Theorem (Advanced composition). Let M be ΣDP . The k-fold composition of M is $(\Sigma_1 S) - DP$ for any $S \in (D_{11}) \ S = K \Sigma_0^2 + \Sigma \sqrt{2K \log V_S}$

Note that this theorem implies that composition of pure DP mechanism satisfies approximate DP; a phenomenon that Can't be obtained from basic composition.

proup.

Comparison with basic composition:

Suppose we want to answer k queries of sensitivity one, with a given

privay guarantee (5.8)-DP.

Basic composition: K Laplace mechanism each of which $\frac{\varepsilon}{\kappa}$ -Op so each adds N~ L(0, 4,5).

Advanced composition: K aplace mechanism each of which

is
$$\left(\frac{2}{\sqrt{2\kappa \log V_{\delta}}}\right)$$
 Dp so each adds $N \sim L(0, \sqrt{\frac{2\kappa \log V_{\delta}}{S^{2}}})$

*why? * The required variance is OWK) smaller! Let M be E Dp.

According to advanced composition, DP pavameter of K fold composition of M is

 $\frac{K}{2} \cdot \left(\frac{S}{\sqrt{2}}\right)^{2} + \left(\frac{S}{\sqrt{2}}\right)^{2}$

In the advanced composition: $\frac{\Sigma^2}{4 \log 16} \approx \frac{\Sigma^2}{60}$

If $K \Sigma^2$ is small, then we can ignore the first term.

$$2 = 2 \sqrt{2 k \log V_S}$$

the previous advanced composition shows that K-fold composition of an S-DP becomes an (5.6) DP mechanism with 2 increasing Sub-linearly in K.

How about K-fold composition of (5,8) - Dp mechanism?

Here is a more general advanced compositions

Theorem (Advanced composition). Let M be
$$(5,8)$$
-DP. The k-fold composition of M is $(5,8)$ -DP where
$$S = \frac{K5^2}{2} + \frac{5}{2} \sqrt{2K \log V_S}$$

$$6 = KS + 6$$
 For any $6 \in (0,1)$.

mechanism:

an (E, S.) - Dp Gaussian mechanism is (5) 2k log/s, KS+8) Dp

$$8 = \frac{5}{2}$$

$$2 \times \log \left(\frac{1}{2} \right) = \frac{5}{2}$$

As you can see, & is a design parameter. We need to pick & in a way that the resulting noise variance is minimal. Another option for 8 in the above example is

$$S'=S$$
.

This choice results in Variance $\frac{4K}{52}$ ($\log \frac{K+1}{5}$) which is

strictly worse than what we had before.

Renyi DP & composition results

Recall that differential privacy was intuitively defined by ensuring M > M, are close, according to HS divergence.

why not other distance measures?

we already know that TV doesn't work.

It turns out that Rényi divergence is a good candidate.

Definition. We say that a randomized mechanism M is

(Mironov 2017) (a, §) - Renyi DP (or (d, §) - RDP) for a>1 & §20

if Da (M/IM) < §

The main motivation [which made RDP extremely ropular in AI]

is the following result:

Theorem. For any $\alpha > 1$, we have D(P|Q) + min D(P|Q) < D(P|Q) + max D(P|Q) + max

In Particular, if Pxy = Px. Py & Qxy = Qx. Qy, then:

D (Pyllax) = D (Pyllax) + D (Pyllax).

Proof. See Lemma 2.2 of "concentrated Pifferential Privacy" by
Bun & Steinke.

this simple result leads to the following simple composition result that resembles basic composition.

Theorem. Let M' be (a,5,)_RDP & M2 be (a,5)_RDP. Then
their composition is (a,5,+52)_RDP.
adaptive

Proof. Simple application the previous theorem.

According to this result, K-fold composition of any (9,5)-RDP is (4, K\$)-RDP.

Limitation of RDP. (4,5)-RDP guarantee doesn't enjoy

any operational interpretation, in terms of hypothesis testing. In fact, it was shown that Da with as 1 has nothing to do with binay hypothesis testing performance.

