

Lecture 9: Optimization in feature-learning regime (contd)

Learning lower bounds

Last lecture: express a model $f(x; \theta) = \frac{1}{m} \sum_{j=1}^m \phi(x; \theta_j)$
recap $= \int \phi(x; \theta) \mu(d\theta)$

- gradient flow wrt parameters (θ_j) \rightarrow wasserstein gradient flow wrt Measure μ .

Today:
→ convergence properties
→ learning lower bounds

→ Recall the continuity equation (Liouville eq) that

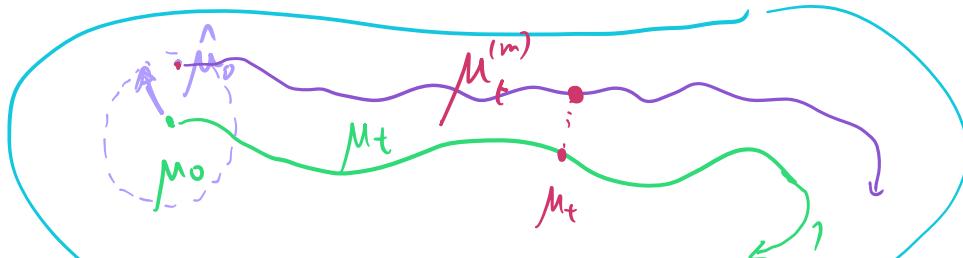
M_t solves:

$$\partial_t \mu_t = \operatorname{div} (\nabla G(\cdot; M_t) \cdot \mu_t) *$$

Natural questions: $\partial_t \mu_t = Q(\mu_t)$

- under what conditions does this PDE converge in time towards the global minimizers of L ? (convergence in t).
- How are the dynamics affected by overparametrization (convergence in width, m)

$$\hat{\mu}_0^{(m)} = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j}$$



$\Theta_i \sim_{iid} \mu_0$
 $P(D)$
 $\arg\min_{\mu} L(\mu)$

$\hookrightarrow \mu_t^{(m)}$ solves $(*)$ with initial condition given by $\mu_0^{(m)}$.

Theorem [CB, R, EVE, MNM, SS'18] For any fixed

$T > 0$, $\mu_T^{(m)}$ converges weakly to μ_T as $m \rightarrow \infty$,
 where μ_T solves $(*)$ with initial condition μ_0 .

\hookrightarrow In other words, Sampling and "learning" commute
 in the limit of $m, t \rightarrow \infty$.

$$\left(\lim_{t \rightarrow \infty} \lim_{m \rightarrow \infty} \mu_t^{(m)} = \lim_{m \rightarrow \infty} \lim_{t \rightarrow \infty} \mu_t^{(m)} \right).$$

The evolution $(\mu_t)_t$ is referred as the mean-field limit of the system. It provides a description of the system at any point in time that "neglects" interactions.

Global convergence in time

Recall that $(*)$ $\partial_t \mu_t = \text{div} (\nabla G(\cdot, \mu_t) \cdot \mu_t)$
 is the WGF associated with the functional

$$L[\mu] = - \int F dm + \frac{1}{2} \iint K(\theta, \theta') dm d\mu'.$$

↳ Stationary point of WGF satisfy

$$0 = \text{div}(\nabla G(\cdot; \hat{\mu}) \cdot \hat{\mu}) \quad \text{Euler-Lagrange equation.}$$

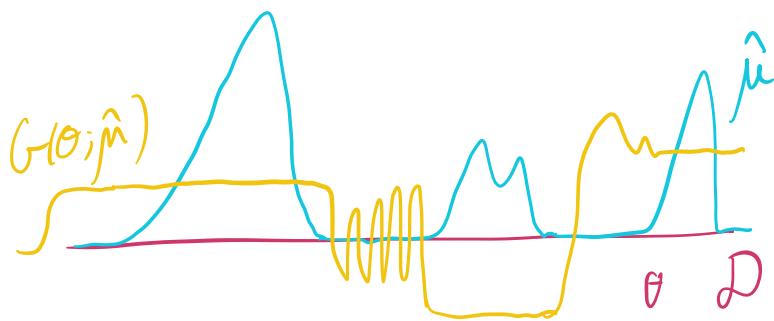
Claim: $\hat{\mu}$ is a stationary point if

$$\nabla G(\theta; \hat{\mu}) = 0 \quad \text{for } \theta \in \text{supp}(\hat{\mu}) \quad (1)$$

