# CS 229br Lecture 2: Dynamics & Bias

## Boaz Barak



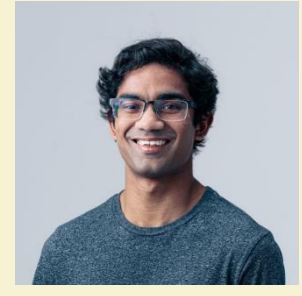**Ankur Moitra**
MIT 18.408

**Yamini Bansal**
Official TF

**Dimitris Kalimeris**
Unofficial TF

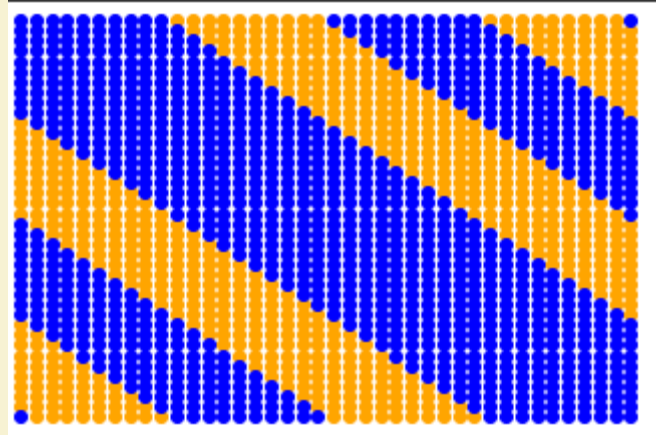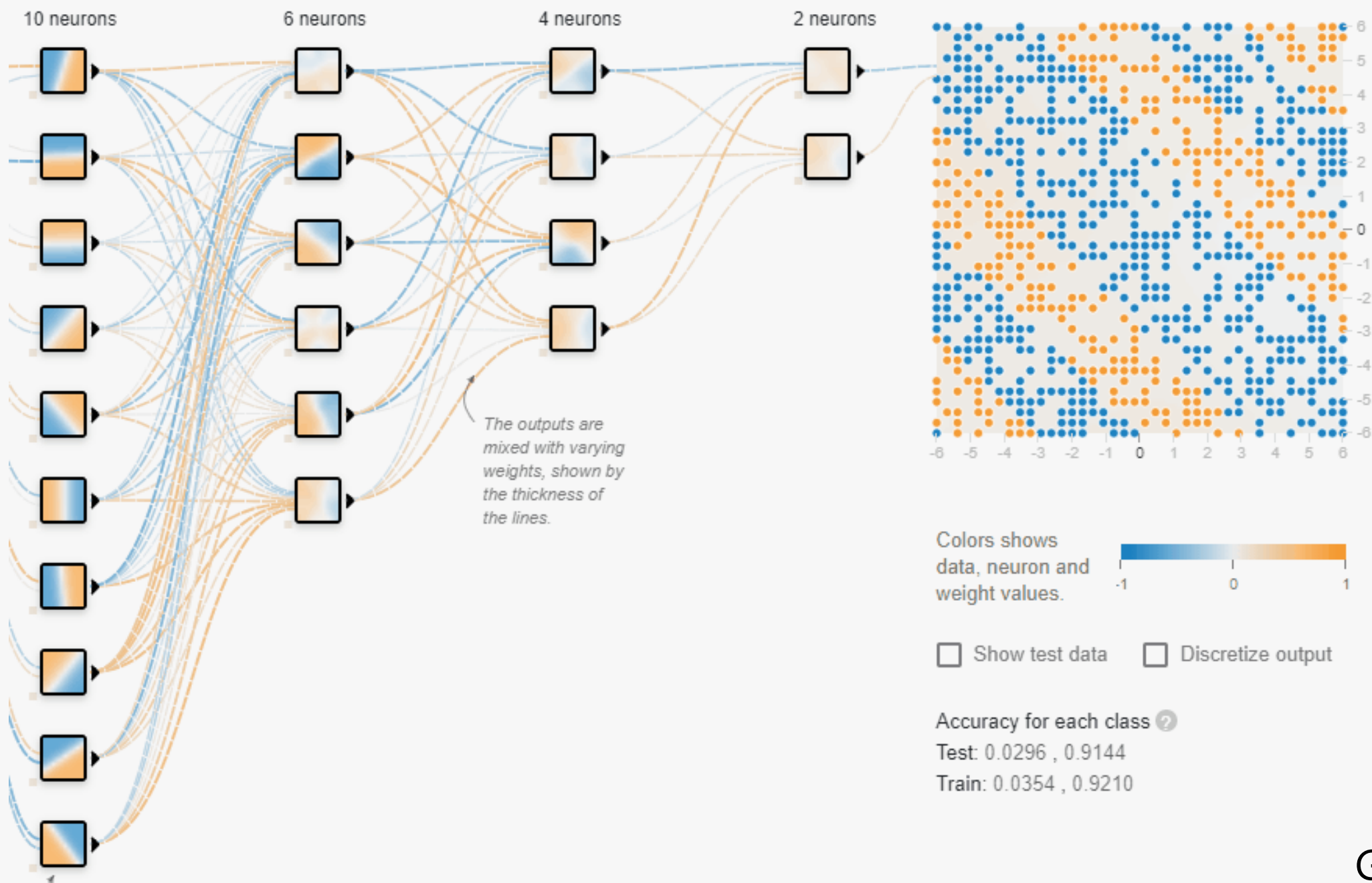**Gal Kaplun**
Unofficial TF

