# 10707
# Deep Learning: Spring 2021

## Andrej Risteski

Machine Learning Department

## Recitation 3:
Brief overview of VC
and Rademacher bounds
in deep learning

# Classical view of generalization

**Meta-"theorem" of generalization:** with probability $1 - \delta$ over the choice of a training set of size m, we have

$$\sup_{f \in \mathcal{F}} | \; \widehat{\mathbb{E}} \; l(f(x), y) \text{ - } \mathbb{E} \; l(f(x), y)| \leq O\left(\sqrt{\frac{\text{complexity}(\mathcal{F}) + \ln 1/\delta}{m}}\right)$$

**Some measures of "complexity":**

(Effective) number of elements in $\mathcal{F}$

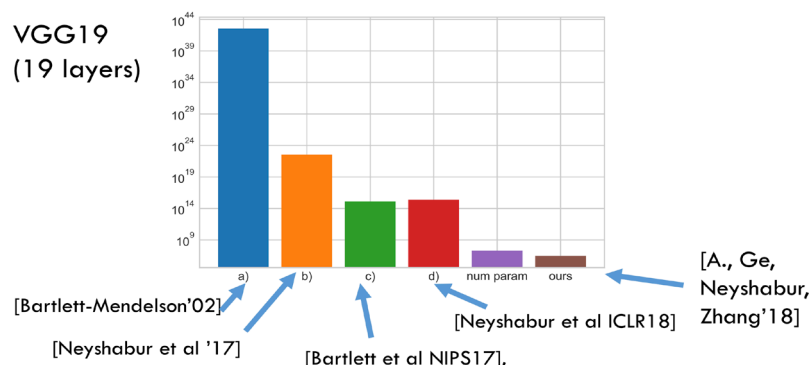VC (Vapnik-Chervonenkis)

Rademacher complexity

PAC-Bayes

# What do these bounds look like?

**VC-dimension** of **fully connected net** with **N** edges in total is O(**N log N**).
*Note*: doesn't depends on magnitude of weights at all.

**Rademacher bounds**: a lot of recent activity [*Neyshabur et al '15, 17; Bartlett-Foster-Telgarsky '18, Golowich-Rakhlin-Shamir '19*]. Roughly, best bounds look like:

$$R_m = O\left(\frac{\sqrt{\Pi_{i=1}^{d}||W_i||_2\mathrm{poly}(d)}}{\sqrt{m}}\right), \text{ where d is the depth}$$

VGG19
(19 layers)



[Bartlett-Mendelson'02]

[Neyshabur et al '17]

[Bartlett et al NIPS17],

[Neyshabur et al ICLR18]

[A., Ge, Neyshabur, Zhang'18]

When plugged in for standard values used in practice, the values these quantities give are **rarely "non-trivial"** (i.e. give a quantity < 1)

(Jiang et al'19): a large-scale investigation of the correlation of extant generalization measures with true generalization.

They show some newer Rademacher bounds have *worse* correlation.

# Recap: VC dimension

Let $\mathcal{F} = \{f: \mathcal{X} \to \{\pm 1\}\}$ be a class of predictors.

*Max # of possible label sequences*

The **growth function** $\Pi_{\mathcal{F}}: \mathbb{N} \to \mathbb{N}$ of $\mathcal{F}$ is defined as

$$\Pi_{\mathcal{F}}(m) = \max_{(x_1, x_2, \dots x_m)} \left|\{(f(x_1), f(x_2), \dots, f(x_m)) \big| f \in \mathcal{F}\}\right|$$

The **VC (Vapnis-Chervonenkis) dimension** of $\mathcal{F}$ is defined as

$$\text{VCdim}(\mathcal{F}) = \max\{m: \Pi_{\mathcal{F}}(m) = 2^m\}$$

Equivalently: the largest m, s.t. $\mathcal{F}$ can **shatter** some set of size m:

that is, for **some** $(x_1, x_2, \dots x_m)$ and **any** labeling $(b_1, b_2, \dots b_m)$, $b_i \in \{\pm 1\}$,

**some** $f \in \mathcal{F}$ can produce that labeling, that is

$$(f(x_1), f(x_2), \dots f(x_m)) = (b_1, b_2, \dots b_m)$$

# Recap: VC dimension

Let $\mathcal{F} = \{f : \mathcal{X} \to \{\pm 1\}\}$ be a class of predictors.

