# MAT280 (ML) - Lecture 1 - 4/4/17

Michael Wolf - Prof. at TIM

- **Mathematic Foundations of Machine Learning**

Not about software/coding/implementation

Grading = 1/3 participation    2/3 research paper presentation (just a paper)
  bi in lecture, comments on lecture notes    (2/3 person team
  online

## What is Machine Learning?

$CS \supseteq AI \supseteq ML \supseteq$ Deep learning

Aim: Data $\rightarrow$ models/algorithms/programs

used when: large amt. of data available, target fn. too difficult to
  implement directly.

Main branches— ✦Supervised ML — training data w/ labels as input
  focus of course    outputs program to predict labels
  of unseen neighbors

Unsupervised ML: input data, output model/pattern in data

Supervised is predictive,    Unsupervised in descriptive

Representation: (structure of programs)  — Decision trees
                                          — Neural networks  ⎤ optimizable
                                          — K-Nearest Neighbors ⎥ structures
                                          — SVMs

## Course Topics:
- Learning Theory
- Neural Networks
- SVMs & Kernel methods

    ... (if time)

# I. Learning Theory　　(supervised ML)
## I.1. Statistical framework

input - "training data" $S = ((x_1, y_1), \ldots, (x_n, y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$

output - "hypothesis" $h: \mathcal{X} \to \mathcal{Y}$

ML algorithm $A: \bigcup_{n \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^n \to \mathcal{Y}^{\mathcal{X}}, \quad S \mapsto h$

range$(A) =: \mathcal{F} \subseteq \mathcal{Y}^{\mathcal{X}}$　↖ could work for any size training set $\forall n \in \mathbb{N}$

assumptions: $(x_i, y_i)$ values of random i.i.d. vars. $(X_i, Y_i)$
(not always valid IRL)　　distributed according to some prob. measure $P$ over $\mathcal{X} \times \mathcal{Y}$

Denote expectation value wrt. $P$ as $\mathbb{E}$. Exp. value of $S$ wrt. $P^n$ as $\mathbb{E}_S$.

Goal: Find a "good" $h$ wrt. "loss function" $L: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$

minimize the "risk": $R(h) := \int_{\mathcal{X} \times \mathcal{Y}} L(y, h(x)) \, dP(x, y)$
　　　　　　　　　　　　　　⌢ don't know this distribut.

Challenge: $P$ is unknown.

Regression: $\mathcal{Y}$ is continuous. If $\mathcal{Y} = \mathbb{R}$, the most common loss fn. is
$$L(y, y') := |y - y'|^2 \quad \text{(quadratic loss)}$$
Then risk is $R(h) = \mathbb{E}\left[ |Y - h(X)|^2 \right]$ ( mean squared error

Classification: $\mathcal{Y}$ is discrete. $h$ called "classifier" ) choice may al be det. by app
Most common loss fn: "0-1 loss" $L(y, y') := 1 - \delta_{yy'}$

Risk $R(h) = P[h(X) \neq Y] = \mathbb{E}\left[ \mathbb{1}_{h(X) \neq Y} \right]$
　　　　　　　　　　　　　　　⌢ indicator fn.

B. classification $\mathcal{Y} = \{-1, 1\}$

Choice of loss fn is determined by problem goals / optimizer constraints

└ E.g. step fn. for goal of $Y \geq 100$, make convex for optimizer

## 1.2 Error Decomposition

· Prior knowledge is encoded in $\mathcal{F}$ (target hyp. space - det. by model choice)

· Cant minimize risk, minimize empirical risk — ERM

$$\hat{R}(h) := \frac{1}{n} \sum_{i=1}^{n} L(y_i, h(x_i))$$

Errors':

$$R_{\mathcal{F}} := \inf_{h \in \mathcal{F}} R(h)$$

$\uparrow$ best we can do w/ model choice

$$R_{yx} := \inf_{h \in yx} R(h)$$

$\uparrow$ best we can do in principle

$$R(h) - R_{yx} = \underbrace{\left(R(h) - R_{\mathcal{F}}\right)}_{\substack{\text{estimation} \\ \text{error}}} + \underbrace{\left(R_{\mathcal{F}} - R_{yx}\right)}_{\substack{\text{approximation} \\ \text{error}}} \begin{array}{l} \leftarrow \text{independent of } S \\ \leftarrow \text{how model choice fits} \end{array}$$

Do ERM, get $\hat{h} \in \mathcal{F}$.

Then $\hat{R}(\hat{h}) \leq \hat{R}(h) \quad \forall h \in \mathcal{F}$.

since $\begin{array}{l} R_{\mathcal{F}} = \inf (sub) \\ -R_{\mathcal{F}} = \sup(-stuff) \end{array}$

Est. error: $R(\hat{h}) - R_{\mathcal{F}} = R(\hat{h}) - \hat{R}(\hat{h}) + \sup_{h \in \mathcal{F}} \left(\hat{R}(\hat{h}) - R(h)\right)$

$$\leq 2 \sup_{h \in \mathcal{F}} \underbrace{\left| R(h) - \hat{R}(h) \right|}_{\substack{\text{generalization} \\ \text{error}}}$$

d: $R(\hat{h}) - R_{\mathcal{F}} \leq 2 \sup_{h \in \mathcal{F}} \left| R(h) - \hat{R}(h) \right|$

$\uparrow \uparrow$
ERM within $\mathcal{F}$

Ex 1 - Linear Regression: $\mathcal{X} \times \mathcal{Y} = \mathbb{R}^d \times \mathbb{R}$

$$\bar{\mathcal{F}} := \{ h : \mathbb{R}^d \to \mathbb{R} \mid \exists v \in \mathbb{R}^d : h(x) = \langle v, x \rangle \}$$

$$\hat{R}(v) := \frac{1}{n} \sum_{i=1}^{n} |\langle v, x_i \rangle - y_i|^2$$

$$\nabla \hat{R}(v) \overset{!}{=} 0 \iff \frac{d}{dv_k} \hat{R}(v) = \frac{1}{n} \sum_{i=1}^{n} (\langle v, x_i \rangle - y_i) x_{i,k} = 0 \quad \forall k$$

$$\hookrightarrow v = A^{-1} b$$

is the ERM

$$\iff Av = b \quad \text{w/} \quad A := \sum_{i=1}^{n} x_i x_i^T$$

$$b := \sum_{i=1}^{n} y_i x_i$$

Ex. 2 - Polynomial Regression: $\mathcal{X} \times \mathcal{Y} = \mathbb{R} \times \mathbb{R}$, $\bar{\mathcal{F}} = \{ h : \mathbb{R} \to \mathbb{R} \mid \exists a \in \mathbb{R}^{m+1} : h(x) = \sum_{k=0}^{m} a_k x^k \}$

$$\Psi : \mathbb{R} \to \mathbb{R}^{m+1}, \quad \Psi(x) := (1, x, x^2, \ldots, x^m)$$

then $\hat{R}(v) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \sum_{k=0}^{m} a_k x^k \right)^2 = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \langle a, \Psi(x_i) \rangle \right)^2$
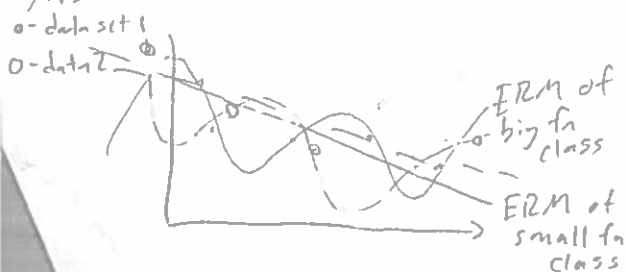
then, do gradient & whatnot (may not get explicit ERM)
Can use this trick for any set of linear basis functions.

$$\text{ERM} \Rightarrow Av = b \quad \text{w/} \quad A = \sum_{i=1}^{n} \Psi(x_i) \Psi(x_i)^T, \quad b = \sum_{i=1}^{n} y_i \Psi(x_i) \leftarrow \text{maybe not invertible}$$

Approx. error $\tilde{R}_{\bar{\mathcal{F}}} - R_{\mathcal{Y}\mathcal{X}}$  ⟵ Inf. risk − Inf. risk in fn. class  overall fns.

$|\bar{\mathcal{F}}| \nearrow \quad |\bar{\mathcal{F}}| \searrow$

Estimation error — $R(h) - R_{\bar{\mathcal{F}}}$ ⟵ Risk of − inf. risk of hypothesis  fn. class

Also known as "bias-variance trade-off"



- data set 1
- data 2

ERM of big fn class

ERM of small fn class

Do high-deg poly fits on random draws of sample, average will have small/zero bias, but sample-to-sample polynomials vary widely.

Linear fits on random draws will have little variance, but bias stays big.

Approaches aiming @ a balanced choice for $\bar{f}$:

- Split data $\longrightarrow$ training data $\to$ optimizing hypothesis $\to h_S$
  $\searrow$ test data $\to$ evaluate performance of $h_S$
  validation data $\to$ tune hyperparameters

- Modification of ERM $\longrightarrow$ structural risk minimization

  Use $\bar{f}_1 \subset \bar{f}_2 \subset \bar{f}_3 \subset \dots$
  & penalize higher levels

  Regularization — minimize
  $$\hat{R}(h) + C_n(h)$$

  e.g. "Tikhonov reg." $C_n(h) = \lambda \|h\|^2$, $\lambda \in \mathbb{R}_+$ chosen by cross-validation.

## 1.3 PAC Learning (Probably Approximately Correct)

$\hookrightarrow$ introduced by Voilant in 1984

Desirable: find uniform bound on $|R(h) - \hat{R}(h)| \leq \dots$

Always have chance of unfair training data! No deterministic bound

Want: $\mathbb{P}_S\left[ |\hat{R}(h) - R(h)| > \varepsilon \right] < \delta$

Take $\mathcal{X}, \mathcal{Y}$ finite, choose 0-1 loss & assume $f: \mathcal{X} \to \mathcal{Y}$ determines "true label" $\in \bar{f}$

so $S = (x_i, f(x_i))$ and $P(x,y) = \delta_{y, f(x)} p(x)$

Lemma: $\forall \varepsilon > 0, \forall h \in \mathcal{Y}^{\mathcal{X}}$, $\overset{\mathbb{P}[h(x) \neq f(x)]}{R(h)} > \varepsilon \implies \mathbb{P}_S\left[ \hat{R}(h) = 0 \right] < e^{-\varepsilon n}$ $\overset{< \text{size of } S}{}$

Pf: $\mathbb{P}_S\left[\hat{R}(h) = 0\right] = \mathbb{P}_S\left[\forall i \in \{1,\dots,n\}: h(x_i) = f(x_i)\right] \overset{\text{i.i.d.}}{=} \prod_{i=1}^n \mathbb{P}\left[h(x_i) = f(x_i)\right] \leq (1-\varepsilon)^n$

$$\leq e^{-\varepsilon n} \qquad \square$$

Thm: $\mathbb{P}_S \left[ |R(h_S) - \bar{R}(h_S)| > \varepsilon \right] < \delta$ if $n \geq \frac{1}{\varepsilon}\left( \ln|\bar{f}| + \ln \frac{1}{\delta} \right)$

under the assumption $\forall S \; \exists h_S \in \bar{f} : \hat{R}(h_S) = 0$.

Pf:

$$\mathbb{P}_S \left[ |R(h_S) - \hat{R}(h_S)| > \varepsilon \right] = \mathbb{P}_S \left[ R(h_S) > \varepsilon \right]$$

$$\leq \mathbb{P}_S \left[ \exists h \in \bar{f} : R(h) > \varepsilon \wedge \bar{R}(h) = 0 \right]$$

union bound $\rightarrow \quad \leq \sum_{h \in \bar{f} : R(h) > \varepsilon} \mathbb{P}_S \left[ \bar{R}(h) = 0 \right]$

Lemma $\rightarrow \quad \leq \sum_{h \in \bar{f} : R(h) > \varepsilon} e^{-\varepsilon n} \leq |\bar{f}| e^{-\varepsilon n} = \delta \qquad$ solve for $n$ get

$\bigg)$

Assumptions:

- $\mathbb{P}[Y=y \mid X=x] = \delta_{y, f(x)}$
- $\forall S = ((x_i, f(x_i))_{i=1}^n : \exists h_s \in \bar{f}: \hat{R}(h_s) = 0$

$\Big\} \Rightarrow \mathbb{P}_s[R(h_s) > \varepsilon] \leq \begin{cases} (1-\varepsilon)^n |\bar{f}| \\ \delta \text{ if} \\ n \geq \frac{1}{\varepsilon} \ln\left(\frac{|\bar{f}|}{\delta}\right) \end{cases}$

Rmk: • These $\mathbb{P}$ bounds assume $\bar{f}$ is finite.

• If $\bar{f} = y^X$, then $n > |X|$ but then there's nothing to generalize to, meaningless!

## 1.4 No Free Lunch

Thm: Let $X, y$ be finite, $|X| > n$. $R_f(h) := \mathbb{P}[h(x) \neq f(x)]$    $(\bar{f} = y^X)$

$\forall S \mapsto h_s \in \bar{f}:$   $\mathbb{E}_f\Big[\mathbb{E}_s[R_f(h_s)]\Big] \geq \left(1 - \frac{1}{|y|}\right)\left(1 - \frac{n}{|X|}\right)$

if uniform dist. over $f$ and $x$ are $x$ used

Pf:

$\mathbb{E}_f \mathbb{E}_s[R_f(h_s)] = \frac{1}{|X|} \mathbb{E}_f \mathbb{E}_s\left[\sum_{x \in X} \mathbb{1}_{h_s(x) \neq f(x)}\right]$

$\geq \frac{1}{|X|} \mathbb{E}_f \mathbb{E}_s\left[\sum_{x \notin X_s} \text{''}\right]$    where $X_s \subset X$ appearing in $S$.

$\vdots$

$\geq \frac{1}{|X|} \sum_{x \notin X_s}\left(1 - \frac{1}{|y|}\right) = \left(1 - \frac{1}{|y|}\right)\left(\frac{|X| - n}{|X|}\right) = \left(1 - \frac{1}{|y|}\right)\left(1 - \frac{n}{|X|}\right)_{\square}$

Random guessing $\longrightarrow \left(1 - \frac{1}{|y|}\right)$. Sophisticated alg. only $\cdot\left(1 - \frac{n}{|X|}\right)$ better this additional factor reflects the fact that training data is already known.

No "better" learners for all data sets — neural nets perform better than decision trees on some data sets but worse on others

— on average, all learners are no better than random guessing!

Need to restrict fn. class $\bar{f}$ a priori.

Cardinality $\langle \bar{f} \rangle$ complexity of $\bar{f}$                      ∠ countable but infly complex!

E.g. $X = \{x_1, x_2, \ldots\}$. $\infty > |y|, y \ni 0$, $\bar{f} := \bigcup_{n \in \mathbb{N}} \{f: X \to y : \forall m > n : f(x_m) = 0\}$

## 1.5 Growth function:

**Def:** $|y| < \infty$, $f \subseteq y^{\chi}$. For every $C \subseteq \chi$ define
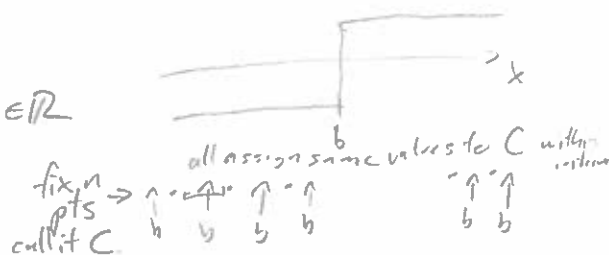
$$\bar{f}_C := \{ f \in y^C \mid \exists F \in \bar{f} \ \forall x \in C : F(x) = f(x) \}$$

The "growth function" $\Gamma : \mathbb{N} \to \mathbb{N}$ assigned to $\hat{f}$ is $\Gamma(n) := \max_{C \subseteq \chi : |C| = n} |\bar{f}_C|$

**Note:** $\Gamma(n) \leq |y|^n$

**Ex:** $\bar{f} \subseteq \{-1, 1\}^{\mathbb{R}}$, $\bar{f} := \{ x \mapsto \text{sgn}[x-b] \}_{b \in \mathbb{R}}$

$\to \Gamma(n) = n+1$

fix $n$ pts $\to$ call it $C$

all assign same values to $C$ with ...

