# Announcements

1. Let's all take a break!

   On 4/21 I will cover no new material. we can just chat about the material so far, or your projects etc

2. No office hours on 4/20 though

3. Many of you have talked to me about your project in office hours or made an appointment to chat some other time. Keep doing that!

# Lower Bounds for SGD

Today we will delve more into SQs and lower bounds

First, consider an idealized version of SGD for fitting the parameters of a model

## Noisy SGD

For $t=1$ to $T$

$$\text{Set } \theta_{t+1} = \theta_t - \eta \nabla_\theta R(\theta) + z \overset{N(0, \sigma^2 I)}{\leftarrow}$$

where $R_\theta = \underset{(x,y)}{\mathbb{E}}\left[(h_\theta(x) - y)^2\right]$

<u>Proposition</u>: when the marginal distribution on X is known, Noisy GD can be approximated by a CSQ algorithm

Recall that a CSQ is of the form

$$\mathbb{E}[y \, q(x)] \pm \tau$$

To formalize what we mean by "approximated" we need the notion of a coupling:

def: A coupling of two r.v.s $X$ and $Y$ with distributions $\mu_X$ and $\mu_Y$ respectively is a joint distribution $\mu_{X,Y}$ on $(X,Y)$ pairs

    ① the marginal on $X$ is $\mu_X$

    ② the marginal on $Y$ is $\mu_Y$

Additionally, when $X$ and $Y$ have the same domain, we call

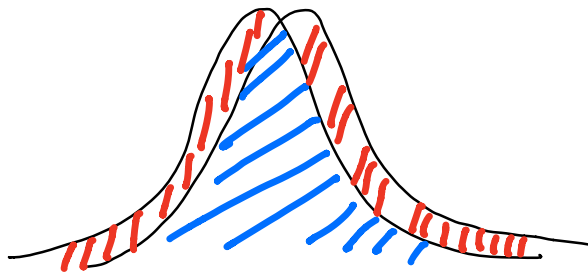$$\mathbb{P}_{\mu_{X,Y}}[X = Y]$$

the <u>coupling probability</u>

We will show:

"If we replace the gradient
computation with CSQs, we
can almost always couple each step"

In particular, we will be interested
in <u>optimal couplings</u> that maximize the
coupling probability

These are easy to visualize pictorially:

<u>Fact 1</u>: The optimal coupling probability between two r.v.s. with p.d.f.s. $f$ and $g$ is

$$\int \min(f(x), g(x)) \, dx$$

$$= 1 - \frac{1}{2} \underbrace{\int |f(x) - g(x)| \, dx}_{\text{total variation distance}}$$

<u>Fact 2</u>: For two spherical Gaussians

$$d_{TV}\left(\mathcal{N}(\mu, \sigma^2 I), \mathcal{N}(\mu', \sigma^2 I)\right) \leq \frac{\|\mu - \mu'\|_2}{2\sigma^2}$$

Now let's return to the proposition, and see how to set $\tau$ along the way

<u>Proof</u>: Let's compute the gradient:

$$\nabla_\theta \mathbb{E}\left[(h_\theta(x) - y)^2\right]$$

$$= \mathbb{E}\left[2(h_\theta(x) - y) \nabla_\theta h_\theta(x)\right]$$

$$= 2\,\mathbb{E}\left[h_\theta(x)\nabla_\theta h_\theta(x)\right] - 2\,\mathbb{E}\left[y\,\nabla_\theta h_\theta(x)\right]$$

$$\underbrace{\qquad\qquad}_{(I)} \qquad\qquad \underbrace{\qquad\qquad}_{(II)}$$

Observe that (I) can be computed just from knowledge of the marginal on X — i.e. without CSQs

Now (II) is a p-dimensional vector and each coordinate can be estimated up to $\pm\tau$ with a CSQ

We can usually couple:

gradient + noise $\Longleftrightarrow$ vector of CSQs + noise

and using Facts 1+2 the failure probability is at most

$$O\left(\frac{\sqrt{p}\,\tau}{\sigma^2}\right)$$

because this is the TV distance of two p-dimensional Gaussians with covariance $\sigma^2 I$ and whose means are at most $\sqrt{p}\,\tau$ in Euclidean distance.

