

Inductive Bias and Optimization in Deep Learning

Nati Srebro (TTIC)

Based on work with **Behnam Neyshabur** (TTIC→Google), **Suriya Gunasekar** (TTIC→MSR), Ryota Tomioka (TTIC→MSR), Srinadh Bhojanapalli (TTIC→Google), **Blake Woodworth**, Pedro Savarese, David McAllester (TTIC), Greg Ongie, Becca Willett (Chicago), **Daniel Soudry**, Elad Hoffer, Mor Shpigel, Itay Sofer (Technion), Ashia Wilson, Becca Roelofs, Mitchel Stern, Ben Recht (Berkeley), Russ Salakhutdinov (CMU), **Jason Lee**, Zhiyuan Li (Princeton), Yann LaCun (NYU/Facebook)

Feed Forward Neural Networks

- Fix architecture (connection graph $G(V, E)$, transfer σ)

$$\mathcal{H}_{G(V, E), \sigma} = \{ f_{\mathbf{w}}(x) = \text{output of net with weights } \mathbf{w} \}$$

- Capacity / Generalization ability / Sample Complexity

- $\tilde{O}(|E|)$ (number of edges, i.e. number of weights)
(with threshold σ , or with RELU and finite precision; RELU with inf precision: $\tilde{\Theta}(|E| \cdot \text{depth})$)



- Expressive Power / Approximation

- Any continuous function with huge network
- Lots of interesting things naturally with small networks
- **Any time T computable function with network of size $\tilde{O}(T)$**



- Computation / Optimization

- NP-hard to find weights even with 2 hidden units
- Even if function exactly representable with single hidden layer with $\Theta(\log d)$ units, even with no noise, and even if we allow a much larger network when learning: no poly-time algorithm always works

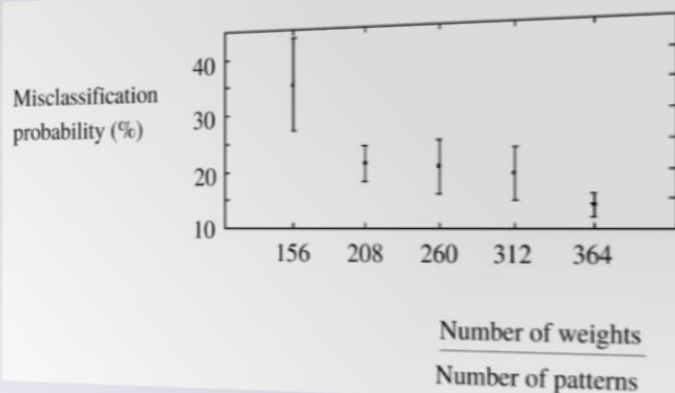


[Kearns Valiant 94; Klivans Sherstov 06; Daniely Linial Shalev-Shwartz '14]

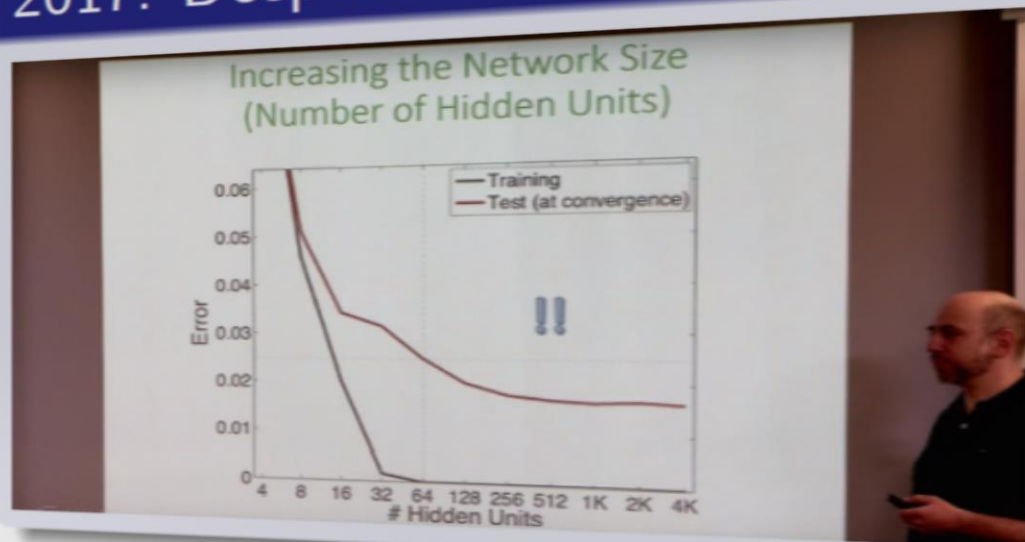
- **Magic property of reality that makes local search “work”**

Generalization: Margins and Size of Parameters

1996: Sigmoid networks



2017: Deep ReLU networks

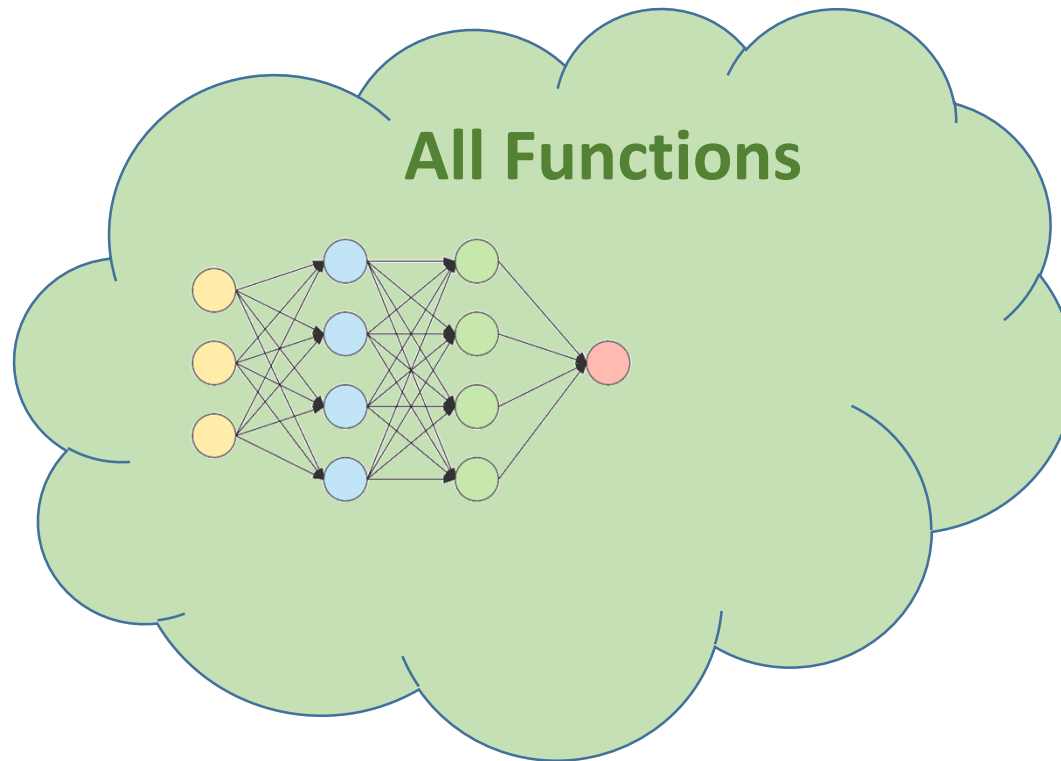


Qualitative behavior explained
small weights theorem.

simons.berkeley.edu

- How to measure the complexity of a ReLU network?

Need to understand optimization alg. not just as reaching *some* (global) optimum, but as reaching a *specific* optimum



