



# When can Deep Networks avoid the Curse of Dimensionality and other theoretical puzzles?

1

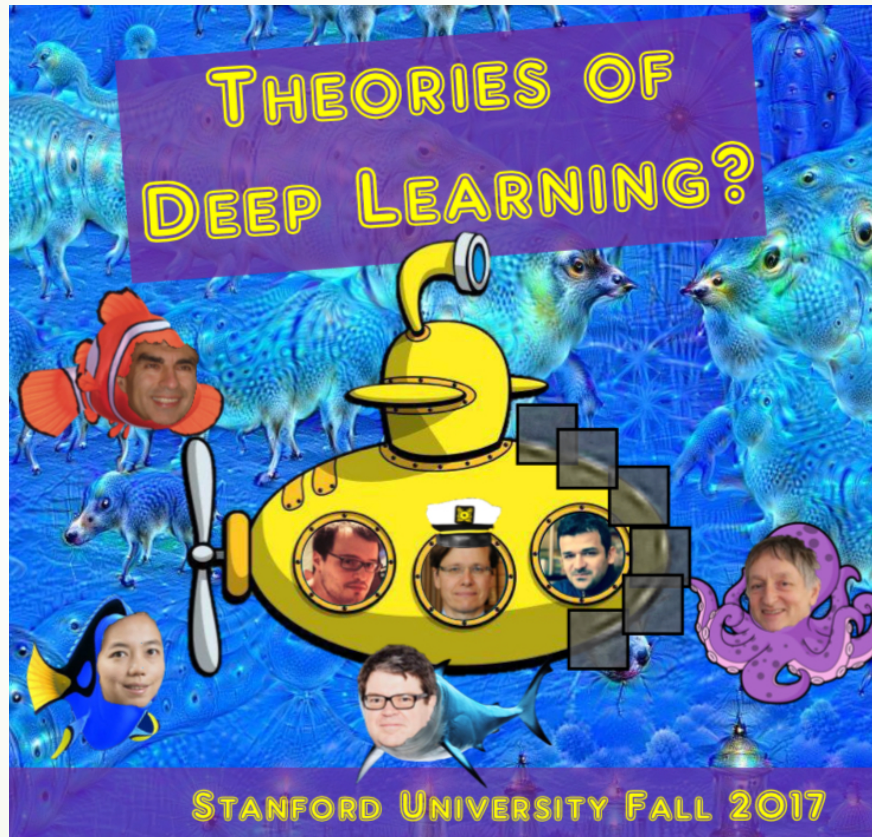
Yuan YAO

HKUST

Based on Tomaso Poggio's talk etc.



# Acknowledgement

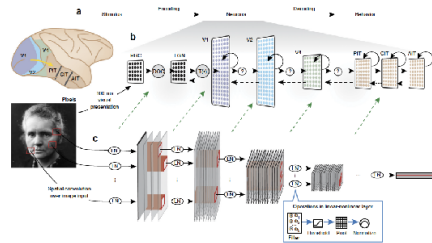


A following-up course at HKUST: <https://deeplearning-math.github.io/>

# Tomaso Poggio's group @CBMM, MIT in ~40 years of pursuit:

## CBMM's focus is the Science and the Engineering of Intelligence

We aim to make progress in understanding intelligence, that is in understanding how the brain makes the mind, how the brain works and how to build intelligent machines. We believe that the science of intelligence will enable better engineering of intelligence.



## CBMM: one of the motivations

Key recent advances  
in the engineering of intelligence  
have their roots  
in basic research on the brain



**It is time for  
a theory of deep learning**



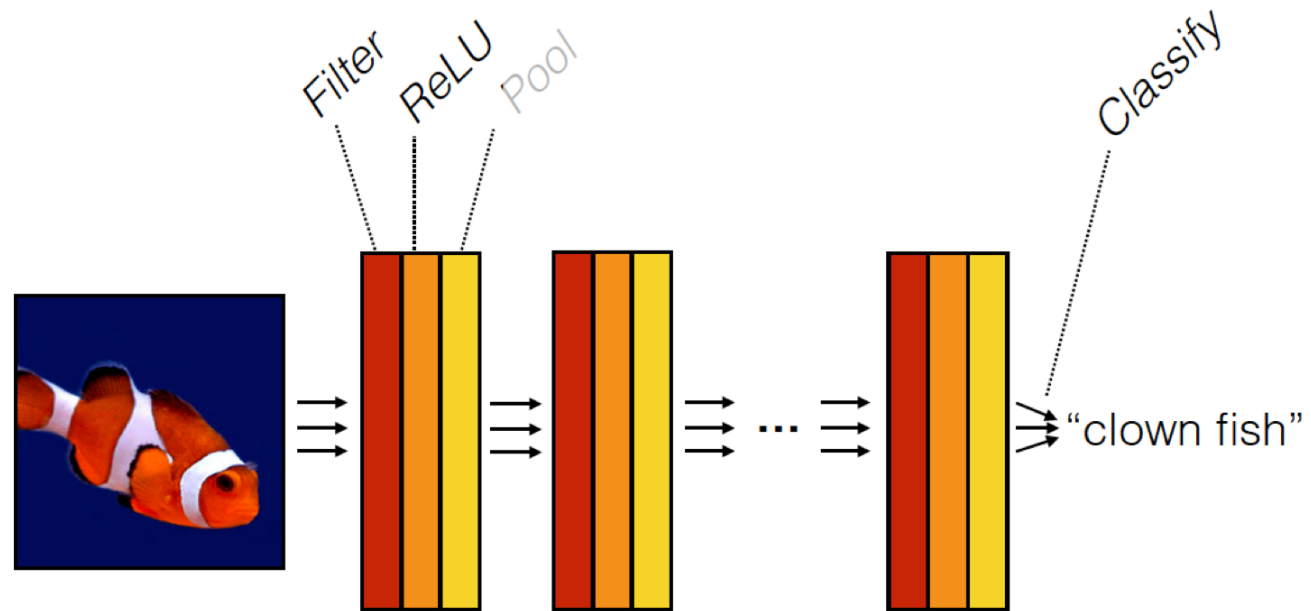


# Deep Networks: Three theory questions

---

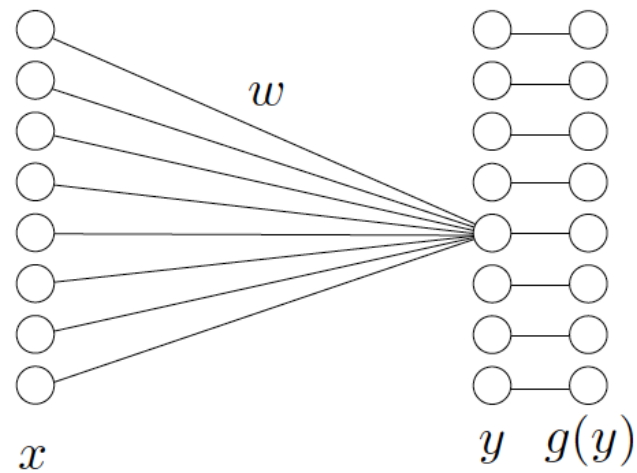
- *Approximation Theory*: When and why are deep networks better than shallow networks?
- *Optimization*: What is the landscape of the empirical risk?
- *Learning Theory*: How can deep learning not overfit?