Finally if we choose

$$\tau \sim \frac{\delta \sigma^2}{\sqrt{p}\,T}$$

with probability $\geq 1-\delta$ we will be able to couple at every step $\implies$ we reach the same model parameters. $\boxtimes$

Thus we have our first lower bound for deep learning:

<u>Corollary</u>: Noisy GD on a polynomial sized deep network must take either exponentially many step or set the noise to be exponentially small to learn parities w.r.t. uniform distribution

These and other lower bounds were shown by Feldman, Guzman, Vempala and Shalev-Schwartz, Shamir, Shamman

What if you don't add your own stochastic noise?

Not only do these proofs break down

Theorem [Abbe, Sandon] SGD is P-complete

i.e. by choosing the appropriate initialization, you can get it to implement any polynomial time algorithm

Proving unconditional lower bounds against SGD $\Rightarrow$ P $\neq$ NP

# Learning Shallow Deep Nets

we just showed a lower bound against learning parities

Main Question: Are there algorithms for learning shallow deep nets under simple distributions?

We will work with deep nets of the following form

$$f(x) = \sum_{i=1}^{m} a_i \sigma(w_i^T x) \quad (*)$$

where $\sigma$ is a sigmoid & $x$ is d-dimensional

we will show superpolynomial lower bounds, following Goel, Gollakota, Jin, Karmalkar, Klivans and Diakonikolas, Kane, Kontonis, Zarifis ← stronger exponential lower bounds

<u>Theorem 1</u>: Any <u>CSQ</u> algorithm for learning nets of the form (*) under the Gaussian distribution must make at least $d^{\Omega(\log m)}$ queries or set $\tau = d^{-\Omega(\log m)}$

Recall the intuition for parities was

$$\begin{array}{c} \text{large} \\ \text{orthogonal} \\ \text{family} \end{array} \implies \begin{array}{c} \text{CSQ} \\ \text{lower bounds} \end{array}$$

And we asserted that for Boolean-valued functions where you know the distribution, you can assume all queries are CSQs wlog

<u>Key Question(s)</u>: But what about for real-valued labels? Is SQ >> CSQ?

For real-valued labels, Szorenyi showed

Theorem 2 (informal) If $\mathcal{H}$ contains a set of $\ell$ functions that are pairwise orth. under $D$ and have squared norm $\Omega(1) \Rightarrow$ any CSQ algorithm must make either $\Omega(\ell)$ queries or have tolerance $\tau = \ell^{-\Omega(1)}$ to get accuracy $\varepsilon > 0$

Main Question: How do we contruct a large orthogonal family of functions of the form (*)?

def: For any set $S \subseteq [d]$ with $|S| = \log m$, let

$$g_S(x) = \sum_{w \in \{\pm 1\}^{\log m}} \chi(w) \, \sigma\left(\frac{w^T x_s}{\sqrt{\log m}}\right)$$

where $\chi(w) = \prod_{i=1}^{\log m} w_i$

coordinates of $x$ restr. to $S$

Let's show these functions are orthogonal:

Lemma 1: For any $S \neq T$, $\mathbb{E}_{X \sim N(0,I)}[g_S(x) g_T(x)] = 0$

The key fact we need is

Fact 3: For any $z \in \{\pm 1\}^d$ and $x$

$$g_S(x \underset{\uparrow}{\circ} z) = X_S(z) g_S(x)$$
pointwise product

Proof: By definition, we have

$$g_S(x \circ z) = \sum_{w \in \{\pm 1\}^{\log m}} X(w) \sigma\left(\frac{w^T (x \circ z)_S}{\sqrt{\log m}}\right)$$

$$= \sum_{w \in \{\pm 1\}^{\log m}} X(w) \sigma\left(\frac{(w \circ z_S)^T x_S}{\sqrt{\log m}}\right)$$

$$= \sum_{w \in \{\pm 1\}^{\log m}} X(w \circ z_S) X(z_S) \sigma\left(\frac{(w \circ z_S)^T x_S}{\sqrt{\log m}}\right)$$

