# Learning Theory

Narayana Santhanam

EE 645

Apr 10, 2024

# This section

PAC Learning

VC dimension and Sauer's lemma

Learnability of
     finite classes
     bounded VC dimension

# PAC learning

Example/Instance space $\mathcal{X}$, label set $\mathcal{Y}$

Hypothesis class $\mathcal{H}$ (set of functions from $\mathcal{X} \to \mathcal{Y}$)

Distribution $D$ over $\mathcal{X}$

Training sample $S$ generated by distribution $D$

Prediction rule $h : \mathcal{X} \to \mathcal{Y}$ that is somehow good

# Loss of a prediction rule

Loss (wrt correct labeling $f$):
$$L_{D,f}(h) = P_{X \sim D}[h(x) \neq f(x)]$$
We cannot observe this in general

Empirical loss on a sample of size $n$,
$$\hat{L}(h) = \frac{1}{n} \sum 1(h(x_i) \neq f(x_i))$$
This we observe in a supervised setting

What can we infer about $L(h)$ from $\hat{L}(h)$?

# IID assumption

Generally expect every example of our training sample to be generated independently

In this case we can expect $\hat{L}(h)$ to concentrate around $L(h)$
    Empirical average $\approx$ real expectation

But by how much? What is the deviation?

# Hoeffding's Inequality

Let $X_1, \ldots, X_n$ be *i.i.d.* variables, the variables bounded in range $X_i \in [a, b]$, and let $\mu = \mathbb{E}X_i$. Then for any $\epsilon > 0$,

$$P\left( \left| \frac{1}{n} \sum_i X_i - \mu \right| > \epsilon \right) \leq 2 \exp\left( -\frac{2n\epsilon^2}{(b-a)^2} \right)$$

# Hoeffding's Inequality

Let $X_1, \ldots, X_n$ be *i.i.d.* variables, the variables bounded in range $X_i \in [a, b]$, and let $\mu = \mathbb{E}X_i$. Then for any $\epsilon > 0$,

$$P\left( \left| \frac{1}{n} \sum_i X_i - \mu \right| > \epsilon \right) \leq 2 \exp\left( -\frac{2n\epsilon^2}{(b-a)^2} \right)$$

If we are working with binary classification (with 0-1 loss), then for each $h$,

$$P\left( \left| \hat{L}(h) - L(h) \right| > \epsilon \right) \leq 2 \exp\left( -2n\epsilon^2 \right)$$

# Hoeffding's Inequality

Let $X_1, \ldots, X_n$ be *i.i.d.* variables, the variables bounded in range $X_i \in [a, b]$, and let $\mu = \mathbb{E}X_i$. Then for any $\epsilon > 0$,

$$P\left(\left|\frac{1}{n}\sum_i X_i - \mu\right| > \epsilon\right) \leq 2\exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

If we are working with binary classification (with 0-1 loss), then for each $h$,

$$P\left(\left|\hat{L}(h) - L(h)\right| > \epsilon\right) \leq 2\exp\left(-2n\epsilon^2\right)$$

In a bad set of training samples $B(h)$, $\hat{L}(h)$ deviates significantly from $L(h)$, but the set of misleading training samples have small probability if $n$ is large enough

# Union bound

If we have finite number of hypothesis, we can argue that collectively, all the bad sets of all $h \in \mathcal{H}$ don't matter: Union bound

$$P\left(\sup_{h \in \mathcal{H}} \left|\hat{L}(h) - L(h)\right| > \epsilon\right) \leq 2|\mathcal{H}| \exp\left(-2n\epsilon^2\right)$$

Here $|\cdot|$ denotes the size of a set

# Union bound

If we have finite number of hypothesis, we can argue that collectively, all the bad sets of all $h \in \mathcal{H}$ don't matter: Union bound

$$P \left( \sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| > \epsilon \right) \leq 2|\mathcal{H}| \exp \left( -2n\epsilon^2 \right)$$

Here $|\cdot|$ denotes the size of a set

This is not artificial—in fact, given we only use finite precision and a finite number of network weights, most deep networks also form finite classes in practice.

# Union bound

If we have finite number of hypothesis, we can argue that collectively, all the bad sets of all $h \in \mathcal{H}$ don't matter: Union bound

$$P\left(\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| > \epsilon\right) \leq 2|\mathcal{H}| \exp\left(-2n\epsilon^2\right)$$

Here $|\cdot|$ denotes the size of a set

This is not artificial—in fact, given we only use finite precision and a finite number of network weights, most deep networks also form finite classes in practice.

