

The Infinite Case: Rademacher Complexity and VC-dimension



Rationale: PAC provides no bounds for the infinite hypothesis class

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{VC}(\epsilon/2, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil$$

Is it still possible to learn if $|\mathcal{H}|$ is infinite?

Είναι άραγε τόσο δύσκολο, ακόμα και για πολύ απλούς ταξινομητές όπως οι axis-aligned rectangles (που όμως έχει άπειρο $|\mathcal{H}|$), να βρούμε ένα άνω φράγμα για το σφάλμα της γενίκευσης;



OXI: έχουμε φράγματα για το σφάλμα της γενίκευσης και στην περίπτωση της απειρομελούς κλάσης υποθέσεων:

φράγμα Rademacher \Leftrightarrow Growth \Leftrightarrow VC-dimension (shattering dimension)

Setup

Nothing new ...

- Samples $S = ((x_1, y_1), \dots, (x_m, y_m))$
- Labels $y_i = \{-1, +1\}$
- Hypothesis $h: X \rightarrow \{-1, +1\}$
- Training error: $\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i]$

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

(2)

(3)

(4)

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i, y_i) == (1, -1) \text{ or } (-1, 1)) \\ 0 & \text{if } (h(x_i, y_i) == (1, 1) \text{ or } (-1, -1)) \end{cases} \quad (2)$$

(3)

(4)

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

(4)

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

Correlation between predictions and labels

An alternative derivation of training error

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

Minimizing training error is thus equivalent to maximizing correlation

$$\arg \max_h \frac{1}{m} \sum_i^m y_i h(x_i) \quad (5)$$

Playing with Correlation

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

Playing with Correlation

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

This gives us Rademacher correlation—what's the best that a random classifier could do?

$$\hat{\mathcal{R}}_S(H) \equiv \mathbb{E}_\sigma \left[\max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] \quad (7)$$

Playing with Correlation

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

This gives us Rademacher correlation—what's the best that a random classifier could do?

$$\hat{\mathcal{R}}_S(H) \equiv \mathbb{E}_\sigma \left[\max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] \quad (7)$$

Notation: $\mathbb{E}_p[f] \equiv \sum_x p(x) f(x)$

Playing with Correlation

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

This gives us Rademacher correlation—what's the best that a random classifier could do?

$$\hat{\mathcal{R}}_S(H) \equiv \mathbb{E}_\sigma \left[\max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] \quad (7)$$

Note: Empirical Rademacher complexity is with respect to a sample.

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$|H| = 2^m$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[\max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right]$$

$$|H| = 2^m$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$h(x_i) \mathbb{E}_\sigma \left[\frac{1}{m} \sum_i^m \sigma_i \right]$$

$$|H| = 2^m$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$h(x_i) \mathbb{E}_\sigma \left[\frac{1}{m} \sum_i^m \sigma_i \right] = 0$$

$$|H| = 2^m$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$h(x_i) \mathbb{E}_\sigma \left[\frac{1}{m} \sum_i^m \sigma_i \right] = 0$$

$$|H| = 2^m$$

$$\mathbb{E}_\sigma \left[\max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right]$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$h(x_i) \mathbb{E}_\sigma \left[\frac{1}{m} \sum_i^m \sigma_i \right] = 0$$

$$|H| = 2^m$$

$$\frac{m}{m} = 1$$

Rademacher Extrema

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$h(x_i) \mathbb{E}_\sigma \left[\frac{1}{m} \sum_i^m \sigma_i \right] = 0$$

$$|H| = 2^m$$

$$\frac{m}{m} = 1$$

- Rademacher correlation is larger for more complicated hypothesis space.
- What if you're right for stupid reasons?

