More Deep Learning: Approximation, Generalization

CPSC 532S: Modern Statistical Learning Theory 16 March 2022 <u>cs.ubc.ca/~dsuth/532S/22/</u>

Admin

- A3 is up; due next Friday the 25th
- There will be an A4 due at the end of the term
 - Going to try to post that soon (before A3 is due) so you can start early
- This week only, Thursday office hours were moved to this morning
- Project proposals due tonight

 - An informal paragraph on Piazza: just tell me what you want to do • Say who your partners are in the post, I'll add them to be able to see it • Again, the scope of these is meant to be small

 - A lit survey doesn't require fully understanding the proofs or anything • An "extension" could be just reading the proofs and talking about when the assumptions hold, etc

Last time

• Deep learning, particularly fully-connected feedforward networks

Last time

- Deep learning, particularly fully-connected feedforward networks
- Universal approximation results:
 - that approximate any continuous function in sup-norm
 - as long as the activation function is non-polynomial

• There exist networks (of depth two, but potentially extremely wide)

SSBD chapter 20:

• 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$



SSBD chapter 20:

- - (remember that computers always represent things as $\{0,1\}^d$...)

• 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$



SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$ • (remember that computers always represent things as $\{0,1\}^d...$) • ...but, it takes exponential width to do that



SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$ • (remember that computers always represent things as $\{0,1\}^d$...)
- ...but, it takes exponential width to do that
- ...but, there's a network of size $\mathcal{O}(T^2)$ that can implement all boolean functions that can be computed in maximum runtime T



SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$ • (remember that computers always represent things as $\{0,1\}^d$...)
- ...but, it takes exponential width to do that
- ...but, there's a network of size $\mathcal{O}(T^2)$ that can implement all boolean functions that can be computed in maximum runtime T



SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$ • (remember that computers always represent things as $\{0,1\}^d$...)
- ...but, it takes exponential width to do that
- ...but, there's a network of size $\mathcal{O}(T^2)$ that can implement all boolean functions that can be computed in maximum runtime T

Circuit Complexity and Neural Networks, Ian Parberry (1994) - <u>UBC access</u>



Proceedings of Machine Learning Research vol 125:1–22, 2020

Universal Approximation with Deep Narrow Networks

Patrick Kidger Terry Lyons Mathematical Institute, University of Oxford

Editors: Jacob Abernethy and Shivani Agarwal

The classical Universal Approximation Theorem holds for neural networks of arbitrary width and bounded depth. Here we consider the natural 'dual' scenario for networks of bounded width and arbitrary depth. Precisely, let n be the number of inputs neurons, m be the number of output neurons, and let ρ be any nonaffine continuous function, with a continuous nonzero derivative at some point. Then we show that the class of neural networks of arbitrary depth, width n + m + 2, and activation function ρ , is dense in $C(K; \mathbb{R}^m)$ for $K \subseteq \mathbb{R}^n$ with K compact. This covers every activation function possible to use in practice, and also includes polynomial activation functions, which is unlike the classical version of the theorem, and provides a qualitative difference between deep narrow networks and shallow wide networks. We then consider several extensions of this result. In particular we consider nowhere differentiable activation functions, density in noncompact domains with respect to the L^p -norm, and how the width may be reduced to just n + m + 1 for 'most' activation functions. 5

KIDGER@MATHS.OX.AC.UK TLYONS@MATHS.OX.AC.UK

mm

Abstract

• Deep networks much better at learning compositional structure

- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

compositional structure an be represented compactly by deep by a compact shallow architecture." – <u>Y. Bengio & LeCun (2007)</u>

- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

– Y. Bengio & LeCun (2007)

• Lots of empirical evidence, but theoretical support pretty limited until recently

- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

 - Telgarsky notes section 5 give a particular such function: shallow net needs huge width to approximate, but narrow not-super-deep net can approximate it efficiently

– <u>Y. Bengio & LeCun (2007)</u>

• Lots of empirical evidence, but theoretical support pretty limited until recently

- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

 - Telgarsky notes section 5 give a particular such function: shallow net needs huge width to approximate, but narrow not-super-deep net can approximate it efficiently

– Y. Bengio & LeCun (2007)

• Lots of empirical evidence, but theoretical support pretty limited until recently

Also proved for a certain class of functions by <u>Mhaskar, Liao, Poggio (2016)</u>

- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

 - Telgarsky notes section 5 give a particular such function: shallow net needs huge width to approximate, but narrow not-super-deep net can approximate it efficiently

– <u>Y. Bengio & LeCun (2007)</u>

• Lots of empirical evidence, but theoretical support pretty limited until recently

Also proved for a certain class of functions by <u>Mhaskar, Liao, Poggio (2016)</u>

• Lu et al. (2017): approximating wide nets with deep nets easier(ish) than vice versa



- Deep networks much better at learning compositional structure
 - "We claim that most functions that can be represented compactly by deep architectures cannot be represented by a compact shallow architecture."

