

Coordinate Descent for Private Composite Empirical Risk Minimization

Paul Mangold

(Aurélien Bellet, Marc Tommasi, Joseph Salmon)

MAGNET SEMINAR

February 17th, 2022

Supervised learning:

$$D = \{d_1, \dots, d_n\} \subseteq \mathcal{X} \times \mathcal{Y}.$$

Learn *good parameters* $w \in \mathbb{R}^p$ for

$$h_w : \mathcal{X} \mapsto \mathcal{Y}.$$

Empirical Risk Minimization

$$\arg \min_{w \in \mathbb{R}^p} F(w) = \frac{1}{n} \underbrace{\sum_{i=1}^n \ell(w; d_i)}_{f(w)}.$$

Assumptions:

- $\ell(\cdot; d)$ convex, component-Lipschitz $\forall d \in \mathcal{X} \times \mathcal{Y}$.
- $\ell(\cdot; d_i)$ component-smooth $\forall d_i \in D$.

$\ell(\cdot; d)$ is component-Lipschitz for $d \in \mathcal{X} \times \mathcal{Y}$:

$$|\ell(w + te_j; d) - \ell(w; d)| \leq L_j |t|.$$

$\ell(\cdot; d)$ is component-Lipschitz for $d \in \mathcal{X} \times \mathcal{Y}$:

$$|\ell(w + te_j; d) - \ell(w; d)| \leq L_j |t|.$$

$$\Rightarrow \text{for } w \in \mathbb{R}^p, |\nabla_j \ell(w; d)| \leq L_j.$$

$\ell(\cdot; d)$ is component-Lipschitz for $d \in \mathcal{X} \times \mathcal{Y}$:

$$|\ell(w + te_j; d) - \ell(w; d)| \leq L_j |t|.$$

$$\Rightarrow \text{for } w \in \mathbb{R}^p, |\nabla_j \ell(w; d)| \leq L_j.$$

for $w, w' \in \mathbb{R}^p$ and $d, d' \in \mathcal{X} \times \mathcal{Y}$,

$$|\nabla_j \ell(w; d) - \nabla_j \ell(w'; d')| \leq 2L_j.$$

$\ell(\cdot; d)$ is component-smooth for $d \in D$:

$$|\nabla_j \ell(w + te_j; d) - \nabla_j \ell(w; d)| \leq M_j |t|.$$

$\ell(\cdot; d)$ is component-smooth for $d \in D$:

$$|\nabla_j \ell(w + te_j; d) - \nabla_j \ell(w; d)| \leq M_j |t|.$$

\Rightarrow for $w, w' \in \mathbb{R}^p$,

$$f(w') \leq f(w) + \langle \nabla f(w), w' - w \rangle + \frac{1}{2} \|w' - w\|_M^2.$$

$$\text{where } \|w\|_M^2 = \sum_{j=1}^p M_j w_j^2.$$

Composite ERM

$$\arg \min_{w \in \mathbb{R}^p} F(w) = \frac{1}{n} \underbrace{\sum_{i=1}^n \ell(w; d_i)}_{f(w)} + \psi(w).$$

Assumptions:

- $\ell(\cdot; d)$ convex, component-Lipschitz $\forall d \in \mathcal{X} \times \mathcal{Y}$.
- $\ell(\cdot; d_i)$ component-smooth $\forall d_i \in D$.
- $\psi(w) = \sum_{j=1}^p \psi_j(w_j)$ convex and separable.

Example: LASSO

$$\arg \min_{w \in \mathbb{R}^p} \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_1.$$

Example: LASSO

$$\arg \min_{w \in \mathbb{R}^p} \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_1.$$

$$M_j = \frac{1}{n} \sum_{i=1}^n |X_{i,j}|^2.$$

$\mathcal{A} : D \mapsto w$ is (ϵ, δ) -Differentially Private

$$\Pr [\mathcal{A}(D) \in \mathcal{S}] \leq e^\epsilon \Pr [\mathcal{A}(D') \in \mathcal{S}] + \delta.$$

(D and D' differ on one element.)

 – C. Dwork, “*Differential Privacy*”, 2006.