**Preetum Nakkiran**
Unofficial TF

Join slack team and sign in on your devices!

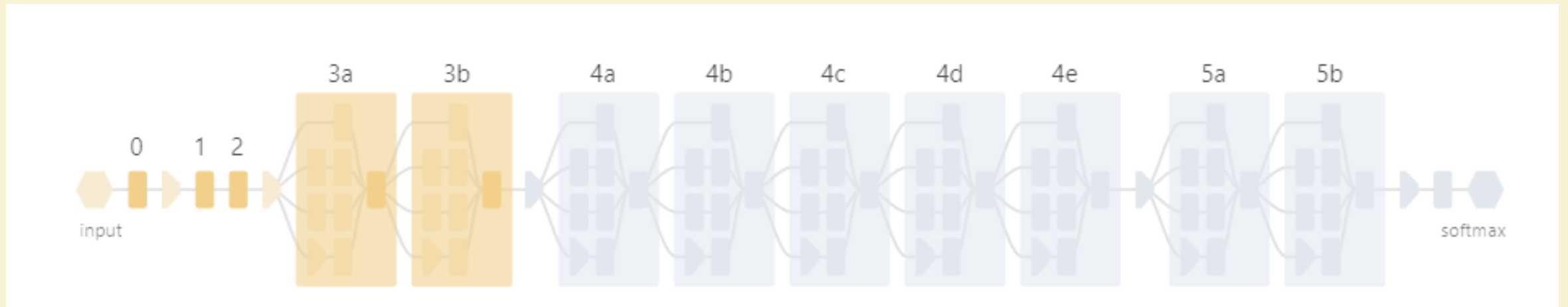# What do networks learn and when do they learn it?

- Simplicity bias

- Learning dynamics – what is learned first

- Different layers – what is learned by which layers?

- Some experimental evidence

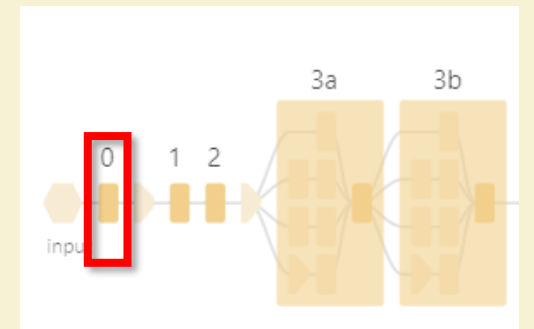- What can we prove?

# What do networks learn?

The outputs are mixed with varying weights, shown by the thickness of the lines.

10 neurons    6 neurons    4 neurons    2 neurons

Colors shows data, neuron and weight values.

