

Lecture 11:

High-dimensional non-convex dynamics

SGD on NNs : if we are free to set architecture and initialization, then we can emulate any polynomial time algorithm.

↳ this is very different than what we do in practice

↳ there is some "hyperparameter tuning" but the space of search is relatively constrained : we only use some standard architectures and "generic" initializations
 (fully connected / CNNs / transformers) (iid initializat°)
 e.g. Gaussian

where we can tune
 the variance

To understand what SGD does on regular NNs, one need to study the dynamics

↳ high-dimensional (many parameters)

+ non-convex

↳ this is hard!

(2)

While it is generically intractable to study these dynamics, there has been a lot of interesting work in the past few years.

Today: we will focus on a very simple setting that presents nonetheless a very interesting phenomenology + is highly instructive.

"Online SGD on non-convex losses from high-dim inference"
Ben Arous, Gheissari, Jagannath, 2021.

Setting: $\alpha \sim N(0, I_d)$

$$y = f(\langle \omega_*, \alpha \rangle) \quad \text{for some } \|\omega_*\|_2 = 1$$

We learn this data using a "single neuron"

$\hat{f}(n) = f(\langle \omega, \alpha \rangle)$ using SGD from a random initialization

$$\omega^* \sim \text{Unif}(\mathbb{S}^{d-1})$$

Test error:

$$\underline{\underline{L}}(\omega) = \frac{1}{2} \mathbb{E} \left[(f(\langle \omega_*, \alpha \rangle) - f(\langle \omega, \alpha \rangle))^2 \right]$$

Hermite polynomials & the information exponent

$\{H_k\}_{k=0}^{\infty}$ the orthogonal basis of Hermite polynomials
 (basis of $L^2(\mathbb{R}, N(0,1))$)

H_k is a degree - k polynomial such that

$$\mathbb{E}_{G \sim N(0,1)} [H_k(G) H_j(G)] = k! \mathbb{1}_{[j=k]}$$

$$H_0(x) = 1 \quad H_1(x) = x \quad H_2(x) = x^2 - 1 \quad H_3(x) = x^3 - 3x$$

Properties: (i) $\frac{d}{dx} H_k(x) = k H_{k-1}(x)$

↙ integration by part

$$(ii) \mathbb{E}[H_k(G) g(G)] = \mathbb{E}[g^{(k)}(G)]$$

$$(iii) H_{k+1}(x) = x H_k(x) - k H_{k-1}(x)$$

Any function $f: \mathbb{R} \rightarrow \mathbb{R}$ s.t. $\mathbb{E}[f(G)^2] < \infty$ can be decomposed as

$$f(x) = \sum_{k=0}^{\infty} \frac{\mu_k(f)}{k!} H_k(x) \text{ s.t. } \mu_k(f) = \mathbb{E}[f(G) H_k(G)]$$

(if k times diff = $\mathbb{E}[f^{(k)}(G)]$)

Def: [Information exponent] We say that $f \in L^2(\mathbb{R}, N(0,1))$ has information exponent k_* if

$$k_* = \min \{ k \in \mathbb{N} : \mathbb{E}[f(G) H_{\mu_k}(G)] \neq 0 \}$$

Lemma: Let k_* be the info exponent of f . Then

$$L(\omega) = \frac{1}{2} \mathbb{E}[(f(\langle \omega_*, x \rangle) - f(\langle \omega, x \rangle))^2] = 1 - O(\langle \omega, \omega_* \rangle^k)$$

Proof: $\frac{1}{2} \|f(\langle \omega_*, \cdot \rangle) - f(\langle \omega, \cdot \rangle)\|_{L^2}^2 = 1 - \mathbb{E}[f(\langle \omega_*, x \rangle) f(\langle \omega, x \rangle)]$

$$= 1 - \sum_{k=k_*}^{\infty} \frac{\mu_k^2}{(k!)^2} \mathbb{E}[H_{\mu_k}(\langle \omega_*, x \rangle) H_{\mu_k}(\langle \omega, x \rangle)]$$

$$\mathbb{E}[H_{\mu_k}(\langle \omega_*, x \rangle) H_{\mu_k}(\langle \omega, x \rangle)] \quad G_1, G_2 \sim N(0,1)$$

$$= \mathbb{E}[H_{\mu_k}(G_1) H_{\mu_k}(\langle \omega, \omega_* \rangle G_1 + \sqrt{1 - \langle \omega, \omega_* \rangle^2} G_2)]$$