Now let $w' \triangleq w \circ z_S$. Then

$$= \chi(z_S) \sum_{w' \in \{\pm 1\}^{\log m}} \chi(w') \, \sigma\left(\frac{w' \circ z_S}{\sqrt{\log m}}\right)$$

$$= \chi(z_S) \, g_S(x) \qquad \blacksquare$$

Now returning to the lemma

Proof of Lemma 1: Because $\mathcal{N}(0, I)$ is sign symmetric, we have

$$\mathop{\mathbb{E}}_{x \sim \mathcal{N}(0,I)}\left[g_S(x) g_T(x)\right] = \mathop{\mathbb{E}}_{z \sim_u \{\pm 1\}^d}\left[\mathop{\mathbb{E}}_{x \sim \mathcal{N}(0,I)}\left[g_S(x \circ z) g_T(x \circ z)\right]\right]$$

And from Fact 3, we have

$$= \mathop{\mathbb{E}}_{z \sim_u \{\pm 1\}^d}\left[\mathop{\mathbb{E}}_{x \sim \mathcal{N}(0,I)}\left[\chi_S(z) \chi_T(z) \, g_S(x) g_T(x)\right]\right]$$

$$= \mathop{\mathbb{E}}_{z \sim_u \{\pm 1\}^d}\left[\chi_S(z) \chi_T(z) \underbrace{\mathop{\mathbb{E}}_{x \sim \mathcal{N}(0,I)}\left[g_S(x) g_T(x)\right]}_{0}\right] \qquad \blacksquare$$

what remains is to show that each
of the functions in the family is
non-trivial

Lemma 2: For any $S$ with $|S| \leq \log m$,
we have
$$\mathbb{E}_{x \sim N(0, I)} \left[ g_S(x)^2 \right] \geq c^{-m}$$

We will omit the proof. But we can now
prove Theorem 1

Proof of Theorem 1: There are $\ell = d^{\Omega(\log m)}$
functions in the orthogonal family

The key point is if we want to learn
within error $\varepsilon \stackrel{\Delta}{=} c^{-m}$ then Lemmas 1+2
imply must identify the correct $g_S$

Now invoking Theorem 2 completes the proof. ☑

Note: There are stronger lower bounds in the probabilistic concept model where

$$x \sim N(0, I) \text{ and } y \in \{\pm 1\}$$
$$\text{with } \mathbb{E}[y|x] = f(x)$$

In this case we get SQ, rather than CSQ, lower bounds

Landscape Design

Next we will study upper bounds for learning
$$y = f(x) = \sum_{i=1}^{m} a_i \sigma(w_i^T x)$$

from Gaussian inputs

Ge, Lee, Ma studied ERM

$$L(\hat{a}, \hat{w}) \triangleq \mathbb{E}\left[\|y - \hat{y}\|^2\right] \quad (*)$$

where $\hat{y} = \hat{a}^T \sigma(\hat{w}^T x)$

Theorem: Assuming the $w_i$'s and $\hat{w}_i$'s are unit vectors, the population risk defined in $(*)$ satisfies

$$L(\hat{a}, \hat{w}) = \sum_{k=1}^{\infty} \hat{\sigma}_k^2 \underbrace{\left\| \sum_{i=1}^{m} \hat{a}_i \hat{w}_i^{\otimes k} - \sum_{i=1}^{m} a_i w_i^{\otimes k} \right\|_F^2}_{f_k} + C$$
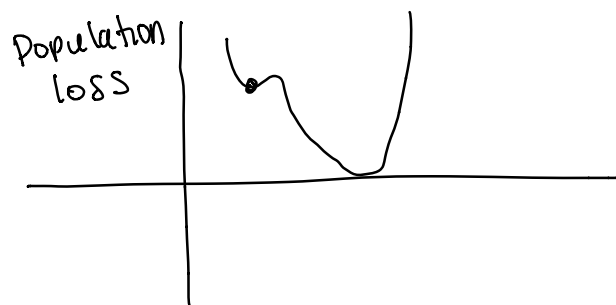
where the inputs come from $N(0, I)$ and the $\hat{\sigma}_k$'s are coefficients of $\sigma$ in the Hermite basis

