Representation, Optimization and Generalization in Deep Learning

Peter Bartlett

UC Berkeley

25 January, 2018

イロト 不同下 イヨト イヨト

Deep neural networks

Game playing



(Jung Yeon-Je/AFP/Getty Images)

Deep neural networks

Image recognition



(Krizhevsky et al, 2012)

Speech recognition



(Graves et al, 2013)

$$h=h_m\circ h_{m-1}\circ\cdots\circ h_1$$

$$h=h_m\circ h_{m-1}\circ\cdots\circ h_1$$

イロン イヨン イヨン イヨン 三日

e.g.,
$$h_i : x \mapsto \sigma(W_i x)$$

 $\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$

$$h=h_m\circ h_{m-1}\circ\cdots\circ h_1$$

イロン イヨン イヨン イヨン 三日

e.g.,
$$h_i : x \mapsto \sigma(W_i x)$$

 $\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$



$$h=h_m\circ h_{m-1}\circ\cdots\circ h_1$$

e.g.,
$$h_i : x \mapsto \sigma(W_i x)$$

 $\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$

 $h_i: x \mapsto r(W_i x)$ $r(v)_i = \max\{0, v_i\}$

イロン イロン イヨン イヨン 三日



$$h=h_m\circ h_{m-1}\circ\cdots\circ h_1$$

e.g.,
$$h_i : x \mapsto \sigma(W_i x)$$

 $\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$

 $h_i: x \mapsto r(W_i x)$ $r(v)_i = \max\{0, v_i\}$





Rich non-parametric family

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions. Nonlinear parameterizations provide better rates of approximation.

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

But...

• Optimization?

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

But...

- Optimization?
 - Nonlinear parameterization.

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

But...

• Optimization?

- Nonlinear parameterization.
- Apparently worse as the depth increases.

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

But...

- Optimization?
 - Nonlinear parameterization.
 - Apparently worse as the depth increases.
- Generalization?

Depth provides an effective way of representing useful features.

Rich non-parametric family

Depth provides parsimonious representions.

Nonlinear parameterizations provide better rates of approximation.

Some functions require much more complexity for a shallow representation.

But...

- Optimization?
 - Nonlinear parameterization.
 - Apparently worse as the depth increases.
- Generalization?
 - What determines the statistical complexity of a deep network?

• Deep residual networks

- Representing with near-identities
- Global optimality of stationary points

- Deep residual networks
 - Representing with near-identities
 - Global optimality of stationary points
- What determines the statistical complexity of a deep network?
 - VC theory: Number of parameters
 - Margins analysis: Size of parameters
 - Understanding generalization failures

Deep residual networks

- Representing with near-identities
- Global optimality of stationary points
- What determines the statistical complexity of a deep network?
 - VC theory: Number of parameters
 - Margins analysis: Size of parameters
 - Understanding generalization failures



(Deep Residual Networks. Kaiming He. 2016)

・ロト <
同 ト <
言 ト <
言 ト ミ のへの
9/59
</p>



(Deep Residual Networks. Kaiming He. 2016)

Revolution of Depth

AlexNet, 8 layers (ILSVRC 2012)



VGG, 19 layers (ILSVRC 2014)



layers (4)

(Deep Residual Networks. Kaiming He. 2016)

Ť

Revolution of Depth

AlexNet, 8 layers (ILSVRC 2012) VGG, 19 layers (ILSVRC 2014) ResNet, 152 layers (ILSVRC 2015)

(Deep Residual Networks. Kaiming He. 2016)



(Deep Residual Networks. Kaiming He. 2016)



(Deep Residual Networks. Kaiming He. 2016)

Advantages

• With zero weights, the network computes the identity.

Advantages

- With zero weights, the network computes the identity.
- Identity connections provide useful feedback throughout the network.

Advantages

- With zero weights, the network computes the identity.
- Identity connections provide useful feedback throughout the network.



(Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. 2016)





(Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. 2016)

Deep Residual Networks: Competition Successes

ImageNet Large Scale Visual Recognition Challenge



(http://image-net.org/)

Deep Residual Networks: Competition Successes

ImageNet Large Scale Visual Recognition Challenge



(http://image-net.org/)

First place:

- Object detection: 16% better than next best
- Object localization: 27% better than next best

COCO (Common Objects in Context)



(http://mscoco.org/)

COCO (Common Objects in Context)



(http://mscoco.org/)

First place:

- Detection: 11% better than next best
- Segmentation: 12% better than next best

Why?

- What is behind the success of residual networks?
- What is important for their performance?
Every invertible* A can be written as

$$A = (I + A_m) \cdots (I + A_1),$$

where $||A_i|| = O(1/m)$.

(Hardt and Ma, 2016)

19/59

イロン イヨン イヨン イヨン 三日

Provided det(A) > 0.