**Lemma:** [Hoeffding's ineq. '63] (Consider $Z_1, \ldots, Z_n$ real indep. rand. vars. w/ range $(Z_i) \subseteq [a_i, b_i]$. Then $\forall \varepsilon > 0$, $\mathbb{P}\left[ \sum_{i=1}^{n} |Z_i - \mathbb{E}[Z_i]| \geq \varepsilon \right] \leq \exp\left[ -\frac{2\varepsilon^2}{\sum_{i=1}^{n}(a_i - b_i)^2} \right]$

$\frac{1}{2} \mathbb{P}\left[ \left| \sum_{i=1}^{n} (Z_i - \mathbb{E}[Z_i]) \right| \geq \varepsilon \right] \leq \quad \cdots$

$\frac{1}{2} \mathbb{P}\left[ \left| \sum_{i=1}^{n} (Z_i - Z_i') \right| \geq \varepsilon \right] \leq \quad \cdots$

↑ independent copies of $Z_i$.

**Thm:** $|y| < \infty$, range $(L) \subseteq [0, c]$, $\delta \in [0,1]$.
With prob. at least $1 - \delta$ wrt. repeated sampling of training data of size $n \in \mathbb{N}$:

$$\forall h \in \bar{f} : |R(h) - \hat{R}(h)| \leq c \sqrt{\frac{8 \ln (\Gamma(2n)\frac{4}{\delta})}{n}}$$
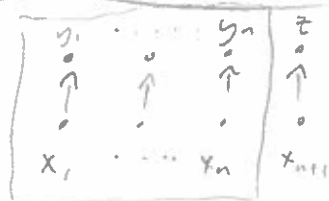
**Pf:** $S, S'$ are i.i.d. rand. vars. over $(\chi \times y)^n$ w/ distr. according to $P^n$.

$|R(h) - \hat{R}(h)| > \varepsilon$, $|\hat{R}'(h) - \hat{R}(h)| < \varepsilon/2 \Rightarrow |R(h) - \hat{R}'(h)| > \varepsilon/2$ by $\Delta$-ineq.

$\mathbb{1}_{|R(h)-\hat{R}(h)|>\varepsilon} \ \mathbb{1}_{|R(h)-\hat{R}'(h)|<\varepsilon/2} \leq \mathbb{1}_{|\hat{R}'(h)-\hat{R}(h)|>\varepsilon/2}$

assume $n \geq \frac{4c^2}{\varepsilon^2} \ln 2$. $\mathbb{E}_{S'}\left[ \mathbb{1}_{|R(h)-\hat{R}'(h)|<\varepsilon/2} \right] \geq 1 - 2\exp\left[ -\frac{\varepsilon^2 n}{2c^2} \right] \geq 1/2$.

**Recall:** $\Gamma(n) := \max\limits_{c \subseteq X, |c|=n} \{|\bar{\mathcal{F}_c}|\}$.

$$
\begin{array}{c|c}
\begin{matrix} y_1 & \cdots & y_n \\ \uparrow & \uparrow & \uparrow \\ x_1 & \cdots & x_n \end{matrix} & \begin{matrix} z \\ \uparrow \\ x_{n+1} \end{matrix}
\end{array}
$$

need restriction on fn class to use info in box !

**Thm:** ... $\forall h \in \bar{\mathcal{F}}:$ $|R(h) - \hat{R}(h)| \leq C \sqrt{\dfrac{8 \ln(\Gamma(2n) \frac{4}{\delta})}{n}}$

**Pf:** Last time, showed $\mathbb{P}_S[\exists h \in \bar{\mathcal{F}}: |R(h) - \hat{R}(h)| > \varepsilon]$

$\leq 2 \mathbb{P}_{SS'}[\exists h \in \bar{\mathcal{F}}: |\hat{R}'(h) - \hat{R}(h)| > \varepsilon/2]$

Note: $2\mathbb{P}_{SS'}[\exists h \in \bar{\mathcal{F}}: |\hat{R}'(h) - \hat{R}(h)| > \varepsilon/2] = 2\mathbb{P}_{SS'}[\exists h \in \bar{\mathcal{F}}: \frac{1}{n}|\sum_{i=1}^{n} L(y_i, h(x_i)) - L(y_i', h(x_i'))| > \varepsilon/2]$

Mult. by $\sigma_i = \{\pm 1\}$, since $S \& S'$ are i.i.d. $= 2\mathbb{P}_{SS'}[\exists h \in \bar{\mathcal{F}}: \frac{1}{n}|\sum_{i=1}^{n}[L(y_i, h(x_i)) - L(y_i', h(x_i'))]\sigma_i| > \varepsilon/2]$

unif. dist. $= 2\mathbb{E}_{SS'}[\mathbb{P}_\sigma[\exists h \in \bar{\mathcal{F}}_{|S \cup S'}: \frac{1}{n}|\sum_{i=1} \cdots| > \varepsilon/2]]$

$\stackrel{\cup\text{-bound}}{\leq} 2\mathbb{E}_{SS'}[\sum\limits_{h \in \bar{\mathcal{F}}_{|S \cup S'}} \mathbb{P}_\sigma[\frac{1}{n}|\cdots| > \varepsilon/2]]$

$\{z_i \leftarrow \text{exp. of } z_i = 0 \text{ b/c of } \sigma_i$

$= 2\mathbb{E}_{SS'}[\sum\limits_{h \in \bar{\mathcal{F}}_{|S \cup S'}} \mathbb{P}_\sigma[|\sum_{i=1}^{n} z_i| > n\varepsilon/2]]$

$\stackrel{\leq}{\text{Hoeffding}} 4\underbrace{\mathbb{E}_{SS'}[|\bar{\mathcal{F}}|_{S \cup S'}]}_{\leq \Gamma(2n)} \exp[-\frac{n\varepsilon^2}{8c^2}] =: \delta \quad \rightarrow \text{solve for } \varepsilon.$

... assumption $\Longleftrightarrow \delta \leq 2\sqrt{2}\,\Gamma(2n) \checkmark$ since $\delta \in (0,1)$ $\square$

---

**Lemma:** $\bar{\mathcal{F}_1} \subseteq \mathcal{Y}^X$, $\bar{\mathcal{F}_2} \subseteq Z^2$. $\bar{\mathcal{F}} := \bar{\mathcal{F}_2} \circ \bar{\mathcal{F}_1} \Rightarrow \Gamma(n) \leq \Gamma_1(n)\Gamma_2(n)$

**1.6 VC-Dimension:**

Consider $|\mathcal{Y}| = 2$.

**Def:** The VC (Vapnik-Chervonenkis) dimension of $\bar{f} \subseteq y^\Lambda$ w/ $|y| = 2$ is
$$VCdim(\bar{f}) := \max\{n \in \mathbb{N} : \bar{f}(n) = 2^n\} \quad (\infty \text{ if max DNE})$$

**Ex:** [sketch] threshold fns. $\bar{f}(n) = n+1 \Rightarrow VCdim(\bar{f}) = 1$

---

**Thm:** Let $d := VCdim(\bar{f})$. Then $\bar{f}(n) \begin{cases} = 2^n & \text{if } n \le d \\ \le \left(\frac{en}{d}\right)^d & \text{if } n > d \end{cases}$

**Pf:** Assume $n > d$ ($n \le d$ pf. by definition)

WTS: $\forall A \subseteq X$ w/ $|A| = n$, $|\bar{f}_{|A}| \overset{?}{\le} |\{B \subseteq A \mid \bar{f}_{|B} = y^B\}|$.

pf. by induction over $A$

Pick $a \in A$ & define
$$\bar{f}' := \{h \in \bar{f}_{|A} \mid \exists g \in \bar{f}_{|A} : h(a) \ne g(a)$$
$$\wedge (h-g)|_{A \backslash a} = 0 \}.$$

$\bar{f}_a := \bar{f}'|_{A \backslash a}$.

Note $|\bar{f}_{|A}| = |\bar{f}_{|A \backslash a}| + |\bar{f}_a|$.

$$\quad \le |\{B \subseteq A \mid |B| \le d\}|$$
$$= \sum_{i=0}^{d} \binom{n}{i} \le \sum_{i=0}^{n} \binom{n}{i}\left(\frac{n}{d}\right)^{d-i}$$
$$= \left(\frac{n}{d}\right)^d \left(1 + \frac{d}{n}\right)^n \le \left(\frac{n}{d}\right)^d e^d.$$

$|\bar{f}_{|A \backslash a}| \le |\{B \subseteq A \mid \bar{f}_{|B} = y^B \wedge a \notin B\}|$ by IH.

$|\bar{f}_a| = |\bar{f}'|_{A \backslash a}| \le |\{B \subseteq A \backslash a \mid \bar{f}'_{|B} = y^B\}|$ by IH.

$\qquad = |\{B \subseteq A \backslash a \mid \bar{f}'_{|B \cup a} = y^{B \cup a}\}|$ by def. of $\bar{f}'$

$\qquad = |\{B \subseteq A \mid \bar{f}'_{|B} = y^B \wedge a \in B\}|$

$\qquad \le |\{B \subseteq A \mid \bar{f}_{|B} = y^B \wedge a \in B\}|$

add together $a \notin B$ cancels/sums with $a \in B$

$\Rightarrow |\bar{f}_{|A}| \le |\{B \subseteq A \mid \bar{f}_B \in y^B\}|$ $\quad\blacksquare$

**Cor:** $(\varepsilon, \delta) \in [0,1]^2$. $R$ = error prob. Then $\forall h \in \bar{f} : |\hat{R}(h) - R(h)| \le \varepsilon$ holds
w/ prob $\ge 1 - \delta$ if $n \ge \frac{32}{\varepsilon^2}\left[d \ln\left(\frac{8d}{\varepsilon^2}\right) + \ln\left(\frac{6}{\delta}\right)\right]$

$\rightsquigarrow$ "$n \ge d/\varepsilon^2$ scaling"

Thm: Let $G$ be an $\mathbb{R}$-vector space of fns. from $X \to \mathbb{R}$

Then $\bar{f} := \{x \mapsto \text{sgn}[g(x)] \mid g \in G\} \subseteq \{\pm 1\}^X$ satisfies

$$VCdim(\bar{f}) = dim(G).$$

Pf: let $k = dim(G) + 1$ & assume $VCdim(\bar{f}) \geq k$.

$\exists \{x_1, \ldots, x_k\} = C \subseteq X$ s.t. $\bar{f}|_C = \{-1, 1\}^C$

$L: G \to \mathbb{R}^k : \quad L(g) := (g(x_1), \ldots, g(x_k))$

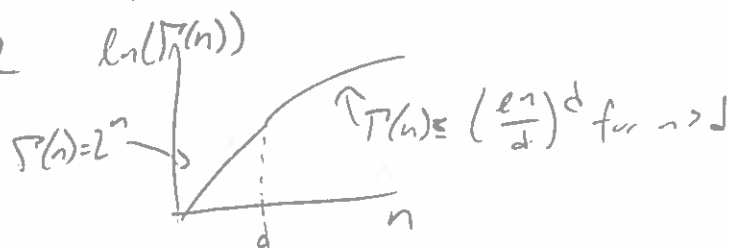$dim(range(L)) \leq dim(G) \implies \exists v \neq 0 \in (range(L))^\perp = ker(L^*)$

Then $\forall g \in G : 0 = \langle L^*(v), g \rangle = \langle v, L(g) \rangle = \sum_{i=1}^{k} v_i \, g(x_i)$

If $\bar{f}|_C = \{-1, 1\}^C$, we can choose $g$ st. $\text{sgn}(g(x_i)) = \text{sgn}(v_i)$

& then $\sum_{i=1}^{k} v_i \, g(x_i) > 0 \quad \text{\Large\char"21AF}.$ $\qquad \blacksquare$

---

<u>MAT 280 - Lec 5 - 4/18/17</u>

Recap: $\bar{f} \subseteq y^X$, $|y| = 2$



$\Gamma(n) = 2^n$

$\Gamma(n) \leq \left(\frac{en}{d}\right)^d$ for $n > d$

WTS

"$n = \frac{1}{\epsilon^2}\left(d \ln\left(\frac{d}{\epsilon^2}\right) + \ln\frac{1}{\delta}\right)$" $\implies \forall h \in \bar{f}: |R(h) - \hat{R}(h)| \leq \epsilon$ w/ prob. $1-\delta$

$\implies$ If $\bar{f} \subseteq$ halfspaces in $\mathbb{R}^m \implies VCdim(\bar{f}) = m+1$.

Ex: $\bar{f}$ is indicator fns of $\cdot$ Euclidean balls in $\mathbb{R}^m \to VCdim(\bar{f}) = m+1$

$\cdot$ axis-aligned boxes in $\mathbb{R}^m \to VCdim(\bar{f}) = 2m$ $\quad \hookleftarrow$ also no. of parameters

counterex: $\bar{f} := \{x \in \mathbb{R} \mapsto \text{sgn}(\sin(\alpha x)) \mid \alpha \in \mathbb{R}\}$. $\quad VCdim(\bar{f}) = \infty$

pf. idea: Choose set of pts. $A = \{2^{-k} : k = 1, \ldots, n\}$ & show can hit each pt w/ opposite sine signs using arbitrarily large frequency $\alpha$.

## 1.7 – Fundamental Theorem of binary classification (qualitative version)

Def: $\text{poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}\right) :=$ fns of the form $(0,1]^2 \ni (\varepsilon, \delta) \mapsto r(\varepsilon, \delta) \in \mathbb{R}^+$
polynomial in both $\frac{1}{\varepsilon}$ & $\frac{1}{\delta}$.

Thm: Consider $\bar{\mathcal{F}} \in \{\pm 1\}^{\mathcal{X}}$, $n := |S|$ - training data set, $S \sim P^n$.

Let $R$ be the error prob. Then TFAE:

(1) $\text{VCdim}(\bar{\mathcal{F}}) < \infty$

(2) $\exists r \in \text{poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}\right)$ s.t. $\forall (\varepsilon, \delta) \in (0,1]^2 \;\; \forall P$:

$\qquad n \geq r(\varepsilon, \delta) \implies \mathbb{P}_S\left[\exists h \in \bar{\mathcal{F}} : |\hat{R}(h) - R(h)| \geq \varepsilon\right] \leq \delta$.

(3) $\exists r \in \text{poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}\right)$ & a learner $S \mapsto h_S \in \bar{\mathcal{F}}$ s.t. $\forall (\varepsilon, \delta) \in (0,1]^2$, $\forall P$:

$\qquad n \geq r(\varepsilon, \delta) \implies \mathbb{P}\left[|R(h_S) - R_{\bar{\mathcal{F}}}| \geq \varepsilon\right] \leq \delta$.

(4) ... as (3) but w/ $S \mapsto h_S$ is ERM.