↳ Global minimizers of L satisfy

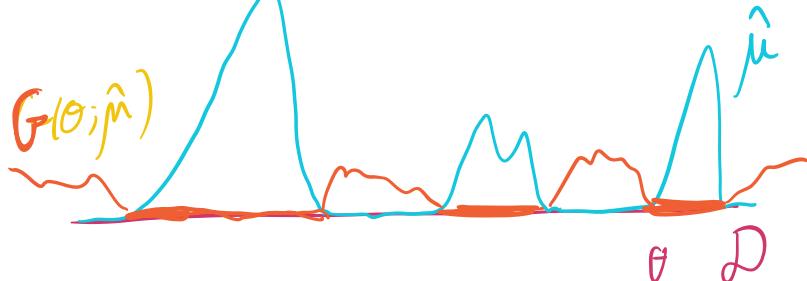
$$\begin{cases} G(\theta; \hat{\mu}^*) = 0 & \text{for } \theta \in \text{supp}(\hat{\mu}^*) \\ G(\theta; \hat{\mu}^*) \geq 0 & \text{for any } \theta. \end{cases} \quad (2)$$

Q): When can we guarantee that $(1) \Rightarrow (2)$?



← stationary point of WGF

U1



← global minimum of L .

Two challenges that we need to address:

- (1) \rightarrow Not "concentrate" mass too soon along dynamics
- (2) \rightarrow From $\nabla \phi(\theta; \mu^*) = 0$ to $\phi(\theta; \mu^*) = 0$?

\hookrightarrow For the first challenge, we can address it by using initial conditions with full support (CB); alternative: add noise into the gradient \rightarrow adds a term Δp_t that enforces full support (MNM).

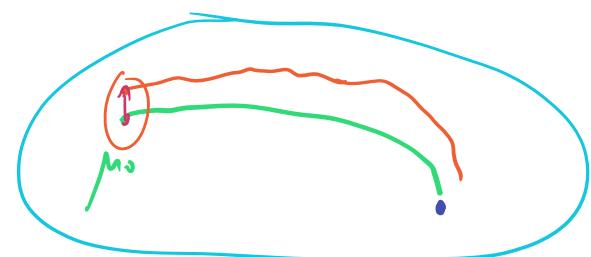
\hookrightarrow The second challenge is addressed thanks to the homogeneity of the neuron. $\phi(\lambda \theta; x) = \lambda^p \phi(\theta; x)$

deux
avec

$$\phi(\theta; x) = P^{-1} \langle \theta, \nabla \phi(\theta; x) \rangle \quad \begin{array}{l} \text{[CB]} \\ \text{[R, EVE]} \end{array}$$

[Euler Rule for homogeneous]

convergence in m



\rightarrow If we have an initial condition μ_0 that leads to global convergence, how many neurons are needed to approximate it uniformly in time?

\rightarrow For $t=0$; $f(x; \theta_1, \dots, \theta_m)$ with $\theta_i \sim_{iid} \mu_0$ approximates the mean function

$$\bar{f}(x; \mu_0) = E_{\Theta \sim \mu_0} [\phi(x; \theta)]$$

→ Approximation error satisfies

$$|E \| f(\cdot; \theta_1 \dots \theta_m) - \bar{f} \|_2^2 = \frac{T(\mu_0)}{m}, \text{ with}$$

$$T(\mu_0) = |E_{\mu_0} \| \phi(\cdot; \theta) \|_2^2 - \| \bar{f} \|_2^2 \quad \text{Monte-Carlo Rate.}$$

Q: How does this error evolve during training?

A: Under mild assumptions in the Mean-Field limit, this error remains bounded (at the MC scale),

uniformly in time:

Theorem [CRBE '20] $\int_0^t \int_n \phi d\mu_t^{(n)} - \int_0^t \int_n \phi d\mu_t$

$$\lim_{t \rightarrow \infty} \lim_{n \rightarrow \infty} m \cdot |E \| \hat{f}_t^{(n)} - f_t \|_2^2 \leq T(\mu_\infty)$$

Remark: CLT analog of the LLN (Law of large numbers)
Mean-field convergence.

→ Extending the finite-time horizon propagation-of-chars results; based on Gronwall's inequality.