*

2 For a linear Gaussian model,

$$y = \mathbf{A}x + \epsilon, \qquad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$

(Hardt and Ma, 2016)

Por a linear Gaussian model,

$$y = \mathbf{A}x + \epsilon, \qquad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

consider choosing A_1, \ldots, A_m to minimize quadratic loss:

 $\mathbb{E}\|(I+A_m)\cdots(I+A_1)x-y\|^2.$

(Hardt and Ma, 2016)

Por a linear Gaussian model,

$$y = \mathbf{A}x + \epsilon, \qquad \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

consider choosing A_1, \ldots, A_m to minimize quadratic loss:

$$\mathbb{E}\|(I+A_m)\cdots(I+A_1)x-y\|^2.$$

If $||A_i|| < 1$, every stationary point of the quadratic loss is a global optimum:

$$\forall i, \ \nabla_{A_i} \mathbb{E} \| (I + A_m) \cdots (I + A_1) x - y \|^2 = 0$$

$$\Rightarrow \qquad \mathbf{A} = (I + A_m) \cdots (I + A_1).$$

(Hardt and Ma, 2016)

Outline

Deep residual networks

- Representing with near-identities
- Global optimality of stationary points
- What determines the statistical complexity of a deep network?



Steve Evans Berkeley, Stat/Math



Phil Long Google

The computation of a smooth invertible map h can be spread throughout a deep network,

$$h_m \circ h_{m-1} \circ \cdots \circ h_1 = h,$$

so that all layers compute near-identity functions:

$$\|h_i - \mathrm{Id}\|_L = O\left(\frac{\log m}{m}\right).$$

The computation of a smooth invertible map h can be spread throughout a deep network,

$$h_m \circ h_{m-1} \circ \cdots \circ h_1 = h,$$

so that all layers compute near-identity functions:

$$\|h_i - \mathrm{Id}\|_L = O\left(\frac{\log m}{m}\right).$$

Definition: the Lipschitz seminorm of f satisfies, for all x, y,

 $||f(x) - f(y)|| \le ||f||_L ||x - y||.$

イロト 不得下 イヨト イヨト 二日

The computation of a smooth invertible map h can be spread throughout a deep network,

$$h_m \circ h_{m-1} \circ \cdots \circ h_1 = h,$$

so that all layers compute near-identity functions:

$$\|h_i - \mathrm{Id}\|_L = O\left(\frac{\log m}{m}\right).$$

Think of the functions h_i as near-identity maps that might be computed as

$$h_i(x) = x + \underbrace{A\sigma(Bx)}_{i}.$$

The computation of a smooth invertible map h can be spread throughout a deep network,

$$h_m \circ h_{m-1} \circ \cdots \circ h_1 = h,$$

so that all layers compute near-identity functions:

$$\|h_i - \mathrm{Id}\|_L = O\left(\frac{\log m}{m}\right).$$

Think of the functions h_i as near-identity maps that might be computed as

$$h_i(x) = x + \underbrace{A\sigma(Bx)}_{i}.$$

As the network gets deeper, the functions $x \mapsto A\sigma(Bx)$ can get flatter.

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

Differentiable,

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$
- Lipschitz inverse: For some M > 0, $||h^{-1}||_L \le M$.

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$
- Lipschitz inverse: For some M > 0, $||h^{-1}||_L \le M$.
- Solution Positive orientation: For some x_0 , $det(Dh(x_0)) > 0$.

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$
- Lipschitz inverse: For some M > 0, $||h^{-1}||_L \le M$.
- Solution Positive orientation: For some x_0 , $det(Dh(x_0)) > 0$.

Then for all *m*, there are *m* functions $h_1, \ldots, h_m : \mathbb{R}^d \to \mathbb{R}^d$ satisfying $\|h_i - \mathrm{Id}\|_L = O(\log m/m)$

Theorem

Consider a function $h : \mathbb{R}^d \to \mathbb{R}^d$ on a bounded domain $\mathcal{X} \subset \mathbb{R}^d$. Suppose that it is

- Differentiable,
- Invertible,
- Smooth: For some $\alpha > 0$ and all x, y, u, $\|Dh(y) - Dh(x)\| \le \alpha \|y - x\|.$
- Lipschitz inverse: For some M > 0, $||h^{-1}||_L \le M$.
- Solution Positive orientation: For some x_0 , $det(Dh(x_0)) > 0$.

Then for all *m*, there are *m* functions $h_1, \ldots, h_m : \mathbb{R}^d \to \mathbb{R}^d$ satisfying $\|h_i - \operatorname{Id}\|_L = O(\log m/m)$ and $h_m \circ h_{m-1} \circ \cdots \circ h_1 = h$ on \mathcal{X} .

Key ideas

1	Assume	h(0)	= 0 and	Dh(0)	$= \mathrm{Id}$
---	--------	------	----------------	-------	-----------------

Key ideas

• Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).

Key ideas

- Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).
- 2 Construct the h_i so that

$$h_1(x) = \frac{h(a_1 x)}{a_1}$$

Key ideas

- Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).
- **2** Construct the h_i so that

$$h_1(x) = rac{h(a_1x)}{a_1}$$

 $h_2(h_1(x)) = rac{h(a_2x)}{a_2}$

h

Key ideas

• Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).