 - Telgarsky notes section 5 give a particular such function: shallow net needs huge width to approximate, but narrow not-super-deep net can approximate it efficiently

 - exponentially smaller deep nets than shallow

– <u>Y. Bengio & LeCun (2007)</u>

• Lots of empirical evidence, but theoretical support pretty limited until recently

Also proved for a certain class of functions by <u>Mhaskar, Liao, Poggio (2016)</u>

Lu et al. (2017): approximating wide nets with deep nets easier(ish) than vice versa Liang and Srikant (2017): can approximate piecewise-constant funcs with



- We have some universal approximation results
- A lot of people use this to say "neural networks can do anything!



- We have some universal approximation results
- A lot of people use this to say "neural networks can do anything!



- We have some universal approximation results
- A lot of people use this to say "neural networks can do anything!

- These kind of approximation results don't tell us:
 - What practically-sized networks can do

n results networks can do anything! 💪"

on't tell us: do

- We have some universal approximation results A lot of people use this to say "neural networks can do anything!

- These kind of approximation results don't tell us:
 - What practically-sized networks can do
 - Gaussian kernels can also do anything (6)...with ridiculously large norm

- We have some universal approximation results A lot of people use this to say "neural networks can do anything!

- These kind of approximation results don't tell us:
 - What practically-sized networks can do

 - Gaussian kernels can also do anything (2)...with ridiculously large norm • Neural nets can do anything... if they're ridiculously large (or large norm)

- We have some universal approximation results A lot of people use this to say "neural networks can do anything!

- These kind of approximation results don't tell us:
 - What practically-sized networks can do
 - Gaussian kernels can also do anything (2)...with ridiculously large norm • Neural nets can do anything... if they're ridiculously large (or large norm)
 - Even if our class approximates, do we generalize? (Does ERM, RLM, ... work?)



- We have some universal approximation results A lot of people use this to say "neural networks can do anything!

- These kind of approximation results don't tell us:
 - What practically-sized networks can do
 - Gaussian kernels can also do anything (6)...with ridiculously large norm • Neural nets can do anything... if they're ridiculously large (or large norm)
 - Even if our class approximates, do we generalize? (Does ERM, RLM, ... work?)
 - Does (S)GD find an approximate ERM / RLM / something that generalizes?



- We have some universal approximation results
- A lot of people use this to say "neural networks can do anything! [] "

- These kind of approximation results don't tell us:
 - What practically-sized networks can do
 - Gaussian kernels can also do anything (2)...with ridiculously large norm • Neural nets can do anything... if they're ridiculously large (or large norm)
 - Even if our class approximates, do we generalize? (Does ERM, RLM, ... work?)
 - Does (S)GD find an approximate ERM / RLM / something that generalizes? • We (pretty much) know it doesn't always find an (approximate) ERM: ERM with deep nets (even for square loss) is NP-hard
- so, if you can prove that it *does*, let me know =)



• and $\Omega(PL\log\frac{P}{I})$, so nearly tight – <u>Bartlett/Harvey/Liaw/Mehrabian (2019</u>)

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $O(PL \log P)$

• and $\Omega(PL\log\frac{P}{L})$, so nearly tight – <u>Bartlett/Harvey/Liaw/Mehrabian (2019)</u> $P = \prod^{L} d_{\ell-1} d_{\ell}$ for fully-connected networks $\ell = 1$

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $O(PL \log P)$

- and $\Omega(PL\log \frac{P}{L})$, so nearly tight <u>Bartlett/Harvey/Liaw/Mehrabian (2019)</u>
 - $P = \prod^{L} d_{\ell-1} d_{\ell}$ for fully-connected networks $\ell = 1$
- For piecewise-constant, e.g. threshold functions, VCdim = $\Theta(P \log P)$

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $\mathcal{O}(PL \log P)$

• and $\Omega(PL\log \frac{P}{L})$, so nearly tight – <u>Bartlett/Harvey/Liaw/Mehrabian (2019</u>)

•
$$P = \prod_{\ell=1}^{L} d_{\ell-1} d_{\ell}$$
 for fully-connected

- For piecewise-constant, e.g. threshold functions, VCdim = $\Theta(P \log P)$
- For piecewise-polynomial, $\mathcal{O}(PL^2 + PL\log P)$, $\mathcal{O}(PU)$ with U units

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $O(PL \log P)$

networks

•
$$P = \prod_{\ell=1}^{L} d_{\ell-1} d_{\ell}$$
 for fully-connected

- For piecewise-constant, e.g. threshold functions, VCdim = $\Theta(P \log P)$
- For piecewise-polynomial, $\mathcal{O}(PL^2 + PL\log P)$, $\mathcal{O}(PU)$ with U units
- For sigmoids/similar, $\mathcal{O}(P^2 U^2)$ and $\Omega(P^2)$

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $O(PL \log P)$ • and $\Omega(PL\log\frac{P}{r})$, so nearly tight – <u>Bartlett/Harvey/Liaw/Mehrabian (2019)</u>

networks

• For ReLU (or general piecewise-linear) nets with P params, VCdim = $\mathcal{O}(PL\log P)$ • and $\Omega(PL\log\frac{P}{L})$, so nearly tight – <u>Bartlett/Harvey/Liaw/Mehrabian (2019)</u>

