \mathbb{R}^d} \left\{ L_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(x_i^\top \theta, y_i) \right\}.$$

Also called *M*-estimation in statistics.

• Special case: conditional quasi maximum likelihood estimator (qMLE):

$$\ell(\eta, y) = -\log p_{\eta}(y)$$

for some parametric family $\{p_{\eta}(y), \eta \in \mathbb{R}\}$, possibly not including the true distribution \mathcal{P}

Goal

Extend the classical asymptotic theory to finite-sample setup.

Motivation 1: Classical asymptotic theory^{*}

• Local regularity assumptions: $L(\theta)$ sufficiently smooth at θ_* , and

$$\mathbf{H} := \nabla^2 L(\theta_*) \succ \mathbf{0}.$$

• Gradient covariance $\mathbf{G} := \mathbf{E}[\nabla_{\theta} \ell(x^{\top} \theta_*, y) \nabla_{\theta} \ell(x^{\top} \theta_*, y)^{\top}]$, and let $M := H^{-1/2}GH^{-1/2}$

 $d_{\text{eff}} := \text{tr}(\mathbf{M})$ is the **effective dimension**. In well-specified models:

$$\mathbf{G}=\mathbf{H} \ \Rightarrow \ \mathbf{M}=\mathbf{I}_d \ \Rightarrow \ d_{\mathrm{eff}}=d.$$

• In the limit $n \to \infty$, Central Limit Theorem & Taylor Expansion give:

$$\begin{split} &\sqrt{n}\mathbf{H}^{-1/2}(\widehat{\theta}_n - \theta_*) \rightsquigarrow \mathcal{N}(0, \mathbf{M}), \\ &n\|\mathbf{H}^{1/2}(\widehat{\theta}_n - \theta_*)\|^2 \rightsquigarrow \mathcal{N}(0, \mathbf{M})^2, \quad 2n(L(\widehat{\theta}_n) - L(\theta_*)) \rightsquigarrow \mathcal{N}(0, \mathbf{M})^2. \\ &\left\{L(\widehat{\theta}_n) - L(\theta_*), \, \|\mathbf{H}^{1/2}(\theta_n - \theta_*)\|^2\right\} = O\left(\frac{d_{\mathsf{eff}}\log(1/\delta)}{n}\right). \end{split}$$

*[Borovkov, 1998; van der Vaart, 1998; Lehmann and Casella, 2006].

n

Motivation 2: Random-design linear regression, I

• Gaussian model $y = \mathcal{N}(x^{\top}\theta, \sigma^2)$ leads to quadratic loss and risk:

$$\ell(x^{\top}\theta, y) = \frac{1}{2\sigma^{2}}(y - x^{\top}\theta)^{2},$$

$$L(\theta) - L(\theta_{*}) = \frac{1}{2} \|\mathbf{H}^{1/2}(\theta - \theta_{*})\|^{2},$$

$$L_{n}(\theta) - L_{n}(\theta_{*}) = \frac{1}{2} \|\mathbf{H}_{n}^{1/2}(\theta - \theta_{*})\|^{2} + \underbrace{\langle \nabla L_{n}(\theta_{*}), \theta - \theta_{*} \rangle}_{\text{Zero-mean}}$$

• In particular, at any θ we have $\nabla^2 L(\theta) \equiv \mathbf{H}$ and $\nabla^2 L_n(\theta) \equiv \mathbf{H}_n$ with

$$\mathbf{H} = \mathbf{E}[xx^{\top}], \quad \mathbf{H}_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^{\top}.$$

Motivation 2: Random-design linear regression, I

• Gaussian model $y = \mathcal{N}(x^{\top}\theta, \sigma^2)$ leads to quadratic loss and risk:

$$\ell(x^{\top}\theta, y) = \frac{1}{2\sigma^2}(y - x^{\top}\theta)^2,$$

$$L(\theta) - L(\theta_*) = \frac{1}{2} \|\mathbf{H}^{1/2}(\theta - \theta_*)\|^2,$$

$$L_n(\theta) - L_n(\theta_*) = \frac{1}{2} \|\mathbf{H}_n^{1/2}(\theta - \theta_*)\|^2 + \underbrace{\langle \nabla L_n(\theta_*), \theta - \theta_* \rangle}_{\text{zero-mean}}$$