Pf: $(1) \Rightarrow (2)$ ✓, $(2) \Rightarrow (4)$ (via Error Decomposition) ✓, $(4) \Rightarrow (3)$ by def. ✓
$\qquad$ by Corr. on §1.3),

WTS: $(3) \Rightarrow (1)$ by no free lunch Thm & contradiction
$\quad$ Choose $\varepsilon = \delta = 1/4$. Let $n = r(\varepsilon, \delta)$ & supp. $\text{VCdim}(\bar{\mathcal{F}}) = \infty$.
$\quad \forall N \in \mathbb{N}$ $\exists C \subseteq \mathcal{X} : |C| = N$ s.t. $|\bar{\mathcal{F}}_{|C}| = 2^N \implies \bar{\mathcal{F}}_{|C} = \{\pm 1\}^C$. (by def. of VCdim

From No-Free-Lunch Thm, $\exists f : C \to \{\pm 1\}$ and $P(x,y) := \left(\mathbb{1}_{x \in C} \cdot \mathbb{1}_{f(x) = y}\right)/N$
over $\mathcal{X} \times \{-1, 1\}$ wrt. which $\mathbb{E}_S[R(h_S)] \geq \frac{1}{2}\left(1 - \frac{n}{N}\right) = \left(1 - \frac{1}{|y|}\right)\left(1 - \frac{n}{|C|}\right)$

but $\mathbb{E}_S[R(h_S)] \leq \mathbb{P}_S[R(h_S) \geq \varepsilon] + \varepsilon\left(1 - \mathbb{P}_S[R(h_S) \geq \varepsilon]\right)$ $\rbrace$ combine

$\implies \mathbb{P}_S\left[R(h_S) \geq \frac{1}{4}\right] \geq \frac{1}{3} - \frac{2n}{3N}$ $\quad \rbrace$ $\delta = 1/4$ b/c $R_{\bar{\mathcal{F}}} = 0$. $\qquad \square$

Quantitative Version yields "sample complexity" of $\bar{\mathcal{F}}$:

$$r_{\bar{\mathcal{F}}}(\varepsilon, \delta) = \Theta\left(\frac{\text{VCdim}(\bar{\mathcal{F}}) + \ln\frac{1}{\delta}}{\varepsilon^2}\right)$$

## 1.8 Rademacher complexities:

**Def:** Consider a set of real-valued fns. $\mathcal{G} \subseteq \mathbb{R}^{Z}$ & a vector $z \in Z^n$. The "empirical Rademacher complexity" of $\mathcal{G}$ w.r.t. $z$

is $\hat{R}(\mathcal{G}) := \mathbb{E}_\sigma \left[ \frac{1}{n} \sup_{g \in \mathcal{G}} \left\{ \sum_{i=1}^n \sigma_i \, g(z_i) \right\} \right]$, $\sigma \in \{-1, 1\}^n$ uniform "Rademacher var"

*(not necc $Z$, any set!)*

If $z_i$'s are iid r.vs. then the "Rademacher complexity" is defined as

$$R_n(\mathcal{G}) := \mathbb{E}(\hat{R}(\mathcal{G})) .$$

with $g(z) = (g(z_1), \ldots, g(z_n))$, $\hat{R}(\mathcal{G}) = \frac{1}{n} \mathbb{E}_\sigma \sup_{g \in \mathcal{G}} \langle \sigma, g(z) \rangle$

↖ will be big if $\mathcal{G}$ is "rich enough" for $g(z)$ to align w/ random signs.

---

**Lemma:** (McDiarmid's Ineq.)

Let $(z_1, \ldots, z_n) =: z$ be indep. rand. vars. w/ values in $Z$, and $\varphi : Z^n \to \mathbb{R}$ s.t. $|\varphi(z) - \varphi(z')| \le r_i$ whenever $z$ & $z'$ differ only in the $i^{th}$ component.

Then $\forall \varepsilon > 0$, $\mathbb{P}\left[ \varphi(z) - \mathbb{E}[\varphi(z)] \ge \varepsilon \right] \le \exp\left\{ -\frac{2\varepsilon^2}{\sum_{i=1}^n r_i^2} \right\}$.

for classification

Want to bound risk by empirical risk: $|R(h) - \hat{R}(h)| \leq$ "growth fn" $\leq$ "VC dim"

$\leq$ "Rademacher compl." $\leq$ "emp. Rademach"

$\leq$ "$L_2$-covering number" $\leq$ "$L_1$-covering #"

$\leq$ "$L_1$-packing"

$\leq$ "$T$" $\leq$ "VC-dim"

beats other VC bound

Recall: $\mathcal{C} \subseteq \mathbb{R}^Z$, $z \in Z$, $Z$ = some set (not necessarily $\mathbb{Z}$)

"Empirical Rad. compl" $\hat{R}(\mathcal{C}) := \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n g(z_i) \sigma_i \right]$

Def: (Rademacher complexity) $\mathbb{E}_Z \left[ \hat{R}(\mathcal{C}) \right] =: R_n(\mathcal{C}) \leftarrow$ Rademacher complexity

Lemma: $\mathcal{C} \subseteq [a,b]^Z$. For any $\varepsilon > 0$, prob. measure $P^n$ on $Z^n$:

$$P_Z \left[ R_n(\mathcal{C}) - \hat{R}(\mathcal{C}) \geq \varepsilon \right] \leq \exp\left[ -\frac{2n\varepsilon^2}{(b-a)^2} \right]$$

Pf: $\varphi: Z^n \to \mathbb{R}$, $\varphi(z) : \hat{R}(\mathcal{C})$. So $\mathbb{E}_Z[\varphi(z)] = R_n(\mathcal{C})$.

Let $z, z' \in Z^n$ differ in only one component. Then for fixed $\sigma_i$,

$\sup_{g \in \mathcal{C}} \sum_i \sigma_i g(z_i)$ changes at most by $|a-b|$.

Thus $|\varphi(z) - \varphi(z')| = |\hat{R}_z(\mathcal{C}) - \hat{R}_{z'}(\mathcal{C})| \leq \frac{|b-a|}{n}$

Then the result follows from McDiarmid's inequality. $\square$

Thm: $\bar{\mathcal{F}} \subseteq Y^X$, $L: Y \times Y \to [0,c]$, $Z := X \times Y$,

$\mathcal{C} := \{ (x,y) \mapsto L(y, h(x)) \mid h \in \bar{\mathcal{F}} \} \subseteq [0,c]^Z$. For any $\delta > 0$

and prob measure $P$ on $Z$, we have w/ prob at least $1-\delta$ and repeated sampling of $S \in (X \times Y)^n$ dist. according to $P^n$:

$\forall h: R(h) - \hat{R}(h) \leq \begin{cases} 2 R_n(\mathcal{C}) + c\sqrt{\dfrac{\ln 1/\delta}{2n}} \\[4mm] 2\hat{R}(\mathcal{C}) + 3c\sqrt{\dfrac{\ln 2/\delta}{2n}} \end{cases} \Big) $ using Lemma

$\to 0$ as $n \to \infty$

**Lemma:** $\mathcal{F} \subseteq \{-1, 1\}^{\mathcal{X}}$, $L = 0\text{-}1$ loss, $S = ((x_i, y_i))_{i=1}^n$, $S_{\mathcal{X}} = (x_i)_{i=1}^n$

$\quad\quad p$ marginal prob of $P$ on $\mathcal{X}$.

$$\hat{R}_S(\mathcal{G}) = \tfrac{1}{2} \hat{R}_{S_{\mathcal{X}}}(\bar{\mathcal{F}}).$$

**Pf:** $L(y, h(x)) = \tfrac{1}{2}(1 - y h(x))$.

$$\hat{R}_S(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{F}} \tfrac{1}{n} \sum \tfrac{\sigma_i}{2}(1 - y_i h(x_i)) \right]$$

$$\underset{\mathbb{E}_\sigma(\sigma_i) = 0}{=} \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{F}} \tfrac{1}{2n} \sum \sigma_i \cdot h(x_i) \right] \quad \Leftarrow \quad \begin{array}{l} \sigma_i \sim -\sigma_i \text{ dist. same} \\ y_i \in \{-1, 1\} \text{ is fixed.} \end{array}$$

$$= \tfrac{1}{2} \hat{R}_{S_{\mathcal{X}}}(\bar{\mathcal{F}}) \qquad \square$$

**Properties:**

- $\hat{R}(c\mathcal{G}) = |c|\hat{R}(\mathcal{G})$ for $c \in \mathbb{R}$.
- $\mathcal{G}_1 \subseteq \mathcal{G}_2 \Rightarrow \hat{R}(\mathcal{G}_1) \le \hat{R}(\mathcal{G}_2)$.
- $\hat{R}(\mathcal{G}_1 + \mathcal{G}_2) = \hat{R}(\mathcal{G}_1) + \hat{R}(\mathcal{G}_2)$.
- $\hat{R}(\text{conv. } \mathcal{G}) = \hat{R}(\mathcal{G})$. (convex hull) $\text{conv}(\mathcal{G}) = \left\{ z \mapsto \sum_{i=1}^n \lambda_i g_i(z) : \begin{array}{l} n \in \mathbb{N}, g_i \in \mathcal{G}, \\ \lambda_i \ge 0, \\ \sum \lambda_i = 1 \end{array} \right\}$
- $\psi : \mathbb{R} \to \mathbb{R} : L\text{-Lipschitz} \Rightarrow \hat{R}(\psi \circ \mathcal{G}) \le L \cdot \hat{R}(\mathcal{G})$

- $\hat{R}(\mathcal{G})$ is NP-hard to compute

## 1.9 Covering Numbers

**Def:** $(M, d)$-pseudometric space (same as metric except $d(x,y) = 0 \not\Rightarrow x = y$)

Let $A, B \subseteq M$, $\varepsilon > 0$. $A$ is an $\varepsilon$-cover of $B$ if $\forall b \in B \; \exists a \in A : d(a, b) \le \varepsilon$.

This is called "internal" if, in addition, $A \subseteq B$.

The "$\varepsilon$-covering number" of $B$, $N(\varepsilon, B)$ is the smallest cardinality of an $\varepsilon$-cover of $B$.

**Def:** $A \subseteq B$ is an "$\varepsilon$-packing" of $B$ if $a, b \in A \wedge a \ne b \Rightarrow d(a, b) > \varepsilon$.

The "$\varepsilon$-packing number" of $B$, $M(\varepsilon, B)$ is largest cardinality of $\varepsilon$-packing of $B$.

**Prop:** $N(\varepsilon/2, B) \ge M(\varepsilon, B) \ge N_{in}(\varepsilon, B) \ge N(\varepsilon, B)$

<u>Covering/Packing # Example:</u> (norm balls in $\mathbb{K}^d$) $\|\cdot\|_{any}$ norm on $\mathbb{K}^d$.

Consider $B_r(x) := \{z \in \mathbb{R}^d \mid \|z - x\| \le r\}$

Supp. $\{x_1, \ldots, x_m\} \subseteq \mathbb{R}^d$ is a max. $\varepsilon$-packing of $B_r(0)$. So $M = M(\varepsilon, B_r(0))$.

Then $i \ne j$, $B_{\varepsilon/2}(x_i) \cap B_{\varepsilon/2}(x_j) = \emptyset$, $B_{\varepsilon/2}(x_i) \subseteq B_{r+\varepsilon/2}(0)$.

$v := \text{vol}(B_1(0))$. Then $\quad M(\varepsilon, B_r(0)) \le \dfrac{\text{vol}(B_{r+\varepsilon/2}(0))}{\text{vol}(B_{\varepsilon/2}(x_i))} = \dfrac{(r + \varepsilon/2)^d \, v}{(\varepsilon/2)^d \, v} \le \left(\dfrac{3r}{\varepsilon}\right)^d$

$\underset{\varepsilon \le r}{\uparrow}$

$\ln M(\varepsilon, B) = O(d \ln \tfrac{1}{\varepsilon})$. essentially same is true when alg. dim $d$ replaced by combinatorial dim - e.g. VC-dim.

Consider $g \in \mathcal{G} \subseteq \mathbb{R}^Z$, $p \in [1, \infty)$, $z \in Z^n$.

Define $\quad \|g\|_{p,z} := \left(\dfrac{1}{n} \sum_{i=1}^{n} |g(z_i)|^p\right)^{1/p} \quad$ seminorm $\quad$ Note: $\|g\|_{p,z} \le \|g\|_{q,z}$ if $p \le q$.

pseudometric $(g_1, g_2) \longmapsto \|g_1 - g_2\|_{p,z}$

Then $\quad M(\varepsilon, \mathcal{G}, \|\cdot\|_{p,z}) \le M(\varepsilon, \mathcal{G}, \|\cdot\|_{q,z})$ if $p \le q$.

<u>Lemma:</u> $\bar{\mathcal{F}} \subseteq \{0,1\}^Z$, $d := \text{VCdim}(\bar{\mathcal{F}})$. (Pf. in Lecture Notes)

$M(\varepsilon, \bar{\mathcal{F}}) \le \left(\dfrac{9}{\varepsilon^p} \ln \dfrac{2e}{\varepsilon^p}\right)^d \overset{\longleftarrow}{\quad}$ bound independent of $n$!

$\underset{\text{dep... } p}{\uparrow}$

## Lec 7 - 4/25/17

$z, z$

<u>Thm</u> (Dudley's chaining thm) Let $z \in Z^n$. $\mathcal{G} \subseteq \mathbb{R}^Z$ equipped w/ $\|\cdot\|_{2,z}$. $\gamma_0 = \sup_{g \in \mathcal{G}} \|g\|_{2,z}$

Then $\quad \widehat{R}(\mathcal{G}) \le \dfrac{12}{\sqrt{n}} \int_0^{\gamma_0} \ln(N(\beta, \mathcal{G}))^{1/2} d\beta$ $\quad$

$\underset{\beta\text{-covering number. if ind. of } n, \text{ this bound } \approx n^{-1/2}!}{\uparrow}$

<u>Pf:</u> $\gamma_j := 2^{-j}\gamma_0, j \in \mathbb{N}$. Let $G_j \subseteq \mathbb{R}^Z$ be minimal $\gamma_j$-cover of $\mathcal{G}$.