Different optimization algorithm

→ Different bias in optimum reached

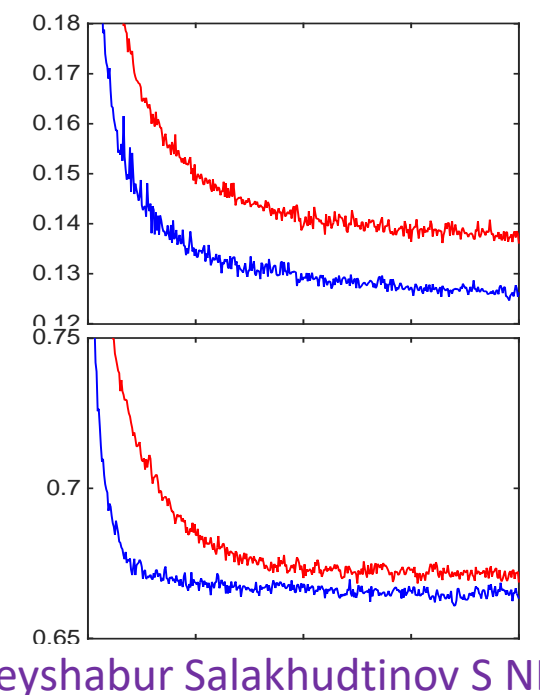
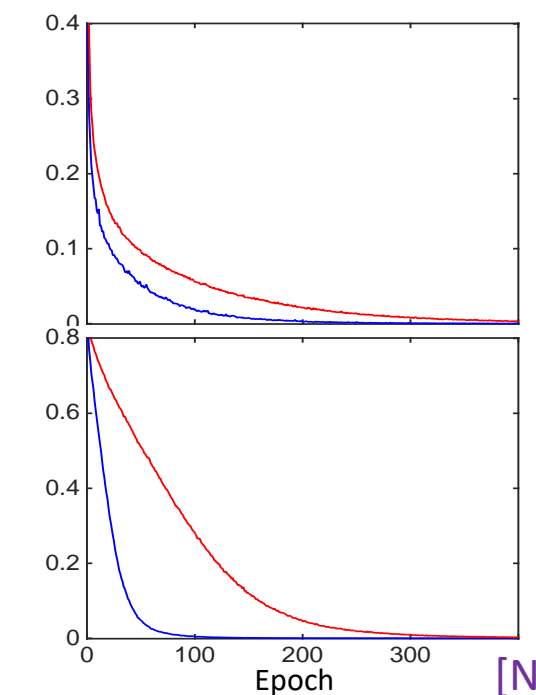
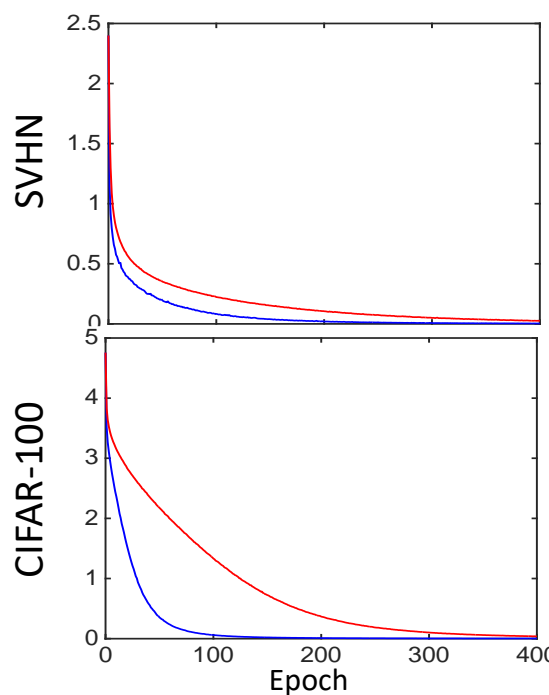
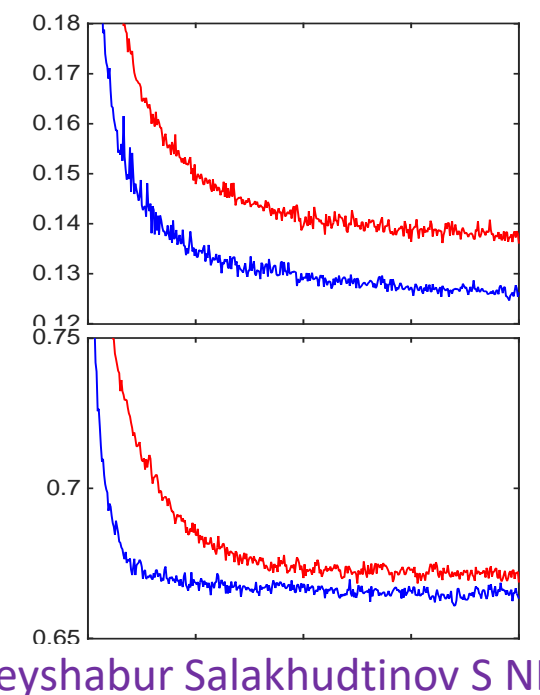
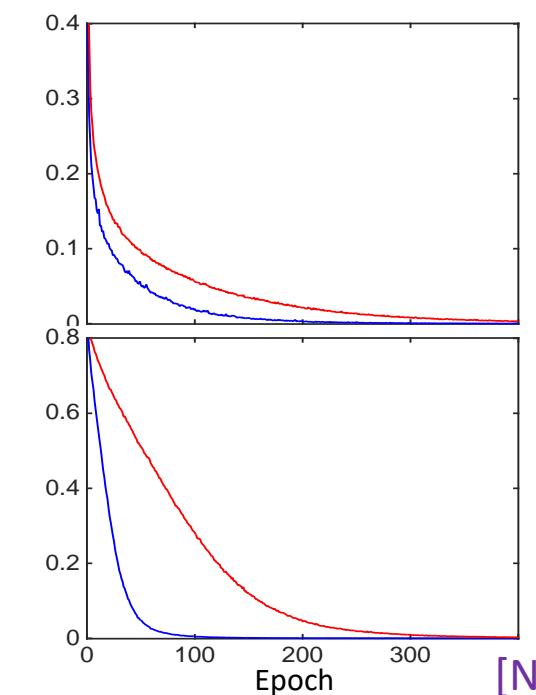
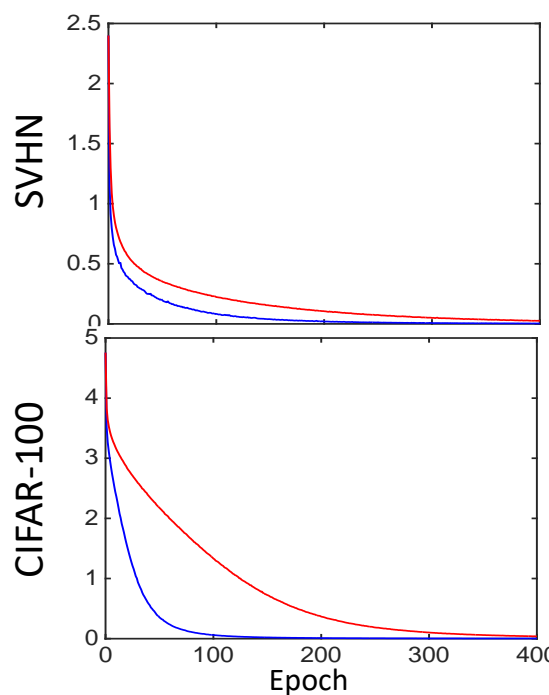
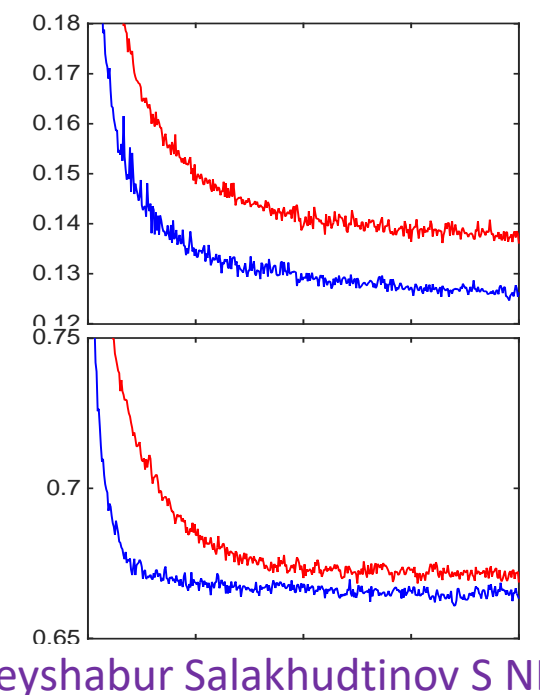
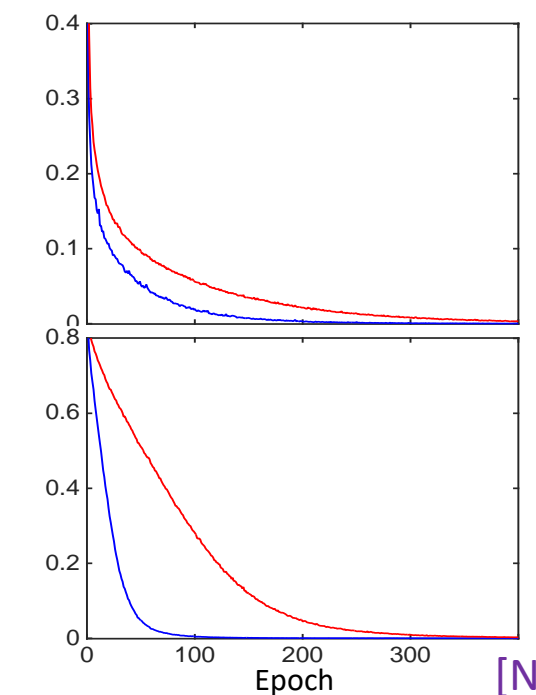
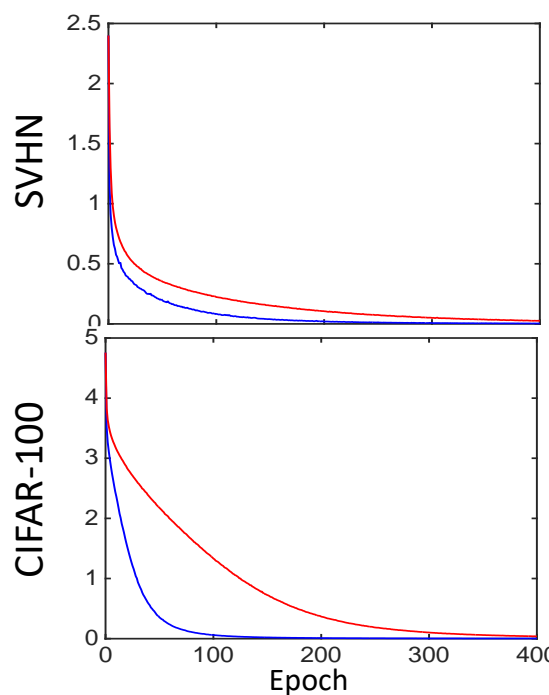
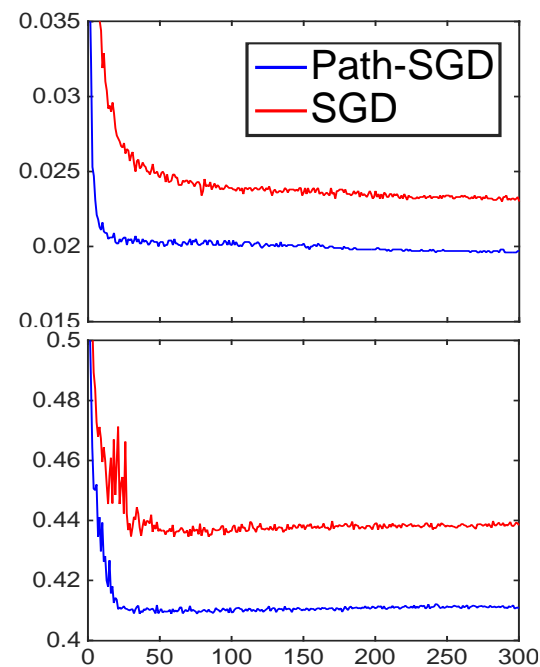
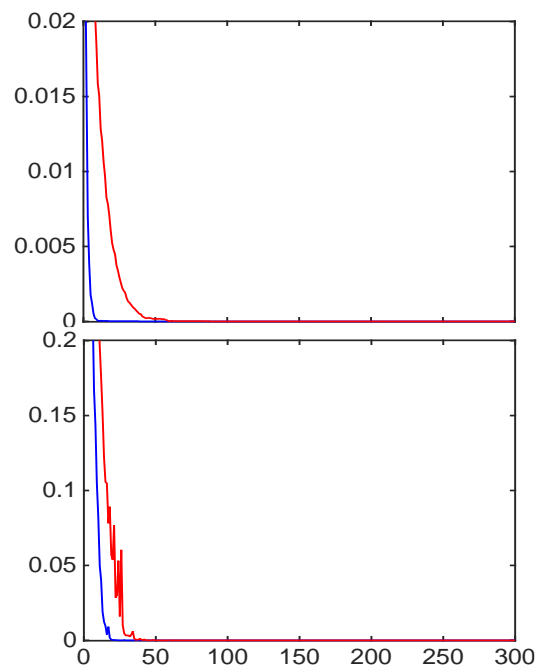
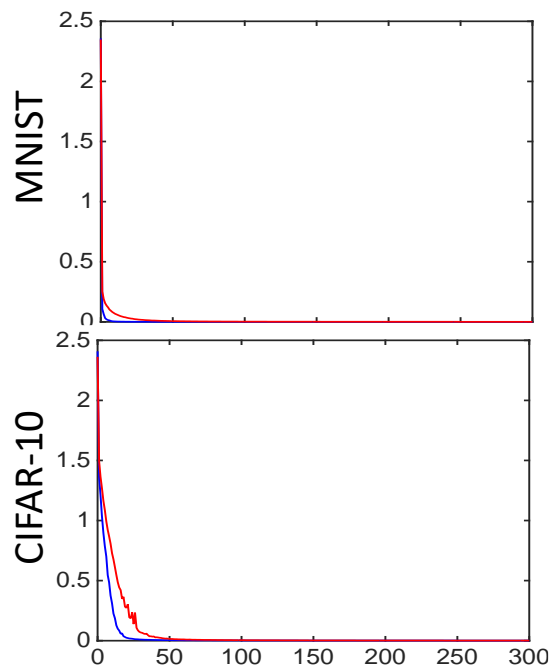
→ Different Inductive bias