•
$$P = \prod_{\ell=1}^{L} d_{\ell-1} d_{\ell}$$
 for fully-connected

- For piecewise-constant, e.g. threshold functions, VCdim = $\Theta(P \log P)$ • For piecewise-polynomial, $\mathcal{O}(PL^2 + PL\log P)$, $\mathcal{O}(PU)$ with U units • For sigmoids/similar, $\mathcal{O}(P^2 U^2)$ and $\Omega(P^2)$
- Theorem 8.13/8.14 of Anthony & Bartlett (1999) textbook <u>UBC access</u>
- networks

- We use networks with a lot of parameters

ResNet-50 has ~25 million parameters and depth 50: VCdim > 1 billion



- We use networks with a lot of parameters
 - ResNet-50 has \sim 25 million parameters and depth 50: VCdim > 1 billion
- We can train our networks to get zero error even for random labels



- We use networks with a lot of parameters
- We can train our networks to get zero error even for random labels
 - Even AlexNet can shatter CIFAR-10, *almost* shatter ImageNet
 - <u>Neyshabur et al. (2015)</u>, <u>Zhang et al. (2017)</u>



training error is 0) under different label corruptions.

• ResNet-50 has \sim 25 million parameters and depth 50: VCdim > 1 billion

Figure 1: Fitting random labels and random pixels on CIFAR10. (a) shows the training loss of various experiment settings decaying with the training steps. (b) shows the relative convergence time with different label corruption ratio. (c) shows the test error (also the generalization error since



- We use networks with a lot of parameters
- We can train our networks to get zero error even for random labels
 - Even AlexNet can shatter CIFAR-10, *almost* shatter ImageNet
 - <u>Neyshabur et al. (2015)</u>, <u>Zhang et al. (2017)</u>

• ResNet-50 has \sim 25 million parameters and depth 50: VCdim > 1 billion

• But these architectures do generalize well – VC of arch. can't explain that


Problems with parameter counting

- We use networks with a lot of parameters
- We can train our networks to get zero error even for random labels
 - Even AlexNet can shatter CIFAR-10, *almost* shatter ImageNet
 - <u>Neyshabur et al. (2015)</u>, <u>Zhang et al. (2017)</u>
 - But these architectures do generalize well VC of arch. can't explain that
 - Uniform stability can't either, since it's data-independent; on-average replace-one stability always can, but hard

• ResNet-50 has \sim 25 million parameters and depth 50: VCdim > 1 billion



Problems with parameter counting

- We use networks with a lot of parameters
- We can train our networks to get zero error even for random labels
 - Even AlexNet can shatter CIFAR-10, *almost* shatter ImageNet
 - <u>Neyshabur et al. (2015)</u>, <u>Zhang et al. (2017)</u>
 - But these architectures do generalize well VC of arch. can't explain that
 - Uniform stability can't either, since it's data-independent; on-average replace-one stability always can, but hard
- Making hidden layers wider can often improve generalization, but worsens parameter counting-based bounds

ResNet-50 has ~25 million parameters and depth 50: VCdim > 1 billion



Problems with parameter counting

- We use networks with a lot of parameters
- We can train our networks to get zero error even for random labels
 - Even AlexNet can shatter CIFAR-10, *almost* shatter ImageNet
 - <u>Neyshabur et al. (2015)</u>, <u>Zhang et al. (2017)</u>
 - But these architectures do generalize well VC of arch. can't explain that
 - Uniform stability can't either, since it's data-independent; on-average replace-one stability always can, but hard
- Making hidden layers wider can often improve generalization, but worsens parameter counting-based bounds
- Remember that \mathscr{H}_k has infinite VCdim for universal kernels, but we can still learn with small-norm predictors

ResNet-50 has ~25 million parameters and depth 50: VCdim > 1 billion



Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\ell}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathfrak{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$.

 $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|$





Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{L}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$. $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{c}$

Base case, L = 0:



Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{L}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$. $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{c}$

Base case, L = 0: $\hat{\mathfrak{R}}_{S}(\{x \mapsto x_{i} : j \in [d]\})$



Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{C}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$. $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{c}$

Base case, L = 0: $\hat{\Re}_{S}(\{x \mapsto x_{j} : j \in [d]\}) \leq \frac{1}{n} \left(\max_{i} ||(x_{1,j}, \dots, x_{n,j})||_{2}\right) \sqrt{2\log d}$



Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{C}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$.