## Computation in a neural net

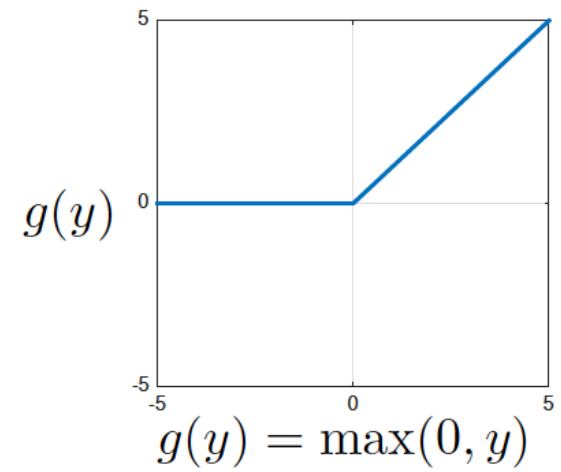


$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Computation in a neural net



Rectified linear unit (ReLU)



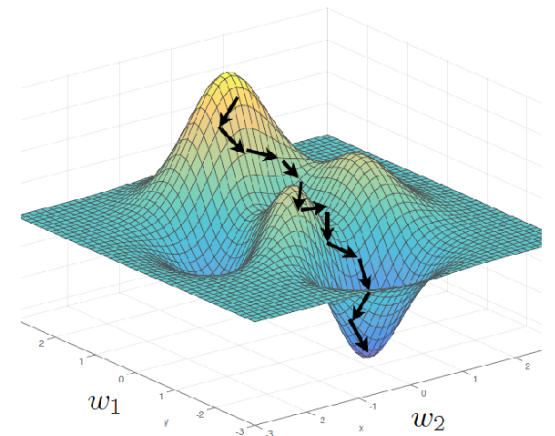
## Gradient descent

$$\operatorname{argmin}_{\mathbf{w}} \sum_i \ell(\mathbf{z}_i, f(\mathbf{x}_i; \mathbf{w})) = L(\mathbf{w})$$

One iteration of gradient de

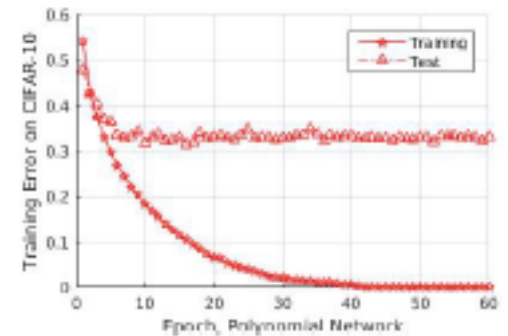
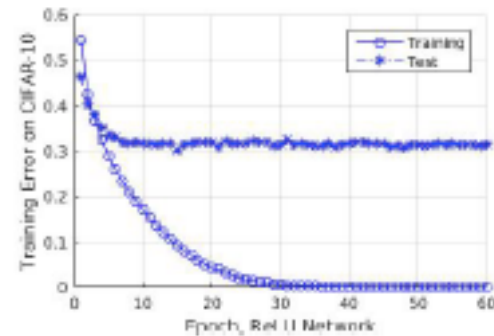
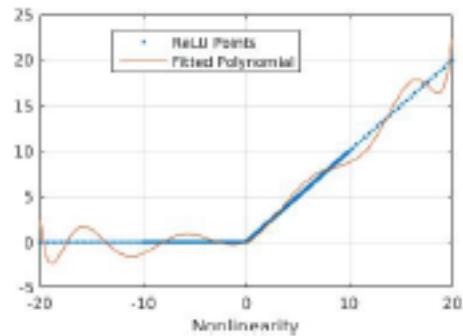
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \frac{\partial L(\mathbf{w}^t)}{\partial \mathbf{w}}$$

learning rate





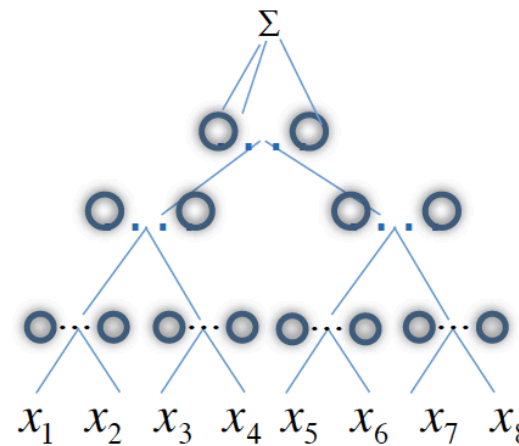
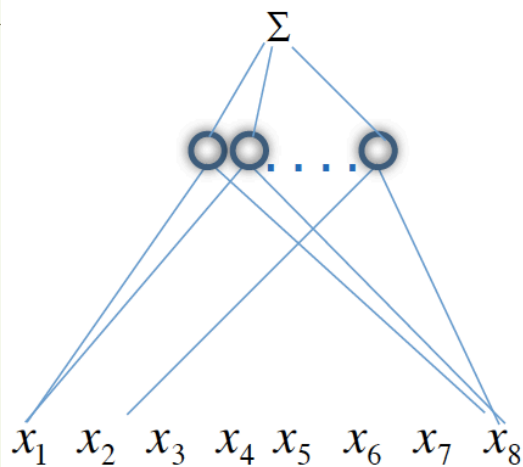
## RELU approximation by univariate polynomial preserves deep nets properties



From now on, we study polynomial networks!

## Deep and shallow networks: universality

**Theorem** *Shallow, one-hidden layer networks with a nonlinear  $\phi(x)$  which is not a polynomial are universal. Arbitrarily deep networks with a nonlinear  $\phi(x)$  (including polynomials) are universal.*



$$\phi(x) = \sum_{i=1}^r c_i |\langle w_i, x \rangle + b_i|_+$$

## Classical learning theory and Kernel Machines (Regularization in RKHS)

---

$$\min_{f \in H} \left[ \frac{1}{\ell} \sum_{i=1}^{\ell} V(f(x_i) - y_i) + \lambda \|f\|_K^2 \right]$$

implies

$$f(\mathbf{x}) = \sum_i^{\ell} \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

**Equation includes splines, Radial Basis Functions and Support Vector Machines (depending on choice of V).**

RKHS were explicitly introduced in learning theory by Girosi (1997), Vapnik (1998).

Moody and Darken (1989), and Broomhead and Lowe (1988) introduced RBF to learning theory. Poggio and Girosi (1989) introduced Tikhonov regularization in learning theory and worked (implicitly) with RKHS. RKHS were used earlier in approximation theory (eg Parzen, 1952-1970, Wahba, 1990).

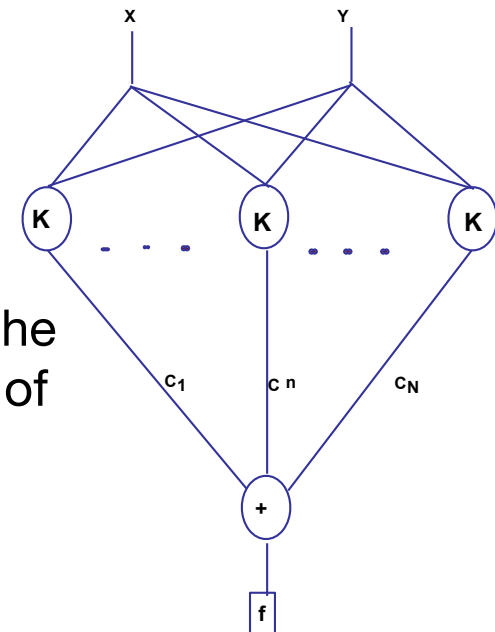
Mhaskar, Poggio, Liao, 2016

## Classical kernel machines are equivalent to shallow networks

Kernel machines...

$$f(\mathbf{x}) = \sum_i^l c_i K(\mathbf{x}, \mathbf{x}_i) + b$$

can be “written” as shallow networks: the value of  $K$  corresponds to the “activity” of the “unit” for the input and the correspond to “weights”



# Curse of dimensionality

---

$$y = f(x_1, x_2, \dots, x_8)$$

## Curse of dimensionality

Both shallow and deep network can approximate a function of  $d$  variables equally well. The number of parameters in both cases depends exponentially on  $d$  as  $O(\epsilon^{-d})$ .