Int by part
on G_1
(property (ii))

$$= \mathbb{E}\left[\frac{d^k}{dG_1^k} H_{\mu_k}(\langle \omega, \omega_* \rangle G_1 + \sqrt{1 - \langle \omega, \omega_* \rangle^2} G_2) \right]$$

$$= \langle \omega, \omega_* \rangle^k \underbrace{\mathbb{E}[H_{\mu_k}^{(k)}(G)]}_{= k!} \quad (\text{property (i)})$$

(5)

Hence $L(\omega) = 1 - \sum_{k=k_*}^{\infty} \frac{\mu_k^2(f)}{k!} \langle \omega_*, \omega \rangle^k$ □

For simplicity, we will set $f(x) = \frac{\text{Ne}_k(x)}{\sqrt{k!}}$ from now on.

[results will hold for more general f of I.E. k]

Test error: $L(\omega) = 1 - \langle \omega, \omega_* \rangle^k = 1 - m^k$ $m = \langle \omega, \omega_* \rangle$

We minimize over $\omega \in \mathbb{S}^{d-1}$. We will consider "spherical gradient descent (gradient & tangent space of sphere manifold)"

$$\begin{aligned} \nabla_{\mathbb{S}^{d-1}} L(\omega) &= (\text{Id} - \omega\omega^\top) \nabla_\omega L(\omega) \quad \rightarrow \quad \text{Diagram showing a point on a sphere with tangent plane and gradient vector } \nabla_\omega L \\ &= (\text{Id} - \omega\omega^\top) (-k \langle \omega, \omega_* \rangle^{k-1} \omega_*) \\ &= -k m^{k-1} (\omega_* - m \omega) \end{aligned}$$

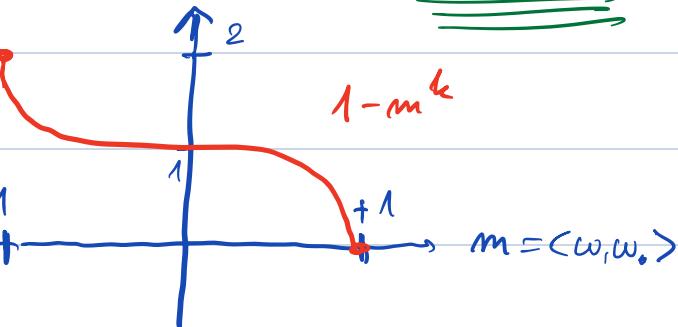
The gradient along the signal direction is

$$\langle \omega_*, \nabla_{\mathbb{S}^{d-1}} L(\omega) \rangle = -k m^{k-1} (1 - m^2)$$