2 Construct the h_i so that

hat

$$h_1(x) = \frac{h(a_1x)}{a_1}$$

$$h_2(h_1(x)) = \frac{h(a_2x)}{a_2}$$

$$\vdots$$

$$m(\cdots(h_1(x))\cdots) = \frac{h(a_mx)}{a_m},$$

Key ideas

• Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).

• Construct the h_i so that $h_1(x) = \frac{h(a_1x)}{a_1}$ $h_2(h_1(x)) = \frac{h(a_2x)}{a_2}$ \vdots $h_m(\cdots(h_1(x))\cdots) = \frac{h(a_mx)}{a_m}$,

9 Pick $a_m = 1$ so $h_m \circ \cdots \circ h_1 = h$.

Key ideas

• Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).

2 Construct the h_i so that

$$h_1(x) = \frac{h(a_1x)}{a_1}$$
$$h_2(h_1(x)) = \frac{h(a_2x)}{a_2}$$
$$\vdots$$
$$h_m(\cdots(h_1(x))\cdots) = \frac{h(a_mx)}{a_m},$$

O Pick a_m = 1 so h_m ◦ · · · ◦ h₁ = h.
O Ensure that a₁ is small enough that h₁ ≈ Dh(0) = Id.

Key ideas

- Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).
- 2 Construct the h_i so that

that

$$h_1(x) = \frac{h(a_1x)}{a_1}$$

$$h_2(h_1(x)) = \frac{h(a_2x)}{a_2}$$

$$\vdots$$

$$h_m(\cdots(h_1(x))\cdots) = \frac{h(a_mx)}{a_m},$$

- Ensure that a_1 is small enough that $h_1 \approx Dh(0) = \text{Id.}$
- **(**) Ensure that a_i and a_{i+1} are sufficiently close that $h_i \approx \text{Id.}$

Key ideas

- Assume h(0) = 0 and Dh(0) = Id (else shift and linearly transform).
- 2 Construct the h_i so that

hat

$$h_1(x) = \frac{h(a_1x)}{a_1}$$

$$h_2(h_1(x)) = \frac{h(a_2x)}{a_2}$$

$$\vdots$$

$$h(\cdots(h_1(x))\cdots) = \frac{h(a_mx)}{a_m},$$

3 Pick
$$a_m = 1$$
 so $h_m \circ \cdots \circ h_1 = h$.

h

- Ensure that a_1 is small enough that $h_1 \approx Dh(0) = \text{Id.}$
- **(5)** Ensure that a_i and a_{i+1} are sufficiently close that $h_i \approx \text{Id.}$
- **(**) Show $||h_i \text{Id}||_L$ is small on small and large scales (c.f. $a_i a_{i-1}$).

The computation of a smooth invertible map h can be spread throughout a deep network,

$$h_m \circ h_{m-1} \circ \cdots \circ h_1 = h,$$

so that all layers compute near-identity functions:

$$\|h_i - \mathrm{Id}\|_L = O\left(\frac{\log m}{m}\right).$$

• Deeper networks allow flatter nonlinear functions at each layer.

- Deep residual networks
 - Representing with near-identities
 - Global optimality of stationary points
- What determines the statistical complexity of a deep network?

Stationary points

Result

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$.

Stationary points

Result

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$.

• e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h^* an empirical risk minimizer.

Stationary points

Result

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$. Consider a function $h = h_m \circ \cdots \circ h_1$, where $||h_i - \mathrm{Id}||_I \le \epsilon < 1$.

• e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h^* an empirical risk minimizer.

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$. Consider a function $h = h_m \circ \cdots \circ h_1$, where $\|h_i - \mathrm{Id}\|_L \le \epsilon < 1$. Then for all *i*,

$$\|D_{h_i}Q(h)\| \geq rac{(1-\epsilon)^{m-1}}{\|h-h^*\|} \left(Q(h)-Q(h^*)
ight).$$

• e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h^* an empirical risk minimizer.

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$. Consider a function $h = h_m \circ \cdots \circ h_1$, where $\|h_i - \mathrm{Id}\|_L \le \epsilon < 1$. Then for all *i*,

$$\|D_{h_i}Q(h)\| \geq rac{(1-\epsilon)^{m-1}}{\|h-h^*\|} \left(Q(h)-Q(h^*)
ight).$$

e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h* an empirical risk minimizer.
D_{hi}Q is a Fréchet derivative; ||h|| is the induced norm. What the theorem says

What the theorem says

• If the composition *h* is sub-optimal and each function *h_i* is a near-identity, then there is a downhill direction in function space: the functional gradient of *Q* wrt *h_i* is non-zero.
- If the composition h is sub-optimal and each function h_i is a near-identity, then there is a downhill direction in function space: the functional gradient of Q wrt h_i is non-zero.
- Thus every stationary point is a global optimum.

- If the composition h is sub-optimal and each function h_i is a near-identity, then there is a downhill direction in function space: the functional gradient of Q wrt h_i is non-zero.
- Thus every stationary point is a global optimum.
- There are no local minima and no saddle points.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

• The theorem does not say there are no local minima of a deep residual network of ReLUs or sigmoids with a fixed architecture.