## Generic functions

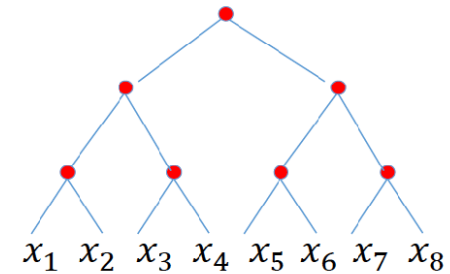
$$f(x_1, x_2, \dots, x_8)$$

## Compositional functions

$$f(x_1, x_2, \dots, x_8) = g_3(g_{21}(g_{11}(x_1, x_2), g_{12}(x_3, x_4))g_{22}(g_{11}(x_5, x_6), g_{12}(x_7, x_8)))$$

## Hierarchically local compositionality

$$f(x_1, x_2, \dots, x_8) = g_3(g_{21}(g_{11}(x_1, x_2), g_{12}(x_3, x_4)), g_{22}(g_{11}(x_5, x_6), g_{12}(x_7, x_8)))$$



### Theorem (informal statement)

Suppose that a function of  $d$  variables is hierarchically, locally, compositional. Both shallow and deep network can approximate  $f$  equally well. The number of parameters of the shallow network depends exponentially on  $d$  as  $O(\epsilon^{-d})$  with the dimension whereas for the deep network it is  $O(d\epsilon^{-2})$ .





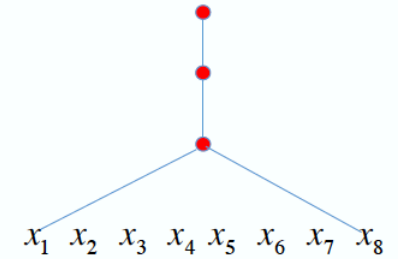
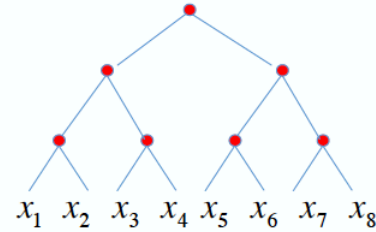
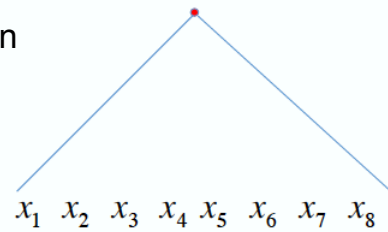
# Historical Results

- ▶ A classical **theorem [Sipser, 1986; Hastad, 1987]** shows that deep circuits are more efficient in representing certain Boolean functions than shallow circuits. Hastad proved that highly-variable functions (in the sense of having high frequencies in their Fourier spectrum) in particular the parity function cannot even be decently approximated by small constant depth circuits
- ▶ The main **result of [Telgarsky, 2016, Colt]** says that there are functions with many oscillations that cannot be represented by shallow networks with linear complexity but can be represented with low complexity by deep networks.
- ▶ Older examples exist: consider a function which is a linear combination of  $n$  tensor product Chui–Wang spline wavelets, where each wavelet is a tensor product cubic spline. It was shown by **Chui and Mhaskar** that is impossible to implement such a function using a shallow neural network with a sigmoidal activation function using  $O(n)$  neurons, but a deep network with the activation function  $(x_+)^2$  do so. In this case, as we mentioned, there is a formal proof of a gap between deep and shallow networks. Similarly, **Eldan and Shamir** show other cases with separations that are exponential in the input dimension.

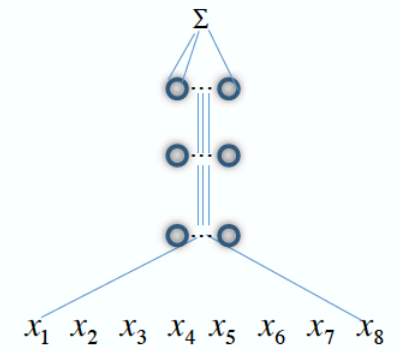
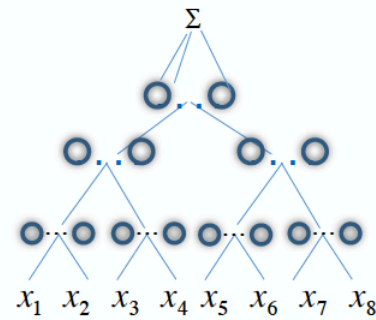
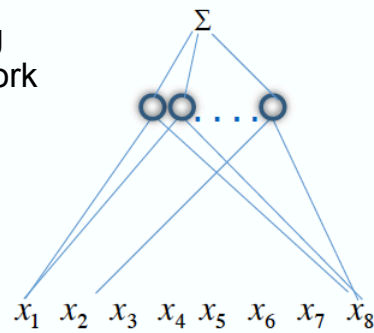


# Microstructure of compositionality

target function



approximating function/network

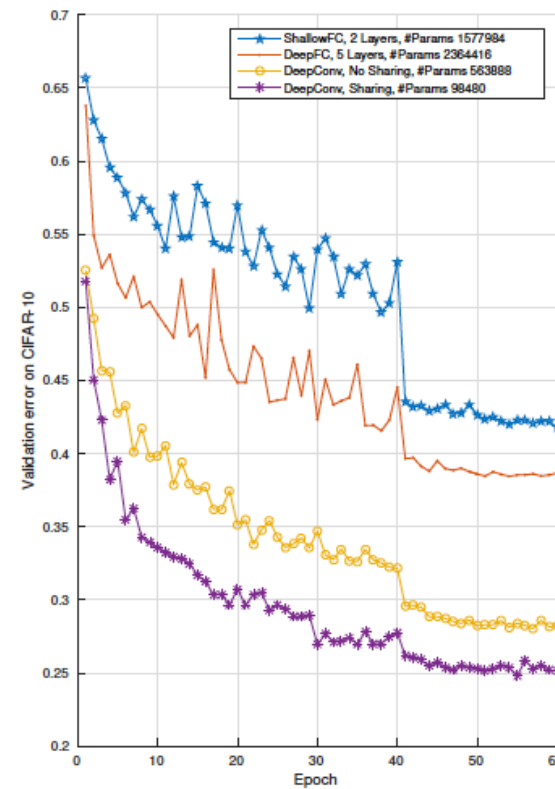
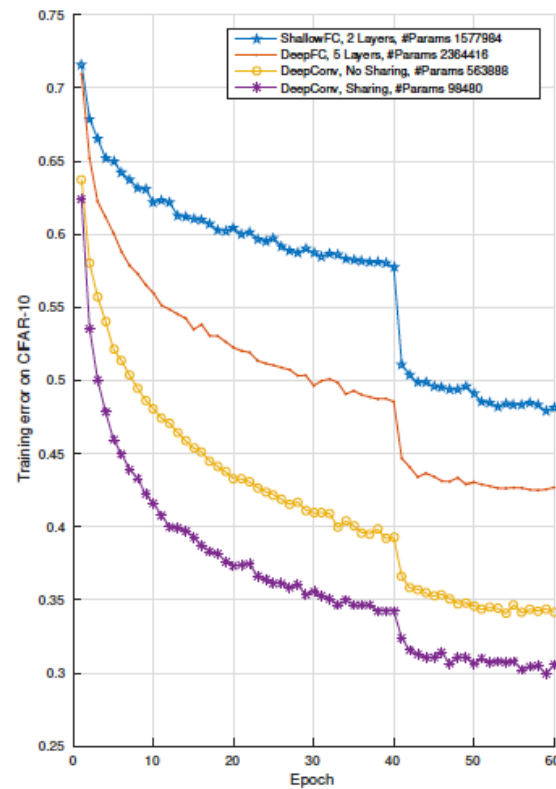


*a*

*b*

*c*

# Locality of constituent functions is key: CIFAR





## Open problem: why compositional functions are important for perception?

- ▶ They seem to occur in computations on text, speech, images...why?
- ▶ **Conjecture (with) Max Tegmark**
- ▶ The locality of the hamiltonians of physics induce compositionality in natural signals such as images
- ▶ or
- ▶ The connectivity in our brain implies that our perception is limited to compositional functions

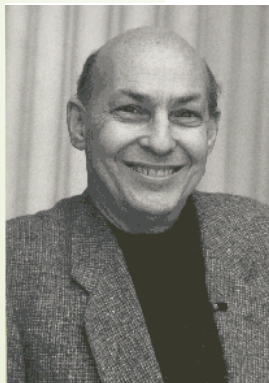
## Why are compositional functions important?