Generalizing Rademacher Complexity

We can generalize Rademacher complexity to consider all sets of a particular size.

$$\mathcal{R}_m(H) = \mathbb{E}_{S \sim D^m} [\hat{\mathcal{R}}_S(H)] \quad (8)$$

Generalizing Rademacher Complexity

Theorem

Convergence Bounds Let F be a family of functions mapping from Z to $[0, 1]$, and let sample $S = (z_1, \dots, z_m)$ where $z_i \sim D$ for some distribution D over Z . Define $\mathbb{E}[f] \equiv \mathbb{E}_{z \sim D}[f(z)]$ and $\hat{\mathbb{E}}_S[f] \equiv \frac{1}{m} \sum_{i=1}^m f(z_i)$. With probability greater than $1 - \delta$ for all $f \in F$:

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (8)$$

Generalizing Rademacher Complexity

Theorem

Convergence Bounds Let F be a family of functions mapping from Z to $[0, 1]$, and let sample $S = (z_1, \dots, z_m)$ where $z_i \sim D$ for some distribution D over Z . Define $\mathbb{E}[f] \equiv \mathbb{E}_{z \sim D}[f(z)]$ and $\hat{\mathbb{E}}_S[f] \equiv \frac{1}{m} \sum_{i=1}^m f(z_i)$. With probability greater than $1 - \delta$ for all $f \in F$:

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (8)$$

f is a surrogate for the accuracy of a hypothesis (mathematically convenient)

Aside: McDiarmid's Inequality

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (9)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp\left\{\frac{-2\epsilon^2}{\sum_i^m c_i^2}\right\} \quad (10)$$

Aside: McDiarmid's Inequality

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (9)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp\left\{\frac{-2\epsilon^2}{\sum_i^m c_i^2}\right\} \quad (10)$$

Proof in Mohri (appendix D.7, p.442) (requires Martingale, constructing $V_k = \mathbb{E}[V | x_1 \dots x_k] - \mathbb{E}[V | x_1 \dots x_{k-1}]$).

Aside: McDiarmid's Inequality

Theorem D.8 (McDiarmid's inequality) *Let $X_1, \dots, X_m \in \mathcal{X}^m$ be a set of $m \geq 1$ independent random variables and assume that there exist $c_1, \dots, c_m > 0$ such that $f: \mathcal{X}^m \rightarrow \mathbb{R}$ satisfies the following conditions:*

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i, \quad (\text{D.15})$$

for all $i \in [m]$ and any points $x_1, \dots, x_m, x'_i \in \mathcal{X}$. Let $f(S)$ denote $f(X_1, \dots, X_m)$, then, for all $\epsilon > 0$, the following inequalities hold:

$$\mathbb{P}[f(S) - \mathbb{E}[f(S)] \geq \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^m c_i^2}\right) \quad (\text{D.16})$$

$$\mathbb{P}[f(S) - \mathbb{E}[f(S)] \leq -\epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^m c_i^2}\right). \quad (\text{D.17})$$

McDiarmid's inequality is used in several of the proofs in this book. It can be understood in terms of stability: if changing any of its argument affects f only in a limited way, then, its deviations from its mean can be exponentially bounded.

Εξαιτίας των δύο ανισοτήτων έχουμε το O(). Βλέπε θεώρημα 3.3 σε Mohri.

Aside: McDiarmid's Inequality

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (9)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp\left\{\frac{-2\epsilon^2}{\sum_i^m c_i^2}\right\} \quad (10)$$

Proof in Mohri (appendix D.7, p.442) (requires Martingale, constructing $V_k = \mathbb{E}[V | x_1 \dots x_k] - \mathbb{E}[V | x_1 \dots x_{k-1}]$).

What function do we care about for Rademacher complexity? Let's define

$$\Phi(S) = \sup_f \left(\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \right) = \sup_f \left(\mathbb{E}[f] - \frac{1}{m} \sum_i f(z_i) \right) \quad (11)$$

Step 1: Bounding divergence from true Expectation

Lemma

Moving to Expectation *With probability at least $1 - \delta$,*

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}$$

Since $f(z_1) \in [0, 1]$, changing any z_i to z'_i in the training set will change $\frac{1}{m} \sum_i f(z_i)$ by at most $\frac{1}{m}$, so we can apply McDiarmid's inequality with

$$\epsilon = \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \text{ and } c_i = \frac{1}{m}.$$

Handwritten derivation of the lemma's proof:

$$\delta \leq \exp\left(\frac{-2\epsilon^2}{\sum \left(\frac{1}{m}\right)^2}\right) = \exp(-2m\epsilon^2)$$
$$\ln(\delta) \leq -2m\epsilon^2$$
$$\frac{-\ln(\delta)}{2m} \geq \epsilon^2 \Rightarrow \sqrt{\frac{\ln(\frac{1}{\delta})}{2m}} \geq \epsilon$$

Step 2: Comparing two different empirical expectations

Define a ghost sample $S' = (z'_1, \dots, z'_m) \sim D$. How much can two samples from the same distribution vary?