- The theorem does not say there are no local minima of a deep residual network of ReLUs or sigmoids with a fixed architecture.
- Except at the global minimum, there is a downhill direction in function space. But this direction might be orthogonal to functions that can be computed with this fixed architecture.

- The theorem does not say there are no local minima of a deep residual network of ReLUs or sigmoids with a fixed architecture.
- Except at the global minimum, there is a downhill direction in function space. But this direction might be orthogonal to functions that can be computed with this fixed architecture.
- We should expect suboptimal stationary points in the ReLU or sigmoid parameter space, but these cannot arise because of interactions between parameters in different layers; they arise only within a layer.

Result

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$. Consider a function $h = h_m \circ \cdots \circ h_1$, where $\|h_i - \mathrm{Id}\|_L \le \epsilon < 1$. Then for all *i*,

$$\|D_{h_i}Q(h)\| \geq rac{(1-\epsilon)^{m-1}}{\|h-h^*\|} \left(Q(h)-Q(h^*)
ight).$$

e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h* an empirical risk minimizer.
D_{hi}Q is a Fréchet derivative; ||h|| is the induced norm.

If $\|f - \operatorname{Id}\|_L \le \alpha < 1$ then

If $\|f - \operatorname{Id}\|_L \le \alpha < 1$ then

31 / 59

1 f is invertible.

- If $\|f \operatorname{Id}\|_L \le \alpha < 1$ then
 - f is invertible.
 - **2** $||f||_L \le 1 + \alpha$ and $||f^{-1}||_L \le 1/(1-\alpha)$.

- If $\|f \operatorname{Id}\|_L \le \alpha < 1$ then
 - f is invertible.
 - 2 $||f||_L \le 1 + \alpha$ and $||f^{-1}||_L \le 1/(1-\alpha)$.
 - For $F(g) = f \circ g$, $||DF(g) Id|| \le \alpha$.

- If $\|f \operatorname{Id}\|_L \le \alpha < 1$ then
 - f is invertible.
 - **2** $||f||_L \le 1 + \alpha$ and $||f^{-1}||_L \le 1/(1 \alpha)$.
 - For $F(g) = f \circ g$, $||DF(g) Id|| \le \alpha$.

• ||f|| denotes the induced norm: $||g|| := \sup\left\{\frac{||g(x)||}{||x||} : ||x|| > 0\right\}$.

・ロ ・ ・ 日 ・ ・ 目 ・ く 目 ・ し 目 ・ の へ (や 31 / 59

- If $\|f \operatorname{Id}\|_L \le \alpha < 1$ then
 - f is invertible.
 - **2** $||f||_L \le 1 + \alpha$ and $||f^{-1}||_L \le 1/(1-\alpha)$.
 - For $F(g) = f \circ g$, $||DF(g) \mathrm{Id}|| \le \alpha$.
 - For a linear map h (such as $DF(g) \mathrm{Id}$), $||h|| = ||h||_L$.
- ||f|| denotes the induced norm: $||g|| := \sup \left\{ \frac{||g(x)||}{||x||} : ||x|| > 0 \right\}.$

イロト 不得下 イヨト イヨト 二日

Proof ideas (2)

Projection theorem implies

$$Q(h)=rac{1}{2}\mathbb{E}\left\|h(X)-h^*(X)
ight\|_2^2+ ext{constant}$$

Proof ideas (2)

Projection theorem implies

$$Q(h) = rac{1}{2}\mathbb{E}\left\|h(X) - h^*(X)
ight\|_2^2 + ext{constant.}$$



$$D_{h_i}Q(h) = \mathbb{E}\left[(h(X) - h^*(X)) \cdot \operatorname{ev}_X \circ D_{h_i}h\right]$$

Proof ideas (2)

Projection theorem implies

$$Q(h) = rac{1}{2}\mathbb{E}\left\|h(X) - h^*(X)
ight\|_2^2 + ext{constant.}$$

$$D_{h_i}Q(h) = \mathbb{E}\left[(h(X) - h^*(X)) \cdot \operatorname{ev}_X \circ D_{h_i}h\right]$$

• ev_x is the evaluation functional, $ev_x(f) = f(x)$.

Proof ideas (2)

Projection theorem implies

$$Q(h) = rac{1}{2}\mathbb{E}\left\|h(X) - h^*(X)
ight\|_2^2 + ext{constant.}$$

2 Then

$$\mathcal{D}_{h_i}Q(h) = \mathbb{E}\left[(h(X) - h^*(X)) \cdot \operatorname{ev}_X \circ \mathcal{D}_{h_i}h
ight].$$

③ It is possible to choose a direction Δ s.t. $\|\Delta\| = 1$ and

 $D_{h_i}Q(h)(\Delta) = c\mathbb{E} \|h(X) - h^*(X)\|_2^2.$

• ev_x is the evaluation functional, $ev_x(f) = f(x)$.