→ Different generalization properties

Cross-Entropy

0/1 Training Error

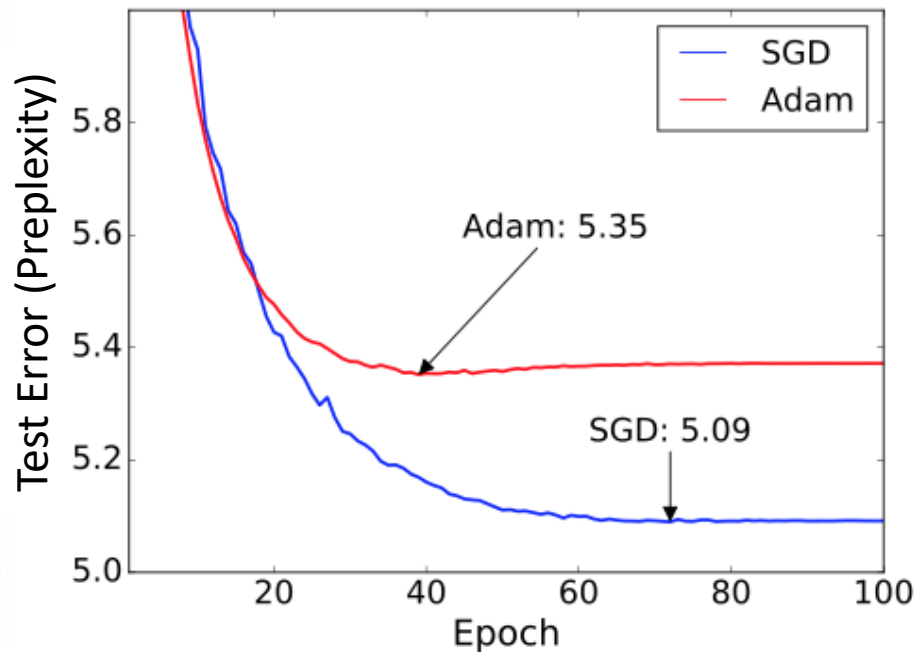
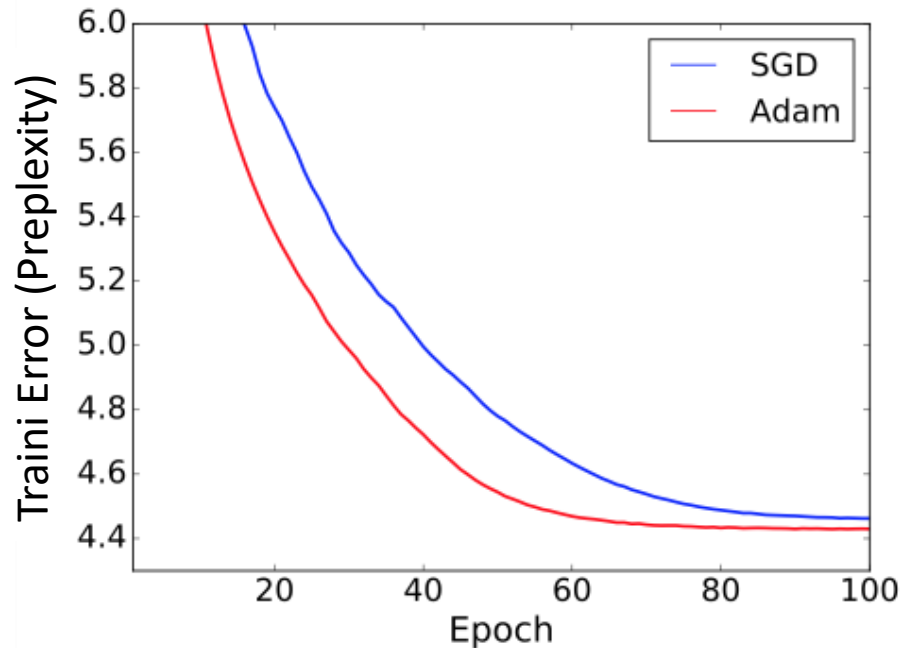
0/1 Test Error



With Dropout

[Neishabur Salakhudtinov S NIPS'15]

SGD vs ADAM



Results on Penn Treebank using 3-layer LSTM

[Wilson Roelofs Stern S Recht, "The Marginal Value of Adaptive Gradient Methods in Machine Learning", NIPS'17]

The Deep Recurrent Residual Boosting Machine

Joe Flow, DeepFace Labs

Section 1: Introduction

We suggest a new amazing architecture and loss function that is great for learning. All you have to do to learn is fit the model on your training data

Section 2: Learning Contribution: our model

The model class h_w is amazing. **Our learning method is:**

$$\arg \min_w \frac{1}{m} \sum_{i=1}^m \text{loss}(h_w(x); y) \quad (*)$$

Section 3: Optimization

This is how we solve the optimization problem (*): [...]

Section 4: Experiments

It works!

Unconstrained Matrix Completion

The diagram illustrates the matrix completion problem. On the left is an 8x8 grid representing a matrix with some observed values (blue cells) and missing values (white cells). The grid is as follows:

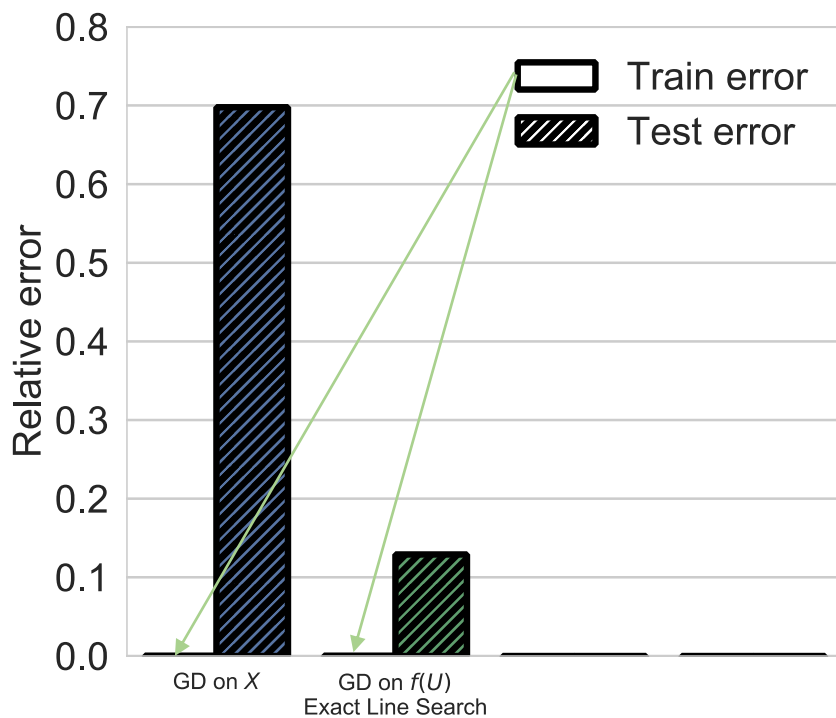
	2		4	5		1	4	2
3	1			2	2		5	4
4		2		4	1		3	1
3			3	4		2		4
2	3		1		4	3		2
	2	2					4	5
	2		4	1	4		2	3
1		3		1	1			4
	4		2	2		5	3	1

This matrix is approximated by a matrix X (red), which is equal to the product of matrix U (blue) and matrix V^T (blue).

$$\min_{X \in \mathbb{R}^{n \times n}} \| \text{observed}(X) - y \|_2^2 \equiv \min_{U, V \in \mathbb{R}^{n \times n}} \| \text{observed}(UV^T) - y \|_2^2$$

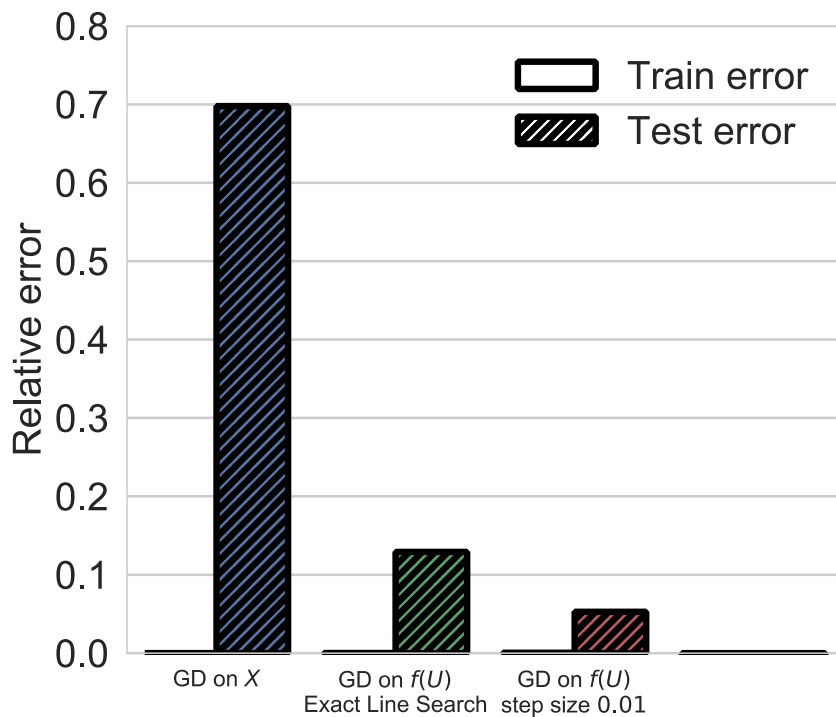
- Underdetermined non-sensical problem, lots of useless global min
- Since U, V full dim, no constraint on X , all the same non-sense global min

What happens when we optimize by gradient descent on U, V ?