Differentially Private Composite ERM:

$$\arg \min_{w \in \mathbb{R}^p} F(w) = f(w) + \psi(w)$$

such that w is (ϵ, δ) -DP.

- ☐ – K. Chaudhuri, C. Monteleoni, and A. D. Sarwate, “*Differentially Private Empirical Risk Minimization*”, 2011.

Solving DP-ERM?

$$\arg \min_{w \in \mathbb{R}^p} F(w) = f(w) + \psi(w)$$

such that w is (ϵ, δ) -DP.

Solving DP-ERM?

$$\arg \min_{w \in \mathbb{R}^p} F(w) = f(w) + \psi(w)$$

such that w is (ϵ, δ) -DP.

The classical: DP-SGD.

Solving DP-ERM?

$$\arg \min_{w \in \mathbb{R}^p} F(w) = f(w) + \psi(w)$$

such that w is (ϵ, δ) -DP.

The classical: DP-SGD.

The challenger: DP-CD.

Stochastic Gradient Descent

$$w^{t+1} = \text{prox}_{\eta\psi} (w^t - \eta\xi),$$

where $\mathbb{E}[\xi] = \nabla f(w^t)$.



- A. Beck and M. Teboulle, “*A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems*”, 2009.
- K. Mishchenko, A. Khaled, and P. Richtárik, “*Proximal and Federated Random Reshuffling*”, 2021.

Private Stochastic Gradient Descent

$$w^{t+1} = \text{prox}_{\eta\psi} \left(w^t - \eta \left(\xi + \mathcal{N}(\sigma^2 \mathbf{1}) \right) \right),$$

where $\mathbb{E}[\xi] = \nabla f(w^t)$.



- R. Bassily, A. Smith, and A. Thakurta, “*Differentially Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds*”, 2014.
- D. Wang, M. Ye, and J. Xu, “*Differentially Private Empirical Risk Minimization Revisited: Faster and More General*”, 2017.

Coordinate Descent

$$w_j^{t+1} = \text{prox}_{\eta_j \psi_j} \left(w_j^t - \eta_j \nabla_j f(w^t) \right),$$

where $j \in \{1, \dots, p\}$ is chosen randomly.



- P. Richtárik and M. Takáč, “Iteration Complexity of Randomized Block-Coordinate Descent Methods for Minimizing a Composite Function”, 2014.

Private Coordinate Descent

$$w_j^{t+1} = \text{prox}_{\eta_j \psi_j} \left(w_j^t - \eta_j \left(\nabla_j f(w^t) + \mathcal{N}(\sigma_j^2) \right) \right),$$

where $j \in \{1, \dots, p\}$ is chosen randomly.

DP-SGD:

$$w^{t+1} = \text{prox}_{\eta\psi} \left(w^t - \eta \left(\xi + \mathcal{N}(\sigma^2 \mathbf{1}) \right) \right),$$

with $\mathbb{E}[\xi] = \nabla f(w^t)$.

DP-CD:

$$w_j^{t+1} = \text{prox}_{\eta_j \psi_j} \left(w_j^t - \eta_j \left(\nabla_j f(w^t) + \mathcal{N}(\sigma_j^2) \right) \right).$$

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth

Step sizes

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth
Step sizes	$\eta \propto 1/\beta$	

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth
Step sizes	$\eta \propto 1/\beta$	$\eta_j \propto 1/M_j$

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth
Step sizes	$\eta \propto 1/\beta$	$\eta_j \propto 1/M_j$
Noise scale		

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth
Step sizes	$\eta \propto 1/\beta$	$\eta_j \propto 1/M_j$
Noise scale	$\sigma^2 = O\left(\frac{\Lambda^2 T}{n^2 \epsilon^2}\right)$	

Since $\sup_{d, d'} \|\nabla \ell(\cdot; d) - \nabla \ell(\cdot; d')\|_2 \leq 2\Lambda$.