Which one of these reasons:

Physics?

Neuroscience? <===

Evolution?

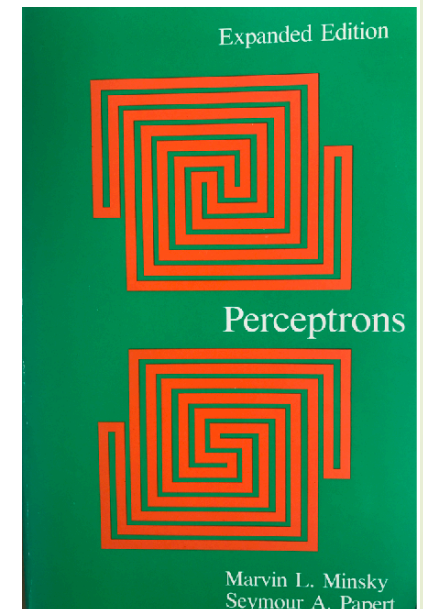


## Locality of Computation

What is special about locality of computation?

Locality in “space”?

Locality in “time”?



**Locality** leads to **Sparsity**.



# Deep Networks: Three theory questions

---

- *Approximation Theory*: When and why are deep networks better than shallow networks?
- *Optimization*: What is the landscape of the empirical risk?
- *Learning Theory*: How can deep learning not overfit?

## Theory II: What is the Landscape of the empirical risk?

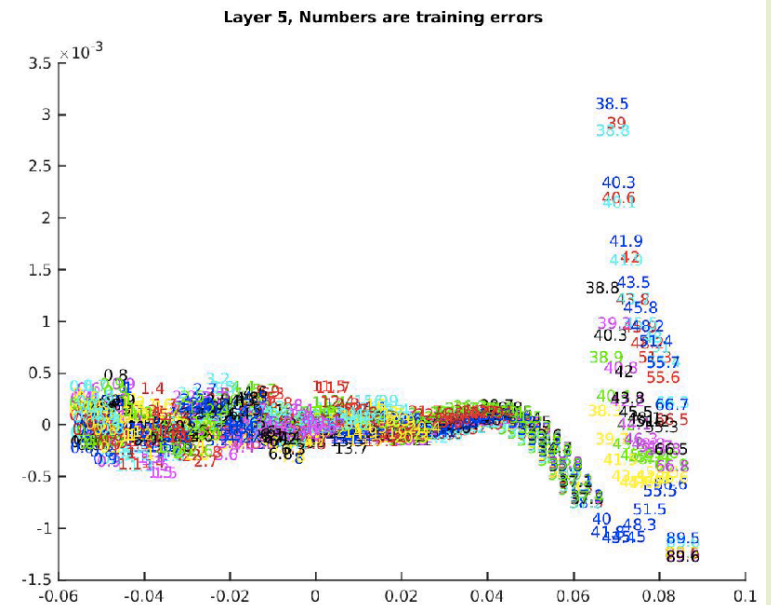
### Observation

Replacing the RELUs with univariate polynomial approximation, Bezout theorem implies that the system of polynomial equations corresponding to zero empirical error has a very large number of degenerate solutions. The global zero-minimizers correspond to flat minima in many dimensions (generically, unlike local minima). Thus SGD is biased towards finding global minima of the empirical risk.



CENTER FOR  
Brains  
Minds+  
Machines

Liao, Poggio, 2017



## Bezout theorem

---

$$p(x_i) - y_i = 0 \quad \text{for } i = 1, \dots, n$$

The set of polynomial equations above with  $k = \text{degree of } p(x)$  has a number of distinct zeros (counting points at infinity, using projective space, assigning an appropriate multiplicity to each intersection point, and excluding degenerate cases) equal to

$$Z = k^n$$

the product of the degrees of each of the equations. As in the linear case, when the system of equations is underdetermined – as many equations as data points but more unknowns (the weights) – the theorem says that there are an infinite number of global minima, under the form of  $Z$  regions of zero empirical error.

## Global and local zeros

---

$$f(x_i) - y_i = 0 \quad \text{for } i = 1, \dots, n$$

$n$  equations in  $W$  unknowns with  $W \gg n$

$$\nabla_w \sum_{i=1}^N (f(x_i) - y_i)^2 = 0$$

$W$  equations in  $W$  unknowns

*There are a very large number of zero-error minima which are highly degenerate unlike the local non-zero minima.*



## Langevin equation

---

$$\frac{df}{dt} = -\gamma_t \nabla V(f(t), z(t)) + \gamma'_t dB(t)$$

$$f_{t+1} = f_t - \gamma_m \nabla V(f_t, z_t) + \gamma'_t W_t.$$

with the Boltzmann equation as asymptotic “solution”

$$p(f) \sim \frac{1}{Z} = e^{-\frac{U(x)}{T}}$$



# DL~SGD is an analogy, NOT a Theorem!

## SGD

---

$$f_{t+1} = f_t - \gamma_t \nabla V(f_t, z_t),$$

$$\nabla V(f_t, z_t) = \frac{1}{|z_t|} \sum_{z \in z_t} \nabla V(f_t, z).$$

We define a noise “equivalent quantity”

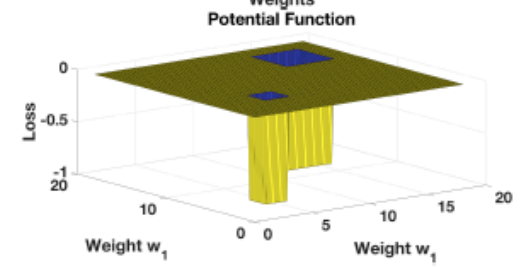
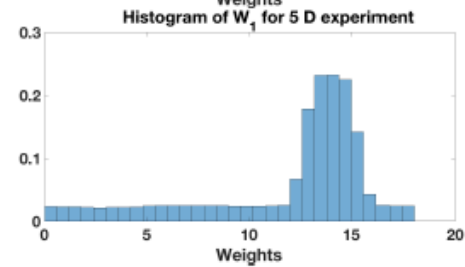
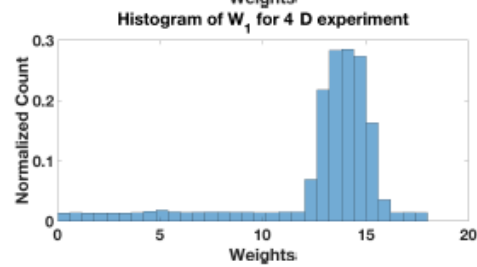
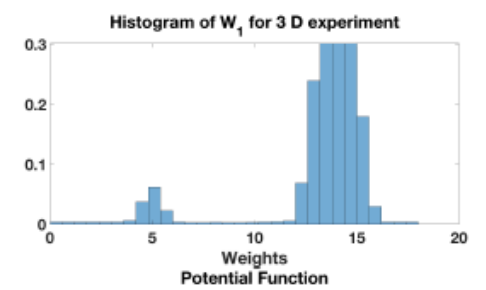
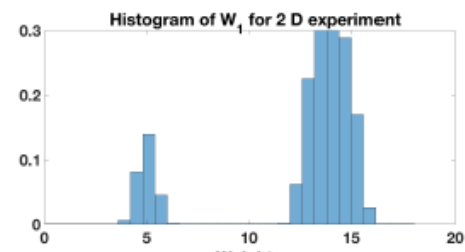
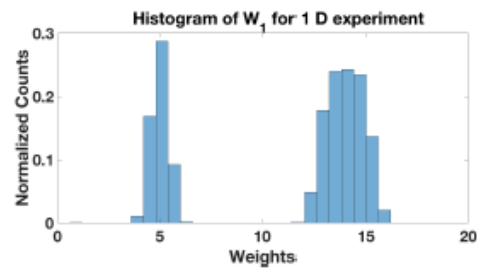
$$\xi_t = \nabla V(f_t, z_t) - \nabla I_{S_n}(f_t),$$

and it is clear that  $\mathbb{E}\xi_t = 0$ .