-1    0    1

☐ Show test data    ☐ Discretize output

Accuracy for each class ⓘ
Test: 0.0296 , 0.9144
Train: 0.0354 , 0.9210

GIF

# Inception V1 (Olah et al, 2020)



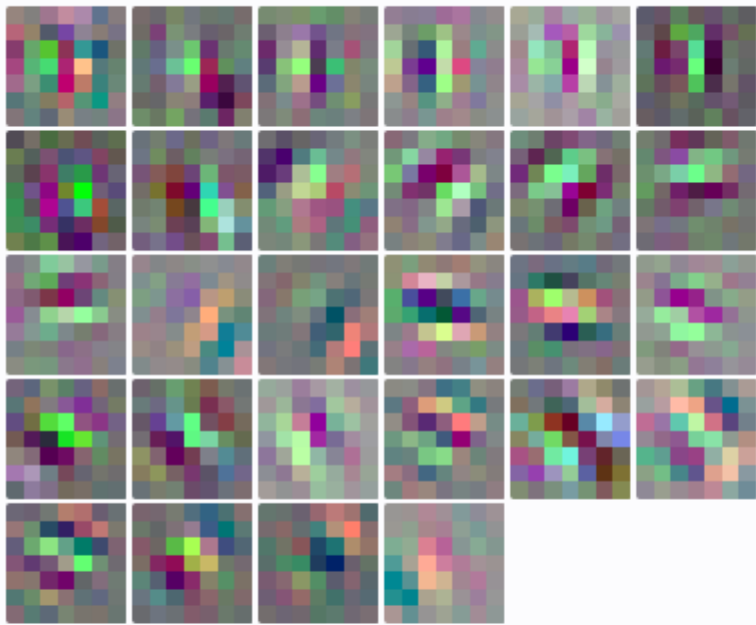Olah, Cammarata, Schubert, Goh, Petrov, Carter 2020
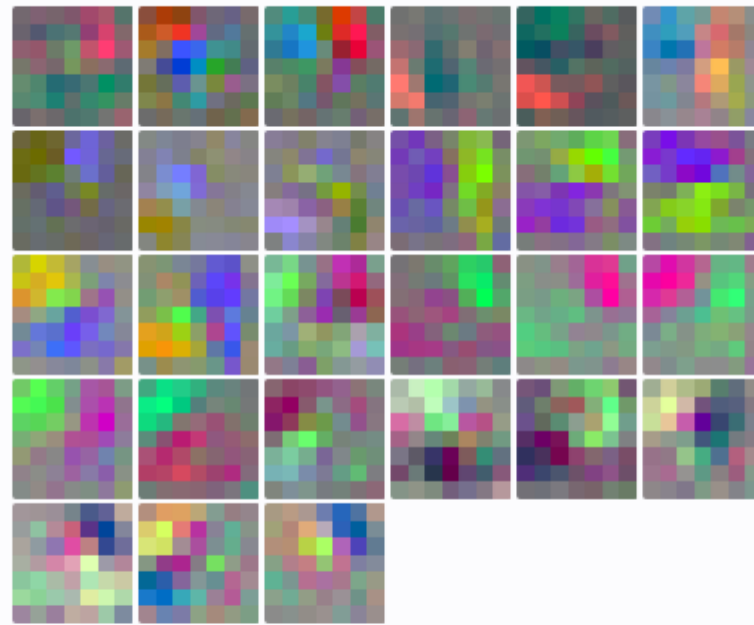
# Inception V1 (Olah et al, 2020)

conv2d0



**Gabor Filters** 44%

Collapse neurons.

Gabor filters are a simple edge detector, highly sensitive to the alignment of the edge. They're almost universally found in the fist layer of vision models. Note that Gabor filters almost always come in pairs of negative reciprocals.
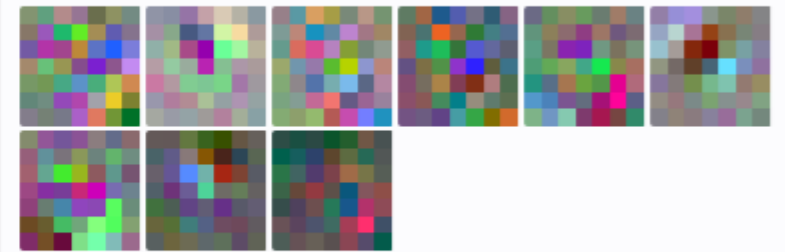
**Color Contrast** 42%

Collapse neurons.

These units detect a color one side of their receptive field, and the opposite color on the other side. Compare to later color contrast (conv2d1, conv2d2, mixed3a, mixed3b).