• In particular, at any θ we have $\nabla^2 L(\theta) \equiv \mathbf{H}$ and $\nabla^2 L_n(\theta) \equiv \mathbf{H}_n$ with

$$\mathbf{H} = \mathbf{E}[xx^{\top}], \quad \mathbf{H}_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^{\top}.$$

Theorem T: Estimation of a sample covariance matrix [Vershynin, 2010]

Assume $\mathbf{H}^{-1/2}X$ is subgaussian, i.e., has tails lighter than $\mathcal{N}(\mu,\mathbf{I}_d)$, and

 $n \gtrsim d + \log(1/\delta).$

Then, with probability at least $1 - \delta$ it holds:

 $0.5\mathbf{H} \preccurlyeq \mathbf{H}_n \preccurlyeq 2\mathbf{H}.$

Motivation 2: Random-design linear regression, II

Theorem 0: Finite-sample risk bound for linear regression [Hsu et al., 2012]

Assume that $\mathbf{H}^{-1/2}x$ and $\mathbf{G}^{-1/2}\nabla \ell_{\theta}(x^{\top}\theta_{*}, y)$ are subgaussian, and

 $n \gtrsim d + \log(1/\delta).$

Then w.p. at least $\geq 1 - \delta$,

 $L(\widehat{\theta}_n) - L(\theta_*) \lesssim \|\mathbf{H}^{1/2}(\widehat{\theta}_n - \theta_*)\|^2 \lesssim \|\mathbf{H}^{-1/2} \nabla L_n(\theta_*)\|^2 \lesssim \frac{d_{\mathsf{eff}} \log(1/\delta)}{n}.$

Motivation 2: Random-design linear regression, II

Theorem 0: Finite-sample risk bound for linear regression [Hsu et al., 2012]

Assume that $\mathbf{H}^{-1/2}x$ and $\mathbf{G}^{-1/2}\nabla \ell_{\theta}(x^{\top}\theta_{*},y)$ are subgaussian, and

 $n \gtrsim d + \log(1/\delta).$

Then w.p. at least $\geq 1 - \delta$,

$$L(\widehat{\theta}_n) - L(\theta_*) \lesssim \|\mathbf{H}^{1/2}(\widehat{\theta}_n - \theta_*)\|^2 \lesssim \|\mathbf{H}^{-1/2} \nabla L_n(\theta_*)\|^2 \lesssim \frac{d_{\mathsf{eff}} \log(1/\delta)}{n}$$

Proof sketch:

- 1. Since $\nabla L_n(\widehat{\theta}_n) = 0$, we have $\|\mathbf{H}_n^{1/2}(\widehat{\theta}_n \theta_*)\|^2 = \|\mathbf{H}_n^{-1/2} \nabla L_n(\theta_*)\|^2$.
- 2. Combining with Theorem T,

$$L(\widehat{\theta}_n) - L(\theta_*) = \frac{1}{2} \|\mathbf{H}^{1/2}(\widehat{\theta}_n - \theta_*)\|^2 \le 2 \|\mathbf{H}^{-1/2} \nabla L_n(\theta_*)\|^2;$$

3. The score $\mathbf{H}^{-1/2} \nabla L_n(\theta_*)$ is the average of *n* i.i.d. subgaussian vectors.

• Generally, risk is not quadratic, and Hessians are not constant:

$$\nabla^2 L(\theta) = \mathbf{H}(\theta), \quad \nabla^2 L_n(\theta) = \mathbf{H}_n(\theta).$$

- To extend the previous argument, we must control the precision of local quadratic approximation of L_n(θ) and L(θ) around θ_{*}.
- We exploit **self-concordance**, a concept introduced in [Nesterov and Nemirovski, 1994] in the theory of interior-point methods, and brought to the statistical learning context in [Bach, 2010] to study logistic regression.

Self-concordant losses

We always assume that $\ell(\eta, y)$ is convex in η .

Definition. $\ell(\eta, y)$ is self-concordant (SC) if for any $(\eta, y) \in \mathbb{R} \times \mathcal{Y}$ it holds $|\ell_{\eta}^{\prime\prime\prime}(\eta, y)| \leq [\ell_{\eta}^{\prime\prime}(\eta, y)]^{3/2}.$

Self-concordant losses

We always assume that $\ell(\eta, y)$ is convex in η .