*Max # of possible label sequences*

The **growth function** $\Pi_{\mathcal{F}} : \mathbb{N} \to \mathbb{N}$ of $\mathcal{F}$ is defined as

$$\Pi_{\mathcal{F}}(m) = \max_{(x_1, x_2, \ldots x_m)} \left| \{ (f(x_1), f(x_2), \ldots, f(x_m)) | f \in \mathcal{F} \} \right|$$

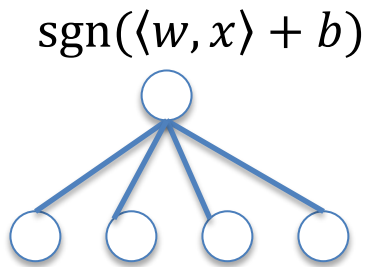The **VC (Vapnis-Chervonenkis) dimension** of $\mathcal{F}$ is defined as

$$\text{VCdim}(\mathcal{F}) = \max\{m : \Pi_{\mathcal{F}}(m) = 2^m\}$$

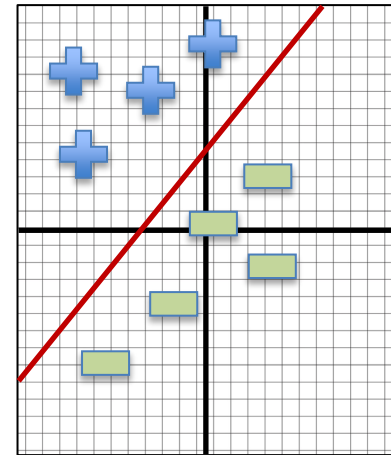The two are closely related (**Sauer's lemma**):

$$\Pi_{\mathcal{F}}(m) = O(m^{\text{VCdim}(\mathcal{F})})$$

# Bounding the VC dimension of a neural network

We'll prove a very simple bound on the VC dimension of a neural network with *__binary__* activation function. (Makes the proof the simplest).

$\text{sgn}(\langle w, x \rangle + b)$

(Each unit calculates a "hyperplane" $\text{sgn}(\langle w, x \rangle + b)$)

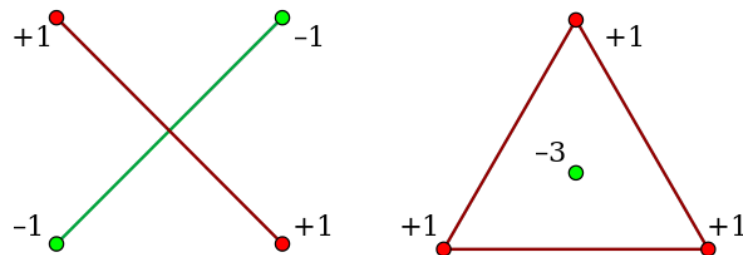For starters, let's bound the VC dimension of a single unit, namely consider

$\mathcal{F} = \{\text{sgn}(\langle w, x \rangle + b) \big| w \in \mathbb{R}^d, b \in \mathbb{R}\}.$

We'll show that $\text{VCdim}(\mathcal{F}) \leq d + 1$.

# Bounding the VC dimension of a neural network

VCdim($\mathcal{F}$) $\leq d + 1$: We want to that for every $(x_1, \ldots, x_{d+2})$, there exists a set of labels $(b_1, \ldots, b_{d+2})$ that cannot be linearly separated.

(*Radon's theorem*): For any $(x_1, \ldots, x_{d+2})$, there exists a partition of the points into sets $S_1, S_2$, s.t. convex hulls of $S_1$ and $S_2$ intersect.



How to use Radon's theorem: label pts in $S_1, S_2$ with +'s and −'s respectively.

**Claim**: no linear hyperplane perfectly separates $S_1, S_2$.

**Pf**: All pts in $S_1$ lie on one side of hyperplane, hence their convex hull does too.

All pts in $S_2$ lie on one side of hyperplane, hence their convex hull does too.

But, intersection is non-empty!

# Bounding the VC dimension of a neural network

We will recursively use this to bound VC dimension of neural nets with binary activation function.