**Other Units** 14%

Units that don't fit in another category.
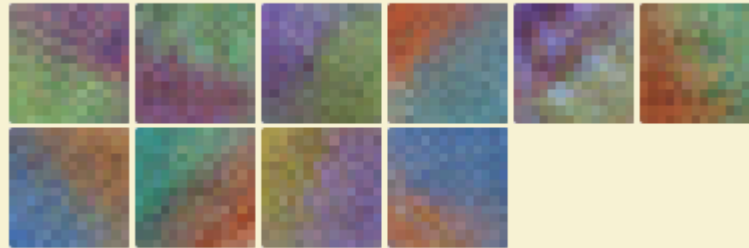
# Inception V1 (Olah et al, 2020)



conv2d1



Gabor Like 17%

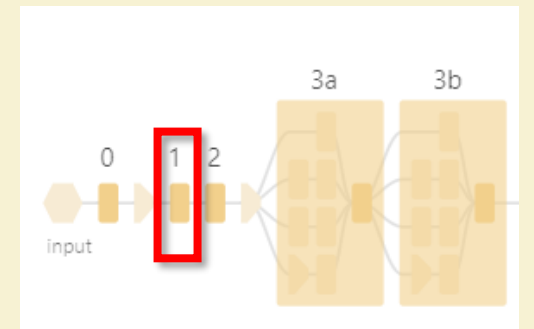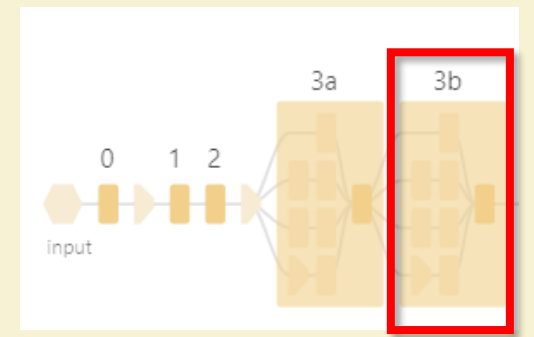Color Contrast 16%

Low Frequency 27%

Complex Gabor 14%

# Inception V1 (Olah et al, 2020)

mixed3b

# SGD Learns simple concepts first



Nakkiran, Kaplun, Kalimeris, Yang, Edelman, Yang, Zhang, B. 2020

# Simplicity bias is a good thing…

A random $f$ fitting $(x_i, y_i)_{i=1..n}$ will never generalize.

# … and a bad thing



| $x$ | $f(x)$ |
|---|---|
| $x_1$ | $y_1$ |
| $x_2$ | $y_2$ |
| $x_3$ | $y_3$ |
| … | … |
| $x_n$ | $y_n$ |
| $x$ | __ |

# Example:

GIF

# The Pitfalls of Simplicity Bias in Neural Networks

**Harshay Shah**
Microsoft Research
harshay.rshah@gmail.com

**Kaustav Tamuly**
Microsoft Research
ktamuly2@gmail.com

**Aditi Raghunathan**
Stanford University
aditir@stanford.edu

**Prateek Jain**
Microsoft Research
prajain@microsoft.com

**Praneeth Netrapalli**
Microsoft Research
praneeth@microsoft.com

# What can we prove?

# (Over-parameterized) Linear Regression

Input: $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^{d+1}$ , $d \gg n$

Goal: Find $w \in \mathbb{R}^d$ s.t. $\langle w, x_i \rangle \approx y_i$

and (more importantly) $\langle w, x \rangle \approx y$ for fresh $(x, y)$



Many solutions for $w$

# (Over-parameterized) Linear Regression

Input: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{d+1}$ , $d \gg n$

Goal: Find $w \in \mathbb{R}^d$ s.t. $\langle w, x_i \rangle \approx y_i$
and (more importantly) $\langle w, x \rangle \approx y$ for fresh $(x, y)$

$$X \quad w \approx y$$

THM: GD / SGD on $\mathcal{L}(w) = \|Xw - y\|^2$ converges to $\arg \min_{w:Xw=y} \|w\|^2$

$$= \lim_{\lambda \to 0} \arg \min_w \|Xw - y\|^2 + \lambda \|w\|^2$$

# Convexity reminders



- $f(x) = x^2$ is convex

- If $f : \mathbb{R}^k \to \mathbb{R}$ convex and $g : \mathbb{R}^d \to \mathbb{R}^k$ linear $f \circ g$ (i.e. $x \mapsto f(g(x))$ ) is convex

- If $f_1, \dots, f_m$ convex and $\alpha_1, \dots, \alpha_m \geq 0$ $\sum \alpha_i f_i$ convex

- If $f$ convex and $\lambda > 0$ $g(x) = f(x) + \lambda \|x\|^2$ strongly convex.