Definition. $\ell(\eta, y)$ is self-concordant (SC) if for any $(\eta, y) \in \mathbb{R} \times \mathcal{Y}$ it holds $|\ell_{\eta}^{\prime\prime\prime}(\eta, y)| \leq [\ell_{\eta}^{\prime\prime}(\eta, y)]^{3/2}.$

• While the above definition is homogeneous in η , the next one is not:

Definition. $\ell(\eta, y)$ is **pseudo self-concordant (PSC)** if instead it holds $|\ell_n'''(\eta, y)| \le \ell_n''(\eta, y).$

- PSC losses are somewhat more common than SC ones.
- However, obtaining optimal rate for **PSC** losses requires larger sample size.

Example 1: Generalized linear models

Conditional negative log-likelihood of y given $\eta = \mathbf{x}^\top \boldsymbol{\theta}$ in the form

$$\ell(\eta, y) = -y\eta + a(\eta) - b(y),$$

where $a(\eta)$ is called the **cumulant**, and is given by

$$a(\eta) = \log \int_{\mathcal{Y}} e^{y\eta + b(y)} \mathrm{d}y.$$

This defines the density $p_{\eta}(y) \propto e^{y\eta + b(y)}$ such that $a(\eta) = \mathbf{E}_{p_{\eta}}[y]$.

PSC: Logistic regression since $(\mathcal{Y} = \{0, 1\})$, and

$$|\mathsf{a}'''(\eta)| = |\mathsf{E}_{\mathsf{p}_\eta}(y - \mathsf{E}_{\mathsf{p}_\eta}[y])^3| \leq \mathsf{E}_{\mathsf{p}_\eta}[(y - \mathsf{E}_{\mathsf{p}_\eta}[y])^2] = \mathsf{a}''(\eta).$$

PSC: Poisson regression: $Y \sim \text{Poisson}(e^{\eta})$, then $a(\eta) = \exp(\eta)$. **SC**: Exponential-response model: $Y \sim \exp(\eta)$, $\eta > 0$, $a(\eta) = -\log(\eta)$.

Example 2: Robust estimation

Loss $\ell(y,\eta) = \varphi(y-\eta)$ with $\varphi(t)$ convex, even, 1-Lipschitz, and $\varphi''(0) = 1$.



PSC: Pseudo-Huber losses: $\varphi(t) = \log \cosh(t)$, $\varphi(t) = \sqrt{1 + t^2} - 1$.

Example 2: Robust estimation

Loss $\ell(y,\eta) = \varphi(y-\eta)$ with $\varphi(t)$ convex, even, 1-Lipschitz, and $\varphi''(0) = 1$.



- **PSC**: Pseudo-Huber losses: $\varphi(t) = \log \cosh(t), \ \varphi(t) = \sqrt{1+t^2} 1.$
 - **SC**: Fenchel dual of the log-barrier $\phi(u) = -\log(1-u^2)/2$ on [-1,1]:

$$arphi(t)=rac{1}{2}\left[\sqrt{1+4t^2}-1+\log\left(rac{\sqrt{1+4t^2}-1}{2t^2}
ight)
ight]$$

Basic result

Recall that in the general case, we have the Hessian process $H(\theta)$, given by

$$\mathbf{H}(\theta) := \mathbf{E}[\ell''(x^{\top}\theta, y)xx^{\top}] = \mathbf{E}[\widetilde{x}(\theta)\widetilde{x}(\theta)^{\top}],$$

where $\widetilde{x}(\theta) := [\ell''(x^{\top}\theta, y)]^{1/2}x$ is the *curvature-scaled design*.

Basic result

Recall that in the general case, we have the Hessian process $H(\theta)$, given by

$$\mathbf{H}(\theta) := \mathbf{E}[\ell''(x^{\top}\theta, y)xx^{\top}] = \mathbf{E}[\widetilde{x}(\theta)\widetilde{x}(\theta)^{\top}],$$

where $\widetilde{x}(\theta) := [\ell''(x^{\top}\theta, y)]^{1/2}x$ is the *curvature-scaled design*.