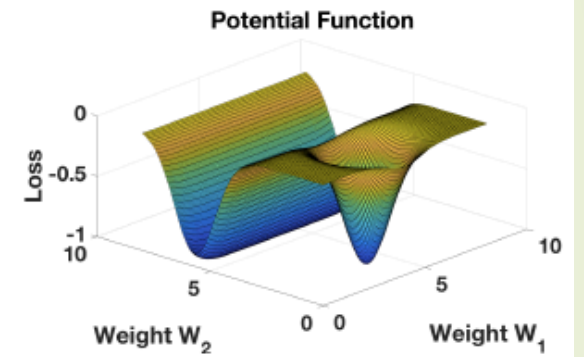
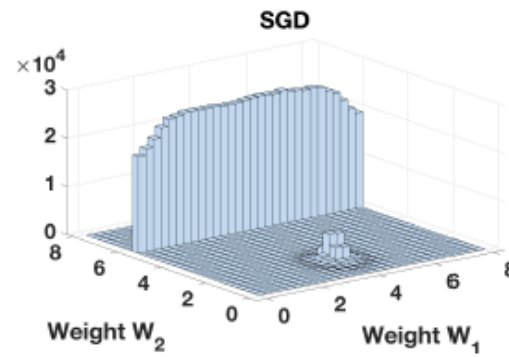
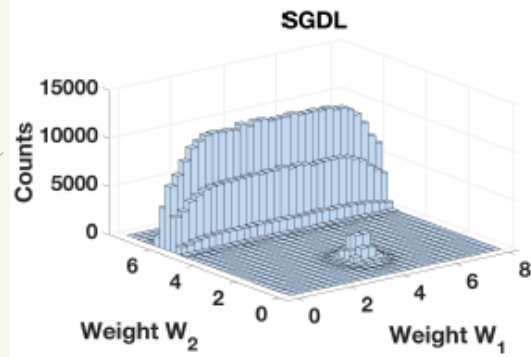
We write Equation 6 as

$$f_{t+1} = f_t - \gamma_t (\nabla I_{S_n}(f_t) + \xi_t).$$

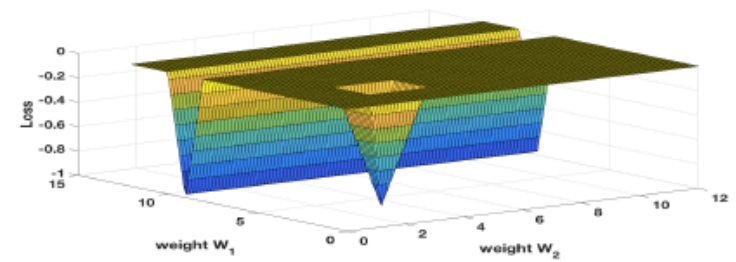
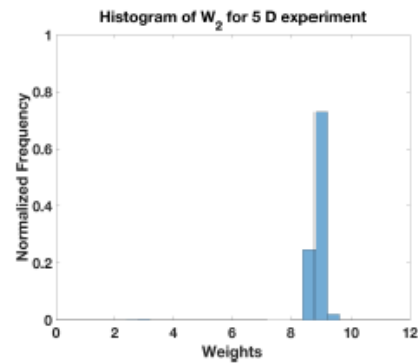
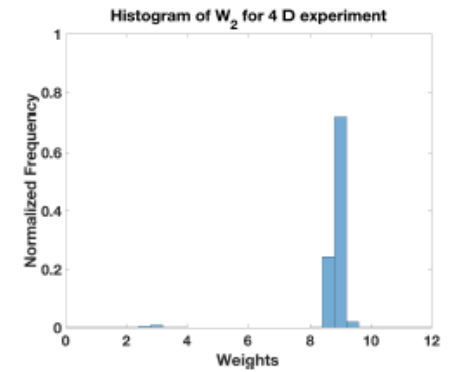
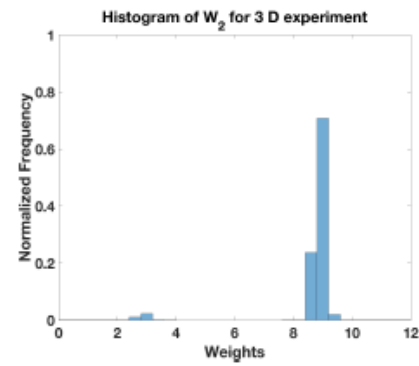
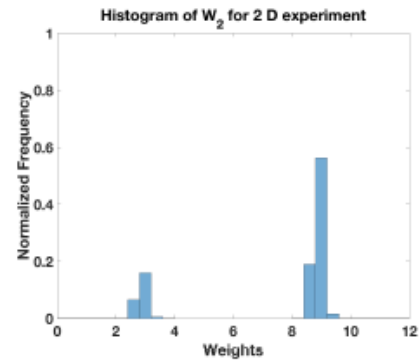
## GDL selects larger volume minima



# GDL and SGD

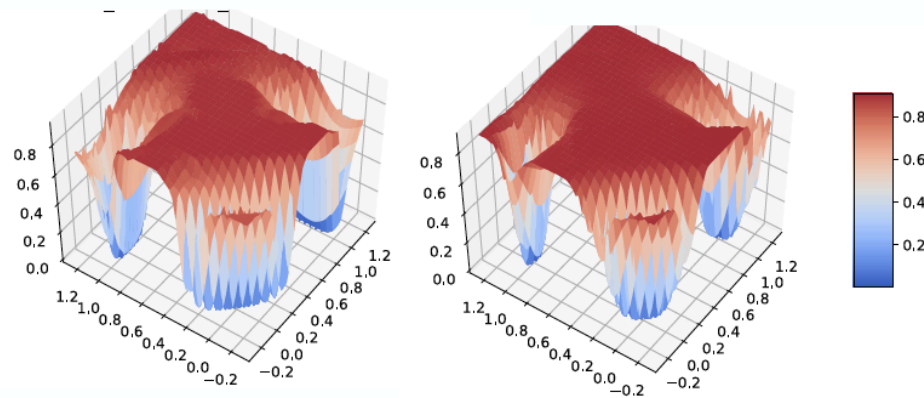


# Concentration because of high dimensionality



## SGDL and SGD observation: summary

- SGDL finds with very high probability large volume, flat zero-minimizers; empirically SGD behaves in a similar way
- Flat minimizers correspond to degenerate zero-minimizers and thus to global minimizers;



CIFAR-10: Natural Labels

Random Labels

Poggio, Rakhlin, Golovits, Zhang, Liao, 2017



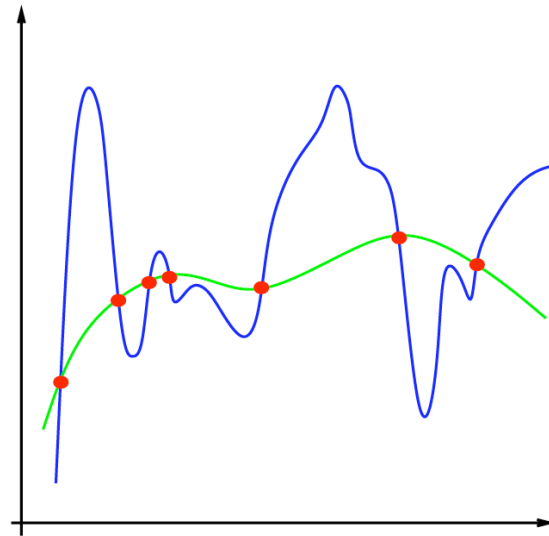


## Deep Networks: Three theory questions

---

- *Approximation Theory*: When and why are deep networks better than shallow networks?
- *Optimization*: What is the landscape of the empirical risk?
- *Learning Theory*: How can deep learning not overfit?