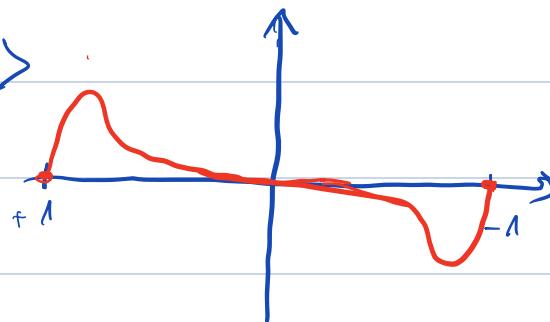
Population landscape

If k is odd

$$L(m) \underset{m}{\rightarrow}$$



$$\langle \omega_*, \nabla_S L(\omega) \rangle$$



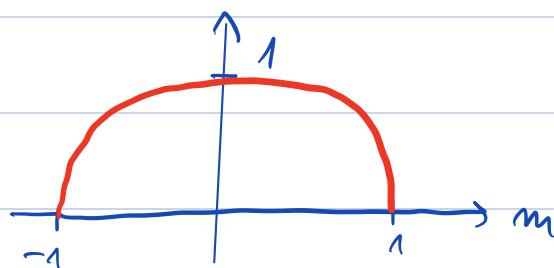
$\omega = \omega_*$ global minimum

$\langle \omega, \omega_* \rangle = 0$ stationary submanifold

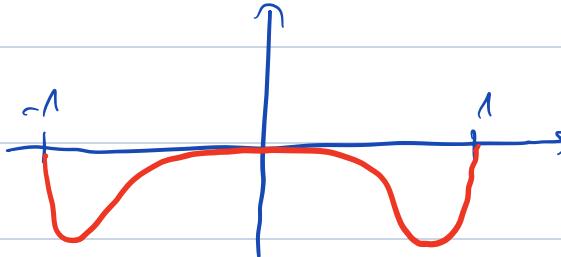
$\omega = -\omega_*$ unstable stationary point (global maxima)

If k is even

$$L(m)$$



$$\langle \omega_*, \nabla_S L(\omega) \rangle$$



$\omega = \pm \omega_*$ 2 global minima

$\langle \omega, \omega_* \rangle = 0$ stationary submanifold

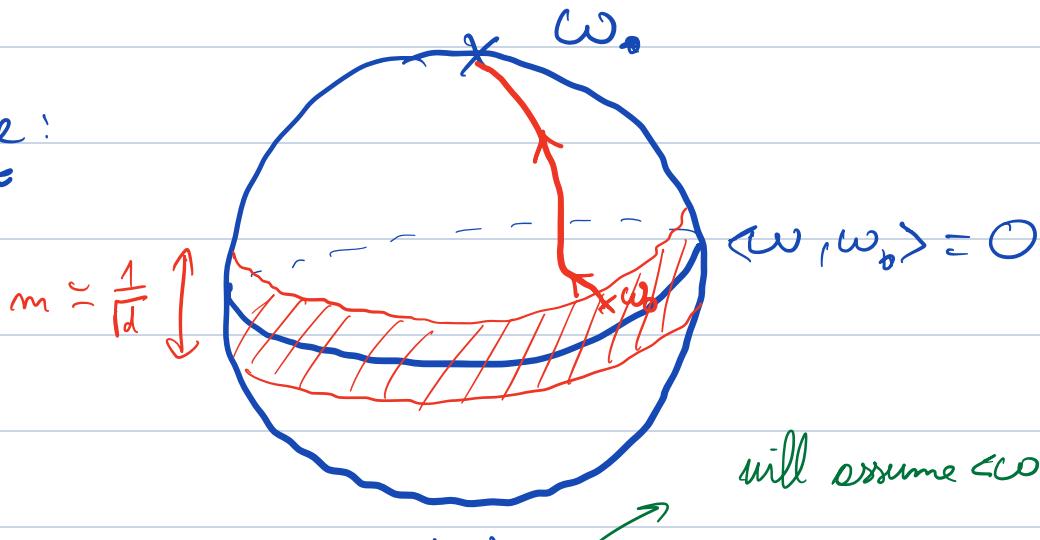
We will assume

$$\langle \omega^*, \omega_* \rangle > 0$$

(prob $1/2$)

and $\omega^t \rightarrow +\omega_*$

Lenscope:



$$\text{w.p. } \geq \frac{1}{2} + \eta$$

will assume $\langle \omega^0, \omega_0 \rangle \gtrsim \frac{c}{\sqrt{d}}$

$$\text{w.h.p. } \langle \omega^0, \omega_0 \rangle \leq \frac{C}{\sqrt{d}}$$

$$\text{At initialization, } \|\nabla_S L(\omega^0)\|_2 \approx \frac{1}{d^{\frac{k-1}{2}}}$$

$$\langle \omega_0, \nabla_S L(\omega^0) \rangle \approx \frac{1}{d^{\frac{k-1}{2}}}$$

very small gradient: dynamic initialized near a saddle-pt.
SGD will take a long time to escape this generic initialization

We want to study online SGD on this problem and in particular establish # of SGD steps to reach the global minimum.