Theorem 1: Finite-sample excess risk bound for self-concordant losses

Assume that the loss is SC, and $\mathbf{G}^{-1/2}\nabla \ell_{\theta}(x^{\top}\theta_{*}, y)$ and $\mathbf{H}(\theta_{*})^{-1/2}\widetilde{x}(\theta_{*})$ are subgaussian. Whenever

 $n \gtrsim d + \log(1/\delta) \lor d_{\mathsf{eff}} \, d \log(1/\delta),$

with probability $1-\delta$ it holds

 $L(\widehat{\theta}_n) - L(\theta_*) \lesssim \|\mathbf{H}^{1/2}(\widehat{\theta}_n - \theta_*)\|^2 \lesssim \|\mathbf{H}^{-1/2} \nabla L_n(\theta_*)\|^2 \lesssim \frac{d_{\mathsf{eff}} \log(1/\delta)}{n}.$

Distribution conditions are local (only at θ_{*});
 arge sample complexity – scaling as the product O(d_{eff} d).

D. M. Ostrovskii

Analysis: Key observation

Given $\mathbf{H}(\theta) = \nabla^2 L(\theta)$, consider **Dikin ellipsoids** of $L(\theta)$ at θ_0 :

$$\Theta(\theta_0, r) := \{\theta : \|\mathbf{H}(\theta_0)^{1/2}(\theta - \theta_0)\|^2 \le r^2\}.$$

Key Observation. Suppose that $\mathbf{H}_n(\theta) \simeq \mathbf{H}_n(\theta_*)$ w.h.p. for any $\theta \in \Theta(\theta_*, r)$. Then, $\widehat{\theta}_n \in \operatorname{Argmin} L_n(\theta)$ belongs to $\Theta(\theta_*, r)$ once

 $\|\mathbf{H}(\theta_*)^{-1/2}\nabla L_n(\theta_*)\|^2 \lesssim r^2,$

Proof sketch:

- Indeed, by definition of $\widehat{\theta}_n$, $L_n(\widehat{\theta}_n) \leq L_n(\theta_*)$. Assume $\widehat{\theta}_n \notin \Theta_n(\theta_*, r)$.
- Pick $\overline{\theta}_{n} \in [\theta_{*}, \widehat{\theta}_{n}]$ on the border of $\Theta_{n}(\theta_{*}, r)$. Still, $L_{n}(\overline{\theta}_{n}) \leq L_{n}(\theta_{*})$. $0 \geq L_{n}(\overline{\theta}_{n}) - L_{n}(\theta_{*}) \approx \langle \nabla L_{n}(\theta_{*}), \overline{\theta}_{n} - \theta_{*} \rangle + \underbrace{\|\mathbf{H}_{n}(\theta_{*})^{1/2}(\overline{\theta}_{n} - \theta_{*})\|^{2}}_{\approx r^{2} \text{ (by Theorem T)}}.$

By Cauchy-Schwarz, we arrive at ||H(θ_{*})^{-1/2}∇L_n(θ_{*})||² ≥ r².
 Contradiction!

D. M. Ostrovskii

Analysis: recap

- Once $\hat{\theta}_n$ has been localized to the neighborhood of θ_* where $L_n(\theta)$ is quadratic, we can mimick the argument for linear regression.
- Localization is guaranteed once

$$\|\mathbf{H}(\theta_*)^{-1/2}\nabla L_n(\theta_*)\|^2 \lesssim r^2,$$

which leads to the second threshold for n:

$$n \gtrsim rac{1}{r^2} d_{ ext{eff}} \log(1/\delta).$$

• Now the question is:

What is the radius r of the Dikin ellipsoid in which $\mathbf{H}_n(\theta) \approx \mathbf{H}_n(\theta_*)$?

- Short answer: we can afford $r^2 \approx 1/d$ using self-concordance.

Analysis: self-concordance at play

What is the radius r of the Dikin ellipsoid in which $\mathbf{H}_n(\theta) \approx \mathbf{H}_n(\theta_*)$?