$|G_j| = N(\gamma_j, \mathcal{G})$. $\forall g \in \mathcal{G} \; \exists g_j \in G_j$ s.t. $\|g - g_j\|_{2,z} \le \gamma_j$. Note: $g_0 := 0$ by def. of

$\mathcal{G} \ni g = g - g_m + \sum_{j=1}^{m} (g_j - g_{j-1}) \quad$ For $m \in \mathbb{N}$

$\widehat{R}(\mathcal{G}) = \dfrac{1}{n} \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n} \sigma_i\left(g(z_i) - g_m(z_i) + \sum_{j=1}^{m} g_j(z_i) - g_{j-1}(z_i)\right) \right]$

$\le \dfrac{1}{n} \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \sum_{i=1}^{n} \sigma_i(g(z_i) - g_m(z_i)) \right] + \dfrac{1}{n} \sum_{j=1}^{m} \mathbb{E}_\sigma \sup_{g \in \mathcal{G}} \sum_{i=1}^{n} \sigma_i(g_j(z_i) - g_{j-1}(z_i)) \quad \leftarrow$ treat as inner products
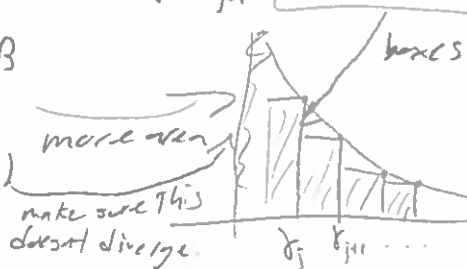
By Cauchy-Schwarz, $\frac{1}{n}\mathbb{E}_\sigma\left[\sup_{g\in G}\sum_{i=1}^{n}\sigma_i(g(z_i)-g_m(z_i))\right] \le \|\sigma\|_{2,z}\|g-g_m\|_{2,z} = 1\cdot\gamma_m$

By Massart's Lemma, $\frac{1}{n}\sum_{j=1}^{m}\mathbb{E}_\sigma\left[\sup_{g\in G}\sum_{i=1}^{n}\sigma_i(g(z_i)-g_{j-1}(z_i))\right] \le \frac{12}{\sqrt{n}}\sum_{j=1}^{m}(\gamma_j-\gamma_{j+1})\underbrace{\sqrt{\ln(N(\gamma_j,G))}}$

Hence $\hat{R}(G) \le \gamma_m + \frac{12}{\sqrt{n}}\int_{\gamma_{m+1}}^{\gamma_0}\sqrt{\ln(N(\beta,G))}\,d\beta$



boxes

more area

make sure this doesn't diverge $\gamma_j \ \gamma_{j+1} \cdots$

as $m\to\infty$, $\gamma_m\to 0$. (be careful $N$ behaves well in this limit)

$\Rightarrow \hat{R}(G) \le \frac{12}{\sqrt{n}}\int_0^{\gamma_0}\ln(N(\beta,G))^{1/2}d\beta.$ $\qquad\square$

**Cor:**
To apply this, note $N(\beta,G)\le M(\beta,G)$ — covering # upper bd by packing #,
then use Lemma $M(\mathcal{E},\mathcal{F})\le\left(\frac{4}{\mathcal{E}^2}\ln\frac{2e}{\mathcal{E}^2}\right)^d$

use $\ln x \le \frac{x}{e}$ $\Rightarrow$ " $\le \left(\frac{16}{\mathcal{E}^4}\right)^d$ & plug into Dudley's

$\hat{R}(G) \le \frac{12}{\sqrt{n}}\int_0^{\gamma_0}\ln\left(\left(\frac{16}{\beta^4}\right)^d\right)^{1/2}d\beta = 12\sqrt{\frac{d}{n}}\int_0^{\gamma_0}\underbrace{\cdots}_{\text{some constant}}d\beta$

$\Rightarrow \hat{R}(G) \le 31\underbrace{\sqrt{d/n}}_{\text{ideal?}}\longleftarrow$ best known scaling!

**Recall:** $R(h)-\hat{R}(h)\le\hat{R}(G)+O\left(\frac{1}{\sqrt{n}}\right)$

So now $R(h)-\hat{R}(h)\le O\left(\sqrt{d/n}\right)+O\left(\frac{1}{\sqrt{n}}\right)$ ! Best bound.

**1.10 Algorithmic Stability:** So far, have only seen learning algorithm in its range, $\mathcal{F}$.
• Small changes in input $\to$ small changes in hypothesis $\Rightarrow$ learning alg. stable.
• Eg. linear regression = stable, high deg poly reg = probably not stable
• Lack of stability = sign of overfitting

$S\to h_S$

**Def:** A learning alg. $\hat{}$ is "uniformly stable" w/ rate $\mathcal{E}:\mathbb{N}\to\mathbb{R}$ w/ loss fn $L$
$\quad$ if $\forall n\in\mathbb{N}$, $S\in(X,Y)^n$, $i\in\{1,\dots,n\}$, $(x',y')\in(X\times Y)$:
$$L(y_i,h_{S^i}(x_i))-L(y_i,h_S(x_i))\le\mathcal{E}(n)$$
$\quad$ w/ $S^i=S$ where $i^{th}$ component is replaced by $(x',y')$.

very strong, implies next def.

**Def:** A learning alg. is "on-average stable" w/ rate $\mathcal{E}:\mathbb{N}\to\mathbb{R}$ if $\forall P$ problems on $X\times Y$.
$$\mathbb{E}_{S\sim P^n}\mathbb{E}_{(x',y')\sim P}\mathbb{E}_{i\sim\text{unif}}[L(\cdots)-L(\cdots)]\le\mathcal{E}(n)$$ w/ all same other assumptions.

Thm: If $S \mapsto h_S$ is $\varepsilon$-on average stable, then        | on-average stable
$$\mathbb{E}_S \left[ R(h_S) - \hat{R}(h_S) \right] \leq \varepsilon(n) .$$        $\Rightarrow$ generalizes well

Pf: Relies heavily on i.i.d. assumption for $S$ & $(x', y')$ $(\sim p^n$ and $\sim p)$

$$\mathbb{E}_S \left[ R(h_S) \right] = \mathbb{E}_S \, \mathbb{E}_{(x', y')} \left[ L(y', h_S(x')) \right]$$
$$= \mathbb{E}_i \, \mathbb{E}_S \, \mathbb{E}_{(x', y')} \left[ L(y_i, h_{S^i}(x_i)) \right] \quad \in Obs. \, 1.$$

$$\mathbb{E}_S \left[ \hat{R}(h_S) \right] = \mathbb{E}_S \, \mathbb{E}_i \left[ L(y_i, h_S(x_i)) \right]$$
$$= \mathbb{E}_S \, \mathbb{E}_{(x', y')} \, \mathbb{E}_i \left[ \text{ ''} \right] \quad \in Obs. \, 2$$

$$\Rightarrow \mathbb{E}_S \left[ R(h_S) - \hat{R}(h_S) \right] = \mathbb{E}_S \, \mathbb{E}_{(x', y')} \, \mathbb{E}_i \left[ L(y_i, h_{S^i}(x_i)) - L(y_i, h_S(x_i)) \right] \leq \varepsilon(n)$$
$$\square$$

$\Bigl($ We will see that the alg. $S \mapsto \underset{h}{\operatorname{argmin}} \, \hat{R}(h) + \lambda \|h\|^2$ is unif. stable $\Bigr)$
Regularization (Tikhonov here) adds "stability" by adding strict convexity.

Def: $f$ is $\alpha$-strongly convex if
$\Updownarrow$
· $x \mapsto f(x) - \frac{\alpha}{2} \|x\|^2$ is convex
· $f(\lambda x + (1-\lambda) y) \leq \lambda f(x) + (1-\lambda) f(y) - \frac{\alpha}{2} \lambda (1-\lambda) \|x - y\|^2$

**Thm:** Assume range $L \subseteq [-c, c]$ & $S \mapsto h_S$ is unif. stable w/ rate $\mathcal{E}_1$.

$\forall \mathcal{E} > 0 \ \forall P:$

$$\mathbb{P}_S \left[ |\hat{R}(h_S) - R(h_S)| \geq \mathcal{E} + \mathcal{E}_1(n) \right] \leq 2 \exp \left[ -\frac{n \mathcal{E}^2}{2(n \mathcal{E}_1(n) + c)^2} \right] \quad \left| \begin{array}{c} \text{PAC} \\ \text{bound} \end{array} \right.$$

**Pf:** Let $\Psi(S) := \hat{R}(h_S) - R(h_S)$. By assumption, $|\mathbb{E}[\Psi(S)]| < \mathcal{E}_1(n)$.

Then $|\Psi(S)| \geq \mathcal{E} + |\mathbb{E}[\Psi(S)]| \Rightarrow |\Psi(S) - \mathbb{E}[\Psi(S)]| \geq \mathcal{E}$.

$\therefore \mathbb{P}\left[ |\hat{R}(h_S) - R(h_S)| \geq \mathcal{E} + \mathcal{E}_1(n) \right] \leq \mathbb{P}\left[ |\Psi(S) - \mathbb{E}[\Psi(S)]| \geq \mathcal{E} \right] \leq 2\exp\left[ -\frac{2\mathcal{E}^2}{n \nu^2} \right]$ by McDiarmid's

where $\nu$ s.t. $\nu \geq |\Psi(S) - \Psi(S^i)|$.

$|\Psi(S) - \Psi(S^i)| \leq \underbrace{\frac{1}{n} \sum_{j \neq i} |L(y_j, h_S(x_j)) - L(y_j, h_{S^i}(x_j))|}_{\downarrow \text{ by unif. stable}} + \frac{2c}{n} + \underbrace{|R(h_S) - R(h_{S^i})|}_{}$

$\leq \mathcal{E}_1(n) + \frac{2c}{n} + \underbrace{|\mathbb{E}_{(x,y)}[L(y, h_S(x)) - L(y, h_{S^i}(x))]|}_{\leq \mathcal{E}_1(n)} \leq 2\left(\mathcal{E}_1(n) + \frac{c}{n}\right) =: \nu.$ $\square$

**Def:** $\Phi$ is an $\alpha$-strongly convex function $(\alpha > 0)$ if $h \mapsto \Phi(h) - \frac{\alpha}{2}\langle h, h \rangle$ is convex

and/or if $\lambda \Phi(h) + (1-\lambda)\Phi(g) \geq \Phi(\lambda h + (1-\lambda)g) + \frac{\alpha}{2}\lambda(1-\lambda)\|h-g\|^2 \quad \forall h, g \in \bar{\mathcal{F}}$

$\forall \lambda \in [0,1].$

**Lemma:** If $\Phi: \bar{\mathcal{F}} \to \mathbb{R}$ is $\alpha$-strongly convex & attains its minimum at $h$, then $\forall g \in \bar{\mathcal{F}}$,

$$\Phi(g) \geq \Phi(h) + \frac{\alpha}{2}\|h-g\|^2.$$

**Pf:** $h$ is minimizer, so $\Phi(h) \leq \Phi(\lambda h + (1-\lambda)g)$

$\Rightarrow \Phi(h) + \frac{\alpha}{2}\lambda(1-\lambda)\|h-g\|^2 \leq \lambda\Phi(h) + (1-\lambda)\Phi(g)$

$(1-\lambda)\Phi(h) + \frac{\alpha}{2}\lambda(1-\lambda)\|h-g\|^2 \leq (1-\lambda)\Phi(g)$

$\Rightarrow \Phi(h) + \frac{\alpha}{2}\lambda\|h-g\|^2 \leq \Phi(g) \quad$ set $\lambda = 1.$ $\square$

**Thm:** $\lambda > 0.$ Let $\mathcal{F}$ be a convex subset of an inner product space.

Assume $h \mapsto L(y, h(x))$ is convex & $\ell$-Lipschitz $\forall x, y$.

Then if $S \mapsto h_S$ minimizes $f_S(h) := \hat{R}(h) + \lambda\langle h, h \rangle$,

$S \mapsto h_S$ is unif. stable w/ rate $\mathcal{E}(h) = \frac{2\ell^2}{\lambda n}.$

Pf. (of Regularization $\Rightarrow$ Unif Stability): Let $h := h_S$, $h' := h_{S^i}$ minimizers of $f_S(h)$ & $f_{S^i}(h)$

$$f_S(h') - f_S(h) = \hat{R}_S(h') - \hat{R}_S(h) + \lambda\left(\|h'\|^2 - \|h\|^2\right)$$

$$= \hat{R}_{S^i}(h') - \hat{R}_{S^i}(h) + \lambda\left(\|h'\|^2 - \|h\|^2\right)$$

$$+ \frac{1}{n}\left[L(y_i, h'(x_i)) - L(y_i, h(x_i)) + L(y_i', h(x_i')) - L(y_i', h'(x_i'))\right]$$

$h'$ min. $f_{S^i}$ $\longrightarrow$ $\leq \underset{\uparrow}{\frac{1}{n}[L\cdots]} \leq \frac{2\ell}{n}\|h - h'\|$

by Lipschitz

Moreover, $\lambda\|h - h'\|^2 \leq f_S(h') - f_S(h)$ since $h$ minimizes $f_S$ which is $2\lambda$-strongly convex

Put them together, get $\|h - h'\| \leq \frac{2\ell}{n\lambda}$

Def. of unif stable: $L(y_i, h'(x_i)) - L(y_i, h(x_i)) \leq \ell\|h - h'\| \leq \frac{2\ell^2}{\lambda n}$. $\quad\square$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Lipschitz

Thm: Let $h^* := \underset{h \in \mathcal{F}}{\arg\min} R(h)$.

If $S \mapsto h_S$ is reg. ERM w/ reg-param. $\lambda$ & $\ell$-Lipschitz loss then

$$\mathbb{E}_S[R(h_S)] \leq R(h^*) + \lambda\|h^*\|^2 + \frac{2\ell^2}{\lambda n}.$$

Pf: $\mathbb{E}_S[\hat{R}(h_S)] \leq \mathbb{E}_S[\hat{R}(h_S) + \lambda\|h_S\|^2] \underset{\uparrow}{\leq} \mathbb{E}_S[\hat{R}(h^*) + \lambda\|h^*\|^2] = R(h^*) + \lambda\|h^*\|^2$.

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ by $h_S$ reg-ERM.

$$\mathbb{E}_S[R(h_S)] = \mathbb{E}_S[\hat{R}(h_S)] + \mathbb{E}_S[R(h_S) - \hat{R}(h_S)]$$

$$\leq R(h^*) + \lambda\|h^*\|^2 + \frac{2\ell^2}{\lambda n} \quad \text{by unif.stable.}$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Future Topics: · PAC-Bayesian · include a priori information about data dist/model

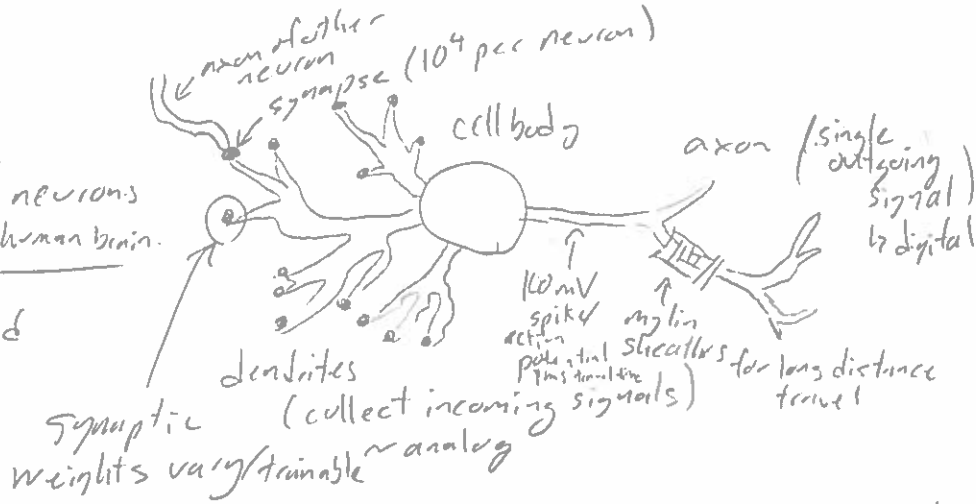$\qquad\qquad\qquad\qquad$ · Ensemble methods — ADA Boost

Next week — start neural networks (for at least 3 weeks)

# II. Neural Networks

## 2.1 Biological NNs:

- grey matter - mostly neurons
- white matter - mostly mylin-covered axons

$10^{11}$ neurons in human brain.