Base case, L = 0: $\hat{\Re}_{S}(\{x \mapsto x_{j} : j \in [d]\}) \leq \frac{1}{n} (\max_{j} ||(x_{1,j}, ..., x_{n,j})||_{2}) \sqrt{2 \log d}$ $= \frac{1}{n} ||X||_{2,\infty} \sqrt{2 \log d}$

 $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{c}$



Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lips Let \mathscr{F}_L be the set of *L*-layer no-inter with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq B$

Base case, L = 0: $\hat{\Re}_{S}(\{x \mapsto x_{j} : j \in [d]\}) \leq \frac{1}{n} (\max_{j} || (x_{1,j}, ..., x_{n,j}) ||_{2}) \sqrt{2 \log d}$ $= \frac{1}{n} ||X||_{2,\infty} \sqrt{2 \log d} = \frac{1}{n} ||X||_{2,\infty} (2\rho B)^{0} \sqrt{2 \log d}$

schitz with
$$\sigma_{\ell}(0) = 0.$$

rcept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)}),$
 $\leq \frac{1}{n} ||X||_{2,\infty} (2\rho B)^L \sqrt{2\log d}.$
 $||M||_{b,c} = \left\| (||M_{\cdot 1}||_b, ..., ||M_{\cdot d}||_{b,c}) \right\|_{2,\infty} (2\rho B)^{0} \sqrt{2\log d}$
 $\sqrt{2\log d} = \frac{1}{2} ||X||_{2,\infty} (2\rho B)^{0} \sqrt{2\log d}$



Theorem: Fix
$$\sigma_1, \ldots, \sigma_L$$
 each ρ -Lipschitz with $\sigma_\ell(0) = 0$.
Let \mathscr{F}_L be the set of L -layer no-intercept nets, $f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)})$,
with $\|W_\ell^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.
 $\|M\|_{b,c} = \left\| \left(\|M_{\cdot 1}\|_b, \ldots, \|M_{\cdot d}\| \right) \right\|_{2,\infty} (1-p) \|W_{\cdot 1}\|_{2,\infty} \|W_{$

Inductive step:



Theorem: Fix
$$\sigma_1, \ldots, \sigma_L$$
 each ρ -Lipschitz with $\sigma_\ell(0) = 0$.
Let \mathscr{F}_L be the set of L -layer no-intercept nets, $f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)})$,
with $\|W_\ell^{\top}\|_{1,\infty} \leq B$. Then $\hat{\Re}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.
Base case, $L = 0$:
 $\hat{\Re}_S(\{x \mapsto x_j : j \in [d]\}) \leq \frac{1}{n} (\max_j \|(x_{1,j}, \ldots, x_{n,j})\|_2) \sqrt{2\log d}$
 $= \frac{1}{n} \|X\|_{2,\infty} \sqrt{2\log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2\log d}$
Inductive step:
 $\hat{\Re}_S(\mathscr{F}_{\ell+1}) = \hat{\Re}_S \left(\left\{ x \mapsto \sigma_{\ell+1} \left(\|W_{\ell+1}^{\top}\|_{1,\infty} g(x) \right) : g \in \operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \right\}$
 $\mathcal{F} \in \mathcal{F}_{\ell,\ell} \Rightarrow \mathfrak{T}_{\ell,\ell} (\mathcal{F}_{\ell,\ell}) + \mathcal{F}_{\ell,\ell} = \mathcal{F}_{\ell,\ell} \in \mathcal{F}_{\ell,\ell}$





Theorem: Fix
$$\sigma_1, \ldots, \sigma_L$$
 each ρ -Lipschitz with $\sigma_\ell(0) = 0$.
Let \mathscr{F}_L be the set of L -layer no-intercept nets, $f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)})$,
with $\|W_\ell^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.
 $\|M\|_{b,c} = \|(\|M_{\cdot 1}\|_b, \ldots, \|M_{\cdot d}\|_{b,c})$
Base case, $L = 0$:
 $\hat{\mathfrak{R}}_S(\{x \mapsto x_j : j \in [d]\}) \leq \frac{1}{n} (\max_j \|(x_{1,j}, \ldots, x_{n,j})\|_2) \sqrt{2\log d}$
 $= \frac{1}{n} \|X\|_{2,\infty} \sqrt{2\log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2\log d}$
Inductive step:
 $\hat{\mathfrak{R}}_S(\mathscr{F}_{\ell+1}) = \hat{\mathfrak{R}}_S \left(\{x \mapsto \sigma_{\ell+1} (\|W_{\ell+1}^{\mathsf{T}}\|_{1,\infty} g(x)) : g \in \operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \}$
 $\leq \widetilde{\rho B \mathfrak{R}}_S(\operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell))$

$$(\tilde{f}_{\ell})$$





$$\begin{aligned} & \text{Theorem: Fix } \sigma_1, \dots, \sigma_L \text{ each } \rho\text{-Lipschitz with } \sigma_\ell(0) = 0. \\ & \text{Let } \mathscr{F}_L \text{ be the set of } L\text{-layer no-intercept nets, } f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)}), \\ & \text{with } \|W_\ell^{\mathsf{T}}\|_{1,\infty} \leq B. \text{ Then } \hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2 \log d}. \\ & \|M\|_{b,c} = \left\| \left(\|M_{\cdot 1}\|_{b}, \dots, \|M_{\cdot d}\|_{b} \right) \right\|_c \\ & \hat{\mathfrak{R}}_S(\{x \mapsto x_j : j \in [d]\}) \leq \frac{1}{n} \left(\max_j \|(x_{1,j}, \dots, x_{n,j})\|_2 \right) \sqrt{2 \log d} \\ & = \frac{1}{n} \|X\|_{2,\infty} \sqrt{2 \log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2 \log d} \\ & \text{Inductive step:} \\ & \hat{\mathfrak{R}}_S(\mathscr{F}_{\ell+1}) = \hat{\mathfrak{R}}_S \left(\left\{ x \mapsto \sigma_{\ell+1} \left(\|W_{\ell+1}^{\mathsf{T}}\|_{1,\infty} g(x) \right) : g \in \operatorname{conv} \left(- \mathscr{F}_\ell \cup \mathscr{F}_\ell \right) \right\} \right) \\ & \leq \rho B \, \hat{\mathfrak{R}}_S(\operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell)) \quad \hat{\mathfrak{R}}_S(\operatorname{conv}(\mathscr{G})) = \hat{\mathfrak{R}}_S(\mathscr{G}) \end{aligned}$$