$n = 50, m = 300, A_i$ iid Gaussian, X^* rank-2 ground truth
 $y = \mathcal{A}(X^*) + \mathcal{N}(0, 10^{-3}), y_{\text{test}} = \mathcal{A}_{\text{test}}(X^*) + \mathcal{N}(0, 10^{-3})$

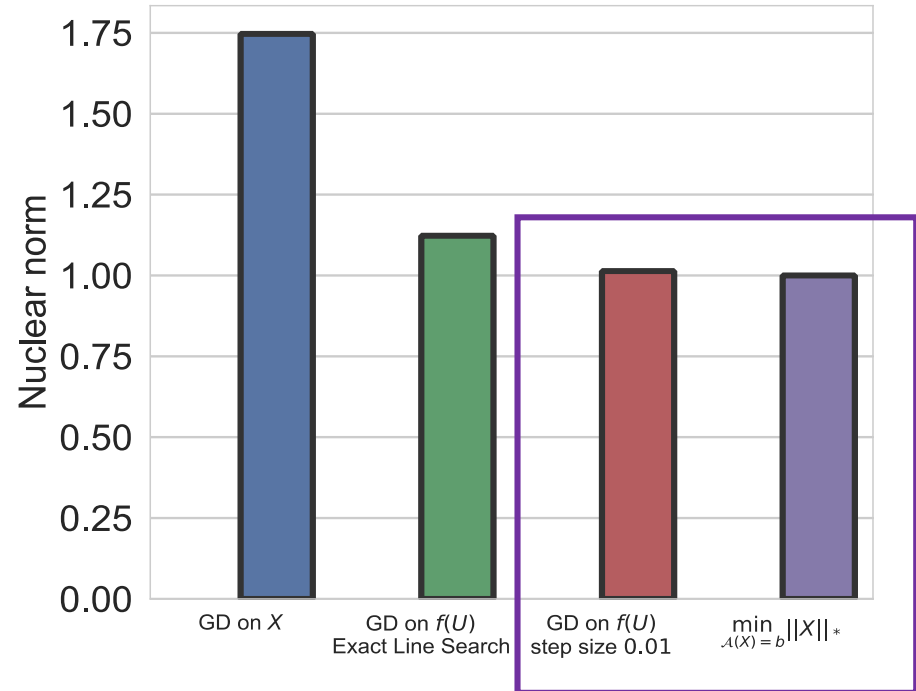
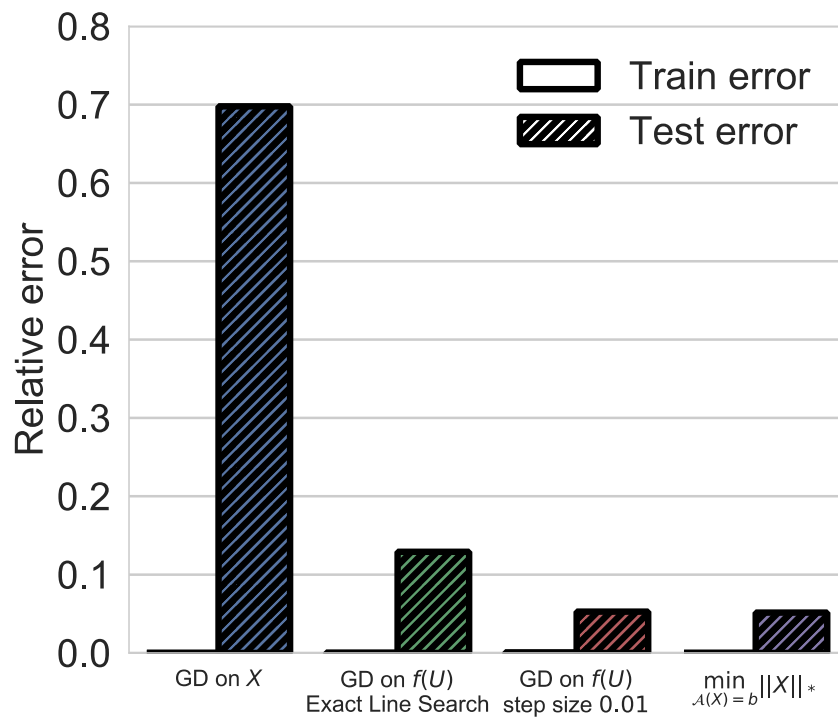
Gradient descent on $f(U, V)$ gets to “good” global minima



$n = 50$, $m = 300$, A_i iid Gaussian, X^* rank-2 ground truth
 $y = \mathcal{A}(X^*) + \mathcal{N}(0, 10^{-3})$, $y_{\text{test}} = \mathcal{A}_{\text{test}}(X^*) + \mathcal{N}(0, 10^{-3})$

Gradient descent on $f(U, V)$ gets to “good” global minima

Gradient descent on $f(U, V)$ generalizes better with smaller step size



Grad Descent on U, V with inf. small stepsize and initialization

→ **min nuclear norm solution**

$$\arg \min \|X\|_* \text{ s.t. } \text{obs}(X) = y$$

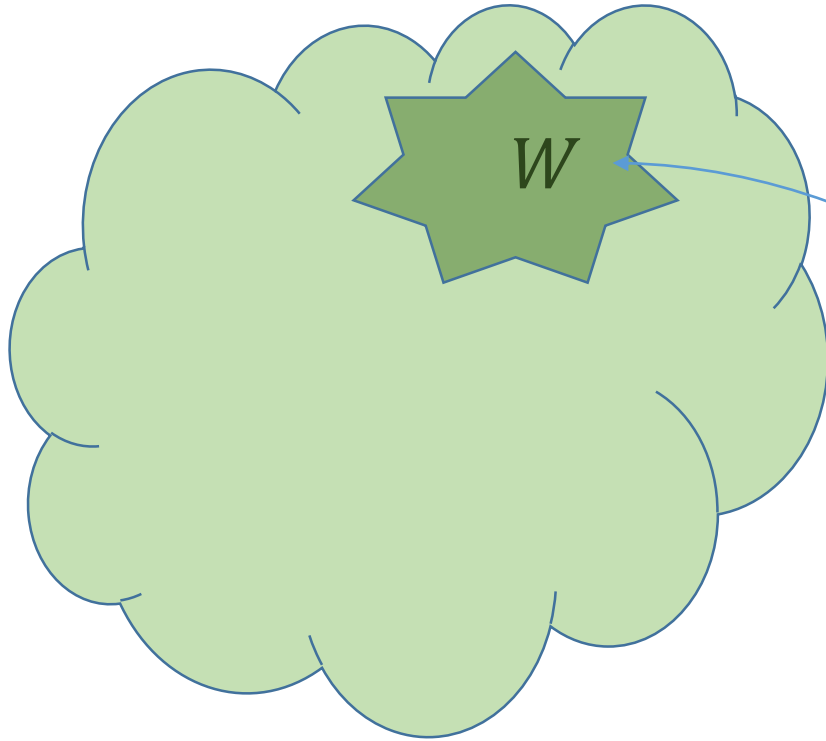
(exact and rigorous only under additional conditions!)

→ good generalization if Y (aprox) low rank

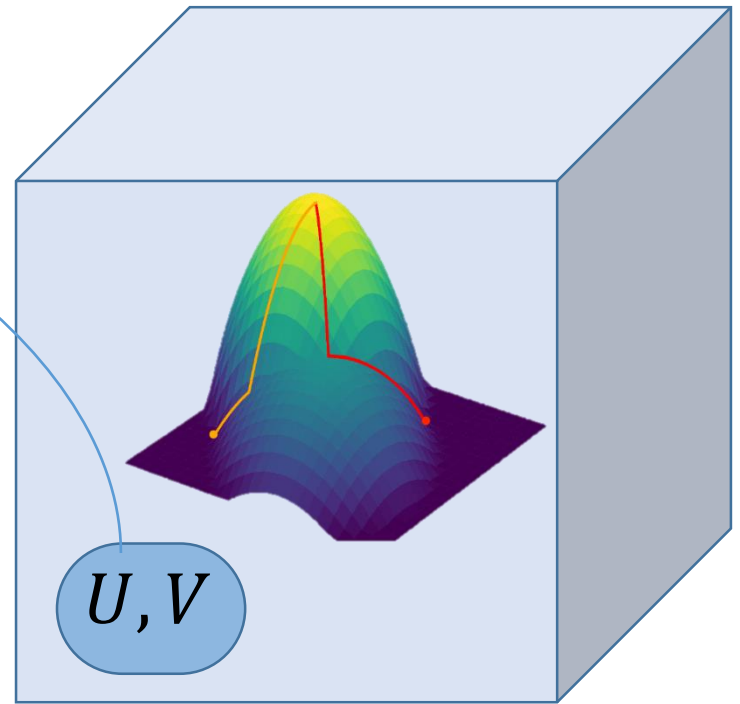
[Gunasekar Woodworth Bhojanapalli Neyshabur S 2017]

[Yuanzhi Li, Hongyang Zhang, Tengyu Ma 2018][Sanjeev Arora, Nadav Cohen, Wei Hu, Yuping Luo 2019]

Predictor Space



Parameter Space



Optimization Geometry and hence Inductive Bias effected by:

- Geometry of local search in parameter space
- Choice of parameterization

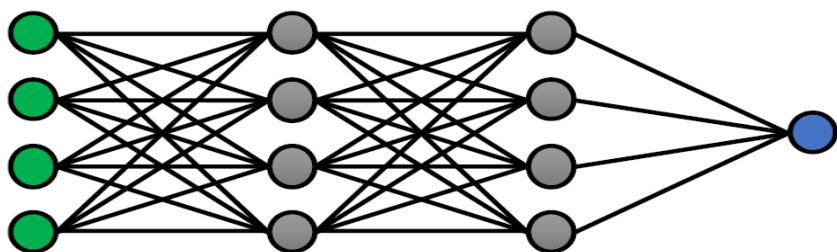
- **Matrix completion** (also: reconstruction from linear measurements)

- $W = UV$ is over-parametrization of all matrices $W \in \mathbb{R}^{n \times m}$
- GD on $U, V \rightarrow$ implicitly minimize $\|W\|_*$