Gradient flow limit

A standard approach is to approximate SGD by its continuous

(8)

$$\text{limit (step size } \eta \rightarrow 0\text{)}: \dot{\omega}^t = -\nabla_{\omega} L(\omega^t)$$

We only need to track $m_t = \langle \omega^t, \omega_* \rangle$

$$m_0 = \frac{c}{\sqrt{d}} \quad m_t = k (m_t)^{k-1} (1 - m_t^2)$$

$$\text{As long as } m_t \text{ small} \quad \dot{m}_t = k (m_t)^{k-1}$$

$$\frac{dm_t}{m_t^{k-1}} = k dt$$

$$\text{If } k=1: \quad m_T \approx m_0 + kT \quad \rightarrow \text{converge in } T = O_d(1)$$

$$\text{If } k=2: \quad m_T \approx m_0 e^{kT} \quad \rightarrow \text{converge in } T = O_d(\log d)$$

$$\text{If } k > 2: \quad \frac{1}{k-2} \left[\frac{1}{m_0^{k-2}} - \frac{1}{m_T^{k-2}} \right] = kT \quad \rightarrow m_T \approx \frac{1}{(d^{\frac{k}{2}-1} - kT)^{\frac{1}{k-2}}} \\ \rightarrow \text{converge in } T = d^{\frac{k}{2}-1}$$

Comparing k steps of SGD with step size η , we can show
the propagation of error: with high proba $\sqrt{\frac{1}{\eta}} \stackrel{\text{steps}}{\times} \eta$

$$\left| \underbrace{L(\omega^k)}_{\text{SGG}} - \underbrace{L(\omega^{k\eta})}_{\text{GF}} \right| \leq e^{ck\eta} \sqrt{d\eta}$$

(9)

For $k=1$, this shows that SGD with step size $\eta = \Theta_d(\frac{1}{d})$ achieves small test error in $\Theta_d(d)$ steps / samples
 [note that $S_2(d)$ samples is information-theoretic optimal]

For $k > 1$, because of propagation of error, we cannot approximate discrete SGD by GF throughout the entire dynamic.

We need a different analysis!

Projected online SGD

We consider the following algorithm: initializing at ω^0

At each step t :

$$\left\{ \begin{array}{l} \tilde{\omega}^{t+1} = \omega^t - \eta \nabla_S l(\omega^t, \alpha^t) \\ \omega^{t+1} = \frac{\tilde{\omega}^{t+1}}{\|\tilde{\omega}^{t+1}\|_2} \end{array} \right.) \text{ project back on the sphere}$$

$$l(\omega^t, \alpha^t) = \frac{1}{2} (f(\langle \omega^t, \alpha^t \rangle) - f(\langle \omega, \alpha^t \rangle))^2$$

$$\epsilon_{\text{iid}} \sim N(0, I_d)$$

(10)

Thm: [Ben Arous, Gheissari, Jagannath, '21]