Proof ideas (2)

Projection theorem implies

$$Q(h) = rac{1}{2}\mathbb{E}\left\|h(X) - h^*(X)
ight\|_2^2 + ext{constant.}$$

2 Then

$$\mathcal{D}_{h_i}Q(h) = \mathbb{E}\left[(h(X) - h^*(X)) \cdot \operatorname{ev}_X \circ \mathcal{D}_{h_i}h
ight].$$

③ It is possible to choose a direction Δ s.t. $\|\Delta\| = 1$ and

$$D_{h_i}Q(h)(\Delta) = c\mathbb{E} \|h(X) - h^*(X)\|_2^2.$$

• Because the h_i s are near-identities,

$$c \geq rac{(1-\epsilon)^{m-1}}{\|h-h^*\|}.$$

• ev_x is the evaluation functional, $ev_x(f) = f(x)$.

Result

For (X, Y) with an arbitrary joint distribution, define the squared error,

$$Q(h) = \frac{1}{2}\mathbb{E} \|h(X) - Y\|_2^2,$$

define the minimizer $h^*(x) = \mathbb{E}[Y|X = x]$. Consider a function $h = h_m \circ \cdots \circ h_1$, where $\|h_i - \mathrm{Id}\|_L \le \epsilon < 1$. Then for all *i*,

$$\|D_{h_i}Q(h)\| \geq rac{(1-\epsilon)^{m-1}}{\|h-h^*\|} \left(Q(h)-Q(h^*)
ight).$$

e.g., if (X, Y) is uniform on a training sample, then Q is empirical risk and h* an empirical risk minimizer.
D_{hi}Q is a Fréchet derivative; ||h|| is the induced norm.

• If the mapping is not invertible?

If the mapping is not invertible?
 e.g., h: ℝ^d → ℝ?

If the mapping is not invertible?
 e.g., h: ℝ^d → ℝ?

If *h* can be extended to a bi-Lipschitz mapping to \mathbb{R}^d , it can be represented with flat functions at each layer.

- If the mapping is not invertible?
 - e.g., $h : \mathbb{R}^d \to \mathbb{R}$?

If *h* can be extended to a bi-Lipschitz mapping to \mathbb{R}^d , it can be represented with flat functions at each layer. What if it cannot?

- If the mapping is not invertible?
 - e.g., $h: \mathbb{R}^d \to \mathbb{R}$?

If *h* can be extended to a bi-Lipschitz mapping to \mathbb{R}^d , it can be represented with flat functions at each layer. What if it cannot?

• Implications for optimization?

- If the mapping is not invertible?
 - e.g., $h : \mathbb{R}^d \to \mathbb{R}$?

If *h* can be extended to a bi-Lipschitz mapping to \mathbb{R}^d , it can be represented with flat functions at each layer. What if it cannot?

 Implications for optimization? Related to Polyak-Łojasiewicz function classes; proximal algorithms for these classes converge quickly to stationary points.

- If the mapping is not invertible?
 - e.g., $h : \mathbb{R}^d \to \mathbb{R}$?

If *h* can be extended to a bi-Lipschitz mapping to \mathbb{R}^d , it can be represented with flat functions at each layer. What if it cannot?

- Implications for optimization? Related to Polyak-Łojasiewicz function classes; proximal algorithms for these classes converge quickly to stationary points.
- Do stochastic gradient methods produce near-identities?

Deep residual networks

- Representing with near-identities
- Global optimality of stationary points

• What determines the statistical complexity of a deep network?

- VC theory: Number of parameters
- Margins analysis: Size of parameters
- Understanding generalization failures

< □ > < □ > < 直 > < 直 > < 直 > < 直 > 36 / 59 • Assume network maps to {-1,1}. (Threshold its output)

- Assume network maps to {-1,1}. (Threshold its output)
- Data generated by a probability distribution P on $\mathcal{X} \times \{-1, 1\}$.

- Assume network maps to {-1,1}. (Threshold its output)
- Data generated by a probability distribution P on $\mathcal{X} \times \{-1, 1\}$.
- Want to choose a function f such that P(f(x) ≠ y) is small (near optimal).

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$. For every prob distribution P on $\mathcal{X} \times \{-1, 1\}$, with probability $1 - \delta$ over n iid examples $(x_1, y_1), \ldots, (x_n, y_n)$, every f in \mathcal{F} satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left(\operatorname{VCdim}(\mathcal{F}) + \log(1/\delta) \right) \right)^{1/2}.$$

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$. For every prob distribution P on $\mathcal{X} \times \{-1, 1\}$, with probability $1 - \delta$ over n iid examples $(x_1, y_1), \ldots, (x_n, y_n)$, every f in \mathcal{F} satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left(\operatorname{VCdim}(\mathcal{F}) + \log(1/\delta) \right) \right)^{1/2}.$$

For uniform bounds (that is, for all distributions and all *f* ∈ *F*, proportions are close to probabilities), this inequality is tight within a constant factor.