We will show: VC-dimension of **fully connected net** with **N** edges in total is O(**N log N**)

**Main idea**: growth function behaves nicely wrt to compositions and concatenations.

# Bounding the VC dimension of a neural network

> **Claim 1**: If $\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2$ (**Cartesian product, i.e. concatenation**), we have $\Pi_{\mathcal{F}}(m) \leq \Pi_{\mathcal{F}_1}(m)\Pi_{\mathcal{F}_2}(m)$

**Pf**: Follows since any $\left(f(x_1), f(x_2), \ldots, f(x_m)\right)$ can be written as

$$\left(\left(f_1(x_1), f_2(x_1)\right), \left(f_1(x_2), f_2(x_2)\right), \ldots, \left(f_1(x_m), f_2(x_m)\right)\right).$$

> **Claim 2**: If $\mathcal{F} = \mathcal{F}_1 \circ \mathcal{F}_2$ (**i.e. compositions**), we have $\Pi_{\mathcal{F}}(m) \leq \Pi_{\mathcal{F}_1}(m)\Pi_{\mathcal{F}_2}(m)$

**Pf**: $\left|\left\{\left(f(x_1), f(x_2), \ldots, f(x_m)\right)\right| f \in \mathcal{F}\right\}\right|$ can be written as

$$\left|\bigcup_{(y_1,\ldots,y_m)\in\{(f_1(x_1),f_1(x_2),\ldots,f_1(x_m))| f_1\in\mathcal{F}\}}\{(f_2(y_1), f_2(y_2), \ldots, f_2(y_m))|f_2 \in \mathcal{F}_2\}\right|$$

$$\leq \sum_{(y_1,\ldots,y_m)\in\{(f_1(x_1),f_1(x_2),\ldots,f_1(x_m))| f_1\in\mathcal{F}\}} |\{(f_2(y_1), f_2(y_2), \ldots, f_2(y_m))|f_2 \in \mathcal{F}_2\}|$$

$$\leq \Pi_{\mathcal{F}_1}(m)\Pi_{\mathcal{F}_2}(m)$$

# Simple bounds on VC dimension of neural networks

**We write a neural net as a sequence of concatenations and compositions.**

Let $\mathcal{F}_{ij}$ be the set of functions (as we vary the input weights) calculated at the j$^{th}$ node on the i$^{th}$ layer.

If we view the **i$^{th}$ layer** as a function, it can be written as the concatenation of the outputs of all the nodes, namely $\mathcal{F}_i = \mathcal{F}_{i1} \times \mathcal{F}_{i2} \times \cdots \times \mathcal{F}_{n_i}$

Furthermore, the **entire network** can be written as $\mathcal{F}_l \circ \mathcal{F}_{l-1} \circ \cdots \circ \mathcal{F}_1$

Using the lemmas on the previous slide, we have $\Pi_{\mathcal{F}}(m) \leq \Pi_{ij} \Pi_{\mathcal{F}_{ij}}(m)$

Note that each $\mathcal{F}_{ij}$ is a **hyperplane**, so by Sauer's lemma and the VC dimension bound we proved, we have $\Pi_{\mathcal{F}_{ij}}(m) \leq m^{d_i - 1}$

Putting together, we have $\Pi_{\mathcal{F}}(m) \leq m^N$.( $N$ is the # of params in net. )

If a size m set is shattered: we have $2^m \leq m^N \Rightarrow$ m = O(N log N)

By definition of VC dimension, the proof follows.

# "Rademacher" bounds

(Result by *Neyshabur-Bhojanapalli-Srebro '18*, following *Arora, Ge, Neyshabur, Zhang '18* )

*Theorem*: For a ReLU network with layers $W^1, W^2, \ldots, W^d$, at most h nodes per layer, output margin $\gamma$ on a training set S, the generalization error can be bounded by

$$\tilde{O}\left(\frac{\sqrt{\text{hd}^2 \max_{x \in S} \|x\| \, \Pi_{i=1}^{d} \|W^i\|^2_2 \sum_i^d \frac{\|W^i\|^2_F}{\|W^i\|^2_2}}}{\gamma \sqrt{m}}\right)$$

"Lipschitz constant"

"Stable rank" (at most rank)

# "Rademacher" bounds

Result by *Bartlett-Foster-Telgarsky '18* was based on covering number and Rademacher arguments.

Follow up by *Neyshabur-Bhojanapalli-Srebro '18* was based on PAC-Bayes arguments.

We present an alternate proof by *Arora, Ge, Neyshabur, Zhang '18* based on compression arguments.

*Golowich-Rakhlin-Shamir '19* sharpen the bounds using further Rademacher arguments.