	DP-SGD	DP-CD
Assumptions on f	Λ -Lipschitz β -smooth	L -comp.-Lipschitz M -comp.-smooth
Step sizes	$\eta \propto 1/\beta$	$\eta_j \propto 1/M_j$
Noise scale	$\sigma^2 = O\left(\frac{\Lambda^2 T}{n^2 \epsilon^2}\right)$	$\sigma_j^2 = O\left(\frac{L_j^2 T}{n^2 \epsilon^2}\right)$

Since $\sup_{d, d'} |\nabla_j \ell(\cdot; d) - \nabla_j \ell(\cdot; d')| \leq 2L_j$.

What about this $O(\cdot)$?

DP-SGD: sampling rate $q \leq 1$,

$$\epsilon \leq \frac{1}{\alpha - 1} \log \left(\frac{1}{2} \sum_{k=0}^{\infty} \binom{\alpha}{k} (1-q)^{\alpha-k} q^k \exp \left(\frac{k^2 - k}{2\sigma^2} \right) \left(\operatorname{erfc} \left(\frac{z_1 - k}{\sqrt{2}\sigma} \right) + \operatorname{erfc} \left(\frac{k - z_1}{\sqrt{2}\sigma} \right) \right) \right) + \frac{\log(1/\delta)}{\alpha - 1}.$$

What about this $O(\cdot)$?

DP-SGD: sampling rate $q \leq 1$,

$$\epsilon \leq \frac{1}{\alpha - 1} \log \left(\frac{1}{2} \sum_{k=0}^{\infty} \binom{\alpha}{k} (1-q)^{\alpha-k} q^k \exp \left(\frac{k^2 - k}{2\sigma^2} \right) \left(\operatorname{erfc} \left(\frac{z_1 - k}{\sqrt{2}\sigma} \right) + \operatorname{erfc} \left(\frac{k - z_1}{\sqrt{2}\sigma} \right) \right) \right) + \frac{\log(1/\delta)}{\alpha - 1}.$$

DP-CD: *not needed!*

- ☰ – I. Mironov, K. Talwar, and L. Zhang, “Rényi Differential Privacy of the Sampled Gaussian Mechanism”, 2019.

Utility?

$$\mathbb{E}[F(w^T) - F^*] \leq ?$$

Ad hoc dual norms.

$$\|w\|_M = \sqrt{\sum_{j=1}^p M_j w_j^2} \quad \|w\|_{M^{-1}} = \sqrt{\sum_{j=1}^p \frac{1}{M_j} w_j^2}.$$

	DP-SGD	DP-CD
	convex	convex
Assumptions on f	Λ -Lipschitz	L -comp.-Lipschitz
	β -smooth	M -comp.-smooth
Utility ($\mathbb{E}[F(w) - F^*] \leq \dots$)		

	DP-SGD	DP-CD
Assumptions on f	convex Λ -Lipschitz β -smooth	convex L -comp.-Lipschitz M -comp.-smooth
Utility ($\mathbb{E}[F(w) - F^*] \leq \dots$)	$O\left(\frac{\Lambda R_I \sqrt{p}}{n\epsilon}\right)$	

Where $R_I^2 = \max(F(w^0) - F(w^*), \|w^0 - w^*\|_2^2)$.

	DP-SGD	DP-CD
Assumptions on f	convex Λ -Lipschitz β -smooth	convex L -comp.-Lipschitz M -comp.-smooth
Utility ($\mathbb{E}[F(w) - F^*] \leq \dots$)	$O\left(\frac{\Lambda R_I \sqrt{p}}{n\epsilon}\right)$	$O\left(\frac{\ L\ _{M-1} R_M \sqrt{p}}{n\epsilon}\right)$

Where $\|L\|_{M-1} = (\sum_{j=1}^p L_j^2 / M_j)^{1/2}$,
 $R_I^2 = \max(F(w^0) - F(w^*), \|w^0 - w^*\|_2^2)$,
 $R_M^2 = \max(F(w^0) - F(w^*), \|w^0 - w^*\|_M^2)$.