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$. For every prob distribution P on $\mathcal{X} \times \{-1, 1\}$, with probability $1 - \delta$ over n iid examples $(x_1, y_1), \ldots, (x_n, y_n)$, every f in \mathcal{F} satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left(\operatorname{VCdim}(\mathcal{F}) + \log(1/\delta) \right) \right)^{1/2}.$$

- For uniform bounds (that is, for all distributions and all *f* ∈ *F*, proportions are close to probabilities), this inequality is tight within a constant factor.
- For neural networks, VC-dimension:

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$. For every prob distribution P on $\mathcal{X} \times \{-1, 1\}$, with probability $1 - \delta$ over n iid examples $(x_1, y_1), \ldots, (x_n, y_n)$, every f in \mathcal{F} satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left(\operatorname{VCdim}(\mathcal{F}) + \log(1/\delta) \right) \right)^{1/2}.$$

- For uniform bounds (that is, for all distributions and all *f* ∈ *F*, proportions are close to probabilities), this inequality is tight within a constant factor.
- For neural networks, VC-dimension:
 - increases with number of parameters

Theorem (Vapnik and Chervonenkis)

Suppose $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$. For every prob distribution P on $\mathcal{X} \times \{-1, 1\}$, with probability $1 - \delta$ over n iid examples $(x_1, y_1), \ldots, (x_n, y_n)$, every f in \mathcal{F} satisfies

$$P(f(x) \neq y) \leq \frac{1}{n} \left| \{i : f(x_i) \neq y_i\} \right| + \left(\frac{c}{n} \left(\operatorname{VCdim}(\mathcal{F}) + \log(1/\delta) \right) \right)^{1/2}.$$

- For uniform bounds (that is, for all distributions and all *f* ∈ *F*, proportions are close to probabilities), this inequality is tight within a constant factor.
- For neural networks, VC-dimension:
 - increases with number of parameters
 - depends on nonlinearity and depth
Consider the class \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Consider the class \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Piecewise constant (linear threshold units):

$$\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(p).$$

(Baum and Haussler, 1989)

Consider the class \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Piecewise constant (linear threshold units):

Piecewise linear (ReLUs):

$$\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(p).$$

(Baum and Haussler, 1989)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(pL).$

(B., Harvey, Liaw, Mehrabian, 2017)

・ロ ・ ・ 一 ・ ・ 注 ・ ・ 注 ・ う へ C
38 / 59

Consider the class \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Piecewise constant (linear threshold units):

Piecewise linear (ReLUs):

Piecewise polynomial:

$$\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(p).$$

(Baum and Haussler, 1989)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(pL).$

(B., Harvey, Liaw, Mehrabian, 2017)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(pL^2).$

(B., Maiorov, Meir, 1998)

Consider the class \mathcal{F} of $\{-1, 1\}$ -valued functions computed by a network with L layers, p parameters, and k computation units with the following nonlinearities:

Piecewise constant (linear threshold units):

Piecewise linear (ReLUs):

Piecewise polynomial:

Generation Sigmoid:

$$\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(p).$$

(Baum and Haussler, 1989)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(pL).$

(B., Harvey, Liaw, Mehrabian, 2017)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(pL^2).$

(B., Maiorov, Meir, 1998)

 $\operatorname{VCdim}(\mathcal{F}) = \tilde{O}(p^2k^2).$

(Karpinsky and MacIntyre, 1994)

Generalization in Neural Networks: Number of Parameters



- Deep residual networks
 - Representing with near-identities
 - Global optimality of stationary points
- What determines the statistical complexity of a deep network?
 - VC theory: Number of parameters
 - Margins analysis: Size of parameters
 - Understanding generalization failures

• Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.
- For a pattern-label pair (x, y) ∈ X × {−1,1}, if yf(x) > 0 then f classifies x correctly.

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.
- For a pattern-label pair (x, y) ∈ X × {−1,1}, if yf(x) > 0 then f classifies x correctly.
- We call yf(x) the margin of f on x.

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.
- For a pattern-label pair (x, y) ∈ X × {−1,1}, if yf(x) > 0 then f classifies x correctly.
- We call yf(x) the margin of f on x.
- We can view a larger margin as a more confident correct classification.

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.
- For a pattern-label pair (x, y) ∈ X × {−1,1}, if yf(x) > 0 then f classifies x correctly.
- We call yf(x) the margin of f on x.
- We can view a larger margin as a more confident correct classification.
- Minimizing a continuous loss, such as

$$\sum_{i=1}^n \left(f(X_i)-Y_i\right)^2,$$

encourages large margins.

イロン 不通 と イヨン イヨン

- Consider a real-valued function $f : \mathcal{X} \to \mathbb{R}$ used for classification.
- The prediction on $x \in \mathcal{X}$ is $sign(f(x)) \in \{-1, 1\}$.
- For a pattern-label pair (x, y) ∈ X × {−1,1}, if yf(x) > 0 then f classifies x correctly.
- We call yf(x) the margin of f on x.
- We can view a larger margin as a more confident correct classification.
- Minimizing a continuous loss, such as

$$\sum_{i=1}^n \left(f(X_i)-Y_i\right)^2,$$

encourages large margins.

• For large-margin classifiers, we should expect the fine-grained details of *f* to be less important.