Catch is, we don't have to wait till we are guaranteed convergence like above: usually our estimators work good well before we need to sample to reduce the right side to within a given confidence

# Vapnik Chervonenkis dimension

Again, binary classification, 0-1 loss.

# Vapnik Chervonenkis dimension

Again, binary classification, 0-1 loss.

A set of points $S$ is shattered by a hypothesis class $\mathcal{H}$ if all $2^{|S|}$ labelings on $S$ are produced by hypothesis in $\mathcal{H}$, namely $|\mathcal{H}(S)| = 2^{|S|}$

Examples

# Vapnik Chervonenkis dimension

The VC dimension of $\mathcal{H}$ is the size of the largest set $S$ of points it shatters.

If the VC dimension of $\mathcal{H}$ is $d$, it doesn't mean every set of $d$ points is shattered by $\mathcal{H}$
      only that some set of $d$ points is

But it does mean *no* set of $d + 1$ points can be shattered by $\mathcal{H}$

Larger VC dimension, more power

# Sauer's lemma

If $\mathcal{H}$ has VC dimension $d$, how many labelings on a sample $S$ of size $n$ can it generate?

Trivially, if $n > d$, then number of labelings is $< 2^n$
But one would imagine $2^n$ is a gross overestimate
Proposed by Erdös, solved (1972) and re-proved several times in other contexts
        including by Vapnik and Chervonenkis

# Sauer's lemma

If $\mathcal{H}$ has VC dimension $d$ and $S$ is a sample of size $n$,

$$|\mathcal{H}(S)| \leq \sum_{i=0}^{d} \binom{n}{i} \overset{\text{def}}{=} L(n, d).$$

Proof (simple, and by induction)

We prove a stronger result that

$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

# Induction argument

To prove:
$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

Proof: When $n = 1$, either both sides are 1 or both are 2.

Induction hypothesis: Assume true for all sets $S$ with size $< n$, will prove for all $S$ of size $n$

Hence qed

# Proof

To prove:
$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

Let $S'$ be the sample with the last example removed (so size $n-1$) and let
$$Y_0 = \{y(S') : y(S') \in \mathcal{H}(S')\}$$

and
$$Y_1 = \{y(S') : (y(S'), 0) \text{ and } (y(S'), 1) \in \mathcal{H}(S)\}$$

Clearly
$$|\mathcal{H}(S)| = |Y_0| + |Y_1|$$

## Proof

To prove:
$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

Recall $S'$ be the sample with the last example removed (so size $n-1$) and that

$$Y_0 = \{y(S') : y(S') \in \mathcal{H}(S')\}$$

From induction hypothesis

$$|Y_0| \leq |\{B \in S : \mathcal{H} \text{ shatters } B \text{ and } y_n \notin B\}|$$

# Proof

To prove:
$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

Recall $S'$ be the sample with the last example removed (so size $n-1$) and

$$Y_1 = \{y(S') : (y(S'), \textcolor{blue}{0}) \text{ and } (y(S'), \textcolor{red}{1}) \in \mathcal{H}(S)\}$$

Let $\mathcal{H}'$ be a subset of $\mathcal{H}$. We put a pair $h, h'$ into $\mathcal{H}'$ if $h, h'$ agree on $S'$ but disagree on the last example.

Now we claim $|Y_1| = |\mathcal{H}'(S')|$ and therefore that

$$|\mathcal{H}'(S')| \leq |\{B \in S : \mathcal{H} \text{ shatters } B \text{ and } y_n \in B\}|$$

# Proof

Therefore
$$|\mathcal{H}(S)| = |Y_0| + |Y_1|,$$

but
$$|Y_0| \leq |\{B \in S : \mathcal{H} \text{ shatters } B \text{ and } y_n \notin B\}|$$

and
$$|Y_1| \leq |\{B \in S : \mathcal{H} \text{ shatters } B \text{ and } y_n \in B\}|$$

and so, the result follows!

# Next steps

How Sauer's lemma gives us learnability results for infinite classes
   log VCDim instead of log #hypothesis

Caveat: not strong enough to explain neural networks
   More refinements to come

# VC dimension and PAC learnability

Let $S$ be a sample of size $n$, generated iid $D$

For each sample, what is the worst deviation between sample error $\hat{L}(h)$ made by some $h \in \mathcal{H}$ and its true generalization error $L(h)$?