If $\| \partial_t \hat{\mu}_t - \partial_t \mu_t \| \leq \beta/t \cdot \| \hat{\mu}_t - \mu_t \|$ then

$$\| \hat{\mu}_T - \mu_T \| \leq \| \hat{\mu}_0 - \mu_0 \| \exp \left(\int_0^T \beta(s) ds \right)$$

$$\cdot Q(\hat{\mu}_t) - Q(\mu_t)$$

→ We use a different technique, based on a Volterra-type Kernel.

→ MF take-aways:

- ⊕ Rich qualitative description of empirically observed phenomena.
- ⊕ Consistent with the "mystery" of over-parametrization.
 - ↳ as m grows, optim is "easier" with generalization under control.
- ⊕ Probably more expressive than NTK.
- ⊖ Qualitative results: no rates of convergence neither in time nor in m (width)!

Q: Can we hope to efficiently learn any target function f^* that can be well approximated in our shallow function class \mathcal{F} ?

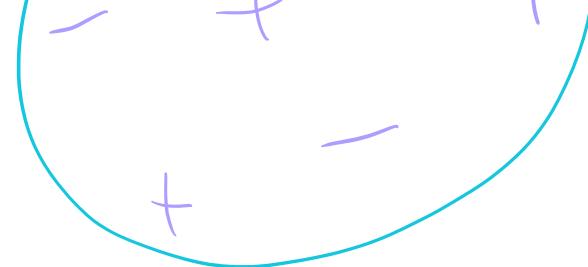
Learning Lower Bounds

↳ Recall from first lecture the example of learning 1-Lipschitz functions

→ Key source of hardness:



there are "too many"
 $\text{exp}(d)$ candidate functions
 in my class.



→ Learning Parities with Gradient Descent [Shalev-Shwartz, Shamir, Shwartz] 17

- Let $v^* \in \{0, 1\}^d$ a point in the d -dim hypercube chosen uniformly at random.

- Now consider the target function

$$f_{v^*} : \{0, 1\}^d \rightarrow \{-1, 1\}$$

$$x \mapsto (-1)^{\langle x, v^* \rangle}$$

In words, $f_{v^*}(x)$ is the parity of x within the subset v^* .

→ Dataset $\{x_i\}_{i=1..n} \sim \text{Unit } \{0, 1\}^d$

$$y_i | x_i = f_{v^*}(x_i)$$

- We wish to learn f_{v^*} from the dataset $\{(x_i, y_i)\}_i$ using a shallow Relu network.

$$f(x; \theta_1 \dots \theta_M) = \sum_{j=1}^M a_j \tau(\langle x, w_j \rangle + b_j)$$

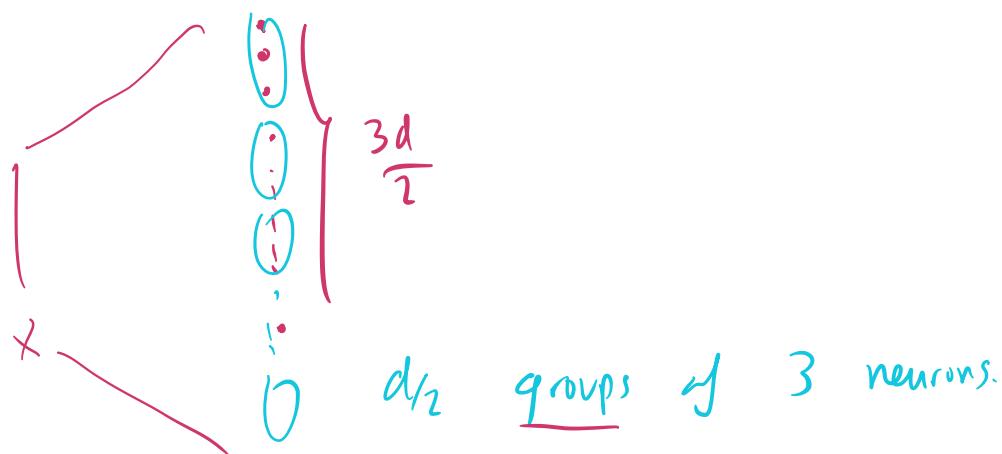
$$\min_{\theta_1 \dots \theta_M} \frac{1}{n} \sum_{i=1}^n (f(x_i; \vec{\theta}) - y_i)^2$$

M: width of net
 n: # of training points.