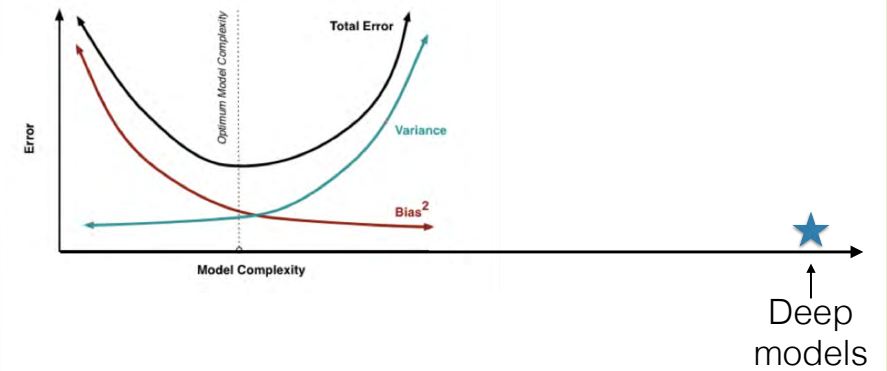
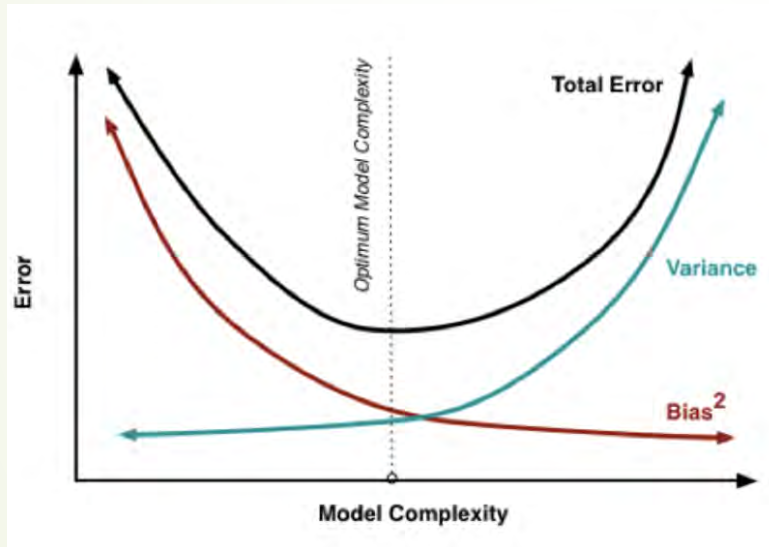
## Problem of overfitting



Regularization or similar to control overfitting

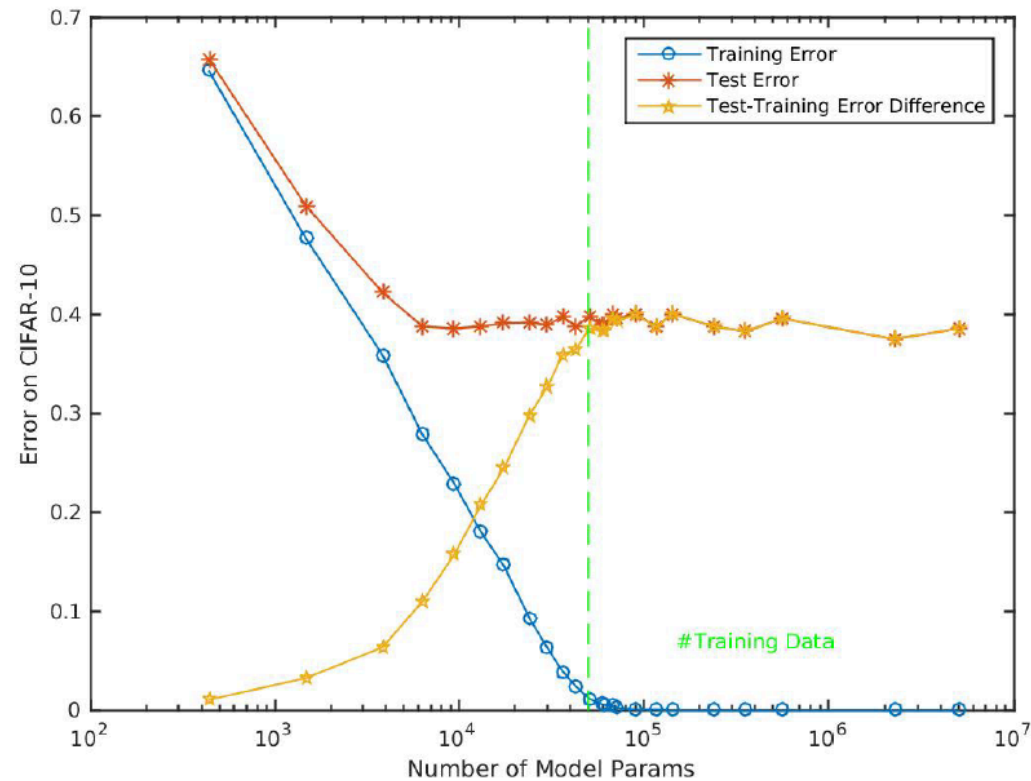


# Deep Learning Puzzle



Models where  $p > 20n$  are common

# Deep Learning Puzzle: No Overfitting



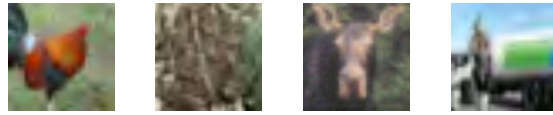
*Explaining this figure is our main goal!*



# How to reduce generalization error?

- Model capacity
- Regularization (norms, dropout, etc.)
- Implicit regularization (early stopping)
- Data augmentation (fake data, crops, shifts, etc.)

All of these are sufficient but by no means necessary!