Namely

$$\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$$

Today' class: bound on this

# VC dimension and PAC learnability

$$S = (z_1 \ldots z_n)$$
$$\quad y_1 \qquad y_n$$

$$\hat{L}(h) = \frac{1}{n} \sum 1(h(z_i) \neq y_i)$$

Let $S$ be a sample of size $n$, generated iid $D$

$$L(h) = \mathbb{E}_{z \sim D} 1(h(z) \neq y)$$

For each sample, what is the worst deviation between sample error $\hat{L}(h)$ made by some $h \in \mathcal{H}$ and its true generalization error $L(h)$? Namely

$$\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$$

## Today' class: bound on this

Full disclosure, we will bound $\mathbb{E}_S \sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$, from which we bound $P(\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| > \epsilon)$ via a Markov inequality

$$\mathbb{E}_S \sup |\hat{L}(h) - L(h)|$$

# Ghost sample

Let $S'$ be a "ghost sample" (an imaginary sample also generated iid $D$)

Let $\hat{L}_{S'}(h)$ be the empirical error of hypothesis $h$ on $S'$

# Ghost sample

$$\mathbb{E}_{S'} \hat{L}_{S'}(h) = \mathbb{E}_{z_1' \ldots z_n'} \frac{1}{n} \sum \mathbb{1}\left(h(z_i') \neq y_i\right)$$

$$= \frac{1}{n} \sum \mathbb{E}_{z_i' \sim D} \mathbb{1}\left(h(z_i') \neq y_i\right) = \underset{z \sim D}{P\left(h(z) \neq y\right)}$$

Let $S'$ be a "ghost sample" (an imaginary sample also generated iid $D$)

Let $\hat{L}_{S'}(h)$ be the empirical error of hypothesis $h$ on $S'$

Still have $L(h) = \mathbb{E}_{S' \sim D} \hat{L}_{S'}(h)$ (expectation over all ghost samples)

Using above, we write for each $h \in \mathcal{H}$,

$$|\hat{L}(h) - L(h)| = |\mathbb{E}_{S'}(\hat{L}(h) - \hat{L}_{S'}(h))| \leq \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

# Ghost sample

Let $S'$ be a "ghost sample" (an imaginary sample also generated iid $D$)

We observed for each $h$,

$$|\hat{L}(h) - L(h)| \leq \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

# Ghost sample

Let $S'$ be a "ghost sample" (an imaginary sample also generated iid $D$)

We observed for each $h$,

$$|\hat{L}(h) - L(h)| \leq \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

So for each training sample,

$$\sup_{h \in \mathcal{H}}|\hat{L}(h) - L(h)| \leq \sup_{h \in \mathcal{H}}\mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

## Ghost sample

$$\max_{0 < x < 1} (x - x) = 0 \leq \max_{0 < x < 1} x + \max_{0 < x < 1} (-x)$$

Let $S'$ be a "ghost sample" (an imaginary sample also generated iid $D$)

We observed for each $h$,

$$|\hat{L}(h) - L(h)| \leq \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

So for each training sample,

$$\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| \leq \sup_{h \in \mathcal{H}} \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)|$$

But

$$\sup_{h \in \mathcal{H}} \mathbb{E}_{S'}|\hat{L}(h) - \hat{L}_{S'}(h)| \leq \mathbb{E}_{S'} \sup_{h \in \mathcal{H}} |\hat{L}(h) - \hat{L}_{S'}(h)|$$

and hence

$$\mathop{\mathbb{E}}_{S} \sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| \leq \mathop{\mathbb{E}}_{S} \mathbb{E}_{S'} \sup_{h \in \mathcal{H}} |\hat{L}(h) - \hat{L}_{S'}(h)|$$

UNIVERSITY of HAWAIʻI MĀNOA

# Symmetrization with ghost sample

Training sample $S = (z_1, \ldots, z_n)$ and ghost sample $S' = (z'_1, \ldots, z'_n)$. Both are drawn *i.i.d.* $D$.

# Symmetrization with ghost sample

Training sample $S = (z_1, \ldots, z_n)$ and ghost sample $S' = (z'_1, \ldots, z'_n)$. Both are drawn *i.i.d.* $D$.

Swap out the first element of ghost with that of train,

# Symmetrization with ghost sample

Training sample $S = (z_1, \ldots, z_n)$ and ghost sample $S' = (z'_1, \ldots, z'_n)$. Both are drawn *i.i.d.* $D$.