If we choose

$$k=1$$

$$\eta = \tilde{\mathcal{O}}\left(\frac{1}{d}\right)$$

$$T = \mathcal{O}(d)$$

$$k=2$$

$$\eta = \tilde{\mathcal{O}}\left(\frac{1}{d}\right)$$

$$T = \mathcal{O}(d \log d)$$

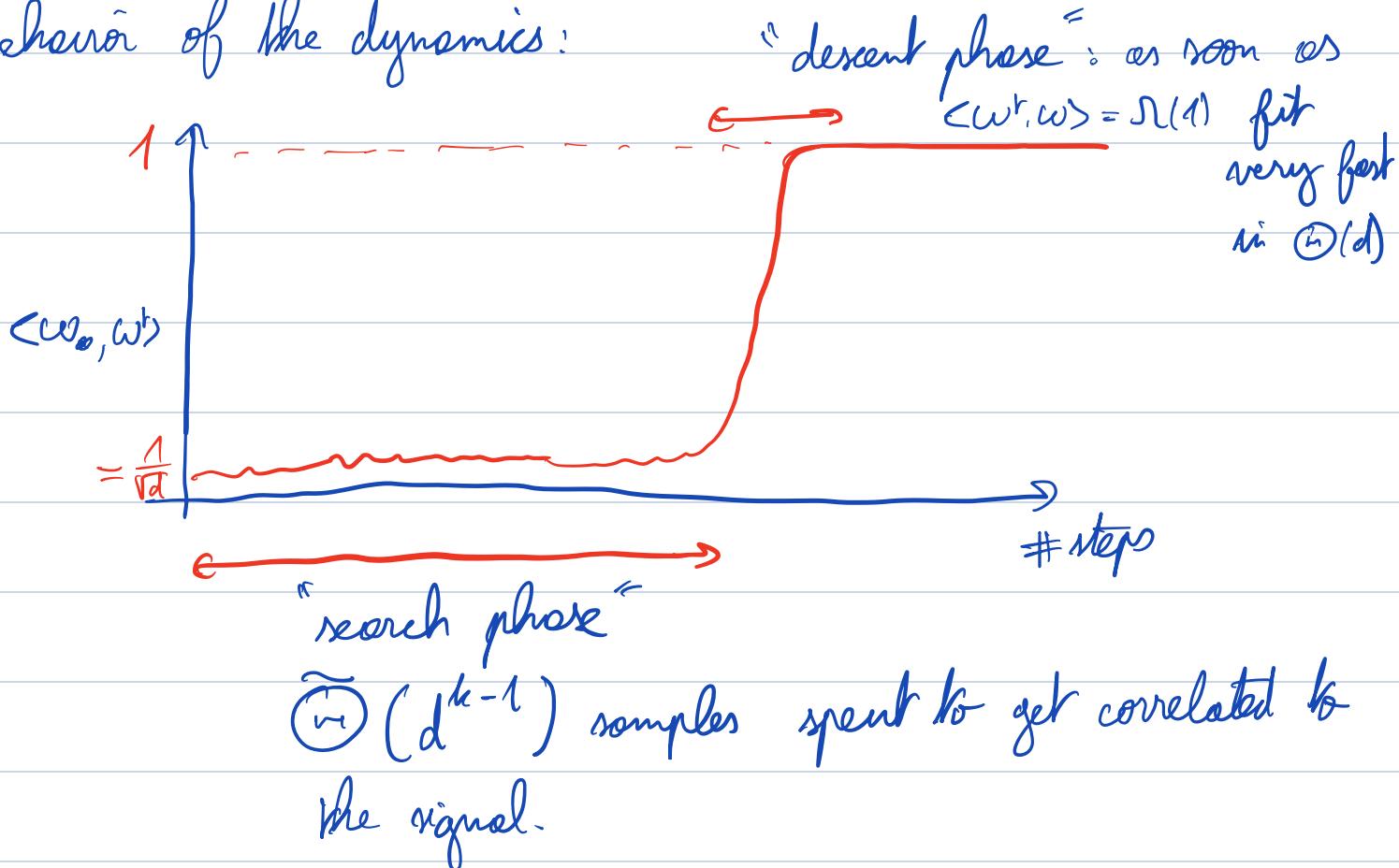
$$k > 2$$

$$\eta = \tilde{\mathcal{O}}\left(d^{-\frac{k}{2}}\right)$$

$$T = \tilde{\mathcal{O}}(d^{k-1})$$

Then online SGD with step size η and T steps will have $\langle w_0, w^T \rangle \approx 1$ w.h.p.

Behavior of the dynamics:



Main difficulty in this problem is leaving the high entropy region of the initialization. As soon as one has left this region, the "fitting" is fast.

Rule: Remember the case of learning k-parities. This is the same phenomenon!