イロン 不通 と イヨン イヨン

Theorem (B., 1996)

↓ □ → ↓ □ → ↓ = → ↓ = → へで ↓ 2 / 59

Theorem (B., 1996)

n training examples

イロン イヨン イヨン イヨン 三日

42 / 59

```
(X_1, Y_1), \ldots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}
```

Theorem (B., 1996)

n training examples

42 / 59

 $(X_1, Y_1), \ldots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$ $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$

Theorem (B., 1996)

n training examples

 $(X_1, Y_1), \ldots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$ $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$

 $\Pr(\operatorname{sign}(f(X)) \neq Y)$

Theorem (B., 1996)

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

• The bound depends on the margin loss plus an error term.

Theorem (B., 1996)

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

- The bound depends on the margin loss plus an error term.
- Minimizing quadratic loss or cross-entropy loss leads to large margins.

Theorem (B., 1996)

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

- The bound depends on the margin loss plus an error term.
- Minimizing quadratic loss or cross-entropy loss leads to large margins.
- fat_F(γ) is a scale-sensitive version of VC-dimension. Unlike the VC-dimension, it need not grow with the number of parameters.

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $\|w\|_1 \leq B$, then

$$\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$$

- The bound depends on the margin loss plus an error term.
- Minimizing quadratic loss or cross-entropy loss leads to large margins.
- fat_F(γ) is a scale-sensitive version of VC-dimension. Unlike the VC-dimension, it need not grow with the number of parameters.

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $||w||_1 \leq B$, then

$$\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$$

• Same ideas used to give rigorous dimension-independent generalization bounds for SVMs (B. and Shawe-Taylor, 1999)

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $||w||_1 \leq B$, then

 $\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$

• Same ideas used to give rigorous dimension-independent generalization bounds for SVMs (B. and Shawe-Taylor, 1999)

• ... and margins analysis of AdaBoost.

(Schapire, Freund, B., Lee, 1998)

◆□ → ◆□ → ◆ 三 → ◆ 三 → ○ Q (0) 42 / 59

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $||w||_1 \leq B$, then

 $\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$

・ロン ・四 と ・ ヨ と ・ ヨ

42 / 59

• The scale of functions $f \in \mathcal{F}$ is important.

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $\|w\|_1 \leq B$, then

 $\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$

- The scale of functions $f \in \mathcal{F}$ is important.
- Bigger fs give bigger margins, so $fat_{\mathcal{F}}(\gamma)$ should be bigger.

1. With high probability over *n* training examples $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathcal{X} \times \{\pm 1\}$, every $f \in \mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ has $\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\sqrt{\frac{\operatorname{fat}_{\mathcal{F}}(\gamma)}{n}}\right).$

2. If functions in \mathcal{F} are computed by two-layer sigmoid networks with each unit's weights bounded in 1-norm, that is, $\|w\|_1 \leq B$, then

 $\operatorname{fat}_{\mathcal{F}}(\gamma) = \tilde{O}((B/\gamma)^2).$

- The scale of functions $f \in \mathcal{F}$ is important.
- Bigger fs give bigger margins, so $fat_{\mathcal{F}}(\gamma)$ should be bigger.
- The output y of a sigmoid layer has ||y||∞ ≤ 1, so ||w||₁ ≤ B controls the scale of f.

イロト イポト イヨト イヨト

43 / 59



• Qualitative behavior explained by small weights theorem.



simons.berkeley.edu

• Qualitative behavior explained by small weights theorem.



simons.berkeley.edu

- Qualitative behavior explained by small weights theorem.
- How to measure the complexity of a ReLU network?

- Deep residual networks
 - Representing with near-identities
 - Global optimality of stationary points
- What determines the statistical complexity of a deep network?
 - VC theory: Number of parameters
 - Margins analysis: Size of parameters
 - Understanding generalization failures

CIFAR10



http://corochann.com/

Stochastic Gradient Training Error on CIFAR10



(Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, 2017) 46 / 59

Training margins on CIFAR10 with true and random labels



Training margins on CIFAR10 with true and random labels



• How does this match the large margin explanation?

Training margins on CIFAR10 with true and random labels



- How does this match the large margin explanation?
- Need to account for the scale of the neural network functions.
Training margins on CIFAR10 with true and random labels



- How does this match the large margin explanation?
- Need to account for the scale of the neural network functions.
- What is the appropriate notion of the size of these functions?

Spectrally-normalized margin bounds for neural networks. B., Dylan J. Foster, Matus Telgarsky, NIPS 2017. arXiv:1706.08498



Dylan Foster Cornell



Matus Telgarsky UIUC

New results for generalization in deep ReLU networks

• Measuring the size of functions computed by a network of ReLUs.

New results for generalization in deep ReLU networks

• Measuring the size of functions computed by a network of ReLUs. (c.f. sigmoid networks: the output y of a layer has $||y||_{\infty} \le 1$, so $||w||_1 \le B$ keeps the scale under control.)

New results for generalization in deep ReLU networks

- Measuring the size of functions computed by a network of ReLUs. (c.f. sigmoid networks: the output y of a layer has $||y||_{\infty} \le 1$, so $||w||_1 \le B$ keeps the scale under control.)
- Large multiclass versus binary classification.