# Linear regression SGD

$$\arg\min\|Xw - y\|^2$$

1. Let $w_0 \leftarrow 0^d$
2. For $t = 0,1,\dots$ :
   - Pick $i \sim [n]$
   - Let $w_{t+1} = w_t - \eta \ \nabla_w(\langle x_i, w\rangle - y_i)^2$

$$x_i^\top(\langle x_i, w\rangle - y_i)$$

Can ignore factor of 2

CLAIM: $\nabla(\langle x_i, w\rangle - y_i)^2 = 2\,x_i^\top x_i w - 2 y_i x_i^\top$

"PF":

i) In one dim $\dfrac{d(xw-y)^2}{dw} = 2x^2 w - 2yx$

ii) Dimensions match

# Linear regression SGD

$$\arg\min\|Xw - y\|^2$$



1. Let $w_0 \leftarrow 0^d$
2. For $t = 0,1,\dots$ :
   - Pick $i \sim [n]$
   - Let $w_{t+1} = w_t - \eta\, x_i^\top(\langle x_i, w\rangle - y_i)$

COR 1: If $w_t \in Span\,\{x_1^\top .. x_n^\top\}$ then $w_{t+1} \in Span\,\{x_1^\top .. x_n^\top\}$

COR 2: If $rank(X) = n$ then
$Xw_\infty = y$ & $w_\infty \in Span\,\{x_1^\top .. x_n^\top\}$

COR 3: $w_\infty = \arg\min\limits_{w:Xw=y}\|w\|^2$

COR 4: $w_\infty = \lim\limits_{\lambda \to 0}\arg\min\limits_{w}\|Xw - y\|^2 + \lambda\|w\|^2$

# GD / SGD dynamics

$$w_{t+1} = w_t - \eta \, \nabla \|Xw_t - y\|^2 \quad = w_t - \eta \nabla (w_t^\top X^\top X w_t - w_t^\top X^\top y)(w_t)$$

$$= w_t - \eta (X^\top X w_t - X^\top y)$$

Let $w_\infty$ s.t. $Xw_\infty = y$ . Then $\quad w_{t+1} = w_t - \eta (X^\top X w_t - X^\top X w_\infty)$

$$w_{t+1} - w_\infty = (I - \eta X^\top X)(w_t - w_\infty)$$

Make progress as long as $0 \prec I - \eta X^\top X \prec 1$: $\quad \eta < \frac{1}{\lambda_1}$ , progress $\approx \frac{\lambda_d}{\lambda_1} = \frac{1}{\kappa}$

$$X^\top X = \begin{pmatrix} \lambda_1 & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \lambda_d \end{pmatrix}$$ Random $X$:



$d < n \qquad\qquad d \approx n \qquad\qquad d > n$

# Actual GD / SGD

$$w_{t+1} = w_t - \eta \, \nabla \| \ldots$$

Let $w_\infty$ s.t. $Xw_\infty = y$



$$w_{t+1} - w_\infty = (I - \eta X^\top X)(w_t - w_\infty)$$

Make progress as long as $0 \prec I - \eta X^\top X \prec 1$:     $\eta < \dfrac{1}{\lambda_1}$ , progress $\approx \dfrac{\lambda_d}{\lambda_1} = \dfrac{1}{\kappa}$

$$X^\top X = \begin{pmatrix} \lambda_1 & \cdots & \\ \vdots & \ddots & \vdots \\ & \cdots & \lambda_d \end{pmatrix}$$     Random $X$:



$d < n$          $d \approx n$          $d > n$

Grosse lecture notes

Hastie-Montanari-Rosset-Tibshirani, 20

# Beyond linear regression

# Implicit regularization in deep networks

## Depth 2 network



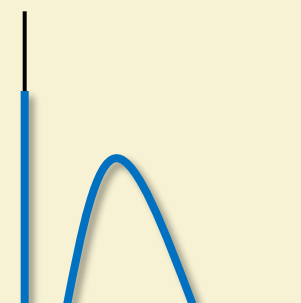Parameter space: $\mathbb{R}^{d \times h + h \times m}$