Proof of the theorem

Note that $\eta T = d^{\frac{T-1}{2}}$ which matches the continuous time guarantee from GF. However GF will not be a good approximation of SGD for $\eta T \gg 1$.

Instead, we will (approximately) decompose the dynamics into a deterministic drift induced by the signal that drives dynamics to $\rightarrow w_*$ + a mean zero noise term.

$$\text{Dynamics: } w^{t+1} = \frac{w^t - \eta \nabla_S l(w^t, x^t)}{\|w^t - \eta \nabla_S l(w^t, x^t)\|_2}$$

Below we will focus on search phase, i.e., getting $\langle w^t, w_* \rangle = \Omega_d(1)$ (descent phase can be studied by GF)

For η small enough, $\omega^{t+1} \approx \omega^t - \eta \nabla_S l(\omega^t, \alpha^t)$ (*)

In these notes, we will ignore the projection step and directly consider (*).

$$\omega^{t+1} = \omega^t - \eta \nabla_S l(\omega^t, \alpha^t) = \omega_0 - \eta \sum_{\delta=0}^t \nabla_S l(\omega^\delta, \alpha^\delta)$$

It is enough to track $m^t = \langle \omega^t, \omega_* \rangle$

$$m_{t+1} = m_0 + \eta \sum_{\delta=0}^t g^\delta$$

$$g^\delta := -\langle \omega_*, \nabla_S l(\omega^\delta, \alpha^\delta) \rangle = \frac{(f(\langle \omega_*, \alpha^\delta \rangle) - f(\langle \omega^\delta, \alpha^\delta \rangle))}{(\alpha^\delta - \omega^\delta \langle \omega^\delta, \alpha^\delta \rangle)}$$

$$\begin{aligned} \text{Denote } \bar{g}^\delta &= \mathbb{E}_{\alpha^\delta} [g^\delta] = -\langle \omega_*, \nabla_S L(\omega^\delta) \rangle \\ &= k m_\delta^{k-1} (1 - m_\delta^2) \end{aligned}$$

Then

$$\begin{aligned} m_{t+1} &= m_0 + \eta \underbrace{\sum_{\delta=0}^t \bar{g}^\delta}_{D^t: \text{"deterministic drift}} + \eta \underbrace{\sum_{\delta=0}^t g^\delta - \bar{g}^\delta}_{M^t: \text{martingale}} \\ &= \frac{c}{d} \end{aligned}$$

(13)

$$M^t \text{ is a martingale} \quad (M^t = \sum_{s=0}^t X_s) \quad \mathbb{E}[X_t | X_{\leq t-1}] = 0$$

① Controlling the martingale:

We can use Doob's maximal inequality: for every $p \geq 1$,

$$\mathbb{P}\left(\max_{0 \leq t \leq T} |M_t| \geq \lambda\right) \leq \frac{p \mathbb{E}[|M_T|^p]}{(p-1) \lambda^p}$$

One can show using hypercontractivity of low-degree polynomials
on Gaussian data $\|f\|_{L^q}^2 \leq (q-1)^k \|f\|_{L^2}^2$
for any degree k polynomial

$$\mathbb{E}[|M_T|^p] \lesssim C_p \eta^p T^{\frac{p}{2}}$$

$$\frac{C_p \eta^p T^{\frac{p}{2}}}{\lambda^p} = \delta$$

Hence, with prob 1- δ ,

$$\max_{0 \leq t \leq T} |M_t| \leq C \sqrt{T} \frac{\delta^{-\frac{1}{2}}}{\eta}$$

dependency can
be improved to
 $\log(\frac{1}{\delta})$

In particular if $C \frac{\gamma \sqrt{T}}{\sqrt{8}} \leq \frac{m_0}{2}$

that is $\boxed{\gamma \sqrt{T} = O(\frac{1}{\sqrt{d}})}$

(A1)

Then: $\frac{m_0}{2} + D^t \leq m_t \leq \frac{3}{2} m_0 + D^t$ for all $t \leq T$.

and we can neglect the Martingale / noise term

② Controlling the drift term

In the search phase (m^t small) $\bar{g}^t \approx k m_d^{k-1}$

Hence for all $t \leq T$ (under (A1))

$$m_t \geq \frac{m_0}{2} + \gamma k \sum_{s=0}^{t-1} m_s^{k-1}$$

We can study this sequence and lower bound their growth. The Lemma below can be seen as a discrete lower bound version of Bihari - LaSalle inequality.