- 6 horizontal layers
- cortical columns (~100 neurons)

| computer | brain |
|---|---|
| 70 W/laptop | 20W |
| #$10^{10}$ transistors | #$10^{11}$ neurons |
| connectedness: ~3 | ~$10^4$ |
| clockrate GHz | 100 Hz - not universal → constraints logical depth |
| deterministic | stochastic |
| 2D | 3D |

axon of other neuron
synapse ($10^4$ per neuron)
cell body
axon (single outgoing signal) ↳ digital

100 mV spike action potential this timeline
mylin sheathes for long distance travel

dendrites (collect incoming signals) ~analog

synaptic weights vary/trainable

**Relevant for Model**
- many dendritic inputs single axon output
- analog in/digital out
- synaptic variability

## Rosenblatt's Perceptron '58

activation fn. $\sigma: \mathbb{R} \to \mathbb{R}$, for perceptron

$$\sigma(z) := \mathbb{1}_{z \geq 0} = \text{sgn}(z).$$

$\mathbb{R} \ni x_1, x_2, \ldots, x_d$ with weights $w_1, w_2, \ldots, w_d$

$$\sigma\left(w_0 + \sum_i w_i x_i\right)$$

bias/threshold   weights

other activation fns
$\sigma(z) = \frac{1}{1+e^{-z}}$ - "Logistic sigmoid"

$\tanh z$

**rectified linear unit** (ReLU)
$\sigma(z) = \max\{0, z\}$
almost exclusively used now

## Networks

architecture - graph $(V, E)$

↳ graph is weighted + directed + acyclic feed-forward)

vertices $V$ = artificial neurons
edges $E$ = connections btwn $V$ w/ weights $w_{ij} \in \mathbb{R}$
$i, j \in V$ and $(i, j) \in E$
+ bias per vertex

often we consider multi-layered Networks/graphs

Feed-Forward Multilayer Acyclic Neural Network.
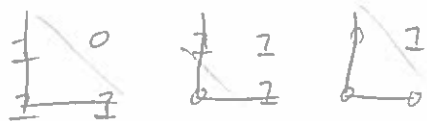
what kind of fns. can we represent by this model?

# Representations & Approximations

single perceptron: $f(x) := \sigma\left(w_0 + \sum_i w_i x_i\right)$, $f: \mathbb{R}^d \to \mathbb{R}$

$\langle w, x \rangle$

Constant on hyperplanes $\to$ can use $\sigma$-step fn. to use this model to separate space via a hyperplane.

$\hookrightarrow$ Can use on linearly (via hyperplane) separable data for binary classification

NAND, OR, AND: $\mathbb{R}^2 \to \mathbb{R}$ $\leftarrow$ can we rep. these via a perceptron?

 $\leftarrow$ can separate via hyperplane. so yes! $\to$ can build any boolean fn. via composition of these
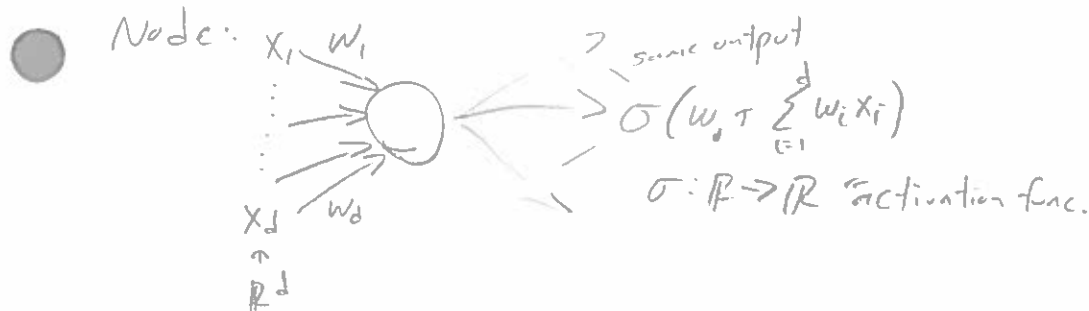
XOR - requires composition of above perceptrons
- can't be done by single perceptron since data isn't linearly separable.



email: mwolf@math.ucdavis.edu
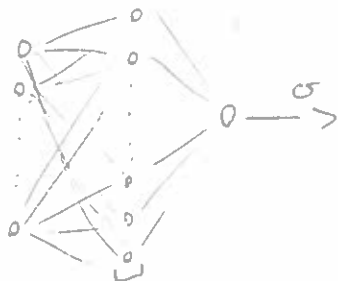for project group 1-3 people.

Neural Network:

Node:



$$\sigma\left(w_0 + \sum_{i=1}^{d} w_i x_i\right)$$

some output

$\sigma: \mathbb{R} \to \mathbb{R}$ "activation func."

$x \uparrow \mathbb{R}^d$

**Thm:** Every $f: \{0,1\}^d \to \{0,1\}$ can be represented by a feed-forward NN w/ a single hidden layer containing at most $2^d$ neurons, if $\sigma(z) := \mathbb{I}_{z \geq 0}$ is used.

**Pf:** If $a, b \in \{0,1\}$ then $2ab - a - b \leq 0$ with "=" iff $a = b$.

Thus $\mathbb{I}_{x=u} = \sigma\left(\sum_{i=1}^{d} 2 x_i u_i - x_i - u_i\right)$. $A := f^{-1}(\{1\})$

$$f(x) = \sigma\left(-1 + \sum_{u \in A} \mathbb{I}_{x=u}\right) = \sigma\left(-1 + \sum_{u \in A} \sigma\left(\sum_{i=1}^{d} 2 x_i u_i - x_i - u_i\right)\right)$$



$A \in$ at most $2^d$ since $A = f^{-1}(\{1\})$ so at most $A = \{0,1\}^d$. $\square$

Assume input space is $\mathcal{X} = \mathbb{R}^d$, output space $\mathcal{Y} = \{0, 1\}$.
Consider $\sigma(z) = \mathbb{I}_{z \geq 0}$. Then every indiv. neuron $j$ is characterized by a halfspace $\mathcal{H}_j \subseteq \mathbb{R}^d$ (looking at the hidden layer $j = 1, \dots, m$)

$$x \overset{\mathbb{R}^d}{\longmapsto} f(x) = \sigma\left(w_0 + \sum_{j=1}^{m} w_j \mathbb{I}_{x \in \mathcal{H}_j}\right)$$

$$\mathcal{A} := \left\{A \subseteq \{1, \dots, m\} \mid \sum_{j \in A} w_j \geq -w_0\right\} \in \text{combinations of neurons in hidden layer who "fire w/ final output fire"}$$

$$\mathbb{R}^d \supseteq f^{-1}(\{1\}) = \bigcup_{A \in \mathcal{A}} \underbrace{\bigcap_{j \in A} \mathcal{H}_j}_{\substack{\text{convex polyregion} \\ \text{for fixed } A \\ \text{(intersection of half-planes)}}}$$

**Thm (Zaslavsky's):** (Given $n$-hyperplanes in $\mathbb{R}^d$, how many regions can we make?)

Let $h_1, \ldots, h_n \subseteq \mathbb{R}^d$ by hyperplanes. The number $N$ of connected components

of $\mathbb{R}^d \setminus \bigcup_{j=1}^{n} h_j$ is $N \le \sum_{i=0}^{d} \binom{n}{i} \le \left(\frac{en}{d}\right)^d$

(↑ tight bound, often the exact #.)

---

**Pf:** Shift hyperplanes away from origin, each can be characterized w/ a single vector.

$$h := \{x \in \mathbb{R}^d \mid \langle w, x \rangle = 1\}, \quad \mathcal{X} \text{ is the set of all such hyperplanes.}$$

↑ $w \in \mathbb{R}^d$

$$\bar{f} := \{h \mapsto \text{sgn}[g(h)] \mid g \in G\}. \quad G = \{h \mapsto \langle x, w \rangle - 1 \mid x \in \mathbb{R}^d\} \subseteq \mathbb{R}^{\mathcal{X}}.$$

∩              ↑ affine space w/ $\dim(G) = d$.

$\{-1, 1\}^{\mathcal{X}}$

**Obs:** $VC\dim(\bar{f}) = d$.

Now assume $A \subseteq \mathcal{X}$ w/ $|A| = n$. separates $\mathbb{R}^d$ into $N$ regions.

Then $\bar{f}|_A$ contains at least $N$ different fns.

Thus $N \le \Gamma(n) \le \sum_{i=0}^{d} \binom{n}{i}$. □

↑ growth fn. of $\bar{f}$

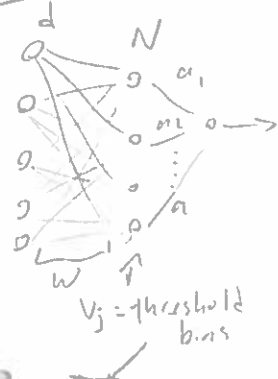*(margin right:)* but still need more neurons ↑ anyway to make... problem: too many neurons = memorization = ↑ overfitting

**Thm:** Let $A = \{x_1, \ldots, x_N\} \subseteq \mathbb{R}^d$ and $f : A \to \mathbb{R}$.

There is a feed-forward NN that implements $F : \mathbb{R}^d \to \mathbb{R}$

w/ a single hidden layer containing $N$ neurons ($\& \ 2N+d$ parameters)

so that $F|_A = f$. We can use $\sigma(z) = \max\{0, z\}$ in the hidden layer

$\& \ \sigma(z) = z$ for the output layer.

**Pf:** $F : \mathbb{R}^d \to \mathbb{R}, \quad F(x) = \sum_{j=1}^{N} a_j \max\{0, \langle w, x \rangle - v_j\}, \quad a, v \in \mathbb{R}^N$



$w \in \mathbb{R}^d$

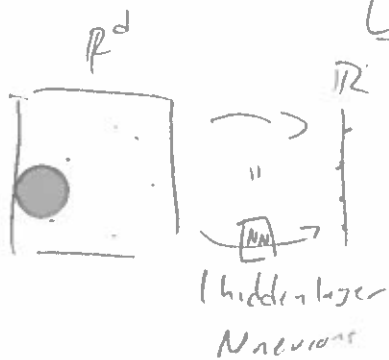All first-layer connections have shared weights $w$.

Let $M_{ij} = \max\{0, \langle w, x_i \rangle - v_j\}$

Then $F(x) = Ma$. If $M^{-1}$ exists,

$\Rightarrow a_j = (M^{-1})_{ji} f(x_i)$ solves the problem exactly.

Assuming $x_i$ distinct, $\exists \ w \in \mathbb{R}^d$ s.t. $\langle w, x_i \rangle = z_i$ distinct as well.

Reorder $z$'s & pick $v$'s s.t. $v_1 < z_1 < v_2 < z_2 < \ldots$

*(bottom margin:)* $v_j$ = threshold bins

1 hidden layer
N neurons

**Thm:** Let $\sigma \in C(\mathbb{R})$. The set $\mathcal{F}_\sigma$ of fns.

representable by a NN w/ a single hidden layer &
activation function $\sigma$ is dense in $C(\mathbb{R})$ wrt.
the topology of unif. conv. on compact sets
**iff** $\sigma$ is not a polynomial.

pf (sketch): Assume $\sigma \in C^\infty(\mathbb{R})$. If $\sigma$ is not a poly, there is a $z \in \mathbb{R}$ s.t.
$$\sigma^{(k)}(z) \neq 0 \quad \forall k \in \mathbb{N}.$$

$$\frac{\sigma((\lambda+\delta)x+z) - \sigma(\lambda x + z)}{\delta} \text{ is in } \mathcal{F}_\sigma \text{ if } \delta \neq 0.$$

$$\frac{d}{d\lambda} \sigma(\lambda x + z) \Big|_{\lambda=0} = x \underbrace{\sigma^{(1)}(z)}_{\neq 0}, \quad \text{Thus } f(x) = x \text{ is in closure of } \overline{\mathcal{F}_\sigma}.$$

Higher derivatives yield $x \mapsto x^k \underbrace{\sigma^{(k)}(z)}_{\neq 0} \Rightarrow$ monomials $\in$ closure of $\overline{\mathcal{F}_\sigma}$

Hence all polynomials $\in$ closure of $\overline{\mathcal{F}_\sigma}$. $\overset{\text{Weierstrauss}}{\Longrightarrow}$ $\overline{\mathcal{F}_\sigma}$ dense in $C(\mathbb{R})$
(since polys are.) $\square$

**Lemma:** Let $U \in \mathbb{R}^d$ compact. Then $\mathcal{E} = \text{span}\{f: U \to \mathbb{R} \mid f(x) = \exp[\sum_{i=1}^{d} w_i x_i], w_i \in \mathbb{R}\}$
is dense in $(C(K), \|\cdot\|_\infty)$.

Pf: Stone-Weierstrauss states $\mathcal{E}$ is dense if (i) $\mathcal{E}$ is an algebra. ✓

to get $f=1$, plug in $w_{i,s}=0 \to$ (ii) $\mathcal{E}$ contains a non-zero const. fn

use $w = x-y$
$\Rightarrow f(x-y) = e^{\|x-y\|^2} \neq 1 \Rightarrow f(x) \neq f(y)$. $\to$ (iii) $\forall x,y \in K: x \neq y, \exists f \in \mathcal{E}$ s.t. $f(x) \neq f(y)$.
"$f(x) \neq f(y)$. $\square$

**Thm:** $d, d' \in \mathbb{N}$, $K \subseteq \mathbb{R}^d$ compact, $\sigma \in C(\mathbb{R})$ nonpolynomial. Then the set of fns.
presentable by a NN w/ a single hidden layer that uses $\sigma$ as an act. fn.
is dense in the set of conts $f: U \subseteq \mathbb{R}^d \to \mathbb{R}^{d'}$.

Pf: WLOG. $d' = 1$. $K \subseteq \mathbb{R}^d$ is compact.

$\forall \varepsilon > 0$, $\exists n \in \mathbb{N}$, $v_1, \ldots, v_n \in \mathbb{R}^d$, $s \in \{\pm 1\}^n$ s.t. $g: \mathbb{R}^d \to \mathbb{R}$ $\Big\}$ from Lemma

$$g(x) := \sum_{i=1}^n s_i e^{v_i \cdot x} \text{ satisfies } \|g - f\|_\infty \leq \varepsilon/2.$$

$K_1 := \bigcup_{i=1}^n \{ v_i \cdot x \mid x \in K \} \subseteq \mathbb{R}$ compact. Then $\exists a_j, b_j, w_j$ s.t.

$\dfrac{1 \leq \varepsilon}{2}$

$$\sup_{y \in K_1} \left| e^y - \sum_{j=1}^L a_j \sigma(w_j y - b_j) \right| \leq \varepsilon/2n$$

$$\left\| f - \underbrace{\sum_{j=1}^L \sum_{i=1}^n s_i a_j \sigma(w_j v_i \cdot x - b_j)}_{\text{expressable by 1-layer NN}} \right\|_\infty \leq \left\| f - \sum_{i=1}^n s_i e^{v_i \cdot x} \right\|_\infty + \sum_{i=1}^n \left\| e^y - \sum_{j=1}^L a_j \sigma(w_j y) \right\|$$

$$\leq \varepsilon/2 + n \, \varepsilon/2n = \varepsilon. \qquad \square$$

Kolmogorov's superposition thm: $\forall n \in \mathbb{N}$ $\exists \, \psi_j \in C([0,1])$, $j \in \{0, 1, \ldots 2n\}$

$\exists \, \lambda \in \mathbb{R}_+^n$ s.t. for every $f \in C([0,1]^n, \mathbb{R})$ $\exists \, \phi \in C([0,1])$ s.t.

$$f(x_1, \ldots, x_n) = \sum_{j=0}^{2n} \phi\left( \sum_{k=1}^n \lambda_k \psi_j(x_k) \right)$$

not of practical use

Can interpret as a 2-hidden layer NN w/ exact #nodes + activation fns.
· 1ˢᵗ hidden layer has $n(2n+1)$ neurons w/ $\psi_j$'s as act. fns.
· 2ⁿᵈ hidden layer has $2n+1$ neurons using $\phi$
· Output layer uses $\sigma(z) = z$.