In particular, for $\sigma \equiv$ ReLU we have

$$\hat{\sigma}_1 = \frac{1}{2} \quad \text{and for } n \geq 2$$

$$\hat{\sigma}_n = \begin{cases} \dfrac{((n-3)!!)^2}{\sqrt{2\pi n}} & n \text{ even} \\ \\ 0 & n \text{ odd} \end{cases}$$

However, experimentally, SGD fails in the sense that it has spurrious local minima — i.e. it gets stuck

Population loss

They gave a fix that they called landscape design

Main Question Can we design a new loss function G that has

① no spurrious local minima

②  every global minimum of $G$
$\Uparrow$
ground truth parameters

③  can estimate the gradient of
$G$ from samples

Concretel, consider
$$\hat{y}' = \hat{a}^T \gamma (\hat{w}^T x)$$

where $\gamma = \hat{\sigma}_2 h_2 + \hat{\sigma}_4 h_4$
         $\uparrow$           $\uparrow$
(normalized prob.'s)  Hermite polynomials

In particular $h_2(t) = \frac{1}{\sqrt{2}} (t^2 - 1)$ and

$$h_4(t) = \frac{1}{\sqrt{24}} (t^4 - 6t^2 + 3t)$$

Now consider

$$G(\hat{a}, \hat{w}) = \mathbb{E}[\|\hat{y}' - y\|^2] \quad (\Delta)$$

Theorem: Under the same assumptions as before, the population risk in $(\Delta)$ satisfies

$$G(\hat{a}, \hat{w}) = \hat{\sigma}_2^2 f_2 + \hat{\sigma}_4^2 f_4 + C$$

Moreover if the $a_i$'s are nonnegative and $W$ is orthogonal then $G(\hat{a}, \hat{w})$ satisfies conditions ①, ②, ③ of landscape design

Diakonikolas, Kane, Kontonis, Zarifis gave PCA-based algorithms that just need $a_i$'s to be nonnegative and Gaussian inputs

**Key Question:** Is nonnegativity of the $a_i$'s actually needed?

The CSQ lower bounds critically need both positive and negative $a_i$'s to get an orthogonal family

## Fixed Parameter Tractability

Chen, Klivans, Meka gave a

$$f(m, \Delta, \varepsilon) \, \text{poly}(d)$$

<div align="center">Lipschitz constant of deep net</div>

time algorithm for learning deep nets with $m$ hidden units from Gaussian inputs

**Key Point:** This is much better than the $d^{\Omega(m)}$ CSQ lower bound

# Filtered PCA

Consider a matrix

$$M_\psi \triangleq \mathbb{E}\left[\psi(y)\left(xx^T - I\right)\right]$$

Let $V = \text{span}\{w_i\}$

$\uparrow$ weight vectors in the first layer

Observe that the output $y$ of the network only depends on $\Pi_V(x)$

__Claim 1__: For any direction $c \perp V$ we have $c \in \ker(M_\psi)$

The main idea is to show for a suitable choice of $\tau$ we have

for any vector $c$ where $|F(c)| \geq \tau$

$\Downarrow$

$(\square)$

$$\|\Pi_V(c)\|^2 \geq 2m$$

Then we could pick $\psi(y) \overset{\Delta}{=} 1_{|y| \geq \tau}$

Now let's run PCA on $M\psi$ to find a direction in $V$:

Claim 2 Under $(\square)$, $M\psi$ has a positive eigenvalue

Note by Claim 1, we know this direction must be in $V$

Proof: We will compute the quantity

"sum of the eigenvalues of $M\psi$"

This is just the trace, and by Claim 1 we can write this as

$$\text{Tr}\left(V^T M_\psi V\right) \qquad (0)$$

where we have abused notation and $V$ is a $d \times m$ matrix whose columns are an orthonormal basis for the subspace $V$

then we have

$$(0) = \text{Tr}\left(V V^T M_\psi\right) = \langle \Pi_V, M_\psi \rangle$$

where $\langle A, B \rangle$ is the matrix inner-product

$$\sum_{ij} A_{ij} B_{ij}$$

Now plugging in the definition for $M_\psi$, we have

$$(0) = \mathbb{E}\left[ 1_{|y| \geq \tau}\left(\langle \Pi_V, xx^T \rangle - \underbrace{\langle \Pi_V, I \rangle}_{m}\right)\right]$$

$$= \mathbb{E}\left[ 1_{|y| \geq \tau}\left(\|\Pi_V x\|^2 - m\right)\right]$$

Now using (□) we get

$$\geq \mathbb{P}\left[|y| \geq \tau\right] m$$

Thus the top eigenvector of $M\psi$ will be in the subspace $V$ ▨

To put it in words

"If we condition on the response being sufficiently large, the resulting distribution is non-Gaussian"

and we can find such a direction by using PCA

At this point (□) seems optimistic; let's do an illustrative example:

(Assume wlog $\mathbb{E}[f(x)] = 0$)

① Let $\sigma(x) = x^2$, quadratic activation)

② Let $V = \text{span}\{e_1, \ldots, e_k\}$

Now let $\psi(y) = y^2$. We are interested in

$$\mathbb{E}\left[y^2(x_1^2 + \ldots + x_m^2 - m)\right]$$

which we can break up into $m$ terms of the form

$$\mathbb{E}\left[y^2(x_1^2 - 1)\right]$$

Now the key is to use Gaussianity

<u>Lemma</u> [Stein] For any differentiable function $g$

$$\mathbb{E}_{x \sim N(0,1)}\left[g(x)x\right] = \mathbb{E}_{x \sim N(0,1)}\left[g'(x)\right]$$

**Proof:** Using integration by parts

$$\mathbb{E}[g'(x)] = \int_{-\infty}^{\infty} g'(x)\left(\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}\right) dx$$

$$= \underbrace{\frac{g(x) e^{-\frac{x^2}{2}}}{\sqrt{2\pi}}\Bigg|_{-\infty}^{\infty}}_{0} + \underbrace{\int_{-\infty}^{\infty} x g(x)\left(\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}\right) dx}_{\mathbb{E}[x g(x)]}$$

$$\boxed{}$$

**Corollary 2:** $\mathbb{E}[g(x)(x^2-1)] = \mathbb{E}[g''(x)]$

**Proof:** we have

$$\mathbb{E}[g(x)(x^2-1)] = \mathbb{E}[g(x)x^2] - \mathbb{E}[g(x)]$$

$$\overset{\text{stein's lemma}}{=} \mathbb{E}\left[\frac{d}{dx}(g(x)x)\right] - \mathbb{E}[g(x)]$$

$$= \mathbb{E}[g'(x)x] \overset{\text{stein's lemma}}{=} \mathbb{E}[g''(x)]$$

$$\boxed{}$$

Now applying Corollary 2

$$\mathbb{E}\left[ f(x)^2 (X_1^2 - 1) \right] = \mathbb{E}\left[ \frac{d^2}{dX_1} f(x)^2 \right]$$

$$= 2 \underbrace{\mathbb{E}\left[ \left( \frac{d}{dX_1} f(x) \right)^2 \right]}_{>0} + 2 \underbrace{\mathbb{E}\left[ f(x) \frac{d^2}{dX_1^2} f(x) \right]}_{0}$$

which again implies the top eigenvector of $M_\varphi$ is in $V$