$$\begin{aligned} \text{Theorem: Fix } \sigma_1, \dots, \sigma_L \text{ each } \rho\text{-Lipschitz with } \sigma_\ell(0) &= 0. \\ \text{Let } \mathscr{F}_L \text{ be the set of } L\text{-layer no-intercept nets, } f^{(\ell)} &= \sigma_\ell(W_\ell f^{(\ell-1)}), \\ \text{with } \|W_\ell^{\mathsf{T}}\|_{1,\infty} &\leq B. \text{ Then } \hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}. \\ \|M\|_{b,c} &= \left\| \left(\|M_{\cdot1}\|_{b}, \dots, \|M_{\cdot d} \|\right) \right\|_{\mathcal{H}} \right\|_{\mathcal{H}} \\ \hat{\mathfrak{R}}_S(\{x \mapsto x_j : j \in [d]\}) &\leq \frac{1}{n} \left(\max_j \|(x_{1,j}, \dots, x_{n,j})\|_2 \right) \sqrt{2\log d} \\ &= \frac{1}{n} \|X\|_{2,\infty} \sqrt{2\log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2\log d} \\ \text{Inductive step:} \\ \hat{\mathfrak{R}}_S(\mathscr{F}_{\ell+1}) &= \hat{\mathfrak{R}}_S \left(\left\{ x \mapsto \sigma_{\ell+1} \left(\|W_{\ell+1}^{\mathsf{T}}\|_{1,\infty} g(x) \right) : g \in \operatorname{conv} \left(-\mathscr{F}_\ell \cup \mathscr{F}_\ell \right) \right\} \\ &\leq \rho B \, \hat{\mathfrak{R}}_S \left(\operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \right) \quad \hat{\mathfrak{R}}_S(\operatorname{conv}(\mathscr{G})) = \hat{\mathfrak{R}}_S(\mathscr{G}) \\ &\leq \rho B \, \hat{\mathfrak{R}}_S \left(-\mathscr{F}_\ell \cup \mathscr{F}_\ell \right) \end{aligned}$$





Theorem: Fix
$$\sigma_1, \ldots, \sigma_L$$
 each ρ -Lipschitz with $\sigma_\ell(0) = 0$.
Let \mathscr{F}_L be the set of L -layer no-intercept nets, $f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)})$,
with $\|W_\ell^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.
 $\|M\|_{b,c} = \|(\|M_{\cdot 1}\|_b, \ldots, \|M_{\cdot d}\|_b)$
Base case, $L = 0$:
 $\hat{\mathfrak{R}}_S(\{x \mapsto x_j : j \in [d]\}) \leq \frac{1}{n} (\max_j \|(x_{1,j}, \ldots, x_{n,j})\|_2) \sqrt{2\log d}$
 $= \frac{1}{n} \|X\|_{2,\infty} \sqrt{2\log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2\log d}$
Inductive step:
 $\hat{\mathfrak{R}}_S(\mathscr{F}_{\ell+1}) = \hat{\mathfrak{R}}_S\left(\left\{x \mapsto \sigma_{\ell+1} \left(\|W_{\ell+1}^{\mathsf{T}}\|_{1,\infty} g(x)\right) : g \in \operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell)\right\}$
 $\leq \rho B \, \hat{\mathfrak{R}}_S(\operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell)) \quad \hat{\mathfrak{R}}_S(\operatorname{conv}(\mathscr{G})) = \hat{\mathfrak{R}}_S(\mathscr{G})$
 $\leq \rho B \, \hat{\mathfrak{R}}_S(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \quad \hat{\mathfrak{R}}_S(A \cup B) \leq \hat{\mathfrak{R}}_S(A) + \hat{\mathfrak{R}}_S(B)$ if $0 \in A, 0$