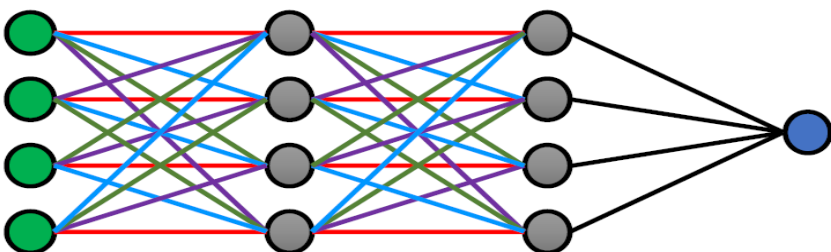
- **Linear Convolutional Network:**

- Complex over-parametrization of linear predictors β
- GD on weight \rightarrow implicitly minimize $\|DFT(\beta)\|_p$ for $p = \frac{2}{depth}$.
(sparsity in frequency domain)

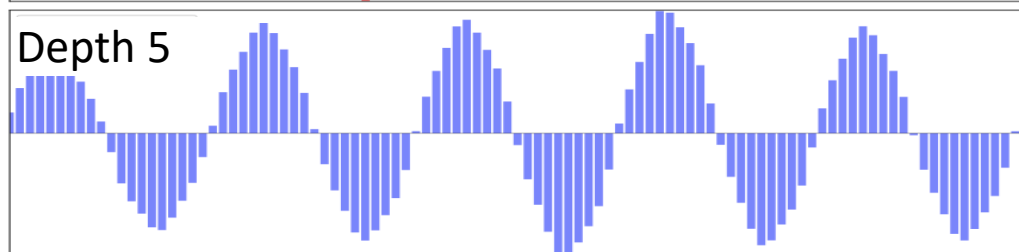
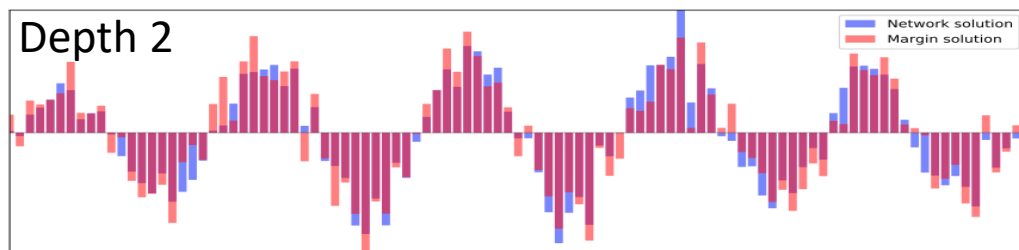
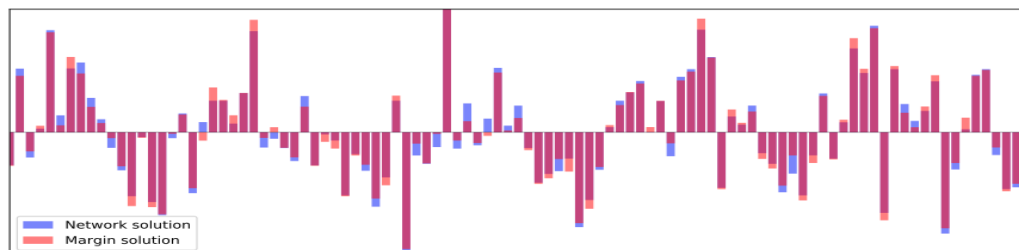
[Gunasekar Lee Soudry S 2018]



$$\min \|\beta\|_2 \text{ s.t. } \forall_i y_i \langle \beta, x_i \rangle \geq 1$$



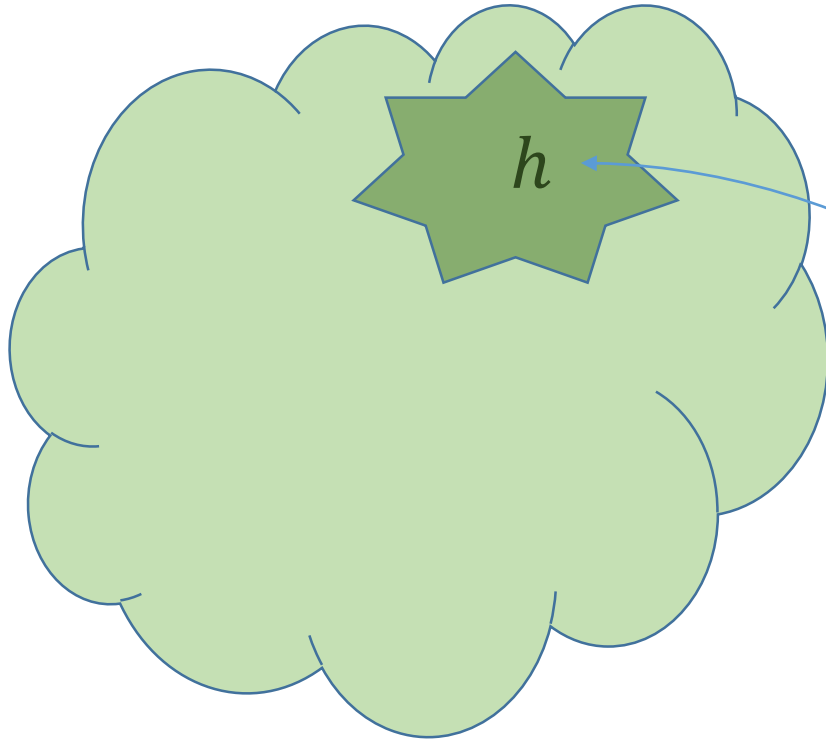
$$\min \|DFT(\beta)\|_{2/L} \text{ s.t. } \forall_i y_i \langle \beta, x_i \rangle \geq 1$$



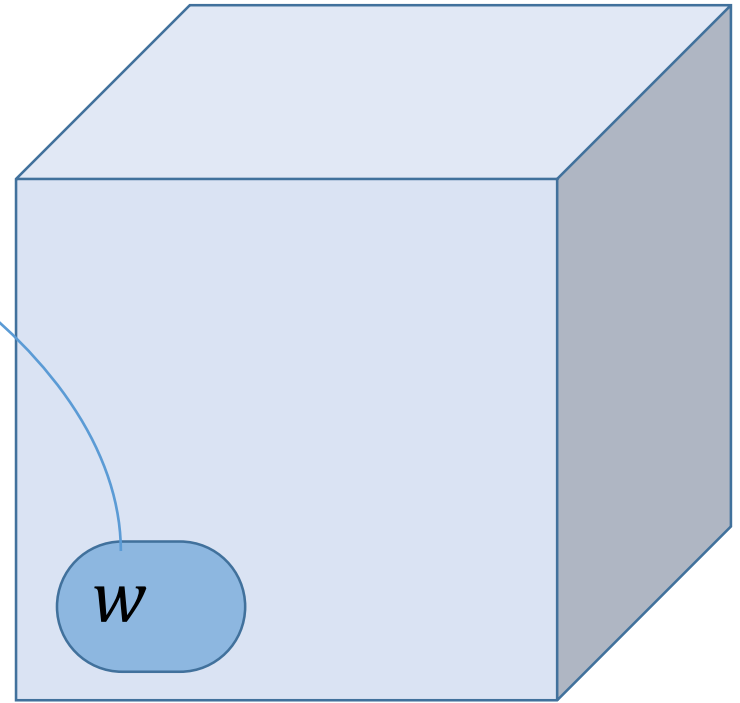
- **Matrix completion** (also: reconstruction from linear measurements)
 - $W = UV$ is over-parametrization of all matrices $W \in \mathbb{R}^{n \times m}$
 - GD on $U, V \rightarrow$ implicitly minimize $\|W\|_*$
- **Linear Convolutional Network:**
 - Complex over-parametrization of linear predictors β
 - GD on weight \rightarrow implicitly minimize $\|DFT(\beta)\|_p$ for $p = \frac{2}{depth}$.
(sparsity in frequency domain) [Gunasekar Lee Soudry S 2018]
- **Infinite Width ReLU Net:**
 - Parametrization of essentially all functions $h: \mathbb{R}^d \rightarrow \mathbb{R}$
 - Weight decay \rightarrow for $d = 1$, implicitly minimize

$$\max\left(\int |h''| dx, |h'(-\infty) + h'(+\infty)|\right)$$
 [Savarese Evron Soudry S 2019]
 - For $d > 1$, implicitly minimize $\int |\partial_b^{d+1} Radon(h)|$
(roughly speaking; need to define more carefully to handle non-smoothness, extra correction term for linear part) [Ongie Willett Soudry S 2019]

All Functions



Parameter Space



Optimization Geometry and hence Inductive Bias effected by:

- Geometry of local search in parameter space
- Choice of parameterization

Doesn't it all boil down
to the NTK?

Is it all just a Kernel?

$$f(\mathbf{w}, \mathbf{x}) \approx f(w^{(0)}, \mathbf{x}) + \langle \mathbf{w}, \boldsymbol{\phi}_{w^{(0)}}(\mathbf{x}) \rangle$$

focus on “unbiased
initialization”: $f(w^{(0)}, \mathbf{x}) = 0$

$$\boldsymbol{\phi}_{\mathbf{w}}(\mathbf{x}) = \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$$

Corresponding to a kernelized linear model with kernel:

$$K_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') = \langle \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}), \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}') \rangle$$

Kernel regime: 1st order approx about $w^{(0)}$ remains valid throughout optimization

$$K_{\mathbf{w}^{(t)}} \approx K_{\mathbf{w}^{(0)}} = K_0$$

➔ GD on squared loss converges to $\arg \min \| \mathbf{h} \|_{K_0} \text{ s.t. } \mathbf{h}(\mathbf{x}_i) = \mathbf{y}_i$

Kernel Regime and Scale of Init

- For D -homogenous model, $f(cw, x) = c^D f(w, x)$, consider gradient flow with:

$$\dot{w}_\alpha = -\nabla L_S(w) \quad \text{and} \quad w_\alpha(0) = \alpha w_0 \quad \text{with unbiased } f(w_0, x) = 0$$

We are interested in $w_\alpha(\infty) = \lim_{t \rightarrow \infty} w_\alpha(t)$

- For squared loss, under some conditions [Chizat and Bach 18]:

$$\lim_{\alpha \rightarrow \infty} \sup_t \left\| w_\alpha \left(\frac{1}{\alpha^{D-1}} t \right) - w_K(t) \right\| = 0$$

Gradient flow of linear least squares
w.r.t tangent kernel K_0 at initialization

and so $f(w_\alpha(\infty), x) \xrightarrow{\alpha \rightarrow \infty} \hat{h}_K(x)$ where $\hat{h}_K = \arg \min \|h\|_{K_0} \text{ s.t. } h(x_i) = y_i$