Thm: (VC-dim of NNs): For arbitrary $n_0, w \in \mathbb{N}$, fix an architecture of a layered feed-forward NN w/ $n_0$ inputs, a single output, $w$ parameters (#weights + biases).

$\bar{\mathcal{F}} \subseteq \{-1, 1\}^{\mathbb{R}^{n_0}}$ that can be implemented by such a NN using $\sigma = \text{sgn}$.

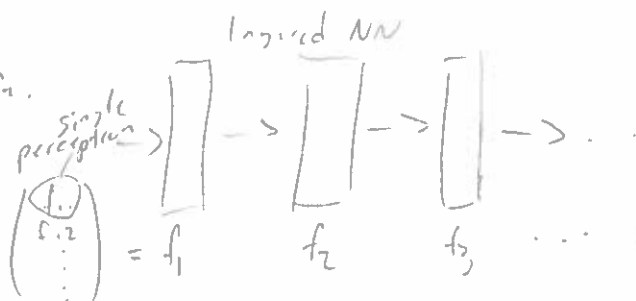Then $\mathrm{VCdim}(\bar{\mathcal{F}}) \leq 2w \log_2(ew)$.

Pf (Sketch): Uses composition property of growth fn.

$$\overline{\Pi_f}(n) \leq \prod_{i=1}^m \Pi_{f_i}(n) \leq \prod_{i=1}^m \prod_{j=1}^{n_i} \underbrace{\Pi_{f_{ij}}(n)}_{\text{perception}}$$

$m = $ #layers  $n_i = $ #nodes in layer $i$

$$\leq \prod_{i=1}^m \prod_{j=1}^{n_i} \left( \frac{en}{w_{ij}} \right)^{w_{ij}} < (en)^w$$

$$\underbrace{\phantom{xx}}_{\frac{1}{w} \leq 1}$$

layered NN

single perception → $\left(\begin{array}{c} 1 \\ 5.2 \\ \vdots \end{array}\right) = f_1$  $\boxed{\phantom{x}} \to \boxed{\phantom{x}} \to \boxed{\phantom{x}} \to \cdots$
$\qquad\qquad\qquad\qquad\qquad f_1 \qquad f_2 \qquad f_3$
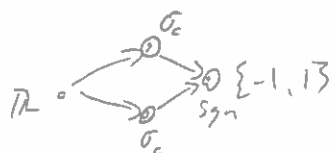
Then use growth fn → VCdim bounds

$\square$

Recall: VC dimension tells us how much data is required in binary classification training for good generalization guarantee.

Last time: w/ activation fn $\int$. $VCdim(\bar{f}) = O(w \log_2 w)$
where $w = $ # of weight parameters in NN.

Recall: $VCdim(\{0,1\}^{\{0,1\}^n}) = 2^n$   $\Rightarrow$ NN has to grow exp. in $n$.

Consider (diagram) $\mathbb{R} \circ \xrightarrow{\sigma_c} \circlearrowright \{-1,1\}$ w/ $\sigma_c$: (graph)   $VCdim(\bar{f}) = \infty$
          $\xrightarrow{\sigma_c}$ sgn

Linear combs + sgn allow NN to amplify crazy properties of $\sigma_c$ that are unobvious

Ex: $\sigma_c(z) = \frac{1}{1+e^{-z}} + c z^3 e^{-z^2} \sin(z)$, $c \geq 0$   ← this class has VCdim = $\infty$.

Consider $\bar{f} = \{f: \mathbb{R}_+ \to \mathbb{R} \mid \exists \alpha \in \mathbb{R}: f(x) = \text{sgn} \sin(\alpha x)\}$. WTS $\bar{f} \subseteq \bar{f}_{NN}$

let $f \in \bar{f}_{NN}$ s.t. $f(x) = \text{sgn}\left[\sigma_c(\alpha x) + \sigma_c(-\alpha x) - 1\right]$
$= \text{sgn}\left[\underbrace{2c(\alpha x)^3 e^{-\alpha^2 x^2}}_{>0} \sin(\alpha x)\right] = \text{sgn}\left[\sin(\alpha x)\right]$.   □

Thus we must require more specific properties on the $\sigma$'s!

Prop: [Warren] Let $\{P_1, \ldots, P_m\}$ be polys in $k$ variables of degree $\leq d$ w/ coeffs $\in \mathbb{R}$.
Let $\gamma(k, d, m) = $ max # of connected components of $\mathbb{R}^k \setminus \bigcup_{i=1}^m P_i^{-1}(\{0\})$
Then   $\gamma(k, d, m) \leq (4 e d m / k)^k$   if $m \geq k$.

Thm: Consider $s$ atomic predicates, each of which is a poly. inequalities of deg. $\leq d$ in $m+k$ vars.
                              └→ truth $\times \mapsto \{0,1\}$
Consider $\Psi: \mathbb{R}^m \times \mathbb{R}^k \to \{0,1\}$ Boolean combs of those predicates. (and/or/xor/etc.)
Define $\bar{f} = \{\Psi(\cdot, w) \mid w \in \mathbb{R}^k\} \subseteq \{0,1\}^{\mathbb{R}^m}$
Then   (i) $T(n) \leq \gamma(k, d, 2ns)$
       (ii) $VCdim(\bar{f}) \leq 2k \log_2(8eds)$

Later, show $VCdim(\bar{f}_{NN})$
is $O(w \log w) \to$ fixed # layers
$O(wL \log w)$ ⎫ p-w linear
$\Omega(wL \log \frac{w}{L})$ ⎬ w/ L layers
↑ lower bd ⎭

Pf: WLOG. All poly inegs. compared to 0 e.g. $p \geq 0$, $p \leq 0$.

Estimate $|\bar{f}|_A|$ where $|A| = n$, $A \subseteq \mathbb{R}^m$.

For $a \in A$, $\mathcal{Y}_a(w) := \mathcal{Y}(a, w)$. $P_a :=$ polys in $\mathcal{Y}_a$..

let $P := \bigcup_{a \in A} P_a$, then $|P| \leq n |P_a| \leq ns$.

Thus
$VCdim(\bar{f}_{NN})$
$= O(w)$

Note $|\bar{f}|_A| \leq |\{Q \in \{-1, 0, 1\}^P \mid Q(p) = \text{sgn}(p(w)), w \in \mathbb{R}^k\}|$
$\hookrightarrow 0 \, m \, 0$

Consider $P' := \{p + \varepsilon \mid p \in P\} \cup \{p - \varepsilon \mid p \in P\}$, $\varepsilon > 0$

→ Then $|P'| \leq 2|P| \leq 2ns$. Now things on bdry $(Q(x) = 0)$ than cells in sit inside cells in $P'$.

Thus $|\bar{f}|_A| \leq \gamma(k, d, 2ns)$ by def. of $\gamma$ from Warren's
$\Rightarrow T(n) \leq ''$ .

Then (ii) follows from Warren's prop. □

Now $VCdim(\bar{s}) = O(tk)$
Adding $(x \mapsto e^x, +, \cdot, /, -) \Rightarrow VCdim(\bar{F}) = O(t^2 k^2)$ by similar proof.
What about $VCdim(\bar{f}_{NN})$? → show $O(wN)$

Thm: Let $N = \#$ neurons, $m = \#$ inputs, $w = \#$ parameters (weights + biases), feed-forward
$\sigma =$ piecewise-poly w/ $p$ pieces & degree $d$.     output uses $\sigma(z) = \mathbb{1}_{z \geq 0}$.

Then $VCdim(\bar{f}_{NN}) \leq 2w[N \log_2 p + \log_2(16e \max\{\delta + 1, 2d\delta\})]$

w/ $\delta =$ depth of network (# of layers).
                                                    $\delta(i) =$ depth of $i^{th}$ neuron

Pf: Regard entire network as $\mathcal{Y}: \mathbb{R}^m \times \mathbb{R}^w \to \{0, 1\}$. Label neurons forwards $i = 1, \dots, N$.
                                                    $a_i =$ output of $i^{th}$ neuron (before $\sigma$)

Define $I: \mathbb{R} \to \{1, \dots, p\}$, $I^{-1}(\{j\})$ interval corr. to $j^{th}$ piece of poly.

$a_1$ is quadratic on $\mathbb{R}^m \times \mathbb{R}^w$. $I(a_1)$ can be det. by $p$ ineqs. quadratic on $\mathbb{R}^m \times \mathbb{R}^w$.
deg(i) of $a_i$ over $\mathbb{R}^m \times \mathbb{R}^w$ is $d^{\delta(i)} + \sum_{j=0}^{\delta(i)} d^j$. Then $I(a_i)$ can be det. by $p$ ineqs. on $\mathbb{R}^m \times \mathbb{R}^w$ w/ deg
$\Rightarrow p^i$ predicates poly of deg $\leq$ deg(i). In total, have $\sum_{j=1}^i p^j$ poly preds. to det. $I(a_i)$.
Last step, predicate is $\mathbb{1}_{a_N \geq 0} \to 1$ ineq. Adding all up, get $\leq 2p^N$ poly predicates of deg
$\Rightarrow \deg = \max\{\delta + 1, 2d\delta\}$. Then use Thm to get VCdim bound.                     □

Recall:



perceptron — VC dim $n$
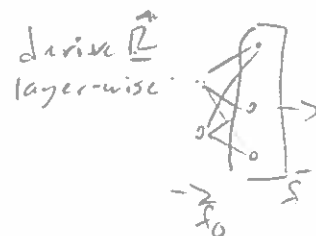
$x_1 \to \bigcirc \to$
$\vdots \quad \sigma\left(\underline{v} + \sum_{i=1}^{n} w_i x_i\right)$
$x_n$

NN — composition of perceptrons
has VC dim $O(w \log w)$

**Thm:** $a, b \in \mathbb{R}$, $\tilde{\sigma} : \mathbb{R} \to \mathbb{R}$, $L$-lipschitz. Assume $\bar{\mathcal{F}}_0 \subseteq \mathbb{R}^{\mathcal{X}}$ includes zero-fn.

Define $\bar{\mathcal{F}} := \left\{ x \mapsto \tilde{\sigma}\left(v + \sum_{j=1}^{m} w_j f_j(x)\right) \mid |v| \le a, \|w\|_1 \le b, f_j \in \bar{\mathcal{F}}_0 \right\} \subseteq \mathbb{R}^{\mathcal{X}}$.

wrt. $x \in \mathcal{X}^n$, $\hat{\mathcal{R}}(\bar{\mathcal{F}}) \le L\left(\frac{a}{\sqrt{n}} + 2b\,\hat{\mathcal{R}}(\bar{\mathcal{F}}_0)\right)$

derive $\hat{\mathcal{R}}$
layer-wise

**Pf:** $\hat{\mathcal{R}}(\bar{\mathcal{F}}) \le \frac{L}{n} \mathbb{E}_\sigma\left[ \sup_{v, w, f_j} \sum_{i=1}^{n} \sigma_i\left(v + \sum_{j=1}^{m} w_j f_j(x_i)\right) \right]$

$\underbrace{\{w_j f_j \in b \, \mathrm{conv}\{\bar{\mathcal{F}}_0 - \bar{\mathcal{F}}_0\}}_{} =: \mathcal{C}_1$

$\{x \mapsto v \mid |v| \le a\} =: \mathcal{C}_2$

$\Rightarrow \hat{\mathcal{R}}(\bar{\mathcal{F}}) \le L\,\hat{\mathcal{R}}(\mathcal{C}_1 + \mathcal{C}_2) = L\left(\hat{\mathcal{R}}(\mathcal{C}_1) + \hat{\mathcal{R}}(\mathcal{C}_2)\right) \le L\left(\frac{a}{\sqrt{n}} + 2b\,\hat{\mathcal{R}}(\bar{\mathcal{F}})\right)$

$\hat{\mathcal{R}}(\mathcal{C}_1) = b\,\hat{\mathcal{R}}(\mathrm{conv}\{\bar{\mathcal{F}}_0 - \bar{\mathcal{F}}_0\}) = b\,\hat{\mathcal{R}}(\bar{\mathcal{F}}_0 - \bar{\mathcal{F}}_0) = b\left(\hat{\mathcal{R}}(\bar{\mathcal{F}}_0) + \hat{\mathcal{R}}(-\bar{\mathcal{F}}_0)\right) = 2b\,\hat{\mathcal{R}}(\bar{\mathcal{F}}_0)$.

$\hat{\mathcal{R}}(\mathcal{C}_2) \le \frac{a}{n} \mathbb{E}_\sigma[|z|] \le \frac{a}{n} \mathbb{E}_\sigma[|z|^2]^{1/2} = \frac{a}{n}\sqrt{n} = \frac{a}{\sqrt{n}}$.

w/ $z := \sum_{i=1}^{n} \sigma_i$

$\square$

**Rmk:** Note that $\|w\|_1 \le b$ plays a large role in bounding $\hat{\mathcal{R}}(\bar{\mathcal{F}})$!
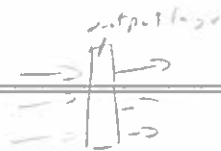keeping weight values low in each layer $\Rightarrow$ generalization.

# Training NNs.

Consider Classification: common choices:

- # of neurons in output layer = # of classes (say $m \in \mathbb{N}$)
- softmax activation $\sigma$ at output

$$\sigma : \mathbb{R}^n \to \mathbb{R}^m \qquad \sigma(z)_j := e^{z_j} / \sum_{k=1}^n e^{z_k}$$

- loss fn: cross-entropy (= neg. log. likelihood)

$$L(h(x), y) = - \sum_{c=1}^m \mathbb{1}_{y=c} \ln[h(x)_c]$$

$\{1, \dots, m\}$

Empirical risk: $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i) +$ "regularizer"

suppose $h$ depends on $w_1, \dots, w_N \in \mathbb{R}$, we minimize $\hat{R}(h)$ via gradient descent

$$\mathbb{R}^N \ni w^{(t+1)} = w^{(t)} - \underset{\hat{o}}{\alpha_t} \nabla_w \hat{R}(h)$$

Naively: # NN evaluations for a single step $\sim nN \to$ bad! can get $nN \gtrsim 2 \cdot 10^6$.

- Reduce $n$ via Stochastic optimization
- Reduce $N$ via backpropagation

constant!

## Backpropagation: (Automatic differentiation)

Bookkeeping: Layers $\ell \in \{0, \ldots, m\}$. $N_\ell = \#$ neurons in $\ell^{\text{th}}$ layer.

$\mathbb{R} \ni w_{jk}^\ell :=$ weight from $k^{\text{th}}$ neuron in layer $\ell-1$ to the $j^{\text{th}}$ neuron in layer $\ell$.



$b_j^\ell =$ threshold value for $j^{\text{th}}$ neuron, layer $\ell$.

$x_j^\ell =$ output of $j^{\text{th}}$ neuron, layer $\ell$.

$x_\ell = \sigma(\underbrace{w^\ell x^{\ell-1} + b^\ell}_{z^\ell})$.

$f: \ldots \to \mathbb{R}$   would like $f: \mathbb{R}^N \to \mathbb{R} \to \nabla_w f$

Define: $\delta_j^\ell := \dfrac{\partial f}{\partial z_j^\ell} = \sum_k \dfrac{\partial f}{\partial z_k^{\ell+1}} \cdot \dfrac{\partial z_k^{\ell+1}}{\partial z_j^\ell} = \sum_k \delta_k^{\ell+1} \cdot w_{kj}^{\ell+1} \cdot \sigma'(z_j^\ell)$

Compute recursively (backwards) from $\ell = m$.        forward to set $x_j^\ell$, backward to set $\delta_j^\ell$.