DP-SGD

DP-CD

	μ_I -strongly-convex	μ_M -strongly-convex
	<i>w.r.t.</i> $\ \cdot\ _2$	<i>w.r.t.</i> $\ \cdot\ _M$
Assumptions on f	Λ -Lipschitz	L -comp.-Lipschitz
	β -smooth	M -comp.-smooth

Utility

$(\mathbb{E}[F(w) - F^*] \leq \dots)$

	DP-SGD	DP-CD
Assumptions on f	μ_I -strongly-convex <i>w.r.t.</i> $\ \cdot\ _2$ Λ -Lipschitz β -smooth	μ_M -strongly-convex <i>w.r.t.</i> $\ \cdot\ _M$ L -comp.-Lipschitz M -comp.-smooth
Utility ($\mathbb{E}[F(w) - F^*] \leq \dots$)	$O\left(\frac{\Lambda^2 p}{\mu_I n^2 \epsilon^2}\right)$	

	DP-SGD	DP-CD
Assumptions on f	μ_I -strongly-convex <i>w.r.t.</i> $\ \cdot\ _2$ Λ -Lipschitz β -smooth	μ_M -strongly-convex <i>w.r.t.</i> $\ \cdot\ _M$ L -comp.-Lipschitz M -comp.-smooth
Utility ($\mathbb{E}[F(w) - F^*] \leq \dots$)	$O\left(\frac{\Lambda^2 p}{\mu_I n^2 \epsilon^2}\right)$	$O\left(\frac{\ L\ _{M^{-1}}^2 p}{\mu_M n^2 \epsilon^2}\right)$

Where $\|L\|_{M^{-1}} = (\sum_{j=1}^p L_j^2 / M_j)^{1/2}$.

CD vs. SGD: Who wins?

Balanced M_j 's:

→ DP-SGD up to p times better.

CD vs. SGD: Who wins?

Imbalanced M_j 's:

→ DP-CD up to $\frac{\max_j M_j}{\min_j M_j}$ times better.

Practical Comments.

- Gradient Clipping.
- Hyperparameters.
- Private Smoothness Constants?

$$|\nabla_j \ell(w; d)| \leq L_j \quad \text{for } d \in \mathcal{X} \times \mathcal{Y}.$$

$$|\nabla_j \ell(w; d)| \leq L_j \quad \text{for } d \in \mathcal{X} \times \mathcal{Y}.$$

$$\Rightarrow \sigma_j^2 = O\left(\frac{L_j^2 T}{n^2 \epsilon^2}\right).$$

Clip!

$$\text{clip}(\nabla_j \ell(w; d), C_j) = \begin{cases} \pm C_j, & \text{if } |\nabla_j \ell(w; d)| > C_j, \\ \nabla_j \ell(w; d), & \text{otherwise.} \end{cases}$$

Clip!

$$\text{clip}(\nabla_j \ell(w; d), C_j) = \begin{cases} \pm C_j, & \text{if } |\nabla_j \ell(w; d)| > C_j, \\ \nabla_j \ell(w; d), & \text{otherwise.} \end{cases}$$
$$|\text{clip}(\nabla_j \ell, C_j)| \leq C_j \Rightarrow \sigma_j = O\left(\frac{C_j^2 T}{n^2 \epsilon^2}\right).$$

	DP-SGD	DP-CD
Clipping	C	
Step sizes	$\eta = \gamma/\beta$	

	DP-SGD	DP-CD
Clipping	C	$C_j = \sqrt{\frac{M_j}{\text{tr}(M)}} C$
Step sizes	$\eta = \gamma/\beta$	$\eta_j = \gamma/M_j$

$$\arg \min_{w \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(w, d_i)}_{f(w)} + \psi(w)$$

$M^{(i)}$ -comp-smooth

$$\arg \min_{w \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(w, d_i)}_{f(w)} + \psi(w)$$

$M^{(i)}$ -comp-smooth

$$M_j = \frac{1}{n} \sum_{i=1}^n M_j^{(i)}.$$

$$\arg \min_{w \in \mathbb{R}^p} \overbrace{\frac{1}{n} \sum_{i=1}^n (X_{i,:}^T w - y_i)^2}^{f(w)} + \lambda \|w\|_1$$

$M^{(i)}$ -comp-smooth

$$M_j = \frac{1}{n} \sum_{i=1}^n |X_{i,j}|^2.$$

$$\arg \min_{w \in \mathbb{R}^p} \overbrace{\frac{1}{n} \sum_{i=1}^n (X_{i,:}^T w - y_i)^2}^{f(w)} + \lambda \|w\|_1$$