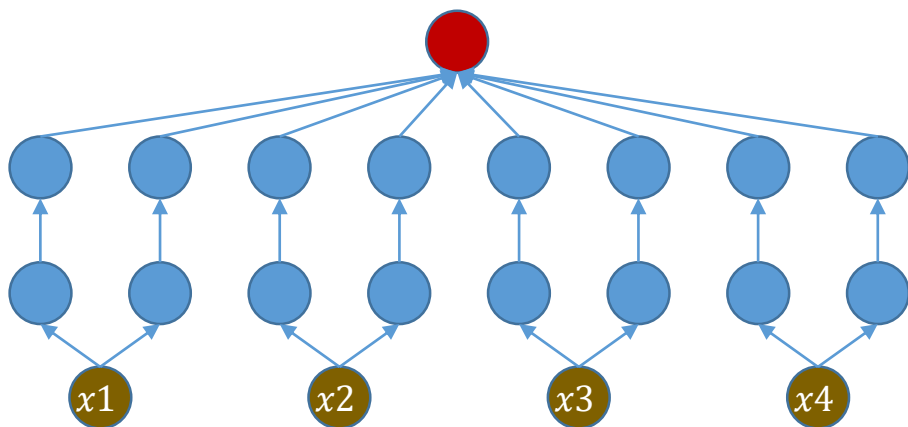
- But: when $\alpha \rightarrow 0$, we got interesting, non-RKHS inductive bias (e.g. nuclear norm, sparsity)

Scale of Init: Kernel vs Rich

Consider linear regression with squared parametrization:

$$f(\mathbf{w}, x) = \sum_j (\mathbf{w}_+[j]^2 - \mathbf{w}_-[j]^2) x[j] = \langle \beta(\mathbf{w}), x \rangle \quad \text{with } \beta(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

And unbiased initialization $\mathbf{w}_\alpha(0) = \alpha \mathbf{1}$ (so that $\beta(\mathbf{w}_\alpha(0)) = 0$).



$$f((U, V), x) = (UV^\top) \cdot \begin{bmatrix} \text{diag}(x) & 0 \\ 0 & -\text{diag}(x) \end{bmatrix}$$
$$U(0) = V(0) = \alpha I$$

Scale of Init: Kernel vs Rich

Consider linear regression with squared parametrization:

$$f(\mathbf{w}, x) = \sum_j (\mathbf{w}_+[j]^2 - \mathbf{w}_-[j]^2) x[j] = \langle \beta(\mathbf{w}), x \rangle \quad \text{with } \beta(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

And unbiased initialization $\mathbf{w}_\alpha(0) = \alpha \mathbf{1}$ (so that $\beta(\mathbf{w}_\alpha(0)) = 0$).

What's the implicit bias of grad flow w.r.t square loss $L_s(\mathbf{w}) = \sum_i (f(\mathbf{w}, x_i) - y_i)^2$?

$$\beta_\alpha(\infty) = \lim_{t \rightarrow \infty} \beta(\mathbf{w}_\alpha(t))$$

In Kernel Regime $\alpha \rightarrow \infty$: $K_0(x, x') = 4\langle x, x' \rangle$ and so

$$\beta_\alpha(\infty) \xrightarrow{\alpha \rightarrow \infty} \hat{\beta}_{L2} = \arg \min_{X\beta=y} \|\beta\|_2$$

In Rich Regime $\alpha \rightarrow 0$: special case of MF with commutative measurements

$$\beta_\alpha(\infty) \xrightarrow{\alpha \rightarrow 0} \hat{\beta}_{L1} = \arg \min_{X\beta=y} \|\beta\|_1$$

For any α :

$$\beta_\alpha(\infty) = ???$$

$$\beta(t) = w_+(t)^2 - w_-(t)^2$$

$$L = \|X\beta - y\|_2^2$$

$$\dot{w}_+(t) = -\nabla L(t) = -2X^\top r(t) \circ 2w_+(t) \quad w_+(t) = w_+(0) \circ \exp\left(-2X^\top \int_0^t r(\tau) d\tau\right)$$

$$\dot{w}_-(t) = -\nabla L(t) = +2X^\top r(t) \circ 2w_-(t) \quad w_-(t) = w_-(0) \circ \exp\left(+2X^\top \int_0^t r(\tau) d\tau\right)$$

$$\beta(t) = \alpha^2 \left(e^{-4X^\top \int_0^t r(\tau) d\tau} - e^{4X^\top \int_0^t r(\tau) d\tau} \right) \quad r(t) = X\beta(t) - y$$

$$s = 4 \int_0^\infty r(\tau) d\tau \in \mathbb{R}^m$$

$$\beta(\infty) = \alpha^2 \left(e^{-X^\top s} - e^{X^\top s} \right) = 2\alpha^2 \sinh X^\top s$$

$$X\beta(\infty) = y$$

$$\min Q(\beta) \quad s.t. \quad X\beta = y$$

$$\nabla Q(\beta^*) = X^\top \nu$$

$$X\beta^* = y$$

$$\beta(\infty) = \alpha^2 \left(e^{-X^\top s} - e^{X^\top s} \right) = 2\alpha^2 \sinh X^\top s$$

$$X\beta(\infty) = y$$

$$\nabla Q(\beta) = \sinh^{-1} \frac{\beta}{2\alpha^2}$$

$$Q(\beta) = \sum_i \int \sinh^{-1} \frac{\beta[i]}{2\alpha^2} = \alpha^2 \sum_i \left(\frac{\beta[i]}{\alpha^2} \sinh^{-1} \frac{\beta[i]}{2\alpha^2} - \sqrt{4 + \left(\frac{\beta[i]}{\alpha^2} \right)^2} \right)$$

$$\min Q(\beta) \quad \text{s.t.} \quad X\beta = y$$

$$\nabla Q(\beta^*) = X^\top \mathbf{v}$$

$$X\beta^* = y$$

$$\sinh^{-1} \frac{\beta(\infty)}{2\alpha^2} = X^\top \mathbf{s}$$

$$X\beta(\infty) = y$$

Scale of Init: Kernel vs Rich

Consider linear regression with squared parametrization:

$$f(\mathbf{w}, x) = \sum_j (\mathbf{w}_+[j]^2 - \mathbf{w}_-[j]^2) x[j] = \langle \beta(\mathbf{w}), x \rangle \quad \text{with } \beta(\mathbf{w}) = \mathbf{w}_+^2 - \mathbf{w}_-^2$$

And unbiased initialization $\mathbf{w}_\alpha(0) = \alpha \mathbf{1}$ (so that $\beta(\mathbf{w}_\alpha(0)) = 0$).

What's the implicit bias of grad flow w.r.t square loss $L_S(\mathbf{w}) = \sum_i (f(\mathbf{w}, x_i) - y_i)^2$?

$$\beta_\alpha(\infty) = \lim_{t \rightarrow \infty} \beta(\mathbf{w}_\alpha(t))$$

In Kernel Regime $\alpha \rightarrow \infty$: $K_0(x, x') = 4\langle x, x' \rangle$ and so

$$\beta_\alpha(\infty) \xrightarrow{\alpha \rightarrow \infty} \hat{\beta}_{L2} = \arg \min_{X\beta=y} \|\beta\|_2$$

In Rich Regime $\alpha \rightarrow 0$: special case of MF with commutative measurements

$$\beta_\alpha(\infty) \xrightarrow{\alpha \rightarrow 0} \hat{\beta}_{L1} = \arg \min_{X\beta=y} \|\beta\|_1$$

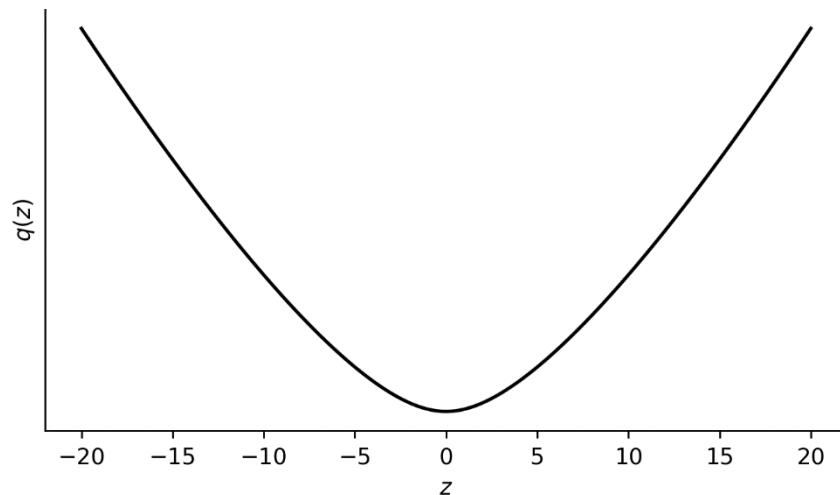
For any α :

$$\beta_\alpha(\infty) = \arg \min_{X\beta=y} Q_\alpha(\beta)$$

where $Q_\alpha(\beta) = \sum_j q\left(\frac{\beta[j]}{\alpha^2}\right)$ and $q(b) = 2 - \sqrt{4 + b^2} + b \sinh^{-1}\left(\frac{b}{2}\right)$

$$\beta_{\alpha}(\infty) = \arg \min_{X\beta=y} Q_{\alpha}(\beta)$$

where $Q_{\alpha}(\beta) = \sum_j q\left(\frac{\beta[j]}{\alpha^2}\right)$ and $q(b) = 2 - \sqrt{4 + b^2} + b \sinh^{-1}\left(\frac{b}{2}\right)$



$$\text{Induced dynamics: } \dot{\beta}_{\alpha} = -\sqrt{\beta_{\alpha}^2 + 4\alpha^4} \odot \nabla L_s(\beta_{\alpha})$$