Thus we only need to run the NN twice for a single computation of $\nabla_w f$

Note: $\dfrac{\partial f}{\partial b_j^\ell} = \sum_k \dfrac{\partial f}{\partial z_k^\ell} \dfrac{\partial z_k^\ell}{\partial b_j^\ell} = \sum_k \delta_k^\ell \cdot \underbrace{\delta_{kj}}_{\text{Kronecker } \delta} = \delta_j^\ell$.        comp. of all gradient pieces w/ only 2 runs of NN.

$\dfrac{\partial f}{\partial w_{jk}^\ell} = \sum_i \underbrace{\dfrac{\partial f}{\partial z_i^\ell}}_{\delta_i^\ell} \cdot \underbrace{\dfrac{\partial z_i^\ell}{\partial w_{jk}^\ell}}_{\text{gives } \delta_{ij}} = \delta_j^\ell x_k^{\ell-1}$.

This technique is useful for network/circuit fns. from $\mathbb{R}^N \to \mathbb{R}^M$

Can compute $\nabla f$ w/ same complexity/#comps. as $f$ itself        can write a computational graph w/ $O(N)$ nodes

Thus   $w \mapsto \langle \nabla f(w), v \rangle$   is linear in $N$.

$\bullet$   $\nabla^2 f(w) v$ is computable in same time as $f$.

Work for one training step $Nu \to$ use backprop to reduce to $\underline{n}$

Recall: Computing $D_w \hat{R}(h_w) = \frac{1}{n} \sum_{i=1}^{n} D_w L(h_w(x_i), y_i) + \text{"""}$

---

We can use Stochastic gradient descent to reduce $n$ to $1$.

Expectation value of random single gradient eval is full gradient.

Pick $X_i$ at random & compute $D_w L(h_w(x_i, y_i)) + \dots$

& step in that (neg.) direction, on average following true $D_w \hat{R}(h_w$

Why use the gradient at all?

Consider diff. $f:$ $f(x+v) \approx f(x) + \langle Df(x), v \rangle$ &larr; Minimize?

Take argind $\{ \langle Df(x), v \rangle : \|v\| = 1 \} = \begin{cases} -\nabla f(x) & \text{if } \|\cdot\|_2 & \text{as step direction} \\ \dfrac{-\rho^{-1} Df(x)}{\|\rho^{-1} Df(x)\|} & \text{if } \|\cdot\| = \langle v, \rho v \rangle^{1/2} & \rho > 0 \\ & & \rho \text{(matrix)} \end{cases}$

(could use new "direction" $\prod_{i=1}^{\delta \leftarrow \text{depth} \#} \|w_i\|$ &larr; product of all weights in layer $i$

for better generalization

Gradient Descent: Randomly choose data point to eval $Df$, $x_t \in \mathbb{R}^d$

$$x_{t+1} = x_t - \alpha Df(x_t), \quad \alpha > 0$$

Lemma: $f \in C^1(\mathbb{R}^d)$ w/ $Df$ $L$-Lipschitz. Then $\forall x, y \in \mathbb{R}^d$,

$$|f(x) - f(y) - \langle Df(x), y-x \rangle| \leq \frac{L}{2} \|y - x\|^2$$

$f$ &larr; every point $f(x)$ is bounded by a lower + upper quadratic.

Pf: $f(x) - f(y) = \int_0^1 \langle Df(x+t(y-x)), y-x \rangle dt$

$|f(x) - f(y) - \langle Df(x), y-x \rangle| \leq \int_0^1 |\langle Df(x+t(y-x)) - Df(x), y-x \rangle| dt$

$\underset{\text{Schwarz}}{\overset{\text{cauchy}}{\leq}} \int_0^1 \|Df(x+t(y-x)) - Df(x)\|_2 \|y-x\|_2 \, dt$

$\leq \int_0^1 tL \|y-x\|^2 dt = \frac{L}{2} \|y-x\|^2$ $\quad \square$

Use this lemma for GD:

Set $x = x_t$, $y = x_{t+1}$: $\boxed{f(x_t) - f(x_{t+1}) \geq \alpha (1 - \frac{\alpha L}{2}) \|\nabla f(x_t)\|^2}$ (*)

w/ $x_{t+1} = x_t - \alpha \nabla f(x_t)$     $> 0$ if $0 < \alpha < 2/L$

& max if $\alpha = 1/L$.

Thm: $f \in C^1$, $\nabla f - L$-Lipschitz, $\alpha \in (0, 2/L)$.

(i) $f(x_{t+1}) \leq f(x_t)$ unless $\nabla f(x_t) = 0$.

(ii) If $f$ is bounded from below, then $\lim_{t \to \infty} \|\nabla f(x_t)\| = 0$.

& $\min_{t \in T} \|\nabla f(x_t)\|^2 \leq \frac{\alpha(2 - \alpha L)}{2T} \underset{\underset{\alpha = 1/L}{\uparrow}}{=} \frac{1}{2LT}$.

Pf: We just showed (i). WTS (ii).

Consider $\sum_{t=0}^{T-1} (*) \quad <= \overset{telescoping}{=}> \quad f(x_0) - f(x_T) \geq \alpha (1 - \frac{\alpha L}{2}) \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2$

Can lower bound by min. of sum elements · #steps:

$$f(x_0) - f(x_T) \geq \alpha (1 - \frac{\alpha L}{2}) \cdot T \cdot \min_{t \in T} \|\nabla f(x_t)\|^2$$

Then use simple algebra, trivial.     □

Lemma: $f \in C^1(\mathbb{R}^n)$ & $\mu$-strongly convex, $\nabla f - L$-Lipschitz, then

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|_2^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

pf: $\phi(x) := f(x) - \frac{\mu}{2} \|x\|^2$ is convex (literally the def. of $f$ $\mu$-strongly convex).

& $\nabla \phi(x) = \nabla f(x) - \mu x$ is $(L - \mu)$-Lipschitz.  . . .     □

Thm: $f \in C^1(\mathbb{R}^n)$, $\mu$-strongly convex, $\nabla f - L$-Lipschitz, $\alpha \in (0, \frac{2}{L + \mu})$.

$\|x_t - x^*\|^2 \leq (1 - \frac{2\alpha \mu L}{\mu + L})^t \|x_0 - x^*\|_2^2$

$\leq \frac{L}{2} (\frac{\kappa - 1}{\kappa + 1})^{2t} \|x_0 - x^*\|_2^2$

(linear convergence)
w/ quadratic rate?

$\underset{\alpha = \frac{1}{\mu + L}}{\uparrow}$, where $\kappa = \frac{L}{\mu}$ "condition number"

Pf: $\|x_{t+1} - x^*\|_2^2 = \|x_t - \alpha \nabla f(x_t) - x^*\|_2^2$

$= \|x_t - x^*\|_2^2 + \alpha^2 \|\nabla f(x_t)\|^2 - 2\alpha \langle \nabla f(x_t), x_t - x^* \rangle$

lemma
$\& \nabla f(x^*) = 0$ $\longrightarrow$ $\leq \left(1 - \frac{2\alpha \mu L}{\mu + L}\right) \|x_t - x^*\|_2^2 + \alpha\left(\alpha - \frac{2}{\mu + L}\right) \|\nabla f(x_t)\|_2^2$

$\leq 0$ if $\alpha \in \left(0, \frac{2}{L + \mu}\right)$

$\leq$ $\quad + 0$

Then use telescoping to get $\|x_t - x\|^2 \leq \left(1 - \frac{2\alpha \mu L}{\mu + L}\right)^t \|x_0 - x^*\|^2$.  □

Thus, if we want to be $\varepsilon$-close to optimum, take $\mathcal{O}(\log 1/\varepsilon)$ steps.
Only have this order work if we have (i) strong convexity of $f$
 (ii) $L$-smoothness of $f$ ($\nabla f$-Lipschitz)
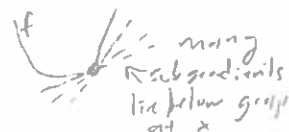 (iii) exact gradient computable.

We will show that relaxing all 3 assumptions brings us to $\mathcal{O}(1/\varepsilon^2)$ steps.

Thm: Assume $f : \mathbb{R}^d \to \mathbb{R}$ convex. $\forall x \in \mathbb{R}^d$, let $g_1(x), \ldots, g_n(x)$ i.id rvs.
w/ values in $\mathbb{R}^d$ s.t. $\mathbb{E}[g_i(x)] = \nabla f(x)$
 a subgradient $\qquad$ (many subgradients lie below graph at $x$)

Assume $\mathbb{E}[\|g_i(x)\|^2] \leq G^2$.

Then consider $x_{t+1} := x_t - \alpha g_{t+1}(x_t)$, $t = 0, \ldots, T-1$, $\bar{x} := \frac{1}{T} \sum_{t=0}^{T-1} x_t$.

Then $\mathbb{E}[f(\bar{x})] - \underbrace{f(x^*)}_{\text{min.}} \leq \frac{2\|x - x^*\| G}{\sqrt{T}}$.

# Lec 15 — 5/23/17

**Polyak '63:** $\frac{1}{2}\|\nabla f(x)\|^2 \geq \mu(f(x) - f(\underset{\uparrow}{x^*})) \quad \forall x \in \mathbb{R}, \ \mu > 0$

$\underset{\text{global min.}}{}$

**Lemma:** If $f \in C'(\mathbb{R}^d)$ is $\mu$-strongly convex, the above holds.

**Pf:** $\phi(x) := f(x) - \frac{\mu}{2}\|x\|^2$ is convex by def. of $f$ $\mu$-strongly convex

$\Rightarrow \phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle \quad \forall x, y$

$\Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2}\|y - x\|^2$

minimize w.r.t. $y$ $\quad f(x^*) \geq f(x) - \frac{1}{2\mu}\|\nabla f(x)\|^2$ $\qquad \square$

**Rmk:** $x \longmapsto |x|$ is convex but doesn't satisfy Polyak's ineq.

$x \longmapsto x^2 + 3(\sin x)^2$ satisfies Polyak w/ $\mu = \frac{1}{32}$ but is not convex →

$\hookrightarrow$ Polyak's ineq. independent of convexity.

## Gradient Descent

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

$$SGD := g(x_t) \text{ where } \mathbb{E}[g(x)] = \nabla f(x)$$

SGD: Fixed step size will not let us stay at min, but rather in an $\alpha$-nbhd of it.
$\hookrightarrow$ but gets close exp-fast!

**Thm:** $f \in C'(\mathbb{R}^d)$, $\mu$-Polyak, L-Lipschitz gradient.
$\forall x \in \mathbb{R}^d$, $g_1, \ldots, g_T$ i.i.d. r.v.s w/ values in $\mathbb{R}^d$ s.t. $\mathbb{E}[g_t(x)] = \nabla f(x)$ & $\mathbb{E}[\|g_t(x)\|^2] \leq \gamma^2$.

Consider $x_{t+1} = x_t - \alpha g_t(x_t)$, $\alpha \in [0, \frac{1}{2\mu}]$. Then

$$\mathbb{E}[f(x_T)] - f(x^*) \leq (1 - 2\mu)^T (f(x_0) - f(x^*)) + \frac{L\alpha\gamma^2}{4\mu}.$$

**Pf:** $f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2}\|y - x\|^2$ from prev. lecture. Use $y = x_{t+1} = x_t - \alpha g_t(x_t)$
$\qquad\qquad x = x_t$

$\hookrightarrow f(x_{t+1}) \leq f(x_t) - \alpha \langle \nabla f(x_t), g_t(x_t) \rangle + \frac{L}{2}\alpha^2 \|g_t(x_t)\|^2$

$\mathbb{E}[f(x_{t+1}) \mid g_1, \ldots, g_{t-1}] \leq f(x_t) - \alpha\|\nabla f(x_t)\|^2 + \frac{L\alpha^2}{2}\mathbb{E}[\|g_t(x_t)\|^2 \mid g_1, \ldots, g_{t-1}]$

$\leq f(x_t) - 2\mu\alpha(f(x_t) - f(x^*)) + \frac{L\alpha^2}{2}\gamma^2.$

$[f(x_{t+1})] - f(x^*) \leq [\mathbb{E}[f(x_t)] - f(x^*)](1 - 2\mu\alpha) + L\alpha^2\gamma^2 \quad \rightarrow$ iterate inequality recursively.

Thm: $C \subseteq \mathbb{R}^d$ compact, convex w/ diam $\delta$ ( $x, y \in C \Rightarrow \|x-y\| \leq \delta$)

$f : C \to \mathbb{R}$ convex w/ global min. at $x^* \in C$.

$\forall x \in \mathbb{R}^d$, $g_1, \ldots, g_T$ are i.i.d. r.v.s w/ $\mathbb{E}[g_t(x)] = \nabla f(x)$

& $\mathbb{E}[\|g_t\|^2] \leq \gamma^2$. Consider $x_t = P_C(x_{t-1} - \alpha_t g_t(x_{t-1}))$

With $\alpha_t \leq \alpha_{t-1}$, $\bar{x} := \frac{1}{T} \sum_{i=0}^{T-1} x_t$, then

$$\mathbb{E}[f(\bar{x})] \leq f(x^*) + \frac{1}{2T}\left(\gamma^2 \|\alpha\|_1 + \frac{\delta^2}{\alpha_T}\right)$$

$$\leq f(x^*) + \sqrt{\tfrac{2}{T}} \gamma \delta \quad \text{for } \alpha_t = \frac{\delta}{\sqrt{2t}\,\gamma}$$

Pf:
$$\mathbb{E}\left[\|x_t - x^*\|^2 \mid g_1, \ldots, g_{t-1}\right]$$

$$\leq \mathbb{E}\left[\|x_{t-1} - \alpha_t g_t(x_{t-1}) - x^*\|^2 \mid \cdots\right]$$

$$= \mathbb{E}\left[\|x_{t-1} - x^*\|^2 - 2\alpha_t \langle g_t(x_{t-1}), x_{t-1} - x^*\rangle + \alpha_t^2 \|g_t(x_{t-1})\|^2\right]$$

$$\leq \|x_{t-1} - x^*\|^2 - 2\alpha_t \langle \nabla f(x_{t-1}), x_{t-1} - x^*\rangle + \alpha_t^2 \gamma^2$$

$$\overset{\text{convexity}}{\leq} \|x_{t-1} - x^*\|^2 - 2\alpha_t \left(f(x_{t-1}) - f(x^*)\right) + \alpha_t^2 \gamma^2$$

$$\to \mathbb{E}[f(x_{t-1})] - f(x^*) \leq \frac{\gamma^2}{2}\alpha_t + \frac{1}{2\alpha_t}\mathbb{E}\left[\|x_{t-1} - x^*\|^2 - \|x_t - x^*\|^2\right]$$

$$\frac{1}{T}\sum_{t=1}^{T} \cdots \quad \leq \quad \mathscr{E} \quad ''$$

$$\left(\text{side calc.}\left[\sum_{t=1}^{T}\frac{1}{2\alpha_t}\left(\|x_{t-1} - x^*\|^2 - \|x_t - x^*\|^2\right) = \frac{1}{2\alpha_1}\|x_0 - x^*\|^2 - \frac{1}{2\alpha_T}\|x_T - x^*\|^2 + \sum_{t=2}^{T}\left(\frac{1}{2\alpha_t} - \frac{1}{2\alpha_{t-1}}\right)\|x_{t-1} - x^*\|\right.\right.$$

$$\underbrace{}_{\leq \delta^2} \qquad \underbrace{}_{\leq 0} \qquad \underbrace{}_{\geq 0} \quad \underbrace{}_{\leq \delta^2}$$

$$\leq \frac{\delta^2}{2\alpha_1} + \frac{\delta^2}{2}\sum_{t=2}^{T}\underbrace{\frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}}}_{\frac{1}{\alpha_T} - \frac{1}{\alpha_1}} \leq \frac{\delta^2}{2\alpha_T}.$$

$$\to \mathbb{E}[f(\bar{x})] \leq f(x^*) + \frac{1}{2T}\left(\gamma^2 \|\alpha\|_1 + \frac{\delta^2}{\alpha_T}\right)$$

$$\alpha_t = \frac{\delta}{\sqrt{2t}\,\gamma} \implies \leq f(x^*) + \sqrt{\tfrac{2}{T}}\gamma\delta. \qquad \blacksquare$$

Stochastic subgradient + convexity $\to$ $O\left(\frac{1}{\varepsilon^2}\right)$ steps (conv. up to $\varepsilon$)

$\qquad$ + Polyak $\qquad \to O\left(\frac{1}{\varepsilon}\right)$ steps

$\qquad$ grad. desc. + Polyak $\qquad \to O(\log 1/\varepsilon)$ steps

$\qquad$ Newton + strong convexity $\to O(\log\log 1/\varepsilon)$ steps

$\qquad$ (→2nd order Hessian inversion method.