Lemme: For a sequence $(u_t)_{t \in \mathbb{N}}$ that satisfies

$$u_t \geq a_0 + a_1 \sum_{s=0}^{t-1} u_s^{k-1}$$

Then u_t is lower bounded by

$$\underline{k=2}: u_t \geq a_0(1+a_1)^t$$

$$\underline{k \geq 2}: u_t \geq \min\left(1, \frac{1}{(a_0^{-(k-2)} - \frac{k-2}{(1+a_1)} a_1 t)^{\frac{1}{k-2}}}\right)$$

$$\text{For } k \geq 2: m_0^{-(k-2)} = d^{\frac{k}{2}-1}$$

Hence $m_T = \mathcal{O}(1)$ as long as

$$qT = \mathcal{O}(d^{\frac{k}{2}-1}) \quad \text{A2}$$

In fact, if $t = o(d^{\frac{k}{2}-1})$ then $m_T \asymp \frac{1}{\sqrt{d}}$

↳ most of the search phase $m_T \asymp \frac{1}{\sqrt{d}}$

Combining A1 and A2

$$\eta \sqrt{T} = \widetilde{\Theta}\left(\frac{1}{\sqrt{d}}\right) \quad \eta T = \widetilde{\Theta}\left(d^{\frac{k}{2}-1}\right)$$

$$\Rightarrow \boxed{T = \widetilde{\Theta}(d^{k-1}) \quad \eta = \widetilde{\Theta}(d^{-\frac{k}{2}})}$$

for $k > 2$

For $k=2$, then $u_T \gtrsim m_0 e^{\eta T}$

$$\rightarrow \eta T = \widetilde{\Theta}(\log d)$$

$$\Rightarrow \boxed{T = \widetilde{\Theta}(\log d) \quad \eta = \widetilde{\Theta}(d^{-1})}$$

Landscape concentration

The above guarantees are for online SGD
 that is # samples = # steps

What about if we use samples?

Batch 1 SGD choose one random sample from the empirical distribution $\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}$

We can do the same drift + martingale decomposition but now

$$\bar{g}^s \rightarrow \hat{g}_m^s = \langle w_*, \nabla_S \hat{E}_m[l(w^s, x)] \rangle$$

$$= \frac{1}{m} \sum_{i=1}^m (y_i - f(\langle w^s, x_i \rangle)) f'(\langle w^s, x_i \rangle) (x_i - w^s \langle w^s, x_i \rangle)$$

If we can show that $\hat{g}_m(w) \geq c \bar{g}(w)$

$$\text{for all } w \in \mathbb{S}^{d-1} \quad \langle w, w_* \rangle \geq \frac{c}{\sqrt{d}}$$

Then the same analysis as for one-pass SGD goes through

For a given $w \in \mathbb{S}^{d-1}$, $\hat{g}_m(w)$ is an empirical average with mean $\bar{g}(w)$

For simplicity take f with, $\|f\|_\infty, \|f'\|_\infty \leq C$.

$$P(|\hat{g}_m(\omega) - \bar{g}(\omega)| \geq \sqrt{\text{Var}(\hat{g}_m(\omega))}) \leq Ce^{-ct^2}$$

Mence $P\left(\min_{\substack{\omega \in S^{d-1} \\ \langle \omega, \omega_0 \rangle \geq d^\gamma}} |\hat{g}_m(\omega) - \bar{g}(\omega)| \geq \sqrt{\frac{c \log d}{m}}\right) \leq \frac{C}{d^C}$