1. Recall that

$$\mathbf{H}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell''(x_i^{\top} \boldsymbol{\theta}, y_i) x_i x_i.$$

2. Integrating $|\ell'''(\eta, y)| \leq [\ell''(\eta, y)]^{\frac{3}{2}}$ from $\eta_* = x^{\top} \theta_*$ to $\eta = x^{\top} \theta$,

$$\frac{1}{(1+[\ell''(\eta_*,y)]^{\frac{1}{2}}|\eta-\eta_*|)^2} \leq \frac{\ell''(\eta,y)}{\ell''(\eta_*,y)} \leq \frac{1}{(1-[\ell''(\eta_*,y)]^{\frac{1}{2}}|\eta-\eta_*|)^2},$$
$$\frac{1}{(1+|\langle \widetilde{x}(\theta_*), \theta-\theta_*\rangle|)^2} \leq \frac{\ell''(x^\top\theta,y)}{\ell''(x^\top\theta_*,y)} \leq \frac{1}{(1-|\langle \widetilde{x}(\theta_*), \theta-\theta_*\rangle|)^2}.$$

3. The ratio is bounded when $|\langle \widetilde{x}(heta_*), heta - heta_*
angle| \le 1/2$, i.e., by Cauchy-Schwarz,

$$\underbrace{\|\mathbf{H}(\theta_*)^{-1/2}\widetilde{x}(\theta_*)\|}_{\approx \sqrt{d}} \cdot \underbrace{\|\mathbf{H}(\theta_*)^{1/2}(\theta - \theta_*)\|}_{r} \leq 1/2 \implies \boxed{r \lesssim \frac{1}{\sqrt{d}}}. \blacksquare$$

Improved result

Theorem 2: Improved sample complexity for self-concordant losses

Assume the loss is SC, $\mathbf{G}^{-1/2}\nabla \ell_{\theta}(x^{\top}\theta_{*}, y)$ is subgaussian, and $\mathbf{H}(\theta)^{-1/2}\widetilde{x}(\theta)$ is subgaussian in the unit Dikin ellipsoid of $L(\theta)$ at θ_{*} :

$$\Theta(\theta_*,1) = \{\theta : \|\mathbf{H}(\theta_*)^{1/2}(\theta-\theta_*)\| \le 1\}.$$

Then, the asymptotic bound holds already when

 $n \gtrsim d \log(d/\delta) \lor d_{\mathsf{eff}} \log(1/\delta),$

Main idea:

- Sample complexity $n \gtrsim d_{\text{eff}} d$ in Theorem 1 is due to Hessian approximation in the small Dikin ellipsoid with $r = O(1/\sqrt{d})$ rather than r = O(1).
- We need to prove that H_n(θ) ≍ H_n(θ_{*}) for θ ∈ Θ(θ_{*}, 1). To do this, we combine self-concordance with a covering argument.

Covering the Dikin ellipsoid



- 1. It is rather easy to prove first that $\mathbf{H}(\theta) \simeq \mathbf{H}(\theta_*)$ on $\Theta(\theta_*, 1)$.
- 2. By **SC**, $\mathbf{H}_n(\theta) \simeq \mathbf{H}_n(\theta_0)$ in a small ellipsoid $\Theta(\theta_0, 1/\sqrt{d})$.
- 3. Now cover $\Theta(\theta_*, 1)$ by $\Theta(\theta_0, 1/\sqrt{d})$ with θ_0 in the epsilon-net $\mathcal{N}_{\varepsilon}$. Control uniform deviations $\mathbf{H}_n(\theta)$ from $\mathbf{H}(\theta)$ on $\mathcal{N}_{\varepsilon}$. Note: $\log |\mathcal{N}_{\varepsilon}| = O(d \log d)$.

• Because of the "incorrect" power of ℓ'' in **PSC**, we need an extra condition: $\mathbf{E}[xx^{\top}] \leq \rho \mathbf{E}[\ell''(x^{\top}\theta_*, y)xx^{\top}].$

for some $\rho > 0$. This condition is standard in logistic regression [Bach, 2010].

- We obtain similar results, but with ho times worse sample complexity.
- Worst-case bounds on ρ can be exponentially bad [Hazan et al., 2014].
- However, this is not the case in practice. E.g., we show that

$$\rho \lesssim \|\theta_*\|_{\Sigma}^{3/2}$$

in logistic regression with Gaussian design $x \sim \mathcal{N}(0, \Sigma)$.

We use **self-concordance** – a concept from optimization – to obtain statistical results – **near-optimal** rates in finite-sample regimes in some statistical models.

Behind the scenes: regularized estimators (ℓ_2 and ℓ_1 -regularization).

Perspectives:

- Heavy-tailed distributions
- Iterative algorithms: stochastic approximation, Quasi-Newton, ...
- Other models: covariance matrix estimation with log det loss, ...

