How to train Private models?

there are 3 broad approaches for solving ERM with prinary guarantee:

1. Output Perturbation: Suppose we have an algorithm that approximates O^+ .

D, & ___ A10. __ O*

To make of private, we can add nother

to it
0*+ (n'.__, nd)

Sensitivity of 0x might be large! Potential drawback: * In SVM, O* is close to a data point a support vector " so changing that datapoint might drastically change 0* = 0 high fensitivity. Appropriate regularization 2 lio, (2,17,1) + 2 Flo) Potential solution: Resulting 0* might have

Resulting 0* might have regularizer
reasonable sensitivity (under some restrictive assumption)

Instead of minimiting 2 (10, 121,7) 2- Objective Perturbation:

we minimite

Canssian
Vandom
Vandom
Vector.

(*) It was shown that the optimizer for this objective can be

thought of on the private version of the optimizer of original ERM. All early works on objective perturbation required the Drawback:

optimization problem to find exact minimum for the objective function (*). This turns out to be very unrealistic ors most popular optimizers one iterative; that is they only approximate the optimum value.

there are <u>Several</u> recent works that relax this requirement, showing objective ferturbention can be quite practical.

3. Gradient Perturbation: This is by far the most Popular approach which is Partially inspired by the fact that gradient-based iterative algorithms are ubiquitous in AI.

This approach is best described in Gradient Descent Alg.

Non-Private GD: := L(O, (X,...))

arg min \(\frac{2}{2} \left(\text{Q}, (X_1,...)_1 \)

\[\text{where } \text{C} \cappa \text{R} \]

is a fet where \(\text{Q} \)

O should live in.

Input: Dataset D={(zinyi)} noss function e, parameter set e

1. Pick $\theta \in \mathbb{C}$ arbitrarily projection operator

2. For t=1 to TSet $\theta_t = T$ $\left(\begin{array}{c} 0 & -1 \\ t-1 & 1 \end{array}\right)$

Output 0,,02,-- OT

It is widely known that moving along the negative of gradient gets as closer to the minimum. Why? You can formalize this observation using Taylor expansion. To assess whitey of this algorithm, we use this folklore theorem

theorem. (Shamire Zhang 2013, "SGD for non-smooth optimization")

Let F be a convex function & let 0* argmin F(0).

Let 9, be an arbitrary point in 8 &

Q= TP (Q-1+G(Q)), where E[G(O)] = PF(O) &

E[$||G_10_1||_2^2$] $||G_2||_2$ the learning rate $||G_1||_2$.

Then for any T70, we have

we can immediately apply this theorem, with $G(Q_{\ell}) = DL(Q_{\ell}, D)$.

To ensure that $E[||G||^2] \le G^2$, we need to be deterministic.

Assumption: loss function (0, vary) is L. Lipschitz

for all (2,13): that is

[(0,1219)] - - (10', (219))] < L. 110-011

P is L-hp =0 1179115 L Note: E[|| G(O) ||2) = E[|| PP (Q,D)||] < (n L)2 thus, Became PL is nL Lip. Thus, for the GD algorithm, me have: 1(QD)-1(0*,D) 5 1(C1) n.L

If we run this algorithm sufficiently large, then £10,0) - £10,0) Something like

$$T=\Omega(n^2)$$

Draw back:

At each iteration, we need to take n gradient, so in total, nT gradients ($\approx v^3$ gradients).

* Is there anyway to reduce computation at each iteration?

perhaps only one gradient at each iteration? Yes!

Stochastic GD (SGD)

Projected	L SGD
Input:	Dataset D={(zi,yi)}; loss function l, parameter set 6
1- Pick	0, E arbitrarily
2- For	t=1 to T
	- Select IE{1,2,, n} uniformly at random
	- Set $\theta_t = \int_{\mathcal{O}} \left(\frac{\partial}{\partial t} - \int_{\mathcal{O}} \mathbf{n} \nabla \mathbf{l} \left(\frac{\partial}{\partial t} , (\mathbf{x}, \mathbf{y}) \right) \right)$
output	0,,021~- OT

claim: $\nabla L(\theta, (x, y_s))$ is a good estimate of $\nabla L(\theta, D)$.

Note that

$$E[\nabla P(0, (x_1, y_2))] = \frac{1}{n} \sum_{i=1}^{n} \nabla P(0, (x_1, y_2)) = \frac{1}{n} \nabla R(0, D)$$

 $\Rightarrow \quad \text{NETPLO,} \ (a_{i}, y_{i})) = \quad \text{DL(0,0)}$

what this means is that $\nabla e(Q, (X, Y))$ is an unbiased estimate of $\nabla f(Q, D)$.

the computation is way less in SGD than in GD:

at each iteration, we only need the gradient, so in total we need only n^2 gradient, down from n^3 in GD.

What is the catch? Convergence is probabilistic!

Using the previous general purpose result (shamir & Zhang 2013):

with
$$G(0) := n \ \nabla C(0, (x, y,)) \rightarrow C[G(0)] = \nabla C(0, D)$$

& $E[||G(0)||^2] < n^2 L^2$ again we assume loss is L-lip

Private Projected SGD (PN-SGD)

9.(D) - [0-70e] TTE 0,-70e - TTE - 0-70e-(ایلایال (هایال

To make projected SGD, differentially private, we need to make each iteration private. To do so, we need to part the guery response through a by mechanism, say Gaussian mechanism.

DR(O, ULY))+Noise

Gaussian

thus, each iteration proceeds as follows:

Private Projected GD Input: Dataset D={(zinyi)}, loss function e, parameter set C 1- Pick O.E C arbitrarily 2- For t=1 to T - Select IEgliz-, n3 uniformly at random $= \frac{1}{8} \frac{$

Output 0,,02, -- OT

At each iteration, grey is Pelo, (x, y,) which is an adaptive query as it depends on the previous iterations output.

Each iteration becomes (ξ, ξ) - Dp with $\xi = \frac{2nL}{5}\sqrt{2\log\frac{1}{\xi}}$ Note that: $\Delta_2^q < 2nL$ for the query $\nabla \ell(Q, \chi_1, \chi_2)$.

guarantee. But, we can do better! At each iteration, we don't use the whole dataset! we only use One data record. 8 not all the records. Thus, the privacy must be significantly better!

Therefore, advanced composition can be used to obtain the privacy

How to quantify this improvement. * Privacy amphification by Sub-Samphing: Let M be an (5,8)-Dp mechanism. If it is run on a uniformly selected subset of dataset of size on (rsi), then it will Provide (log(1+r(e-1)), 88) -DP. 9(D) - M - Z Zo provides more privacy m depends compared to the case only the first part when Q(D) depends on datalet

the whole dataset.

In the SGD: $\gamma = n$ (as the size of dataset =1)

Gaussian mechanism coupled with this sub-sampling is typically called: Sub-sampled Gaussian mechanism.