* Read "Hypothesis testing interpretation & RDP" by Balle et al., AISTAT

How to fix it? RDP Privacy gnarantee needs to be converted back to approximate DP.

Theorem. If M is $(a_1 \xi) - RDP$ with a>1, then it is $(\epsilon_1 \xi) - DP$ Mironov with $\delta \in (011)$ & $\epsilon = \xi + \frac{1}{a-1} \log \frac{1}{2} \epsilon$.

Moments accountant: Was the state-of-the-art technique for studying composition of iterative algorithms used for training A1 models (Such as SGD).

We discuss it further later.

Statistical learning

Let D= {(x,y,),..., (x,yn)} be a dataset of features 2. 2 labels Di. For instance feature might be images, texts or medical records of individuals & labels are the type of animal in the image (cat vr. dog), or the diagnostic. the main goal in machine learning is to discover the relationship between features & labels. This allows us to come up with a classifier. Think of classifier as a conditional distribution Pyix. Any such classifiers are modelled by parameters OEIRd.

To find the best classifier that fits a dataset, we must find 10 GIRd that describes the relationship between is 2 y. For all i G1/2-7n ?.

To mathematically formulate this goal, we take a froblem-specific 1055 function 1:

L(0, (x,y)) quantifies the loss associated

with representing the relationship

between feature x x label. y

between feature x 2 label y using a classier parameterized by 0.

Example: 1- Linear

1- linear regression: Think of linear regression as a continuous

regressor is parameter, Zed by OEIRd

regressor is parameter, Zed by OEIRdas < 0,2>
inner product
loss Function: (LO, 12181) = [y - <0,25]

* Goal is to find 0 that

minimizer 10ss for each point: Zelo, (dig.)

2_ Logistic regression: Here, label 3 \(\(\frac{1}{2} \), +13 binary classification

classifier in farametrized by
$$O \in \mathbb{R}^d$$
 as
$$P_{Y|X=2L} = \frac{1}{1+e^{-\langle O_i X \rangle}}.$$
the goal is to find O such that

1 is maximited when
$$J_i = +1$$

1 $e^{-\langle Q_i \chi_i \rangle}$ is maximited when $J_i = +1$

for all data points in the dataset.

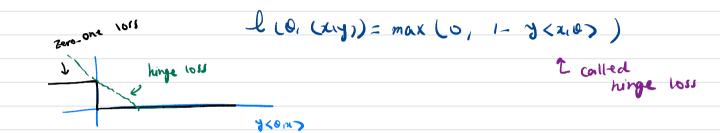
thus, the loss function is defined to be lo, (x,y) = log (1+ e- + < 0,x)

3_ Support Vector machine (SVM): Binary label 3: Ex-1,+13 Prixax (+1)= 1 if <0,x7 >0 4 deterministic classifier, The good is to find OER & bEIR such that 1: < 0,x17 > 0 so any reasonable loss increases with -4 LO(2). One potential (011: - PLO, (xxy)) = 122<0,x><0}

Called Zero-one loss

Since this function is not differentiable

& not convex, we usually consider this loss:



4. Geometric median: Unsupervited learning, so label is not

given.

We are interested in median:

thus, the minimizer of is the geometric median of all feature. In all these examples, we are interested in the following Optimi Zation problem: $\Theta^*:=\underset{\Theta\in\mathbb{R}^{d}}{\operatorname{arg min}}\sum_{i=1}^{\infty}\mathbb{L}\left(O,\left(x_{i},y_{i}\right)\right)$ ERM (Empirical risk minimite tion) Revealing 0* might compromise Privacy ! - For instance, if loc(xry) = 10-x112, then if the number of datagoints is odd, then ox is equal

of \(\frac{1}{2} \cdot \(\text{10} \(\text{12} \(\text{13} \) \)

to one of the point.

_ In NM, the optimal 0 * 11 close to one data point

_ Model inversion attack