Thank you!

- Bach, F. (2010). Self-concordant analysis for logistic regression. *Electronic Journal* of Statistics, 4:384–414.
- Blondel, M., Martins, A. F. T., and Niculae, V. (2018). Learning classifiers with Fenchel-Young losses: Generalized entropies, margins, and algorithms. arXiv preprint arXiv:1805.09717.
- Borovkov, A. A. (1998). *Mathematical statistics*. Gordon and Breach Science Publishers.
- Defazio, A., Bach, F., and Lacoste-Julien, S. (2014). Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654.
- Grigoriadis, M. D. and Khachiyan, L. G. (1995). A sublinear-time randomized approximation algorithm for matrix games. *Operations Research Letters*, 18(2):53–58.
- Hazan, E., Koren, T., and Levy, K. Y. (2014). Logistic regression: tight bounds for stochastic and online optimization. In *Conference on Learning Theory*, pages 197–209.

References II

- Hsu, D., Kakade, S. M., and Zhang, T. (2012). Random design analysis of ridge regression. *The Journal of Machine Learning Research*, 23(9):1–24.
- Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive variance reduction. In Advances in Neural Information Processing Systems, pages 315–323.
- Juditsky, A. and Nemirovski, A. (2011). First order methods for nonsmooth convex large-scale optimization, ii: Utilizing problem's structure. Optimization for Machine Learning, pages 149–183.
- Lacoste-Julien, S., Jaggi, M., Schmidt, M., and Pletscher, P. (2012). Block-coordinate frank-wolfe optimization for structural svms. *arXiv preprint arXiv:1207.4747*.
- Lehmann, E. L. and Casella, G. (2006). *Theory of point estimation*. Springer Science & Business Media.
- Nemirovski, A. and Yudin, D. (1983). Problem complexity and method efficiency in optimization.
- Nesterov, Y. and Nemirovski, A. S. (1994). *Interior-point polynomial algorithms in convex programming*. Society of Industrial and Applied Mathematics.

References III

- Schmidt, M., Le Roux, N., and Bach, F. (2017). Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1-2):83–112.
- Shalev-Shwartz, S. and Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
- Shalev-Shwartz, S., Singer, Y., Srebro, N., and Cotter, A. (2011). Pegasos: Primal estimated sub-gradient solver for SVM. *Mathematical programming*, 127(1):3–30.
- Shalev-Shwartz, S. and Zhang, T. (2013). Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14(Feb):567–599.
- van der Vaart, A. W. (1998). Asymptotic statistics. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv:1011.3027*.
- Xiao, L., Yu, A. W., Lin, Q., and Chen, W. (2017). DSCOVR: randomized primal-dual block coordinate algorithms for asynchronous distributed optimization. arXiv preprint arXiv:1710.05080.

Part III. Covariance Estimation for Heavy-Tailed Distributions

Joint work with Alessandro Rudi



D. M. Ostrovskii, A. Rudi. Affine-Invariant Covariance Estimation for Heavy-Tailed Distributions. *arXiv:1902.03086*.

Covariance Estimation Problem

Estimate the covariance matrix $\mathbf{S} = \mathbf{E}[XX^{\top}]$ from i.i.d. copies $X_1, ..., X_n$ of $X \in \mathbb{R}^d$.

• Sample covariance estimator:

$$\widetilde{\mathbf{S}} = \frac{1}{n} \sum_{i=1}^{n} X_i X_i^{\top}.$$

• Relative spectral-norm guarantee: when X is light-tailed,

$$\frac{\|\widetilde{\mathbf{S}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \sqrt{\frac{\mathbf{r}(\mathbf{S})\log(d/\delta)}{n}} \quad \text{with probability} \quad \geq 1 - \delta,$$

where $\mathbf{r}(\mathbf{S}) = \frac{\operatorname{tr}(\mathbf{S})}{\|\mathbf{S}\|}$ is effective rank (Lounici & Kolchinskii 2014). • Due to affine equivariance, this gives the guarantee

$$\left(1-\sqrt{\frac{d\log(d/\delta)}{n}}\right)\mathbf{S}\preccurlyeq \widetilde{\mathbf{S}}\preccurlyeq \left(1+\sqrt{\frac{d\log(d/\delta)}{n}}\right)\mathbf{S}.$$