# Lecture 16 — 5/30/17

○ Ideas to improve (S)GD:  — momentum (introduces short term memory)
  — exploit $2^{nd}$ order info "$D^2 f \cdot v$"
  — hypergradient descent
  — . . .

## Deep Neural Networks

= NN w/ more than a single hidden layer

• until ~2006, almost all NNs in practice were "shallow"
  — more computing power & GPUs now
  — more data    — more tricks → different initialization
      ↳ ·· activation fn. · sigmoid e tanh · ReLU
      ↳ training changed

○ Problem w/ "vanishing gradient in backpropagation" ✳



$\delta^{L-1} = W \delta^L$ , $\delta^{L-2} = W W' \delta^L$ , . . .

eigenvals of $W$'s kill $\delta$ in deep nets!
Taking partial deriv. of weight far away from output
gives a nearly 0 value → weight updates slowly
✰ Initialization of these far-away weights is crucial.

✗ Neural networks unpopular in academia since even if it works,
nobody understands the representation or what its doing.

Why and when should one use a deep net?
  ↳ Representational efficiency

Results that indicate that deep NNs can be more efficient (regarding representation/approximation) than shallow nets.

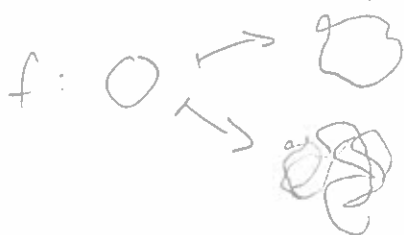- 2013 - number of affine regions is larger when using ReLUs.
- 2014 - Bianchini & Scarselli : $f : \mathbb{R}^d \to \mathbb{R}$

$$A := \{ x \in \mathbb{R}^d \mid f(x) \geq 0 \} \leq \mathcal{O}((\# \text{neurons})^d) \quad \text{for shallow NNs}$$

Consider $\sum_{i=0}^{d-1} b_i(A) \geq$ can be $\Omega(2^\delta)$ for depth $\delta$

$\underbrace{b_i}_{\text{Betti number}}$

(is number of topological "holes")

- 2016 - Poole et al. $f : \mathbb{R}^d \to \mathbb{R}^{d'}$, where $d, d' \geq 2$

shallow:

deep: $f : O \rightrightarrows$

consider : $\dfrac{\text{length of image}}{\text{diam}}$

$\mathbb{E}$ of ↗ w.r.t. randomly chosen network weights (untraining)

• Hence decision boundaries can be more complex in deep!

- 2015 - Telgarsky: $\bar{f}(m, L) \subseteq \mathbb{R}^{\mathbb{R}}$, fns rep. by NN using ReLU.
  
  (with superscript "layer", subscript "neurons/layer")

For every $f \in \bar{f}(m, L)$, define $\hat{f}(x) := \mathbb{1}_{f(x) \geq 1/2}$

Emp. risk of training set $S$: $\hat{R}(f) = \dfrac{1}{|S|} \sum_{(x,y) \in S} \mathbb{1}_{\hat{f}(x) \neq y}$

For $k \in \mathbb{N}$, $n = 2^k$, choose $S := (x_i, y_i)_{i=0}^{n-1}$, with $x_i = \dfrac{i}{n}$, $y_i = i \bmod 2$
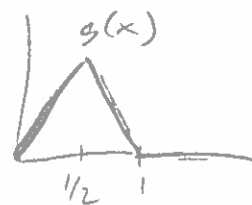
Thm - 1. There is a $h \in \bar{f}(2, 2k)$ for which $\hat{R}(h) = 0$.

2. If $m, L \in \mathbb{N}$ and $m \leq 2^{(\frac{k-2}{L} - 1)}$, then $\forall f \in \bar{f}(m, L)$, $\hat{R}(f) \geq 1/6$.

Rmk: To represent $h \in \bar{f}(2, 2k)$ w/ $\sim \sqrt{k}$ layers still requires $\sim 2^{\sqrt{k}}$ neurons/layer

Telgarsky's thm, proof:

○ Let $g : \mathbb{R} \to \mathbb{R}$, $g(x) := \begin{cases} 2x, & 0 \le x \le 1/2 \\ 2(1-x), & 1/2 \le x \le 1 \\ 0, & o/w. \end{cases}$
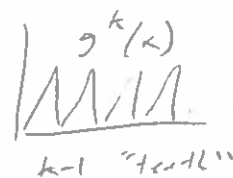


ReLU:
$\sigma(x) := \max\{0,x\}$ $= \sigma(2\sigma(x) - 4\sigma(x - 1/4))$

Thus $g \in \bar{f}(2,2)$.

Then $h(x) := g^k(x) = g \circ \cdots \circ g(x) \in \bar{f}(2, 2k)$



This easily represents the training data set $S$.

1. $\hat{R}(h) = 0$.

2. Every $f \in \bar{f}(m, l)$ is piecewise linear w/ at most $t = (2m)^l$ pieces

If $f_1$ has $t_1$ pieces $\longrightarrow$ $f_1 + f_2$ has $\le (t_1 + t_2)$ pieces
& $f_2$ has $t_2$ pieces $\longrightarrow$ $f_1 \circ f_2$ has $\le t_1 \cdot t_2$ pieces

○ $\rightarrow$ graph of $f$ can cross $1/2$ at most $2t - 1$ times



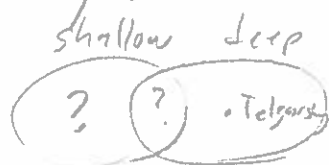$\rightarrow \tilde{f}$ is piecewise constant w/ at most $2t$ intervals

$\rightarrow h - 2t$ points are in intervals containing $\ge 1$ point
$\rightarrow$ out of these, at least $1/3$ misclassified

This example shows deep nets can represent fns.
much more simply than shallow!

& Doesn't show that there are no fns. more simply
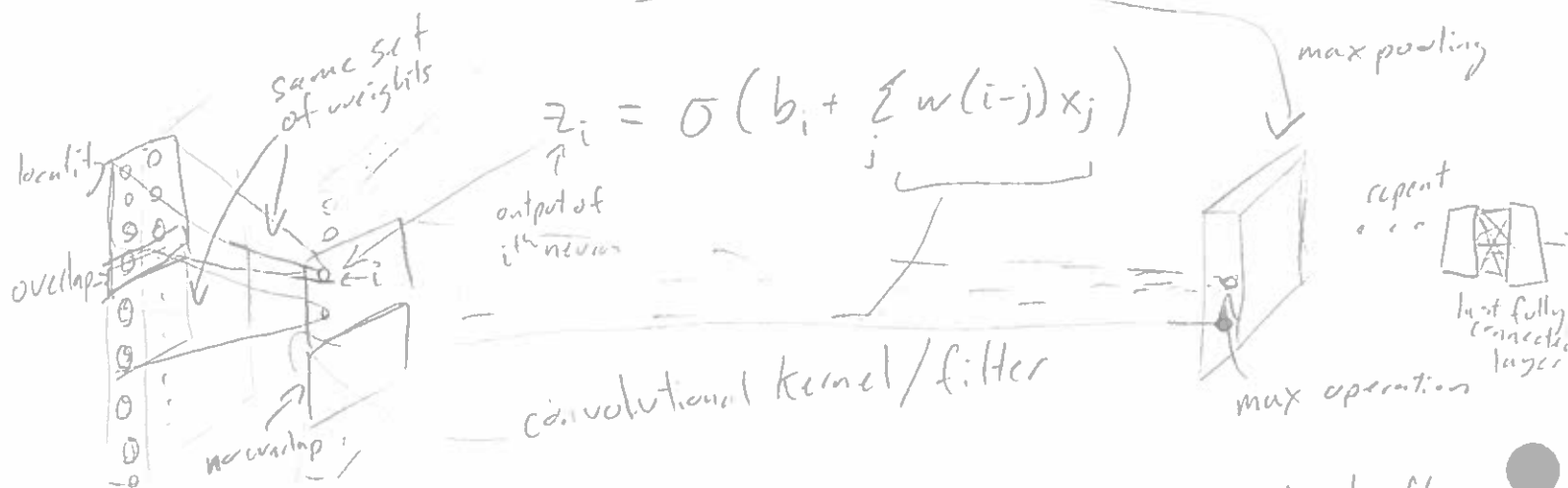○ represented by shallow than deep

shallow deep

• Telgarsky

# Lecture 17 - 6/1/17

## Convolutional Nets

, 2012 ImageNet breakthrough using conv. nets

- repeating ingredients
  - locality - small support of conv. kernel
  - weight sharing - same conv. kernel used throught layer
    - ↳ translation - invariant
  - pooling & down-sampling

$$z_i = \sigma\left(b_i + \sum_j w(i-j)x_j\right)$$

same set of weights

locality

overlap

no overlap

↑ output of $i^{th}$ neuron

$z_i$

convolutional kernel/filter

max pooling

repeat • • •

last fully connected layer

max operation

⊛ Convolution reflects structure in data - 2D like image, local influences
  ↳ Translational invariance of network ⟺ tr. inv. of images - dog is a dog when on left or right

⊛ Conv. is run in parallel!

☆ Depth is usually 10~20 layers

## Optimization — How does all this work?

NP-hardness of ERM: graph $G = (V, E)^{\text{edges}}$, $V = \{1, \ldots, d\}$
$\uparrow$ vertices
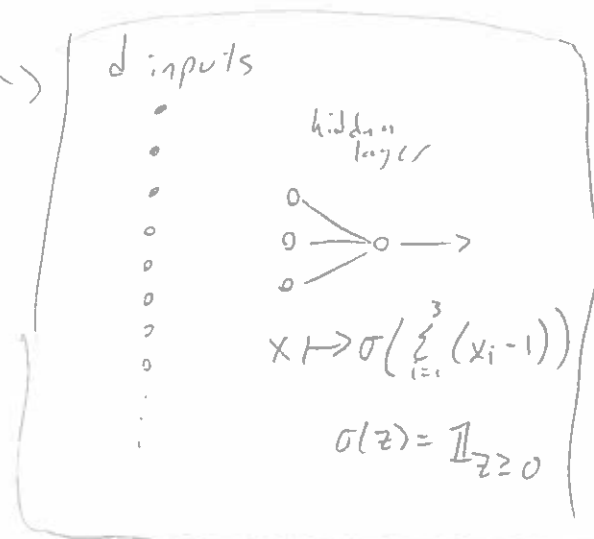
assign $S_G := \{(e_i, 0), (e_i + e_j, 1), (0, 1)\}_{i \in V, (i,j) \in E}$

where $e_i \in \{0, 1\}^d$ s.t. $(e_i)_k = \delta_{ik}$, $\hookrightarrow$ labels vertices 0, labels edges 1

Recall: $G$ is 3-colorable iff $\exists \chi: V \to \{1, 2, 3\}$ s.t. $(i,j) \in E \implies \chi(i) \neq \chi(j)$

Prop: For any $G$ w/ $d$ vertices, $\exists h \in \bar{\mathcal{F}}_d$ $\longrightarrow$
that correctly classifies $S_G$ iff $G$ is 3-colorable.

(embedding an NP-hard problem in a simple NN)

$d$ inputs
⋮

hidden layer



$x \mapsto \sigma\left(\sum_{i=1}^{3} (x_i - 1)\right)$

$\sigma(z) = \mathbb{1}_{z \geq 0}$

Pf: Assume $G$ 3-colorable.

let $w_{\ell, i} = \begin{cases} -1 & \text{if } \chi(i) = \ell \\ 1 & \text{o/w} \end{cases}$
$\ell \in \{1,2,3\}$ ↑ input

$h(x) = 1$ iff $\forall \ell \in \{1, 2, 3\}: \sum_k w_{\ell, k} x_k \geq -1/2$

This classifies $S_G$ correctly since

· $h(0) = 1$

· $h(e_i) = 0$ since if $\chi(i) = \ell$ then $w_{\ell, i} = -1$ so $\sum_k w_{\ell, k} (e_i)_k = -1 < -1/2$.

· $h(e_i + e_j) = 1$ since $\sum_k w_{\ell, k} (e_i + e_j)_k = w_{\ell, i} + w_{\ell, j} \geq 0$ since $G$ is 3-colorable
$(\implies \chi(i) = \ell$ or $\chi(j) = \ell$.

(Converse?)

Converse pf: Assume he $\bar{f}_d$ that correctly classifies $S_G$.

$$h^{-1}(\{1\}) = H_1 \cap H_2 \cap H_3 := H \ni 0 \left( \begin{array}{l} \text{since output is only 1 iff all 3} \\ \text{hidden nodes output } \underline{1} \end{array} \right)$$

$\forall (i,j) \in E, \quad e_i + e_j \in H$

---

$0, e_i + e_j \in H \implies \dfrac{e_i + e_j}{2} \in H$ by convexity.

$\chi(i) := \min\{d \mid e_i \notin H_d\}$

If $(i,j) \in E$, WTS $\chi(i) \neq \chi(j)$. Assume $(i,j) \in E$ but $\chi(i) = \chi(j) = d$.

Then $e_i, e_j \notin H_d$ so by $\underset{\text{convexity}}{}$ $\dfrac{e_i + e_j}{2} \notin H_d$ ↯

Thus $(i,j) \in E \implies \chi(i) \neq \chi(j) \implies G$ is 3-colorable.  □

Hence even a simple $NN$ can embed an NP-hard problem. But that's a discrete combinatorial problem, what about smooth targets?

NP hardness of classifying stationary points

Consider $Q \in \mathbb{Z}^{d \times d}$, $f : \mathbb{R}^d \to \mathbb{R}$. $f(x) = \sum\limits_{i,j=1}^{d} Q_{ij} x_i^2 x_j^2$.

At $x = 0$, $\nabla f = 0$ & $\nabla^2 f = 0 \implies$ stationary pt. but no info on if max, min, saddle

If at $x = 0$ $f$ has local min, then it's a global min.

Suppose $\exists x$ st. $f(x) < 0$, then $\mathbb{R} \ni \lambda \mapsto f(\lambda x) = \lambda^4 \underset{<0}{\underline{f(x)}} \implies$ close to $x$, fn. descents in some directions it cannot be a local min. directions

Def: $Q$ is "copositive" iff $\langle z, Qz \rangle \geq 0 \ \forall z \in \mathbb{R}_+^d$

The question "Is $Q$ $\underline{not}$ copositive?" is NP-complete !!