# CIFAR10

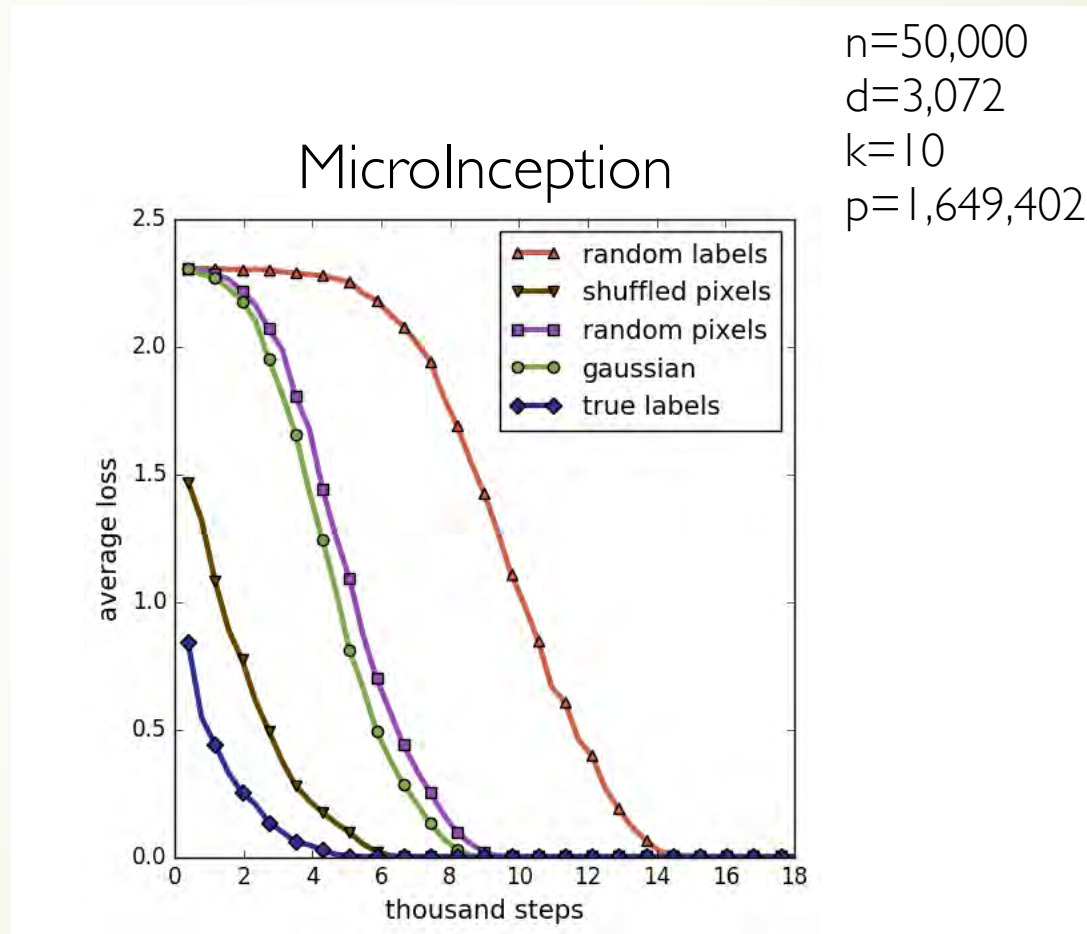
$n=50,000$   
 $d=3,072$   
 $k=10$

What happens when I turn off the regularizers?

<u>Model</u>	<u>parameters</u>	<u>p/n</u>	Train <u>loss</u>	Test <u>error</u>
CudaConvNet	145,578	2.9	0	23%
CudaConvNet (with regularization)	145,578	2.9	0.34	18%
MicroInception	1,649,402	33	0	14%
ResNet	2,401,440	48	0	13%

Ben Recht

# Zhang et al. 2016 on Cifar10

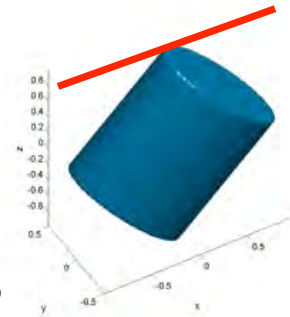


# Why do big models generalize?

- This is true in the linear case too!

$$\underset{w}{\text{minimize}} \quad \|y - Xw\|^2$$

$$X \text{ } n \times p, n < p$$



- *Infinite* number of *global* minima.
- All global minima have the *same* Hessian.
- *At least*  $p-n$  of the Hessian eigenvalues are zero.
- Which one should we pick?
- Regularize to leverage structure.

Sparsity

Rank

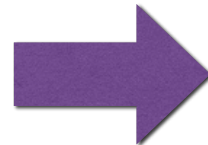
Smoothness

Architecture

- Why do we generalize when fitting the labels exactly?
- Happens for linear models!  $f(x) = w^T x$

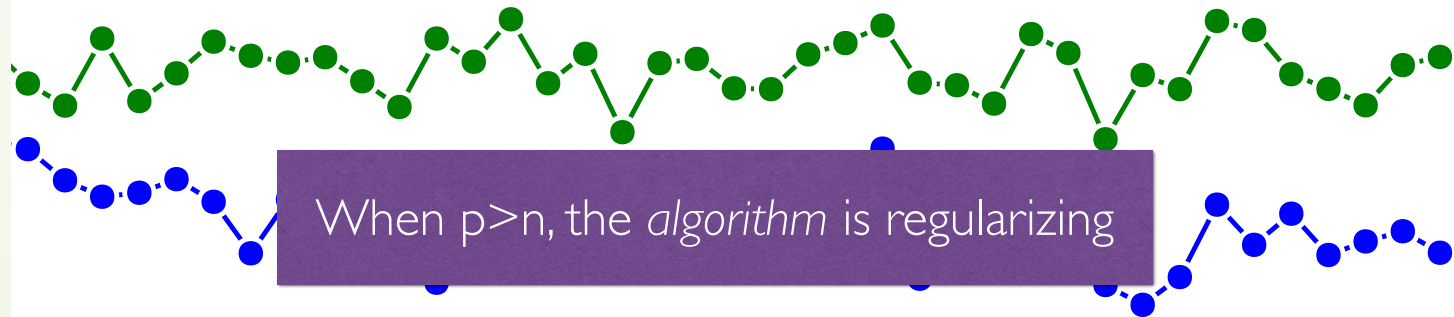
$$\text{minimize } \sum_{i=1}^n (w^T x_i - y_i)^2$$

SGD solution



$$\begin{array}{l} \text{minimize } \|w\| \\ \text{subject to } Xw = y \end{array}$$

If you run SGD you find the minimum norm solution



## Implicit regularization by GD: #iterations controls $\lambda$

**Theorem 3.1** *In the setting of Section 2, let Assumption 1 hold. Let  $\gamma \in ]0, \kappa^{-1}]$ . Then the following hold:*

(i) *If we choose a stopping rule  $t^* : \mathbb{N}^* \rightarrow \mathbb{N}^*$  such that*

$$\lim_{n \rightarrow +\infty} t^*(n) = +\infty \quad \text{and} \quad \lim_{n \rightarrow +\infty} \frac{t^*(n)^3 \log n}{n} = 0 \quad (9)$$

*then*

$$\lim_{n \rightarrow +\infty} \mathcal{E}(\hat{w}_{t^*(n)}) - \inf_{w \in \mathcal{H}} \mathcal{E}(w) = 0 \quad \mathbb{P}\text{-almost surely.} \quad (10)$$

(ii) *Suppose additionally that the set  $\mathcal{O}$  of minimizers of (1) is nonempty and let  $w^\dagger$  be defined as in (2). If we choose a stopping rule  $t^* : \mathbb{N}^* \rightarrow \mathbb{N}^*$  satisfying the conditions in (9) then*

$$\|\hat{w}_{t^*(n)} - w^\dagger\|_{\mathcal{H}} \rightarrow 0 \quad \mathbb{P}\text{-almost surely.} \quad (11)$$

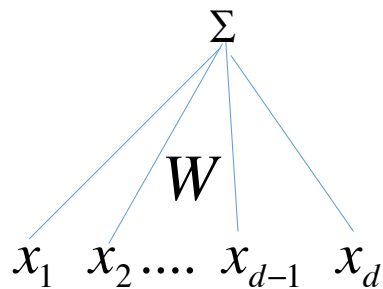


Rosasco, Villa, 2015

For early stopping in general RKHSs: Yao, Rosasco, Caponnetto, 2007.