Heavy-Tailed Distributions

$$\frac{\|\widetilde{\mathbf{S}} - \mathbf{S}\|}{\|\mathbf{S}\|} \lesssim \sqrt{\frac{\mathbf{r}(\mathbf{S})\log(d/\delta)}{n}}$$
$$\left(1 - \sqrt{\frac{d\log(d/\delta)}{n}}\right) \mathbf{S} \preccurlyeq \widetilde{\mathbf{S}} \preccurlyeq \mathbf{S} \left(1 + \sqrt{\frac{d\log(d/\delta)}{n}}\right).$$

- The second guarantee is more useful in some applications (random-design linear regression, noisy PCA).
- Both require light-tailed assumptions on X, i.e. **S** is not robust.
- Minsker (2014) proposes an estimator with a spectral-norm guarantee for heavy-tailed distributions (4th moment):

$$\widehat{\mathbf{S}}^{\mathsf{Min}} = rac{1}{n} \sum_{i=1}^{n} \tau(\|X_i\|) X_i X_i^{\top}.$$

where $\tau(x)$ is the truncation map. Breaks affine equivariance!

Main Idea

• Minsker (2014) proposes an estimator with a spectral-norm guarantee for **heavy-tailed distributions** (4th moment):

$$\widehat{\mathbf{S}}^{\mathsf{Min}} = \frac{1}{n} \sum_{i=1}^{n} \tau(\|X_i\|) X_i X_i^{\top}.$$

where $\tau(x)$ is the truncation map.

• In fact, the desired \preccurlyeq guarantee would hold for

$$\widehat{\mathbf{S}}^* = \frac{1}{n} \sum_{i=1}^n \tau(\|\mathbf{S}^{-1/2} X_i\|) X_i X_i^{\top},$$

but it is unavailable, as $S^{-1/2}X_i$'s are not observed.

Start with $\widehat{\mathbf{S}}_0 = \widehat{\mathbf{S}}^{\text{Min}}$, and imitate $\widehat{\mathbf{S}}^*$ iteratively: $\widehat{\mathbf{S}}_{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^n \tau(\|\widehat{\mathbf{S}}_t^{-1/2} X_i\|) X_i X_i^{\top}$,

Actual Estimator

$$\widehat{\mathbf{S}}_{t+1} \leftarrow \frac{1}{n} \sum_{i=1}^{n} \tau(\|\widehat{\mathbf{S}}_{t}^{-1/2} X_{i}\|) X_{i} X_{i}^{\top},$$

- Separate the sample X₁, ..., X_n into batches, and use the new batch to compute S
 [°]
 [°]
- Iterative regularization: replace $\widehat{\mathbf{S}}_t^{-1/2}$ with $(\widehat{\mathbf{S}}_t + \lambda_t \mathbf{I})^{-1/2}$, where $\lambda_t = 2^{-t} ||\mathbf{S}||$. Convergence in

 $O(\log(cond(\mathbf{S})))$ iterations,

where cond(S) is the condition number of S.

• Similar complexity as for the sample covariance estimator!

- Noisy PCA: convergence of the noisy power method depends on the eigenvalue ratios which are controlled by the ≼ guarantee.
- Random-design linear regression: we achieve optimal convergence rates in the setting with heavy-tailed design.

Convergence: deterministic term

Recall that
$$\Omega_{\mathcal{U}}=\widetilde{O}(R^{*2})$$
 and $\Omega_{\mathcal{V}}=\widetilde{O}(n).$

Lemma.

$$\mathcal{L}_{\mathscr{U},\mathscr{V}} = \frac{1}{n} \max_{j \in [d]} \|X_j\|_2,$$

where $X_j \in \mathbb{R}^n$ is the histogram of the *j*-th feature.

Corollary. Deterministic mirror descent converges with the rate

$$\frac{\mathcal{L}_{\mathscr{U},\mathscr{V}}\sqrt{\Omega_{\mathcal{U}}\Omega_{\mathcal{V}}}}{\sqrt{T}} = \widetilde{O}\left(\frac{R^*}{\sqrt{T}}\max_{j\in[d]}\sqrt{\frac{\|X_j\|_2^2}{n}}\right) = \widetilde{O}\left(\frac{R^*}{\sqrt{T}}\max_{j\in[d]}\mathbf{E}_n^{1/2}[\phi_j^2]\right),$$

where $\mathbf{E}_n^{1/2}[\phi_i^2]$ is the empirical 2nd moment of the *j*-th feature.