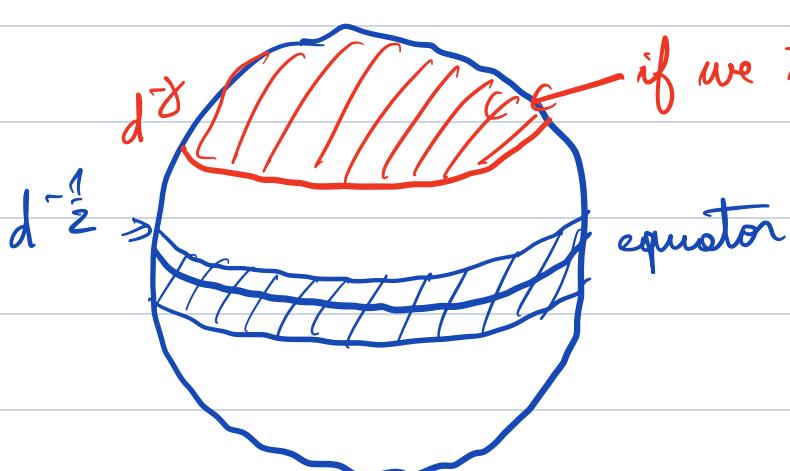
[using a standard uniform convergence bound]

Hence for $\frac{1}{2} \leq \gamma \leq 0$

then if we take $\sqrt{\frac{1}{m}} \leq \frac{1}{d^{\gamma(k-1)}}$

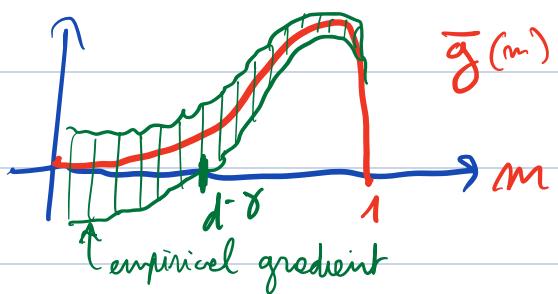
$$\Rightarrow m \geq d^{2\gamma(k-1)+1}$$

then the landscape is "benign" for any $\langle \omega, \omega_0 \rangle \geq d^\gamma$



if we initialize on this cap, then

$$\omega^* \rightarrow \omega_*$$



If $\gamma = \frac{1}{2}$, $n \geq d^k$ random initialization falls in this cap w.h.p.

and $w_t \rightarrow w_*$ in d^{k-1}

↳ so in fact see each sample only once during the dynamics w.h.p.

and we didn't gain anything compared to online SGD

Open question: can we show that batch-1 multi-pass SGD still succeeds with high-probability for $n \ll d^{k-1}$ despite the landscape being not benign?

How tight are these guarantees?

Online SGD requires $n = T \asymp d^{k-1}$ samples

and runtime $Td \asymp d^k$ (each step compute a d -dim gradient)

* for parity fits: SQ lower bound is $\Omega(d^k)$
 → match best SQ runtime

* for Hermite-hc: CSQ lower bound is $\Omega(d^{k/2})$

↳ modifying the dynamics, online SGD succeeds
 with $m = T \asymp d^{k/2}$ (runtime $d^{\frac{k}{2}+1}$)

[Dannan, Nichani, Ge, Lee, '23]

* for Hermite-hc: SQ lower bound is much smaller
 $\Omega(d)$

↳ taking L^1 loss $l(y, \hat{y}) = |y - \hat{y}|$

online SGD succeeds in $m \asymp T \asymp d$

↳ reusing samples with L^2 loss

full-batch GD with $m = \Theta(d)$ and $T = \Theta_1(1)$
 succeeds

* There exists single index functions that seem to have a fundamental statistical/computational trade-off

$n = \Theta(d)$ is info-theoretically optimal

$n = \tilde{\Omega}(d^{k/2})$ samples is necessary (and sufficient) for a polynomial-time algo to exist.

(under low-degree poly conjecture)

[See Kortes and Lede presentation !!!]