Swap out the first element of ghost with that of train, (so we now consider the function in place of $\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$

$$\sup_h |\ell(h, z'_1) - \ell(h, z_1) + \sum_{i=2}^n (\ell(h, z_i) - \ell(h, z'_i)|$$

Above function different, its expectation (over $S, S'$) is not, that is:

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h |\hat{L}(h) - \hat{L}_{S'}(h)|$$

$$= \mathbb{E}_S \mathbb{E}_{S'} \sup_h |\ell(h, z_1) - \ell(h, z'_1) + \sum_{i=2}^n (\ell(h, z_i) - \ell(h, z'_i)|$$

$$= \mathbb{E}_S \mathbb{E}_{S'} \sup_h |\ell(h, z'_1) - \ell(h, z_1) + \sum_{i=2}^n (\ell(h, z_i) - \ell(h, z'_i)|$$

# Symmetrization with ghost sample

In fact, can swap as many examples between train/ghost as we want to get new functions with the same expectation

We could even pick $\sigma_1, \ldots, \sigma_n$ to be independent random variables taking values in $\{-1, 1\}^n$, with equal probabilities and

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h |\hat{L}(h) - \hat{L}_{S'}(h)|$$

$$= \mathbb{E}_S \mathbb{E}_{S'} \sup_h |\sum_{i=1}^n (\ell(h, z_i) - \ell(h, z_i'))|$$

$$= E_\sigma \mathbb{E}_S \mathbb{E}_{S'} \sup_h |\sum_{i=1}^n \sigma_i (\ell(h, z_i) - \ell(h, z_i'))|$$

# Reviewing so far

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h |\hat{L}(h) - \hat{L}_{S'}(h)|$$

$$= E_\sigma \mathbb{E}_S \mathbb{E}_{S'} \sup_h |\sum_{i=1}^n \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|$$

$$= \mathbb{E}_S \mathbb{E}_{S'} E_\sigma \sup_h |\sum_{i=1}^n \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|$$

We now fix a train and ghost sample combination and examine

$$\sup_{h \in \mathcal{H}} |\sum_{i=1}^n \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|$$

The quantity above is very close to the Rademacher complexity, which will give us another way to deal with generalization error.
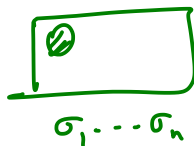
# Better Heoffding

In fact, our earlier version of Hoeffding's inequality wasn't the full version. We don't need the random variables to be iid, only independence suffices.

Let $W_1, \ldots, W_n$ be independent variables, the variables bounded in range $X_i \in [a, b]$, with $\mathbb{E}W_i = \mu$. Then for any $\epsilon > 0$,

$$P\left(|\frac{1}{n}\sum_i W_i - \mu| > \epsilon\right) \leq 2\exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

# Completing the proof for classes with small VC dimension

We were examining for a fixed train/ghost sample:

$$\sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|$$

Now given a fixed train and ghost sample, and a fixed $h$, let
$W_i = \sigma_i(\ell(h, z_i) - \ell(h, z_i'))$.
$\mathbb{E}W_i = 0$, $-1 \leq W_i \leq 1$, and $W_i$ independent (not identical!)

Therefore for each $h \in \mathcal{H}$,

$$P(|\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$

But doesn't $\mathcal{H}$ have infinitely many hypotheses?

# Recall: Sauer's lemma

If $\mathcal{H}$ has VC dimension $d$ and $S$ is a sample of size $n$,

$$|\mathcal{H}(S)| \leq \sum_{i=0}^{d} \binom{n}{i} \stackrel{\text{def}}{=} L(n, d).$$

Proof (simple, and by induction)

We prove a stronger result that

$$|\mathcal{H}(S)| \leq |B \subset S : \mathcal{H} \text{ shatters } B|.$$

# Now for Sauer's lemma

Therefore for each $h \in \mathcal{H}$,

$$P(|\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon) \leq 2\exp\left(-\frac{n\epsilon^2}{2}\right)$$

Now since we have fixed train/ghost, there are only $L(2n, d)$ labelings from $\mathcal{H}$ if it has VC dimension $d$. So effectively only $L(2n, d)$ hypotheses! Therefore

$$P\left(\sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon \Big| S, S'\right)$$

$$\leq 2L(2n, d)\exp\left(-\frac{n\epsilon^2}{2}\right)$$

Remember: this is probability over choice of $\sigma_i$ (the training and ghost samples are being held fixed)