$Bx = A_2 A_1 x$:
Same expressiveness /
functional space

Different parameter space

## Depth 2 linear network



Parameter space: $\mathbb{R}^{d \times h + h \times m}$

## Linear model



Parameter space: $\mathbb{R}^{d \times m}$

# Linear model



$$x \quad B = A_2 A_1$$

Parameter space: $\mathbb{R}^{d \times m}$

For every loss function $\mathcal{L}$:
$$\min \; \mathcal{L}(B)$$

**BUT**
SGD/GD on ↑

Potentially convex function in $B \in \mathbb{R}^{d \times m}$

# Depth 2 linear network



$$x \quad A_1 \quad A_2$$

Parameter space: $\mathbb{R}^{d \times h + h \times m}$

$$\min \; \mathcal{L}(A_1, A_2)$$

=

SGD/GD on ↑

≠

Non-convex function in $(A_1, A_2) \in \mathbb{R}^{d \times h + h \times m}$

# Gradient flow on deep linear nets

Linear model



$$B = A^2$$

Parameter space: $\mathbb{R}^{d \times m}$

Depth 2 linear network



$$A \qquad A$$

Parameter space: $\mathbb{R}^{d \times h + h \times m}$

Simplifying assumptions: $A_1 = A_2$ symmetric

$$\Rightarrow \ B = A^2, A = \sqrt{B}$$

Analyze GD with $\eta \to 0$ on $\min \tilde{\mathcal{L}}(A)$ where $\tilde{\mathcal{L}}(A) = \mathcal{L}(A^2)$

# Gradient flow on deep linear nets

$$\tilde{\mathcal{L}}(A) = \mathcal{L}(A^2)$$
$$B = A^2$$

$$\frac{dA(t)}{dt} = -\nabla\tilde{\mathcal{L}}(\mathrm{A}(t))$$

GF on linear model:

$$\frac{dB(t)}{dt} = -\nabla\mathcal{L}(\mathrm{B}(t))$$

$$\widetilde{\nabla} \qquad \nabla \qquad \widetilde{\nabla} = A\nabla$$

By chain rule $\quad \nabla\tilde{\mathcal{L}}(A) = \nabla\mathcal{L}(A^2)A = \mathrm{A}\nabla\mathcal{L}(A^2)$

Hence $\quad \dfrac{dA^2(t)}{dt} = \dfrac{dA(t)}{dt} \cdot A \quad = -\widetilde{\nabla} \cdot A = -A \cdot \nabla \cdot A$

GF on deep linear net $B = A^2$:

$$\frac{dB(t)}{dt} = -A\,\nabla\mathcal{L}(B(t))A = -\sqrt{B}\,\nabla\mathcal{L}(B(t))\sqrt{B}$$

"The big get bigger"

GF on deep linear net $B = A^2$:

$$\frac{dB(t)}{dt} = -A\, \nabla\mathcal{L}\big(B(t)\big)A = -\sqrt{B}\, \nabla\mathcal{L}\big(B(t)\big)\sqrt{B}$$

Generally GF on deep linear net $B$ evolves* by

$$\frac{dB(t)}{dt} = -\psi_{B(t)}\big(\nabla\mathcal{L}\big(B(t)\big)\big)$$

$$\psi_B(\nabla) =^* \sum B^\alpha\, \nabla B^{1-\alpha}$$

Gradient flow on a Riemannian Manifold     * not equivalent to $\min \mathcal{L}(B) + \lambda\, R(B)$

Saxe, McClelland, Ganguli 2013
Arora, Cohen, Hazan, 2018
Bah, Rauhut, Terstiege, Westdickenberg, 2019

# Riemannian Manifolds



External description:  A smooth subset $\mathcal{M} \subseteq \mathbb{R}^N$

Intrinsic description:  Set $\mathcal{M}$ with "local geometry" at each $x \in \mathcal{M}$

For every $x \in \mathcal{M}$, tangent space $T_x$ - set of directions we can move in

(Gradient of $f(x)$: shortest direction from $x$ to increase $f$)

local inner product on $T_x$ - defined via PSD matrix $M_x$ on $T_x$

# Learning in different layers

# Cartoon



"dog on the beach"

High level

depends on task/data

~ independent of task/data

Low level

# Non-convexity & symmetry breaking

$\frac{1}{2}$  $+$ $\frac{1}{2}$  $=$ *JUNK*

## Intuition:

Initial weights:     0.49  $+$ 0.51      0.51  $+$ 0.49 

Final weights:     0.01  $+$ 0.99      0.99  $+$ 0.01 

(Too) strong hypothesis: All nets are similar up to transformations, depending on initialization and data

# Linear mode connectivity



$$\mathcal{L}\left(\frac{1}{2}f_{w_\infty} + \frac{1}{2}f_{w'_\infty}\right) \gg \frac{1}{2}\mathcal{L}(f_{w_\infty}) + \frac{1}{2}\mathcal{L}(f_{w'_\infty})$$