Theorem 2. For any $0 < \epsilon < d$,

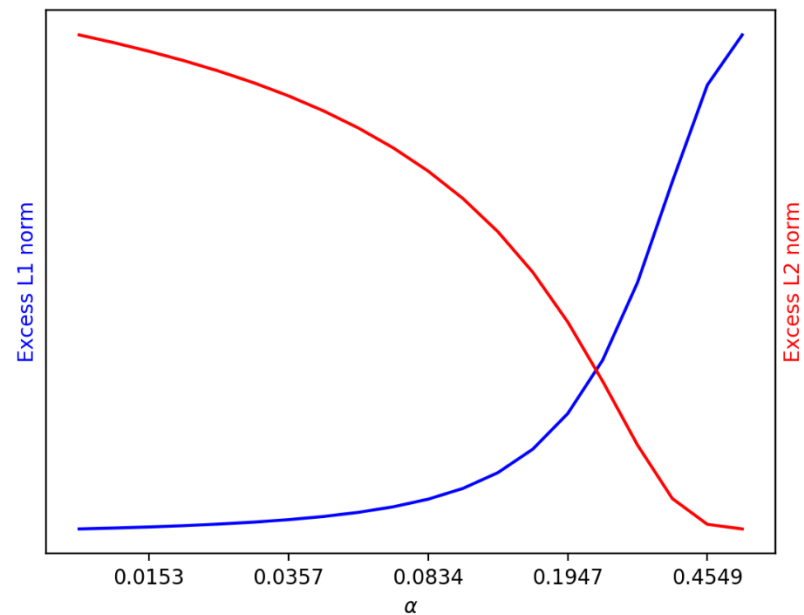
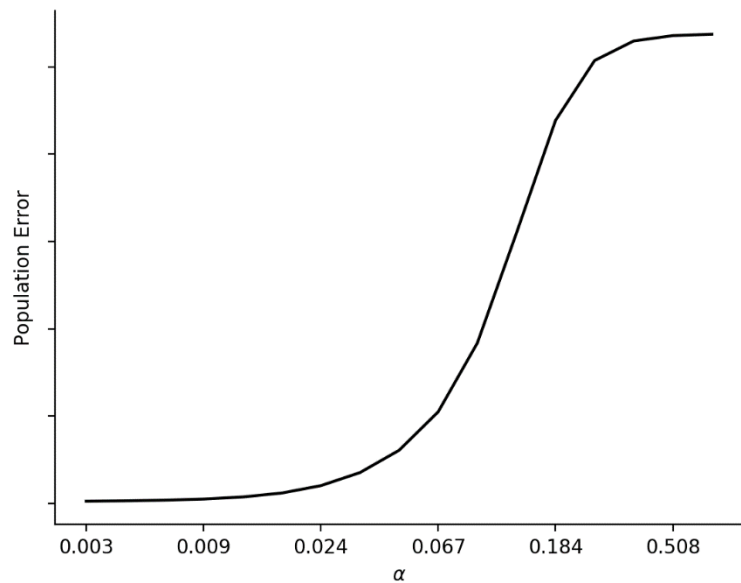
$$\alpha \leq \min \left\{ (2(1 + \epsilon) \|\beta_{L1}^*\|_1)^{-\frac{2+\epsilon}{2\epsilon}}, \exp \left(-\frac{d}{\epsilon \|\beta_{L1}^*\|_1} \right) \right\} \implies \|\hat{\beta}_{\alpha}\|_1 \leq (1 + \epsilon) \|\beta_{L1}^*\|_1$$

Theorem 3. For any $\epsilon > 0$

$$\alpha \geq \sqrt{2(1 + \epsilon) \left(1 + \frac{2}{\epsilon}\right) \|\beta_{L2}^*\|_2} \implies \|\hat{\beta}_{\alpha}\|_2^2 \leq (1 + \epsilon) \|\beta_{L2}^*\|_2^2$$

Sparse Learning

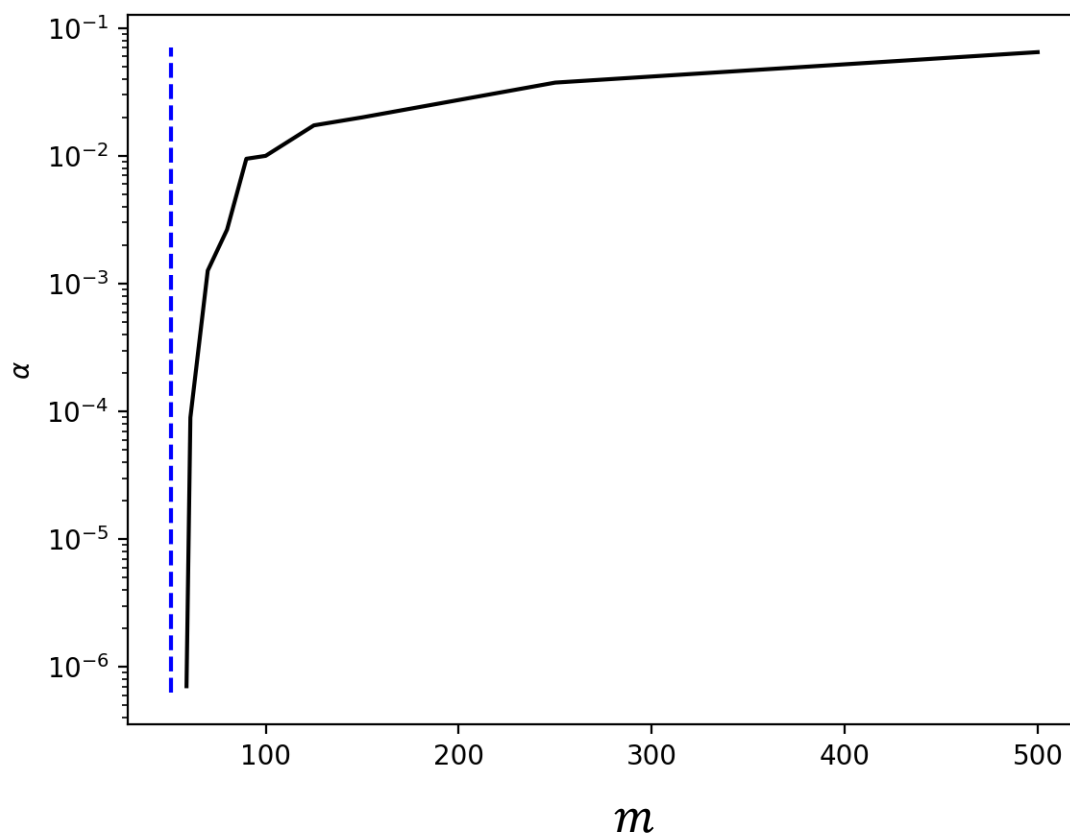
$$y_i = \langle \beta^*, x_i \rangle + N(0, 0.01)$$
$$d = 1000, \quad \|\beta^*\|_0 = 5, \quad m = 100$$



Sparse Learning

$$y_i = \langle \beta^*, x_i \rangle + N(0, 0.01)$$
$$d = 1000, \quad \|\beta^*\|_0 = k$$

How small does α need to be to get $L(\beta_\alpha(\infty)) < 0.025$

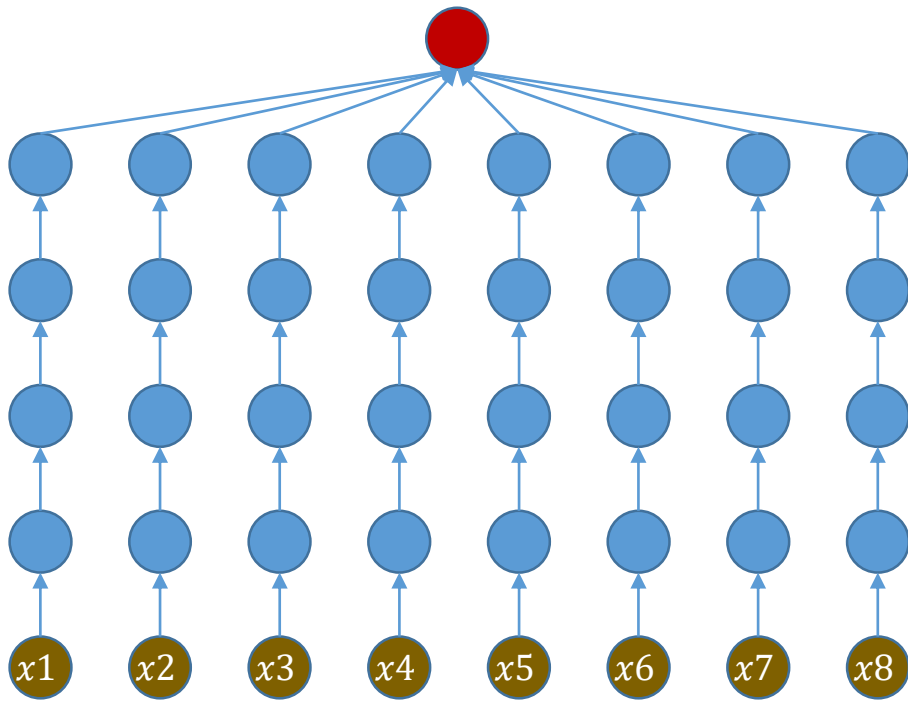


More controlling parameters

- Depth
- Width
- Optimization accuracy
- Stepsize, batchsize, ??

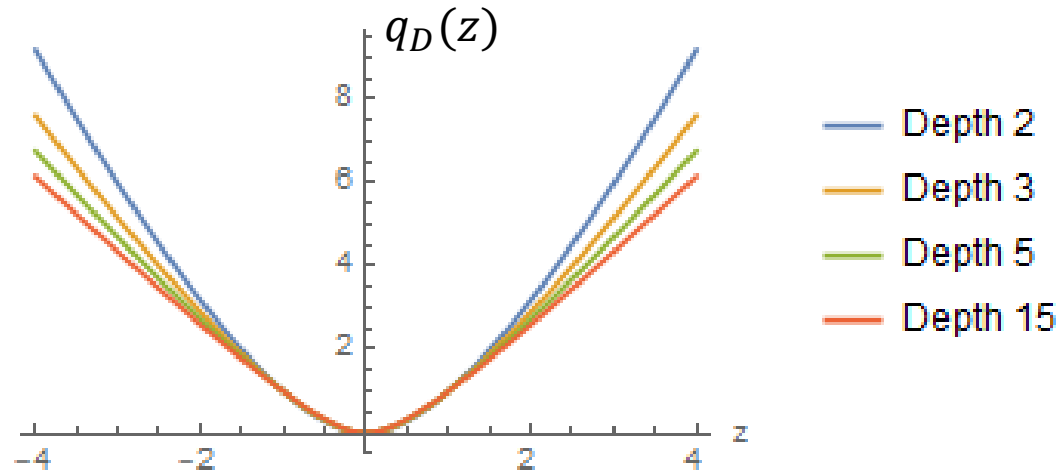
Depth

$$\beta(t) = w_+(t)^D - w_-(t)^D$$



Depth

$$\beta(t) = w_+(t)^D - w_-(t)^D \quad \beta(\infty) = \arg \min Q_D \left(\beta / \alpha^D \right) \text{ s.t. } X\beta = y$$