Theorem: Fix
$$\sigma_1, ..., \sigma_L$$
 each ρ -Lipschitz with $\sigma_\ell(0) = 0$.
Let \mathscr{F}_L be the set of L -layer no-intercept nets, $f^{(\ell)} = \sigma_\ell(W_\ell f^{(\ell-1)})$,
with $\|W_\ell^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.
 $\|M\|_{b,c} = \|(\|M_{\cdot 1}\|_b, ..., \|M_{\cdot d}\|_b)$
Base case, $L = 0$:
 $\hat{\mathfrak{R}}_S(\{x \mapsto x_j : j \in [d]\}) \leq \frac{1}{n} (\max_j \|(x_{1,j}, ..., x_{n,j})\|_2) \sqrt{2\log d}$
 $= \frac{1}{n} \|X\|_{2,\infty} \sqrt{2\log d} = \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^0 \sqrt{2\log d}$
Inductive step:
 $\hat{\mathfrak{R}}_S(\mathscr{F}_{\ell+1}) = \hat{\mathfrak{R}}_S \left(\left\{ x \mapsto \sigma_{\ell+1} \left(\|W_{\ell+1}^{\mathsf{T}}\|_{1,\infty} g(x) \right) : g \in \operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \right\}$
 $\leq \rho B \, \hat{\mathfrak{R}}_S(\operatorname{conv}(-\mathscr{F}_\ell \cup \mathscr{F}_\ell)) \quad \hat{\mathfrak{R}}_S(\operatorname{conv}(\mathscr{G})) = \hat{\mathfrak{R}}_S(\mathscr{G})$
 $\leq \rho B \, \hat{\mathfrak{R}}_S(-\mathscr{F}_\ell \cup \mathscr{F}_\ell) \quad \hat{\mathfrak{R}}_S(A \cup B) \leq \hat{\mathfrak{R}}_S(A) + \hat{\mathfrak{R}}_S(B) \text{ if } 0 \in A, 0$
 $\leq 2\rho B \, \hat{\mathfrak{R}}_S(\mathscr{F}_\ell)$









$\hat{\mathfrak{R}}_{S}(\operatorname{conv}(\mathscr{G})) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sup_{g_{1}, \dots, g_{k} \in \mathscr{G}} \left\langle \sigma, \sum_{j=1}^{k} \alpha_{j}(g_{j})_{S} \right\rangle$

$\hat{\mathfrak{R}}_{S}(\operatorname{conv}(\mathscr{G})) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sup_{g_{1}, \dots, g_{k} \in \mathscr{G}} \left\langle \sigma, \sum_{j=1}^{k} \alpha_{j}(g_{j})_{S} \right\rangle$ $= \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sum_{j=1}^{k} \alpha_{j} \sup_{g_{j} \in \mathscr{G}} \left\langle \varepsilon, (g_{j})_{S} \right\rangle$

 $\hat{\mathfrak{R}}_{S}(\operatorname{conv}(\mathscr{G})) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sup_{g_{1}, \dots, g_{k} \in \mathscr{G}} \left\langle \sigma, \sum_{j=1}^{k} \alpha_{j}(g_{j})_{S} \right\rangle$ $= \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sum_{j=1}^{k} \alpha_{j} \sup_{g_{j} \in \mathscr{G}} \left\langle \varepsilon, (g_{j})_{S} \right\rangle$ $= \frac{1}{n} \mathbb{E}_{\sigma} \left(\sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sum_{j=1}^{k} \alpha_{j} \right) \sup_{g \in \mathcal{G}} \langle \varepsilon, g_{S} \rangle$

 $\hat{\mathfrak{R}}_{S}(\operatorname{conv}(\mathscr{G})) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sup_{g_{1}, \dots, g_{k} \in \mathscr{G}} \left\langle \sigma, \sum_{j=1}^{k} \alpha_{j}(g_{j})_{S} \right\rangle$ $= \frac{1}{n} \mathbb{E}_{\sigma} \sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sum_{j=1}^{k} \alpha_{j} \sup_{g_{j} \in \mathscr{G}} \left\langle \varepsilon, (g_{j})_{S} \right\rangle$ $= \frac{1}{n} \mathbb{E}_{\sigma} \left(\sup_{k \ge 1} \sup_{\alpha \in \Delta_{k}} \sum_{j=1}^{k} \alpha_{j} \right) \sup_{g \in \mathscr{G}} \langle \varepsilon, g_{S} \rangle$ $= \Re_{S}(\mathscr{G})$

$\hat{\mathfrak{R}}_{S}(\mathscr{G}\cup\mathscr{H}) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{g \in (\mathscr{G}\cup\mathscr{H})} \left\langle \sigma, g_{S} \right\rangle$

 $\hat{\mathfrak{R}}_{S}(\mathscr{G}\cup\mathscr{H}) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{g \in (\mathscr{G}\cup\mathscr{H})} \left\langle \sigma, g_{S} \right\rangle$

$\leq \frac{1}{n} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathscr{G}} \left\langle \sigma, g_{S} \right\rangle + \sup_{g \in \mathscr{H}} \left\langle \sigma, g_{S} \right\rangle \Big] \quad \text{if } 0 \in \mathscr{G}, 0 \in \mathscr{H}$

 $\hat{\mathfrak{R}}_{S}(\mathscr{G}\cup\mathscr{H}) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{g \in (\mathscr{G}\cup\mathscr{H})} \left\langle \sigma, g_{S} \right\rangle$ $= \hat{\Re}_{\mathcal{S}}(\mathscr{G}) + \hat{\Re}_{\mathcal{S}}(\mathscr{H})$