$M^{(i)}$ -comp-smooth

$$M_j = \frac{1}{n} \sum_{i=1}^n |X_{i,j}|^2.$$

Remark: standardized data $\rightarrow M_j = 1$.

$$\arg \min_{w \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(w, d_i)}_{f(w)} + \psi(w)$$

$M^{(i)}$ -comp-smooth

Let $\epsilon' \leq \epsilon$ (e.g., $\epsilon' = 0.1\epsilon$).

$$M_j^{priv} = \frac{1}{n} \sum_{i=1}^n M_j^{(i)} + \text{Lap} \left(\frac{p \cdot \max_i M_j^{(i)}}{n\epsilon'} \right).$$

$$\arg \min_{w \in \mathbb{R}^p} \overbrace{\frac{1}{n} \sum_{i=1}^n f_i(w)}^{f(w)} + \psi(w)$$

Let $\epsilon' \leq \epsilon$ (e.g., $\epsilon' = 0.1\epsilon$).

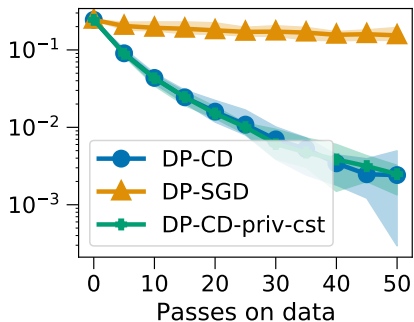
$$M_j^{priv} = \frac{1}{n} \sum_{i=1}^n \text{clip}(M_j^{(i)}, \mathbf{b}_j) + \text{Lap} \left(\frac{p \cdot \mathbf{b}_j}{n\epsilon'} \right).$$

Experiments.

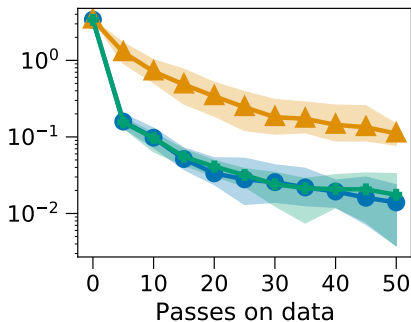
- On logistic regression and LASSO.
- Tune $\left\{ \begin{array}{l} \text{step size,} \\ \text{clipping threshold,} \\ \text{number of iterations.} \end{array} \right.$
- Average over 10 runs.