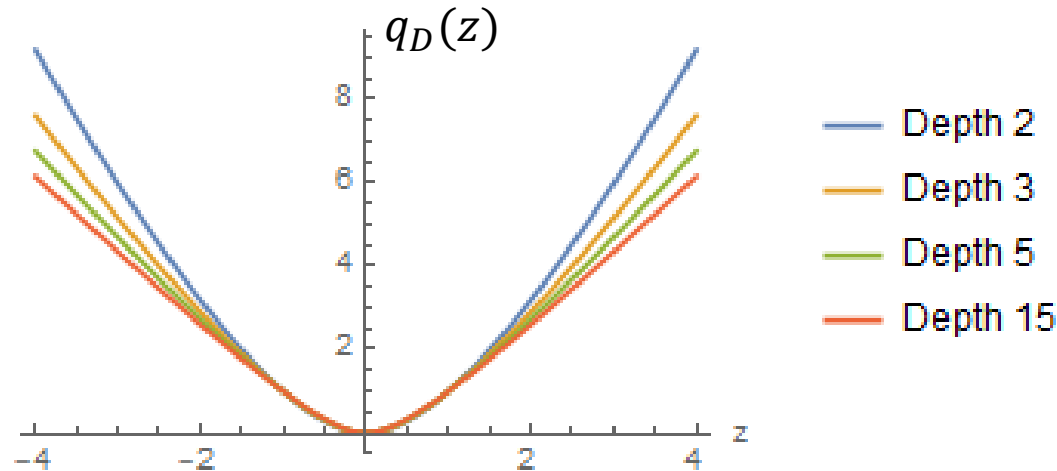
$$h_D(z) = \alpha^D \left((1 + \alpha^{D-2} D(D-2)z)^{\frac{-1}{D-2}} - (1 - \alpha^{D-2} D(D-2)z)^{\frac{-1}{D-2}} \right)$$

$$q_D = \int h_D^{-1}$$

$$Q_D(\beta) = \sum_i q_D \left(\frac{\beta[i]}{\alpha^D} \right)$$

Depth

$$\beta(t) = w_+(t)^D - w_-(t)^D \quad \beta(\infty) = \arg \min Q_D \left(\beta / \alpha^D \right) \text{ s.t. } X\beta = y$$



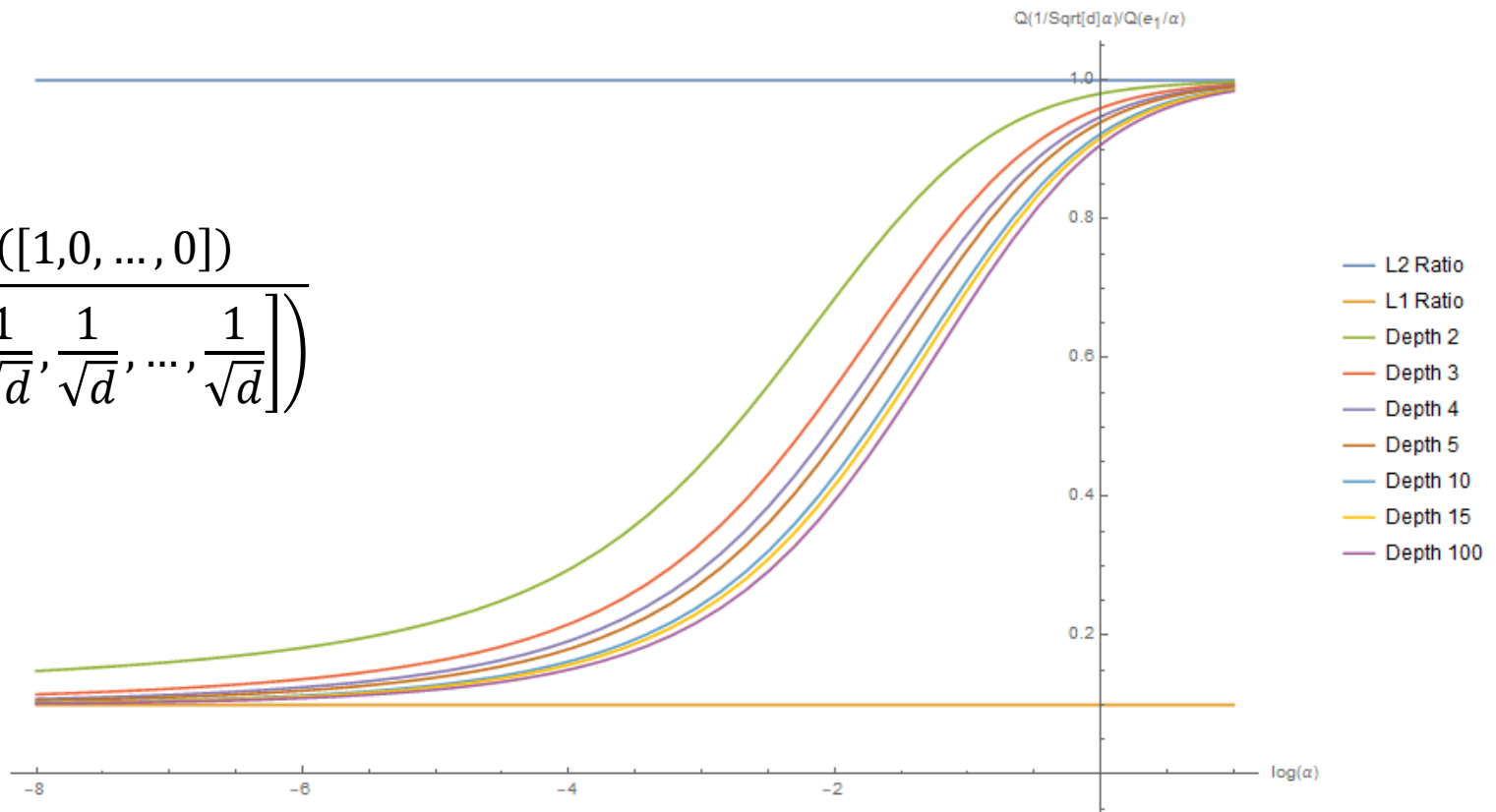
For all depth $D \geq 2$, $\beta(\infty) \xrightarrow{\alpha \rightarrow 0} \arg \min_{X\beta=y} \|\beta\|_1$

- Contrast with explicit reg: For $R_\alpha(\beta) = \min_{\beta=w_+^D - w_-^D} \|w - \alpha \mathbf{1}\|_2^2$, $R_\alpha(\beta) \xrightarrow{\alpha \rightarrow 0} \|\beta\|_{2/D}$
also observed by [Arora Cohen Hu Luo 2019]
- Also with logistic loss, $\beta(\infty) \xrightarrow{\alpha \rightarrow 0} \propto \text{SOSP of } \|\beta\|_{2/D}$ [Gunasekar Lee Soudry Srebro 2018]
- With sq loss, always $\|\cdot\|_1$, but for deep D , we get there quicker

Depth

$$\beta(t) = w_+(t)^D - w_-(t)^D \quad \beta(\infty) = \arg \min Q_D \left(\beta / \alpha^D \right) \text{ s.t. } X\beta = y$$

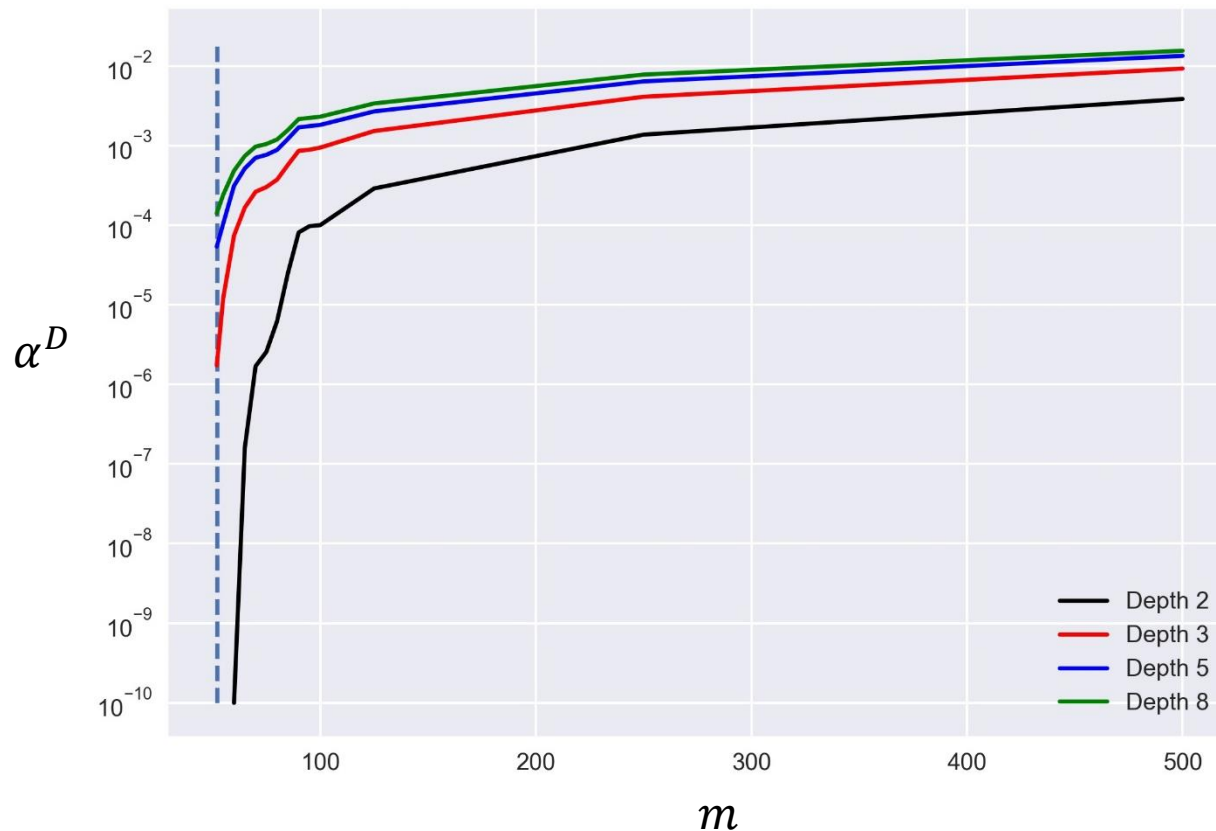
$$\frac{Q([1,0,\dots,0])}{Q\left(\left[\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}, \dots, \frac{1}{\sqrt{d}}\right]\right)}$$



Sparse Learning with Depth

$$y_i = \langle \beta^*, x_i \rangle + N(0, 0.01)$$
$$d = 1000, \quad \|\beta^*\|_0 = k$$

How small does α need to be to get $L(\beta_\alpha(\infty)) < 0.025$



Deep Learning

- Expressive Power
 - We are searching over the space of all functions...
... but with what inductive bias?
 - How does this bias look in function space?
 - Is it reasonable/sensible?
- Capacity / Generalization ability / Sample Complexity
 - What's the true complexity measure (inductive bias)?
 - How does it control generalization?
- Computation / Optimization
 - How and where does optimization bias us? Under what conditions?
 - Magic property of reality under which deep learning “works”