$\leq \frac{1}{n} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathscr{G}} \left\langle \sigma, g_{S} \right\rangle + \sup_{g \in \mathscr{H}} \left\langle \sigma, g_{S} \right\rangle \Big] \quad \text{if } 0 \in \mathscr{G}, 0 \in \mathscr{H}$

 $\hat{\mathfrak{R}}_{S}(\mathcal{G}\cup\mathcal{H}) = \frac{1}{n}\mathbb{E}_{\sigma}\sup_{g\in(\mathcal{G}\cup\mathcal{H})}\left\langle\sigma,g_{S}\right\rangle$ $= \hat{\Re}_{\mathcal{S}}(\mathscr{G}) + \hat{\Re}_{\mathcal{S}}(\mathscr{H})$

$\leq \frac{1}{n} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathscr{G}} \left\langle \sigma, g_{S} \right\rangle + \sup_{g \in \mathscr{H}} \left\langle \sigma, g_{S} \right\rangle \Big] \quad \text{if } 0 \in \mathscr{G}, 0 \in \mathscr{H}$ or if both sets are **symmetric**: for all $g \in \mathcal{G}$, also have $-g \in \mathcal{G}$





 $\hat{\mathfrak{R}}_{S}(\mathcal{G}\cup\mathcal{H}) = \frac{1}{n} \mathbb{E}_{\sigma} \sup_{g \in (\mathcal{G}\cup\mathcal{H})} \left\langle \sigma, g_{S} \right\rangle$ $= \hat{\Re}_{\mathcal{S}}(\mathscr{G}) + \hat{\Re}_{\mathcal{S}}(\mathscr{H})$

$\leq \frac{1}{n} \mathbb{E}_{\sigma} \Big[\sup_{g \in \mathscr{G}} \left\langle \sigma, g_{S} \right\rangle + \sup_{g \in \mathscr{H}} \left\langle \sigma, g_{S} \right\rangle \Big] \quad \text{if } 0 \in \mathscr{G}, 0 \in \mathscr{H}$ or if both sets are **symmetric**: for all $g \in \mathcal{G}$, also have $-g \in \mathcal{G}$

...or if we otherwise know that $\sup \langle \sigma, g_S \rangle \ge 0$, $\sup \langle \sigma, g_S \rangle \ge 0$ $g \in \mathcal{G}$ $g \in \mathcal{H}$ for any assignment of σ







Rademacher of deep nets

Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{C}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^L \sqrt{2\log d}$.

 $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|$

 $\mathcal{O}_{\mathcal{U}}(W_{\mathcal{U}} \mathcal{O}_{\mathcal{U}-1}(\dots \mathcal{O}_{\mathcal{U}}(W_{\mathcal{U}} \mathbf{X})))))))))$ $\mathcal{O}_{\mathcal{U}}(W_{\mathcal{U}} \mathcal{O}_{\mathcal{U}-1}(\dots \mathcal{O}_{\mathcal{U}}(\mathbf{X})))))))))(k)$ Take $\mathcal{O}_{\mathcal{U}-1}(\dots \mathcal{O}_{\mathcal{U}}(\mathbf{X}))))(k)$ $\mathcal{O}_{\mathcal{U}}(W_{\mathcal{U}} \mathcal{O}_{\mathcal{U}-1}(\mathbf{X}))))(k)$



Rademacher of deep nets

Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\mathcal{C}}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$.

Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\mathscr{C}}\|_F \leq B$. Then $\hat{\mathfrak{R}}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_F B^L \left(1 + \sqrt{2L\log 2}\right)$.

- $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{a}$
- **Theorem:** Fix $\sigma_1, \ldots, \sigma_L$ each 1-Lipschitz, positive homogenous ($\sigma_{\ell}(ax) = a\sigma_{\ell}(x)$ for a > 0).
- (More complicated proof: Golowich/Rakhlin/Shamir, COLT 2018 / Telgarsky's 14.2.)



Rademacher of deep nets

Theorem: Fix $\sigma_1, \ldots, \sigma_L$ each ρ -Lipschitz with $\sigma_{\ell}(0) = 0$. Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}^{\mathsf{T}}\|_{1,\infty} \leq B$. Then $\hat{\mathfrak{R}}_{n}(\mathscr{F}) \leq \frac{1}{n} \|X\|_{2,\infty} (2\rho B)^{L} \sqrt{2\log d}$.

Let \mathscr{F}_L be the set of L-layer no-intercept nets, $f^{(\ell)} = \sigma_{\ell}(W_{\ell}f^{(\ell-1)})$, with $\|W_{\ell}\|_F \leq B$. Then $\hat{\Re}_n(\mathscr{F}) \leq \frac{1}{n} \|X\|_F B^L \left(1 + \sqrt{2L\log 2}\right)$.

Can get a slightly better rate via covering numbers: see <u>Telgarsky's section 16.2</u>.

- $\|M\|_{b,c} = \|(\|M_{.1}\|_{b}, \dots, \|M_{.d}\|_{b})\|_{a}$
- **Theorem:** Fix $\sigma_1, \ldots, \sigma_L$ each 1-Lipschitz, positive homogenous ($\sigma_{\ell}(ax) = a\sigma_{\ell}(x)$ for a > 0).
- (More complicated proof: Golowich/Rakhlin/Shamir, COLT 2018 / Telgarsky's 14.2.)

e de je



So, does this solve it?