New results for generalization in deep ReLU networks

- Measuring the size of functions computed by a network of ReLUs. (c.f. sigmoid networks: the output y of a layer has $||y||_{\infty} \le 1$, so $||w||_1 \le B$ keeps the scale under control.)
- Large multiclass versus binary classification.

Definitions

• Consider operator norms: For a matrix A_i,

$$||A_i||_* := \sup_{||x|| \le 1} ||A_ix||_*$$

New results for generalization in deep ReLU networks

- Measuring the size of functions computed by a network of ReLUs. (c.f. sigmoid networks: the output y of a layer has $||y||_{\infty} \le 1$, so $||w||_1 \le B$ keeps the scale under control.)
- Large multiclass versus binary classification.

Definitions

• Consider operator norms: For a matrix A_i ,

$$\|A_i\|_* := \sup_{\|x\| \le 1} \|A_i x\|.$$

• Multiclass margin function for $f : \mathcal{X} \to \mathbb{R}^m$, $y \in \{1, \dots, m\}$:

$$M(f(x), y) = f(x)_y - \max_{i \neq y} f(x)_i.$$

With high probability, every f_A

With high probability, every f_A

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_{\mathcal{A}}(x) := \sigma_{\mathcal{L}}(A_{\mathcal{L}}\sigma_{\mathcal{L}-1}(A_{\mathcal{L}-1}\cdots\sigma_{1}(A_{1}x)\cdots)).$$

<ロ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ 50 / 59

With high probability, every f_A

satisfies

```
\Pr(M(f_A(X), Y) \le 0) \le
```

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_{\mathcal{A}}(x) := \sigma_{L}(A_{L}\sigma_{L-1}(A_{L-1}\cdots\sigma_{1}(A_{1}x)\cdots)).$$

With high probability, every f_A

satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma]$$

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_{\mathcal{A}}(x) := \sigma_L(A_L \sigma_{L-1}(A_{L-1} \cdots \sigma_1(A_1 x) \cdots)).$$

With high probability, every f_A

satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_A(x) := \sigma_L(A_L\sigma_{L-1}(A_{L-1}\cdots\sigma_1(A_1x)\cdots)).$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

With high probability, every f_A with $R_A \leq r$ satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_A(x) := \sigma_L(A_L\sigma_{L-1}(A_{L-1}\cdots\sigma_1(A_1x)\cdots)).$$

.

Scale of f_A : $R_A := \prod_{i=1}^{L} ||A_i||_*$

With high probability, every f_A with $R_A \leq r$ satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Definitions

Network with L layers, parameters A_1, \ldots, A_L :

$$f_A(x) := \sigma_L(A_L\sigma_{L-1}(A_{L-1}\cdots\sigma_1(A_1x)\cdots)).$$

Scale of f_A : $R_A := \prod_{i=1}^{L} \|A_i\|_* \left(\sum_{i=1}^{L} \frac{\|A_i\|_{2,1}^{2/3}}{\|A_i\|_{2,1}^{2/3}} \right)^{3/2}$.

(Assume σ_i is 1-Lipschitz, inputs normalized.)

Stochastic Gradient Training Error on CIFAR10



(Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyats, 2017) 51/59

Training margins on CIFAR10 with true and random labels



• How does this match the large margin explanation?

If we rescale the margins by R_A (the scale parameter):



If we rescale the margins by R_A (the scale parameter):



With high probability, every f_A with $R_A \leq r$ satisfies

$$\Pr(M(f_A(X), Y) \le 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}[M(f_A(X_i), Y_i) \le \gamma] + \tilde{O}\left(\frac{rL}{\gamma\sqrt{n}}\right).$$

Network with *L* layers, parameters A_1, \ldots, A_L :

$$f_A(x) := \sigma(A_L \sigma_{L-1}(A_{L-1} \cdots \sigma_1(A_1 x) \cdots)).$$

Scale of f_A : $R_A := \prod_{i=1}^{L} \|A_i\|_* \left(\sum_{i=1}^{L} \frac{\|A_i\|_{2,1}^2}{\|A_i\|_*^{2/3}} \right)^{3/2}$.





• With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.
- Lower bounds?

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.
- Lower bounds?
- Regularization: explicit control of operator norms?

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.

イロト イポト イヨト イヨト

58 / 59

- Lower bounds?
- Regularization: explicit control of operator norms?
- Role of depth?

- With appropriate normalization, the margins analysis is qualitatively consistent with the generalization performance.
- Margin bounds extend to residual networks.
- Recent work by Golowich, Rakhlin, and Shamir give bounds with improved dependence on depth.
- Lower bounds?
- Regularization: explicit control of operator norms?
- Role of depth?
- Interplay with optimization?

Deep residual networks

- Representing with near-identities
 - Deeper networks allow flatter functions at each layer.
- Global optimality of stationary points
 - With flat functions, stationary points are global minima.
- What determines the statistical complexity of a deep network?
 - VC theory: Number of parameters
 - Margins analysis: Size of parameters
 - Understanding generalization failures