# Preliminaries

Some preliminaries:

A multiclass classifier f incurs can be associated with a <span style="color:#c0622a">margin loss</span>

$$L_\gamma(f) = \mathbb{P}_{(x,y)}[f(x)[y] \leq \gamma + \max_{i \neq y} f(x)[i]]$$

**Idea**: "compress" neural network and count compressed networks.

*Compressibility*: Let $G_{\mathcal{A}} = \{g_A : A \in \mathcal{A}\}$ be a discrete set of classifiers. A classifier $f$ is $(\gamma, S)$-compressible with respect to $G_{\mathcal{A}}$ if on a training set $S = \{(x_i, y_i)\}$, we have:

$$\exists A \in \mathcal{A}, \quad s.t. \ \forall x_i \in S, \quad ||f(x_i) - g_A(x_i)||_\infty \leq \gamma$$

# Generalization from compressibility

**Idea**: "compress" neural network and count compressed networks.

*Compressibility*: Let $G_{\mathcal{A}} = \{g_A : A \in \mathcal{A}\}$ be a discrete set of classifiers. A classifier $f$ is $(\gamma, S)$-compressible with respect to $G_{\mathcal{A}}$ if on a training set $S = \{(x_i, y_i)\}$, we have:

$$\exists A \in \mathcal{A}, \qquad s.t. \ \forall x_i \in S, \qquad ||f(x_i) - g_A(x_i)||_{\infty} \leq \gamma$$

*Theorem*: Suppose $A$ is a set of $q$ parameters each of which can have at most $r$ discrete values and let $|S| = m$.

If the trained classifier f is $(\gamma, S)$- compressible with respect to $G_{\mathcal{A}}$, there exists an $A \in \mathcal{A}$, s.t. with high probability over S,

$$|L_0(g_A) - \hat{L}_{\gamma}(f)| \leq O\left(\sqrt{\frac{q \log r}{m}}\right)$$

**Pf**: Essentially Chernoff + union bound.

# How do you compress?

*How do you compress a neural network?*

*Lemma*: For any matrix $A$, let $\hat{A}$ be the truncated version of $A$ where the singular values of $A$ smaller than $\delta\left\Vert A\right\Vert_2$ are removed. Then, $\left\Vert \hat{A} - A \right\Vert_2 \leq \delta\left\Vert A\right\Vert_2$ and $\hat{A}$ has rank at most $\dfrac{\left\Vert A\right\Vert_F^2}{\delta^2\left\Vert A\right\Vert_2^2}$.

*Implied compression*: write a rank r matrix $\hat{A}^{m\times n}$ as $\hat{U}^{m\times r}\hat{U}^{r\times n}$ which has (m+n)r parameters. (As opposed to mn.)

*Proof*: Denote by r the rank of $\hat{A}$. Max singular value of $\hat{A} - A$ is at most $\delta\left\Vert A\right\Vert_2$.

Since the remaining singular values are at least $\delta\left\Vert A\right\Vert_2$, we have $\left\Vert A\right\Vert_F \geq \left\Vert \hat{A}\right\Vert_F \geq \sqrt{r}\,\delta\left\Vert A\right\Vert_2$.

# Perturbation after compression

*Why spectral norm?* Roughly captures "Lipschitzness" of network. Namely:

*Lemma*: Let $f_{\boldsymbol{w}}$ be a d-layer network with ReLU activations, and let $\boldsymbol{u} = \{U_i\}_{i=1}^{d}$ be perturbations, s.t. $||U_i||_2 \leq \frac{1}{d}||W_i||_2$. Then:

$$|f_{\boldsymbol{w}+\boldsymbol{u}}(x) - f_{\boldsymbol{w}}(x)| \leq e\, ||x||\, (\Pi_i ||W_i||_2) \sum_i \frac{||U_i||_2}{||W_i||_2}$$

*Proof sketch:* Induction on the error at layer i.

$$
\begin{aligned}
|f_{\boldsymbol{w}+\boldsymbol{u}}^i(x) - f_{\boldsymbol{w}}^i(x)| \ &= |(W_i + U_i)\sigma(f_{\boldsymbol{w}+\boldsymbol{u}}^{i-1}(x)) - W_i\sigma(f_{\boldsymbol{w}}^{i-1}(x))| \\[2mm]
&= |(W_i + U_i)(\sigma(f_{\boldsymbol{w}+\boldsymbol{u}}^{i-1}(x)) - \sigma(f_{\boldsymbol{w}}^{i-1}(x))) + U_i\, \sigma(f_{\boldsymbol{w}}^{i-1}(x))| \\[2mm]
&\leq \left\|W_i + U_i\right\|_2 |f_{\boldsymbol{w}+\boldsymbol{u}}^{i-1}(x)) - f_{\boldsymbol{w}}^{i-1}(x)))| + ||U_i||_2 |f_{\boldsymbol{w}}^{i-1}(x))|
\end{aligned}
$$