 Experiment by <u>Dziugaite/Roy (2017)</u>: training a small network on MNIST (0-4 vs 5-9), plotting a Rademacher-based margin bound using a different (but similarly[?] tight) upper bound on the Rademacher complexity



What's left?

- We've shown some bounds on approximation and generalization (each with significant limitations)
- If we could run ERM, this combination could be enough ...but ERM is NP-hard even for square loss even with one ReLU
- What's the optimization error for SGD/similar?

• Neural nets are not convex

- Neural nets are not convex
- Even deep linear networks are not convex

- Neural nets are not convex
- Even deep linear networks are not convex

• But we do know that SGD converges to a *critical point* under fairly mild conditions



- Neural nets are not convex
- Even deep linear networks are not convex
- But we do know that SGD converges to a critical point under fairly mild conditions • e.g.: if $f \ge f^{\text{inf}}$ is differentiable and β -smooth, and there are A, B, C s.t. for all x, $\mathbb{E}[\|\hat{g}(x)\|^2] \le 2A(f(x) - f^{\inf}) + B\|\nabla f(X)\|^2 + C$,
 - then the best iterate has $\mathbb{E}\left[\|\nabla f(x)\|^2\right] \leq \varepsilon^2$ in $\mathcal{O}(\varepsilon^{-4})$ steps (Khaled/Richtárik 2020)



- Neural nets are not convex
- Even deep linear networks are not convex
- But we do know that SGD converges to a critical point under fairly mild conditions • e.g.: if $f \ge f^{inf}$ is differentiable and β -smooth, and there are A, B, C s.t. for all x, $\mathbb{E}[\|\hat{g}(x)\|^2] \le 2A(f(x) - f^{\inf}) + B\|\nabla f(X)\|^2 + C$, then the best iterate has $\mathbb{E}\left[\|\nabla f(x)\|^2\right] \leq \varepsilon^2$ in $\mathcal{O}(\varepsilon^{-4})$ steps (Khaled/Richtárik 2020)
- In deep linear nets, local minima are global minima (Kawaguchi 2016, Laurent/von Brecht 2019)



- Neural nets are not convex
- Even deep linear networks are not convex
- But we do know that SGD converges to a critical point under fairly mild conditions • e.g.: if $f \ge f^{\text{inf}}$ is differentiable and β -smooth, and
 - there are A, B, C s.t. for all x, $\mathbb{E}[\|\hat{g}(x)\|^2] \le 2A(f(x) f^{\inf}) + B\|\nabla f(X)\|^2 + C$, then the best iterate has $\mathbb{E}\left[\|\nabla f(x)\|^2\right] \leq \varepsilon^2$ in $\mathcal{O}(\varepsilon^{-4})$ steps (Khaled/Richtárik 2020)
- In deep linear nets, local minima are global minima (Kawaguchi 2016, Laurent/von Brecht 2019) • ...but there are saddle points, including "bad" ones where $\lambda_{\min}(\nabla^2 f) = 0$


Nonconvex optimization

- Neural nets are not convex
- Even deep linear networks are not convex
- But we do know that SGD converges to a *critical point* under fairly mild conditions • e.g.: if $f \ge f^{inf}$ is differentiable and β -smooth, and
 - there are A, B, C s.t. for all x, $\mathbb{E}[\|\hat{g}(x)\|^2] \le 2A(f(x) f^{\inf}) + B\|\nabla f(X)\|^2 + C$, then the best iterate has $\mathbb{E}\left[\|\nabla f(x)\|^2\right] \leq \varepsilon^2$ in $\mathcal{O}(\varepsilon^{-4})$ steps (Khaled/Richtárik 2020)
- In deep linear nets, local minima are global minima (Kawaguchi 2016, Laurent/von Brecht 2019) • ...but there are saddle points, including "bad" ones where $\lambda_{\min}(\nabla^2 f) = 0$...but gradient descent doesn't get stuck, under some conditions (Arora et al. 2019)





Nonconvex optimization

- Neural nets are not convex
- Even deep linear networks are not convex
- But we do know that SGD converges to a *critical point* under fairly mild conditions • e.g.: if $f \ge f^{inf}$ is differentiable and β -smooth, and
 - there are A, B, C s.t. for all x, $\mathbb{E}[\|\hat{g}(x)\|^2] \le 2A(f(x) f^{\inf}) + B\|\nabla f(X)\|^2 + C$, then the best iterate has $\mathbb{E}\left[\|\nabla f(x)\|^2\right] \leq \varepsilon^2$ in $\mathcal{O}(\varepsilon^{-4})$ steps (Khaled/Richtárik 2020)
- In deep linear nets, local minima are global minima (Kawaguchi 2016, Laurent/von Brecht 2019) • ...but there are saddle points, including "bad" ones where $\lambda_{\min}(\nabla^2 f) = 0$...but gradient descent doesn't get stuck, under some conditions (Arora et al. 2019)
- **Next time:** optimization "acts convex" in the neural tangent kernel regime 17