- When $n \to \infty$, $\max_{j \in [d]} \mathbf{E}_n^{1/2}[\phi_j^2] \to \max_{j \in [d]} \mathbf{E}_*^{1/2}[\phi_j^2]$ by LLN.
- We have equivalence in finite-sample when data is light-tailed.

Variance control (full)

Ideally, we are interested in solving

$$\inf_{\mathbf{p}\in\Delta_d,\mathbf{P}\in\Delta_k^{\otimes d}}\left\{\sup_{U\in\mathcal{U}}\mathbf{E}[\|XU-\xi_U(\mathbf{p})\|_{\mathscr{V}^*}^2]\right\},\quad\inf_{\mathbf{q}\in\Delta_n,\mathbf{Q}\in\Delta_k^{\otimes n}}\left\{\ldots\right\}.$$

Lemma. Second moments $\sup_{U \in \mathcal{U}} \mathbf{E}[\|\xi_U(\mathbf{p}, \mathbf{P})\|_{\mathcal{V}^*}^2]$, ... minimized by: $\mathbf{p}_j^* \propto \|X(:,j)\|_2 \cdot \|U(j,:)\|_1, \qquad \mathbf{P}_{j\kappa}^* \propto |U_{j\kappa}|,$ $\mathbf{q}_i^* \propto \|X(i,:)\|_\infty \cdot \|V(i,:) - Y(i,:)\|_1, \qquad \mathbf{Q}_{i\kappa}^* \propto |V_{i\kappa} - Y_{i\kappa}|.$

Moreover, the corresponding variance proxies satisfy:

$$\sigma_{\mathcal{U}}^2 \leqslant 4R^{*2}\mathcal{L}^2_{\mathscr{U},\mathscr{V}}, \quad \sigma_{\mathcal{V}}^2 \leqslant 8n\mathcal{L}^2_{\mathscr{U},\mathscr{V}} + 8\mathbf{E}_n^2\left[\max_{j\in[d]}|\phi_j|\right]$$

Final bound

$$\mathsf{E}[\operatorname{Gap}^{\mathsf{T}}] = \frac{R^*}{\sqrt{T}} \cdot \widetilde{O}\left(\max_{j \in [d]} \mathsf{E}_n^{1/2}[\phi_j^2] + \mathsf{E}_n\left[\max_{j \in [d]} |\phi_j|\right]\right)$$

For light-tailed data distribution, the two terms are of the same order.

D. M. Ostrovskii

Example 1: Generalized linear models (full)

Conditional negative log-likelihood of y given $\eta = \mathbf{x}^\top \boldsymbol{\theta}$ in the form

$$\ell(\eta, y) = -y\eta + a(\eta) - b(y),$$

where $a(\eta)$ is called the **cumulant**, and is given by

$$a(\eta) = \log \int_{\mathcal{Y}} e^{y\eta + b(y)} \mathrm{d}y.$$

This defines the density $p_{\eta}(y) \propto e^{y\eta + b(y)}$ such that $a(\eta) = \mathsf{E}_{p_{\eta}}[y]$, and

$$\ell^{(s)}_\eta(\eta,y)=a^{(s)}(\eta)=\mathsf{E}_{
ho_\eta}[(y-\mathsf{E}_{
ho_\eta}y)^s],\quad s\geq 2$$

SC/**PSC** relate 2nd and 3rd central moments of $p_{\eta}(\cdot)$. **PSC**: Logistic regression since ($\mathcal{Y} = \{0, 1\}$), and

$$|a'''(\eta)| = |\mathsf{E}_{
ho_\eta}(y - \mathsf{E}_{
ho_\eta}[y])^3| \le \mathsf{E}_{
ho_\eta}[(y - \mathsf{E}_{
ho_\eta}[y])^2] = a''(\eta).$$

PSC: **Poisson regression:** $Y \sim \text{Poisson}(e^{\eta})$, then $a(\eta) = \exp(\eta)$.

SC: Exponential-response

model: $Y \sim \text{Exp}(\eta)$, $\eta > 0$, $a(\eta) = -\log(\eta)$.