# Now for Sauer's lemma

Therefore for each $h \in \mathcal{H}$,

$$P(|\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$

Now since we have fixed train/ghost, there are only $L(2n, d)$ labelings from $\mathcal{H}$ if it has VC dimension $d$. So effectively only $L(2n, d)$ hypotheses! Therefore

$$P\left(\sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon \Big| S, S'\right)$$
$$\leq 2L(2n, d) \exp\left(-2n\epsilon^2\right)$$

Remember: this is probability over choice of $\sigma_i$ (the training and ghost samples are being held fixed)$\,$from which we get a straightforward bound on $\mathbb{E}_S \mathbb{E}_{S'} \mathbb{E}_\sigma \sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|,$

# Now for Sauer's lemma

Therefore for each $h \in \mathcal{H}$,

$$P(|\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z'_i))| \geq \epsilon) \leq 2\exp\left(-\frac{n\epsilon^2}{2}\right)$$

Now since we have fixed train/ghost, there are only $L(2n, d)$ labelings from $\mathcal{H}$ if it has VC dimension $d$. So effectively only $L(2n, d)$ hypotheses! Therefore

$$P\left(\sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z'_i))| \geq \epsilon \Big| S, S'\right)$$
$$\leq 2L(2n, d)\exp\left(-2n\epsilon^2\right)$$

Remember: this is probability over choice of $\sigma_i$ (the training and ghost samples are being held fixed) from which we get a straightforward bound on $\mathbb{E}_S\mathbb{E}_{S'}\mathbb{E}_\sigma \sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z'_i))|$, which upper bounds $\mathbb{E}_S \sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$,

# Now for Sauer's lemma

Therefore for each $h \in \mathcal{H}$,

$$P(|\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon) \leq 2\exp\left(-\frac{n\epsilon^2}{2}\right)$$

Now since we have fixed train/ghost, there are only $L(2n, d)$ labelings from $\mathcal{H}$ if it has VC dimension $d$. So effectively only $L(2n, d)$ hypotheses! Therefore

$$P\left(\sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))| \geq \epsilon \Big| S, S'\right)$$
$$\leq 2L(2n, d)\exp\left(-2n\epsilon^2\right)$$

Remember: this is probability over choice of $\sigma_i$ (the training and ghost samples are being held fixed) from which we get a straightforward bound on $\mathbb{E}_S \mathbb{E}_{S'} \mathbb{E}_\sigma \sup_{h \in \mathcal{H}} |\sum_{i=1}^{n} \sigma_i(\ell(h, z_i) - \ell(h, z_i'))|$, which upper bounds $\mathbb{E}_S \sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)|$, which upper bounds $P(\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| > \epsilon)$ via a Markov inequality

# Putting everything together

For a class $\mathcal{H}$ with VC dimension $d$,

$$P\left(\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| \geq \frac{\sqrt{\log L(2n, d)}}{\eta\sqrt{2n}}\right) \leq \eta.$$

# Putting everything together

For a class $\mathcal{H}$ with VC dimension $d$,

$$P\left(\sup_{h \in \mathcal{H}} |\hat{L}(h) - L(h)| \geq \frac{\sqrt{\log L(2n, d)}}{\eta\sqrt{2n}}\right) \leq \eta.$$

For a given accuracy $\epsilon$ and confidence $\eta$, we need a training sample of size

$$n \geq 4\frac{2d}{(\eta\epsilon)^2} \log\left(\frac{2d}{(\eta\epsilon)^2}\right) + \frac{4d\log(2e/d)}{(\eta\epsilon)^2}.$$

# But even this isn't enough

We need to do beter. For kernel methods, we lift up the points to very high-d space. Even linear classifiers in this high-d space have VC dimension equal to high-d $+1$, which is too large.

Similar line of argument works, but we focus on Rademacher complexity

# But even this isn't enough

We need to do beter. For kernel methods, we lift up the points to very high-d space. Even linear classifiers in this high-d space have VC dimension equal to high-d $+1$, which is too large.

Similar line of argument works, but we focus on Rademacher complexity

Let $\mathcal{F}$ be a set of functions on $S = (z_1, \ldots, z_n)$. Then

$$\mathcal{R}(\mathcal{F}(S)) = \frac{1}{m} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} f(z_i) \right]. \quad \frac{1}{m} \mathbb{E}_\sigma \left[ \sup_f \sum \sigma_i f(z_i) \right]$$

Note that we don't think of the worst case training sample or an average. The Rademacher complexity is defined per sample.