## Implicit regularization by GD+SGD (linear case, no hidden layer)

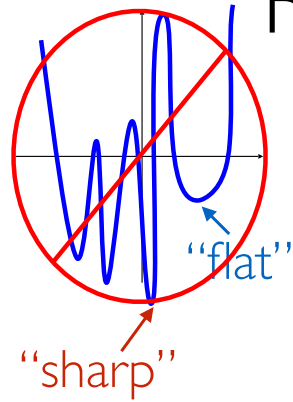


$$W = YX^\dagger$$

**Corollary 1.** *When initialized with zero, both GD and SGD converges to the minimum-norm solution.*

*Min norm solution is the limit for  $\lambda \rightarrow 0$  of regularized solution*

## Resolving “flat” vs “sharp”



$$\text{minimize } \sum_{i=1}^n (w^T x_i - y_i)^2$$

When  $p > n$ , all local minima have the same curvature.  
“flat minimizers?”

$$\begin{array}{l} \text{minimize } \|w\| \\ \text{subject to } Xw = y \end{array}$$

$\|w\|^{-1}$  is the *margin* of the classifier

Small norm  $\Rightarrow$  loss is stable to perturbations in parameters  
“flat minimizer”

Large norm  $\Rightarrow$  loss fluctuates with small perturbations to parameters  
“sharp minimizer”



# Margin in Classifications

...margin all over again

- In statistical learning, when all population points are classified correctly, one can show

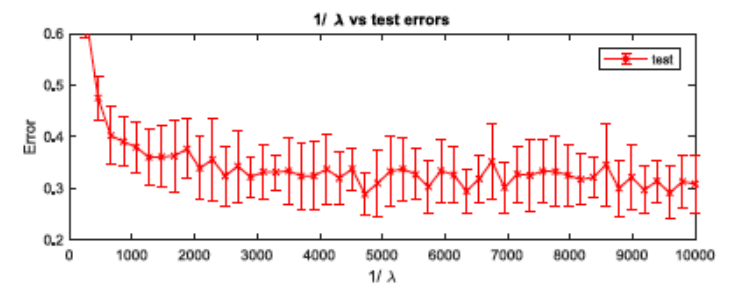
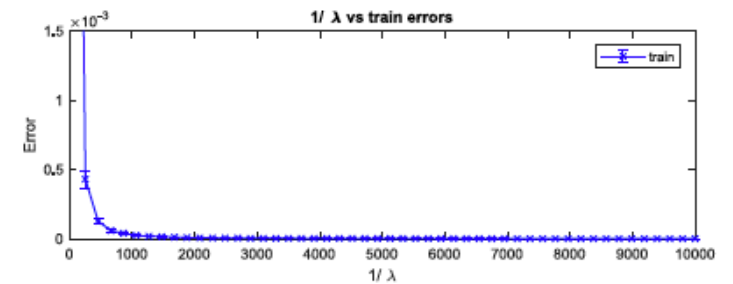
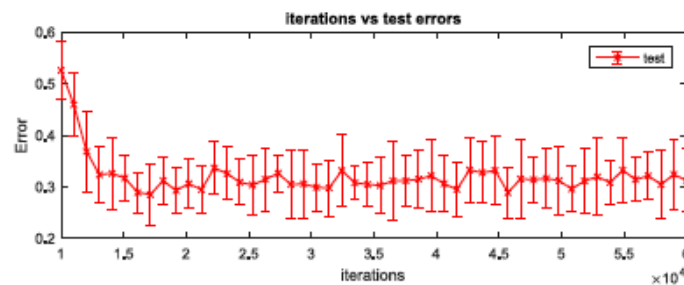
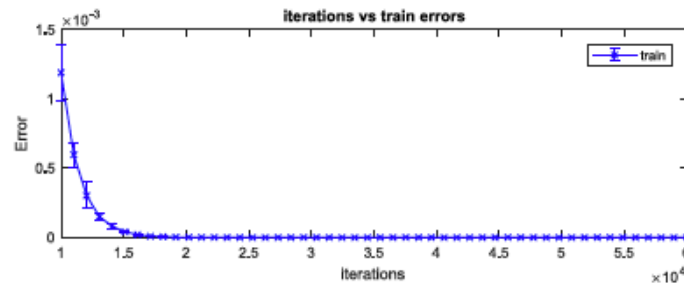
$$\mathbb{E}[\text{test error}] \leq 4 \frac{\|f_\star\|_k}{\sqrt{n}}$$

Inverse margin divided by  $\sqrt{n}$

- Better rates achievable with worse constants.

$$\mathbb{E}[\text{test error}] \leq \tilde{O}\left(\frac{\|f_\star\|_k^2}{n}\right)$$

# Deep linear network: GD as regularizer



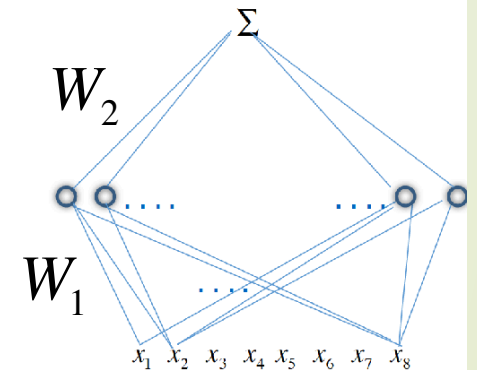
*GD regularizes deep linear networks as it does for linear networks*

# Deep linear networks

**Lemma 3.** For gradient descent and stochastic gradient descent with any mini-batch size,

- any number of the iterations adds no element in  $\text{Null}(X^\top)$  to the rows of  $W_1$ , and hence
- if the rows of  $W_1$  has no element in  $\text{Null}(X^\top)$  at anytime (including the initialization), the sequence converges to a minimum norm solution if it converges to a solution.

**Lemma 4.** If  $W_2 \neq 0$ , every stationary point w.r.t.  $W_1$  is a global minimum.



Remark:  $W_2 W_1 = A$  implies redundant parameters that are controlled if null space is empty

# Deep nonlinear (degree 2) networks

*Dynamical polynomial multilayer systems, training* We now discuss an extension of the above argument to the nonlinear activation case. Consider a polynomial second order (for simplicity and w.l.g) activation function  $h(z) = az + bz^2$ . The dynamical system (see for notation SI) is given by

$$\dot{W}_1 = -2(aW_2^\top E + 2b[(W_1 X) \circ (W_2^\top E)])X^\top \quad (7)$$

and

$$\dot{W}_2 = -2[aEX^\top W_1^\top + bE(((W_1 X)^2)^\top)]. \quad (8)$$

# Linearized dynamics to study stable solutions

$$\dot{\delta}_{W_1} = -2\delta_{W_2^\top} Y X^\top + 2\delta_{W_2^\top} W_2^* W_1^* X X^\top + 2W_2^{*\top} \delta_{W_2} W_1^* X X^\top + 2W_2^{*\top} W_2^* \delta_{W_1} X X^\top$$

and similarly

$$\dot{\delta}_{W_2} = -2Y X^\top \delta_{W_1^\top} + 2\delta_{W_2} W_1^* X X^\top W_1^{*\top} + 2W_2^* \delta_{W_1} X X^\top W_1^{*\top} + 2W_2^* W_1^* X X^\top \delta_{W_1^\top}$$

$$\dot{\delta}_{W_1} = -2\delta_{W_2^\top} Y X^\top$$

If  $W^*$  small

$$\dot{\delta}_{W_2} = -2Y X^\top \delta_{W_1^\top}$$



## Deep nonlinear networks: conjecture

---

The conclusion about the extension to multilayer networks with polynomial activation is thus similar to the linear case and can be summarized as follows:

*For low-noise data and a degenerate global minimum  $W^*$ , GD on a polynomial multilayer network avoids overfitting without explicit regularization, despite overparametrization.*





## Three theory questions: summary

---

- *Approximation theorems*: for hierarchical compositional functions deep but not shallow networks avoid the curse of dimensionality because of locality of constituent functions
- *Optimization remarks*: Bezout theorem suggests many global minima that are found by SGD with high probability wrt local minima
- *Learning Theory results and conjectures*: Unlike the case for a linear network the data dictate - because of the regularizing dynamics of GD - the number of effective parameters, which are in general fewer than the number of weights.

Thank you!