Frankle, Dziugaite, Roy, Carbin, 2019

# Linear mode connectivity



$f_{w_0}$

$S$

$\tilde{S}$

$S$

$f_{w_\infty}$

$f_{\tilde{w}_\infty}$

$f_{w'_\infty}$

$\frac{1}{2} f_{w_\infty} + \frac{1}{2} f_{\tilde{w}_\infty}$

$$\mathcal{L}\left(\frac{1}{2} f_{w_\infty} + \frac{1}{2} f_{w'_\infty}\right) \gg \frac{1}{2}\mathcal{L}(f_{w_\infty}) + \frac{1}{2}\mathcal{L}(f_{w'_\infty})$$

$$\mathcal{L}\left(\frac{1}{2} f_{w_\infty} + \frac{1}{2} f_{\tilde{w}_\infty}\right) \approx \frac{1}{2}\mathcal{L}(f_{w_\infty}) + \frac{1}{2}\mathcal{L}(f_{\tilde{w}_\infty})$$

Frankle, Dziugaite, Roy, Carbin, 2019

\* After pruning / from $w_k$ for $k > 0$

# Layers



similarity to final state

Convnet, CIFAR-10
layer (during training)

0% trained    35% trained    75% trained    100% trained

stage1.block1  stage1.block2  stage2.block1  stage2.block2  stage3.block1  stage3.block2  stage3.block3  stage3.block4  stage4.block1  stage4.block2  stage4.block3  stage4.block4  stage5.block1  stage5.block2  stage5.block3  stage5.block4  fc1  fc2  final_linear  fullmodel

ReRnd
0
1
10
100

Raghu, Gilmer, Yosinski, Sohl-Dickstein, 2017
Zhang, Bengio, Singer 2019

randomness just for symmetry breaking!

# Theoretical insights

Intuition: If data doesn't contain "local correlations" then "can't get off the ground" – learning will not succeed.

## Is Deeper Better only when Shallow is Good?

Eran Malach and Shai Shalev-Shwartz

## Failures of Gradient-Based Deep Learning

Shai Shalev-Shwartz[1], Ohad Shamir[2], and Shaked Shammah[1]

# Canonical "hard" example: parities

For $I \subseteq [d], D_I$ defined as: $x \sim \{\pm 1\}^d, y = \prod_{i \in I} x_i$

$$= \begin{cases} -1, & \text{num}_{-1}(x_I) \text{ odd} \\ +1, & \text{num}_{-1}(x_I) \text{ even} \end{cases}$$

Example: $d = 7, I = \{1,3,6,7\}$

| **1** | 2 | **3** | 4 | 5 | **6** | **7** | |
|---|---|---|---|---|---|---|---|
| +1 | +1 | −1 | +1 | +1 | +1 | −1 | +1 |
| −1 | −1 | −1 | +1 | −1 | −1 | +1 | −1 |
| −1 | +1 | +1 | −1 | −1 | +1 | +1 | −1 |

CLAIM: Given $2d$ samples $(x_i, y_i)_{i=1..2d} \sim D_I$ can recover $I$

# Canonical "hard" example: parity

$$= \begin{cases} -1, & \mathrm{num}_{-1}(x_I) \text{ odd} \\ +1, & \mathrm{num}_{-1}(x_I) \text{ even} \end{cases}$$

For $I \subseteq [d]$, $D_I$ defined as: $x \sim \{\pm 1\}^d$, $y = \prod_{i \in I} x_i$

CLAIM: Given $2d$ samples $(x_i, y_i)_{i=1..2d} \sim D_I$ can recover $I$

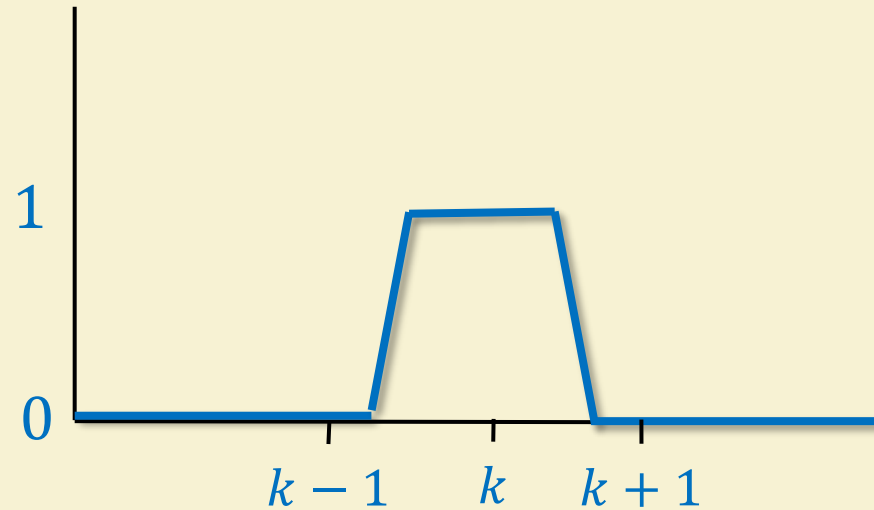PROOF: Let $Z_{i,j} = (1 - x_{i,j})/2$ and $b_i = (1 - y_i)/2$