Imbalanced Dataset

Imbalanced Dataset



ELECTRICITY (RAW)
 $n = 45,312$ $p = 8$
Logistic Regression

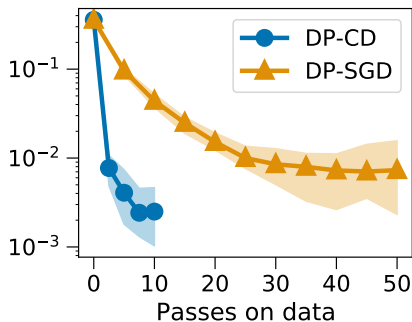


CALIFORNIA (RAW)
 $n = 20,640$ $p = 8$
LASSO

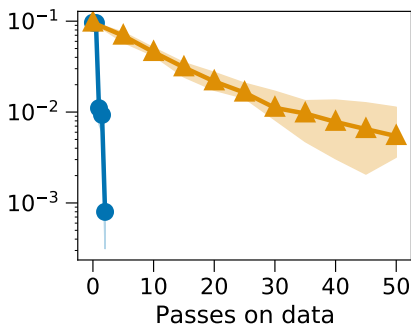
$$\epsilon = 1, \delta = 1/n^2.$$

Balanced Dataset

Balanced Dataset



ELECTRICITY (SCALED)
 $n = 45,312$ $p = 8$
Logistic Regression

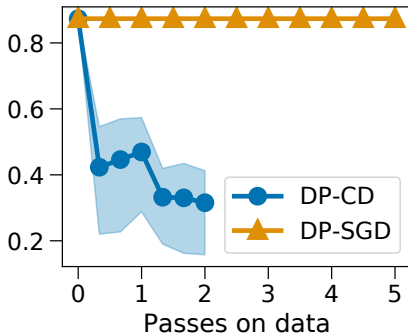


CALIFORNIA (SCALED)
 $n = 20,640$ $p = 8$
LASSO

$$\epsilon = 1, \delta = 1/n^2.$$

Higher dimension

Higher dimension



SPARSE SOLUTION
 $n = 1,000$ $p = 1,000$
LASSO

$$\epsilon = 10, \delta = 1/n^2 \text{ and } \|w^*\|_0 = 10. \quad 32 / 40$$

Partial Conclusion.

- Partial gradients without amplification.
- Large learning rates.
- Good practical performance.

Is DP-CD optimal?

Is DP-CD optimal? Kind of.

Loss	Convex L -comp.-Lipschitz	Strongly-Convex L -comp.-Lipschitz
Utility ($\mathbb{E}[F(w) - F^*] \geq \dots$)	$\Omega\left(\frac{L_{\min}}{L_{\max}} \frac{\ L\ _2 \ w^*\ _2 \sqrt{p}}{n\epsilon}\right)$	$\Omega\left(\frac{L_{\min}^2}{L_{\max}^2} \frac{\ L\ _2^2 p}{\mu_I n^2 \epsilon^2}\right)$

Is DP-CD optimal? Yes, if...

Loss	Convex L -comp.-Lipschitz	Strongly-Convex L -comp.-Lipschitz
Utility ($\mathbb{E}[F(w) - F^*] \geq \dots$)	$\Omega\left(\frac{\ L\ _2 \ w^*\ _2 \sqrt{p}}{n\epsilon}\right)$	$\Omega\left(\frac{\ L\ _2^2 p}{\mu_I n^2 \epsilon^2}\right)$

If $\sum_{j \in \mathcal{S}} L_j^2 = \Omega(\|L\|_2^2)$ for \mathcal{S} with $\text{card } \mathcal{S} \geq \frac{p}{75}$.

Perspectives!

Choose j wisely. (Greedy?)

- Match lower bounds?
- Reduce dependence on p .



- J. Nutini et al., “*Coordinate Descent Converges Faster with the Gauss-Southwell Rule Than Random Selection*”, 2015.
- S. P. Karimireddy et al., “*Efficient Greedy Coordinate Descent for Composite Problems*”, 2019.
- H. Fang et al., “*Greedy Meets Sparsity: Understanding and Improving Greedy Coordinate Descent for Sparse Optimization*”, 2020.

Sparse Approximation?

- After T iterations: $\|w^T\|_0 \leq T$.
- Active sets?



- M. Massias, A. Gramfort, and J. Salmon, “*Celer: A Fast Solver for the Lasso with Dual Extrapolation*”, 2018.
- K. L. Clarkson, “*Coresets, Sparse Greedy Approximation, and the Frank-Wolfe Algorithm*”, 2010.

DP-CD as a subroutine.

e.g., in Iteratively Reweighted Least Squares

$$w^{t+1} = \arg \min_{w \in \mathbb{R}^p} \sum_{i=1}^n \alpha_i^t |x_i \beta - y_i|^2.$$



- P. W. Holland and R. E. Welsch, “*Robust Regression Using Iteratively Reweighted Least-Squares*”, 1977.
- E. J. Candès, M. B. Wakin, and S. P. Boyd, “*Enhancing Sparsity by Reweighted L1 Minimization*”, 2008.

Coordinate-Wise Clipping.

- Even in DP-SGD!
- Could improve Fairness?



- V. Pichapati et al., “*AdaCliP: Adaptive Clipping for Private SGD*”, 2019.
- D. Xu, W. Du, and X. Wu, “*Removing Disparate Impact of Differentially Private Stochastic Gradient Descent on Model Accuracy*”, 2020.

Thank you! Questions? :)

See our paper:

- ☰ – P. Mangold et al., “*Differentially Private Coordinate Descent for Composite Empirical Risk Minimization*”, 2021.