# Perturbation after compression

*Why spectral norm?* Roughly captures "Lipschitzness" of network. Namely:

*Lemma*: Let $f_{\boldsymbol{w}}$ be a d-layer network with ReLU activations, and let $\boldsymbol{u} = \{U_i\}_{i=1}^d$ be perturbations, s.t. $||U_i||_2 \leq \frac{1}{d}||W_i||_2$. Then:

$$|f_{\boldsymbol{w}+\boldsymbol{u}}(x) - f_{\boldsymbol{w}}(x)| \leq e\,||x||\,(\Pi_i||W_i||_2)\sum_i \frac{||U_i||_2}{||W_i||_2}$$

*Proof sketch:* Induction on the error at layer i.

$$|f_{\boldsymbol{w}+\boldsymbol{u}}^i(x) - f_{\boldsymbol{w}}^i(x)| \leq \left||W_i + U_i\right||_2 |f_{\boldsymbol{w}+\boldsymbol{u}}^{i-1}(x)) - f_{\boldsymbol{w}}^{i-1}(x)))| + ||U_i||_2\,|f_{\boldsymbol{w}}^{i-1}(x))|$$

$$\leq \left||W_i + U_i\right||_2 |f_{\boldsymbol{w}+\boldsymbol{u}}^{i-1}(x)) - f_{\boldsymbol{w}}^{i-1}(x)))| + ||U_i||_2||x||\,\Pi_i||W_i||_2$$

Using $||U_i||_2 \leq \frac{1}{d}||W_i||_2$, the lemma follows.

# Putting things together

*Lemma*: Let $f_{\boldsymbol{w}}$ be a d-layer network with ReLU activations, and let $\boldsymbol{u} = \{U_i\}_{i=1}^{d}$ be perturbations, s.t. $||U_i||_2 \leq \frac{1}{d}||W_i||_2$. Then:

$$|f_{\boldsymbol{w}+\boldsymbol{u}}(x) - f_{\boldsymbol{w}}(x)| \leq e\,||x||\,(\Pi_i ||W_i||_2) \sum_i \frac{||U_i||_2}{||W_i||_2}$$

*Putting things together:*

Applying compression lemma for $\delta = \dfrac{\gamma}{3d||x||\Pi_i \left||W^i\right||_2}$ gives $|f_{\boldsymbol{w}+\boldsymbol{u}}(x) - f_{\boldsymbol{w}}(x)| \leq \gamma$

Hence $|L_0(f_{\boldsymbol{w}+\boldsymbol{u}}) - L_\gamma(f_{\boldsymbol{w}})| \leq \gamma$

# Putting things together

*Lemma*: For any matrix $A$, let $\hat{A}$ be the truncated version of $A$ where the singular values of $A$ smaller than $\delta\|A\|_2$ are removed. Then, $\left\|\hat{A} - A\right\|_2 \le \delta\|A\|_2$ and $\hat{A}$ has rank at most $\dfrac{\|A\|_F^2}{\delta^2\|A\|_2^2}$.

*Putting things together:*

Applying compression lemma for $\delta = \dfrac{\gamma}{3d\|x\|\Pi_i\|W^i\|_2}$ gives $|f_{\boldsymbol{w}+\boldsymbol{u}}(x) - f_{\boldsymbol{w}}(x)| \le \gamma$

Hence $|L_0(f_{\boldsymbol{w}+\boldsymbol{u}}) - L_\gamma(f_{\boldsymbol{w}})| \le \gamma$

By truncation lemma, rank in layer i is at most $\dfrac{9d^2\|x\|^2\Pi_i\|W^i\|_2^2}{\gamma^2}\dfrac{\|W^i\|_F^2}{\|W^i\|_2^2}$

Number of nodes is at most h per layer, so compressed encoding has num of params

$\dfrac{1}{\gamma^2}\mathrm{hd}^2\max_{x\in S}\|x\|\,\Pi_{i=1}^d\|W^i\|^2{}_2\sum_i^d\dfrac{\|W^i\|_F^2}{\|W^i\|_2^2}$. Applying compression lemma, proof follows.