Fact: When $M > \frac{3d}{2}$, a shallow ReLU net can express f_v for any v . (no error).

Proof

Set $\begin{cases} w_j = v & \text{for } j = 1 \dots 3d/2 \\ w_j = 0 & \text{for } j > 3d/2 \end{cases}$



For $i = 1 \dots d/2$ set biases and output weights as
 $j = 3i + \{1, 2, 3\}$

$$b_j = \begin{cases} -(2i - \frac{1}{2}) \\ -2i \\ -(2i + \frac{1}{2}) \end{cases} \quad a_j = \begin{cases} 1 \\ -2 \\ 1 \end{cases}$$

$$\sum_{j=3i+1,2,3} a_j \sigma(\underbrace{\langle x, w_j \rangle}_{\text{if } \langle x, w_j \rangle \neq 2i} + b_j) = \begin{cases} 0 & \text{if } \langle x, w_j \rangle \neq 2i \\ \frac{1}{2} & \text{if } \langle x, w_j \rangle = 2i \end{cases}$$

$$\rightarrow \sum_i \left(\sum_{j=3i+1,2,3} a_j \sigma(\langle x, w_j \rangle + b_j) \right) = \begin{cases} \frac{1}{2} & \text{if } \langle x, v \rangle \text{ even} \\ 0 & \text{otherwise.} \end{cases}$$

So we are expressing f_v (up to linear rescaling).

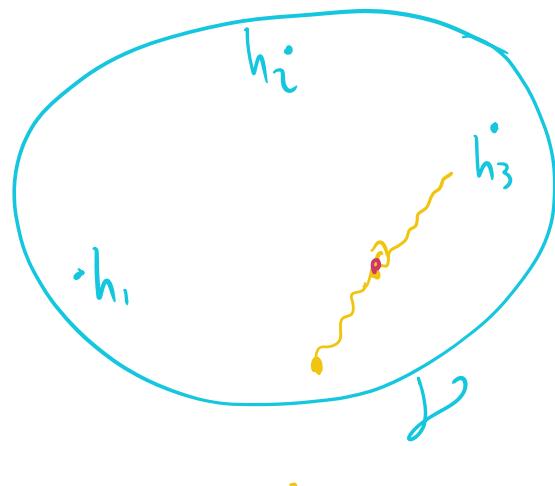
- Can we learn such function with gradient descent?
- Consider the optimisation given by

$$\min_{\theta} L_h(\theta) = \mathbb{E}_x [l(\tilde{f}_{\theta}(x), h(x))]$$

\tilde{f}_{θ} : is a differentiable class (eg NNs)

h : is a target hypothesis picked from a set \mathcal{V} .

- We consider how sensitive the gradient of L_h wrt θ is to h (target function). We measure this through the Variance



$$\text{Var}(V; L; \theta) = \mathbb{E}_h \| \nabla_{\theta} L_h(\theta) - \mathbb{E}_h \nabla_{\theta} L_h(\theta) \|^2$$

Var measures how much "signal" from target ($h+V$) is contained in the gradient at point θ .

Theorem [SSS'17] . Suppose \mathcal{V} is a discrete set

of target functions satisfying $\langle h, h' \rangle_V = 0$ if $h \neq h'$ and $\|h\|_V^2 \leq 1$

Then \tilde{f}_{θ} is a universal function with

- Assume f_θ is differentiable wrt θ , with
$$E_x [\| \nabla_\theta \tilde{f}_\theta(x) \|^2] \leq G(\theta)^2$$

- Pick l : squared loss.

Then $\text{Var}(V; L; \theta) \leq \frac{G(\theta)^2}{|D|}$

Prof: [SSS'17] \rightarrow very easy.

\rightarrow When $D =$ Random parities f_{v_0}

↳ Fact: $\begin{cases} \langle f_v, f_{v'} \rangle = 0 & \text{if } v \neq v' \\ \|f_v\|^2 = 1 \end{cases}$

↳ We conclude that

$$\text{Var}(V; L; \theta) \leq \frac{G(\theta)^2}{2^d}$$

↳ The gradient at any point will be concentrated at a value independent of target. (bad!).

① \rightarrow Learning parities is GD-hard. Is it hard in general?

$$f_v(x) = (-1)^{\langle x, v \rangle}$$

$$\forall i, \quad \langle x_i, v \rangle = f_v(x_i) \pmod 2$$

↳ Gaussian elimination solves parities in linear time!