Let $s_i = 1$ if $i \in I$ and $s_i = 0$ otherwise

Then for every $i$, $\sum_j Z_{i,j} s_j = b_i \pmod 2$

$2d$ linear equations modulo 2 in $d$ variables $s_1, \ldots, s_d$! ∎

# Parities can be expressed by few ReLUs

# ..but are hard to learn

THM: For every* NN architecture $f_w(x)$, SGD on $\min \|f_w(x) - \prod_{i \in I} x_i\|^2$ will require $\exp(\Omega(d))$ steps.

Key fact: For fixed $w$ define r.v. $D_w = \nabla \|f_w(x) - \prod_{i \in I} x_i\|^2(w)$ over the choice of $x \sim \{\pm 1\}^d$, $I \subseteq [d]$.  Then

$$Var(D_w) \leq \frac{poly(d)}{2^d}$$

Possibly large but independent of $I$

Exponentially tiny

Key fact: For fixed $w$ define r.v. $D_w = \nabla\|f_w(x) - \prod_{i \in I} x_i\|^2(w)$ over the choice of $x \sim \{\pm 1\}^d, I \subseteq [d]$ . Then

$$Var(D_w) \leq \frac{poly(d)}{2^d}$$

PF: Fix $w$ & coordinate $i$, and let $D_x = \frac{d}{di}\|f_w(x) - \prod_{i \in I} x_i\|^2(w)$

Then $D_x = \underbrace{2 f_w(x) \frac{d}{di} f_w(x)}_{\text{Independent of } I} - \underbrace{2 \frac{d}{di} f_w(x) \prod_{i \in I} x_i}_{\text{Depends on } I}$

LEMMA: For every $g: \mathbb{R}^d \to \mathbb{R}$, $\left(\mathbb{E}_{x \sim \{\pm 1\}^d} \mathbb{E}_{I \subseteq [d]} g(x) \prod_{i \in I} x_i\right)^2 \leq \frac{\mathbb{E}_x g(x)^2}{2^d}$

LEMMA $\Rightarrow$ FACT $\Rightarrow$ THM

LEMMA: For every $g\colon \mathbb{R}^d \to \mathbb{R}$, $\left(\mathbb{E}_{x\sim\{\pm1\}^d}\mathbb{E}_{I\subseteq[d]}\, g(x)\prod_{i\in I} x_i\right)^2 \leq \dfrac{\mathbb{E}_x g(x)^2}{2^d}$

PF: $\left(\mathbb{E}_{x\sim\{\pm1\}^d}\mathbb{E}_{I\subseteq[d]}\, g(x)\prod_{i\in I} x_i\right)^2 \leq \left(\mathbb{E}_x g(x)^2\right)\cdot\left(\mathbb{E}_x\left(\mathbb{E}_I \prod_{i\in I} x_i\right)^2\right)$

$\mathbb{E}_x\left(\mathbb{E}_I \prod_{i\in I} x_i\right)^2 = \mathbb{E}_x\left(\mathbb{E}_I \prod_{i\in I} x_i\right)\left(\mathbb{E}_J \prod_{j\in J} x_j\right)$

$= \mathbb{E}_I\mathbb{E}_J\,\mathbb{E}_{x\sim\{\pm1\}^d} \prod_{i\in I} x_i \prod_{j\in J} x_j$

$\mathbb{E}_{x\sim\{\pm1\}^d} \prod_{i\in I} x_i \prod_{j\in J} x_j = \prod_{i=1}^{d} \mathbb{E}_{\sigma\in\{\pm1\}}\sigma^{n_i\in\{0,1,2\}} = \begin{cases} 1, & I = J \\ 0, & I \neq J \end{cases}$