# But even this isn't enough

We need to do beter. For kernel methods, we lift up the points to very high-d space. Even linear classifiers in this high-d space have VC dimension equal to high-d $+1$, which is too large.

Similar line of argument works, but we focus on Rademacher complexity
Let $\mathcal{F}$ be a set of functions on $S = (z_1, \ldots, z_n)$. Then

$$\mathcal{R}(\mathcal{F}(S)) = \frac{1}{m} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} f(z_i) \right].$$

Note that we don't think of the worst case training sample or an average. The Rademacher complexity is defined per sample.
If we have a hypothesis class $\mathcal{H}$, and let $f(z_i) = \ell(h, z_i)$, we pretty much have something very similar to what we have been seeing.

# Rademacher complexity

For a training sample $S = (z_1, \ldots, z_n)$, using the ghost sample idea

$$\sup_h L(h) - \hat{L}(h)$$

$$= \sup_h \mathbb{E}_{S'}\big[\hat{L}_{S'}(h) - \hat{L}(h)\big]$$

$$\leq \mathbb{E}_{S'} \sup_h \big[\hat{L}_{S'}(h) - \hat{L}(h)\big]$$

$$= \mathbb{E}_{S'} \sup_h \big[\frac{1}{n}\sum_i (\ell(h, z') - \ell(h, z))\big]$$

We can now swap between the training/ghost samples just like before to ge new functions with the same expectation (under $S$ and $S'$)

# Rademacher central argument

By arguments very similar to before, where $\sigma_i \sim$ iid $B(1/2)$

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h \left[ \hat{L}_{S'}(h) - \hat{L}(h) \right]$$

$$= \mathbb{E}_S \mathbb{E}'_S \mathbb{E}_\sigma \sup_h \left[ \frac{1}{n} \sum_i \sigma_i (\ell(h, z'_i) - \ell(h, z_i)) \right]$$

$$\longrightarrow \leq 2 \mathbb{E}_S \mathbb{E}_\sigma \sup_h \frac{1}{n} \left[ \sum_i \sigma_i \ell(h, z) \right]$$

$$= 2 \mathbb{E}_S \mathcal{R}(\ell(\mathcal{H}(S)))$$

$$\mathbb{E} \, Sup \left( \frac{1}{n} \sum \sigma_i \, \ell(h, z'_i) - \frac{1}{n} \sum \sigma_i \, \ell(h, z_i) \right).$$

$$\leq \; 2 \mathbb{E} \, Sup \; \frac{1}{n} \sum \sigma_i \, \ell(h, z_i)$$

# Rademacher central argument

By arguments very similar to before, where $\sigma_i \sim$iid B$(1/2)$

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h \big[ \hat{L}_{S'}(h) - \hat{L}(h) \big]$$

$$= \mathbb{E}_S \mathbb{E}'_S \mathbb{E}_\sigma \sup_h \big[ \frac{1}{n} \sum_i \sigma_i (\ell(h, z') - \ell(h, z)) \big]$$

$$\leq 2 \mathbb{E}_S \mathbb{E}_\sigma \sup_h \frac{1}{n} \big[ \sum_i \sigma_i \ell(h, z) \big]$$

$$= 2 \mathbb{E}_S \mathcal{R}(\ell(\mathcal{H}(S)))$$

where $\ell(\mathcal{H}(S))$ is the set of loss sequences obtained on $S$ by various labelings in $\mathcal{H}$.

# Rademacher central argument

By arguments very similar to before, where $\sigma_i \sim$iid B(1/2)

$$\mathbb{E}_S \mathbb{E}_{S'} \sup_h \left[ \hat{L}_{S'}(h) - \hat{L}(h) \right]$$

$$= \mathbb{E}_S \mathbb{E}'_S \mathbb{E}_\sigma \sup_h \left[ \frac{1}{n} \sum_i \sigma_i (\ell(h, z') - \ell(h, z)) \right]$$

$$\leq 2 \mathbb{E}_S \mathbb{E}_\sigma \sup_h \frac{1}{n} \left[ \sum_i \sigma_i \ell(h, z) \right]$$

$$= 2 \mathbb{E}_S \mathcal{R}(\ell(\mathcal{H}(S)))$$

where $\ell(\mathcal{H}(S))$ is the set of loss sequences obtained on $S$ by various labelings in $\mathcal{H}$.

In the VC bound, we just used Sauer's lemma for number of labelings on $S$. Sadly, for high dimensional liftings like in kernel methods, this doesn't yield good results. Can we do better?