Lemma

Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[\sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (12)$$

(13)

Step 2: Comparing two different empirical expectations

Define a ghost sample $S' = (z'_1, \dots, z'_m) \sim D$. How much can two samples from the same distribution vary?

Lemma

Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[\sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (12)$$

$$= \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} (\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[f]] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$(14)$$

The expectation is equal to the expectation of the empirical expectation of all sets S'

Step 2: Comparing two different empirical expectations

Define a ghost sample $S' = (z'_1, \dots, z'_m) \sim D$. How much can two samples from the same distribution vary?

Lemma

Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[\sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (12)$$

$$= \mathbb{E}_S \left[\sup_{f \in F} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[f]] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$= \mathbb{E}_S \left[\sup_{f \in F} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]]) \right] \quad (14)$$

$$(15)$$

S and S' are distinct random variables, so we can move inside the expectation

Step 2: Comparing two different empirical expectations

Define a ghost sample $S' = (z'_1, \dots, z'_m) \sim D$. How much can two samples from the same distribution vary?

Lemma

Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[\sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (12)$$

$$= \mathbb{E}_S \left[\sup_{f \in F} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$\leq \mathbb{E}_{S, S'} \left[\sup_f (\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (14)$$

The expectation of a max over some function is at least the max of that expectation over that function

Step 3: Adding in Rademacher Variables

From S, S' we'll create T, T' by swapping elements between S and S' with probability .5. This is still independent, identically distributed (iid) from D . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (15)$$

Step 3: Adding in Rademacher Variables

From S, S' we'll create T, T' by swapping elements between S and S' with probability .5. This is still independent, identically distributed (iid) from D . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (15)$$

Let's introduce σ_i :

$$\hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] = \frac{1}{m} \begin{cases} f(z_i) - f(z'_i) & \text{with prob .5} \\ f(z'_i) - f(z_i) & \text{with prob .5} \end{cases} \quad (16)$$

$$= \frac{1}{m} \sum_i \sigma_i (f(z'_i) - f(z_i)) \quad (17)$$

Step 3: Adding in Rademacher Variables

From S, S' we'll create T, T' by swapping elements between S and S' with probability .5. This is still independent, identically distributed (iid) from D . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (15)$$

Let's introduce σ_i :

$$\hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] = \frac{1}{m} \begin{cases} f(z_i) - f(z'_i) & \text{with prob .5} \\ f(z'_i) - f(z_i) & \text{with prob .5} \end{cases} \quad (16)$$

$$= \frac{1}{m} \sum_i \sigma_i (f(z'_i) - f(z_i)) \quad (17)$$

Thus:

$$\mathbb{E}_{S, S'} \left[\sup_{f \in F} \left(\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \right) \right] = \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in F} \left(\sum_i \sigma_i (f(z'_i) - f(z_i)) \right) \right].$$

Step 4: Making These Rademacher Complexities

Before, we had $\mathbb{E}_{S, S', \sigma} \left[\sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

Step 4: Making These Rademacher Complexities

Before, we had $\mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (18)$$

(19)

Taking the sup jointly must be less than or equal the individual sup.

Step 4: Making These Rademacher Complexities

Before, we had $\mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (18)$$

$$\leq \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) \right] + \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (19)$$

$$(20)$$

Linearity

Step 4: Making These Rademacher Complexities

Before, we had $\mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (18)$$

$$\leq \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i \sigma_i f(z'_i) \right] + \mathbb{E}_{S,S',\sigma} \left[\sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (19)$$

$$= \mathcal{R}_m(F) + \mathcal{R}_m(F) \quad (20)$$

Definition

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\Phi(\mathcal{S}) \leq \mathbb{E}_{\mathcal{S}}[\Phi(\mathcal{S})] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Step 1

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Definition of Φ

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Drop the sup, still true

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq \mathbb{E}_{S,S'} \left[\sup_f (\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Step 2

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in F} \left(\sum_i \sigma_i (f(z'_i) - f(z_i)) \right) \right] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Step 3

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Step 4

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Recall that $\hat{\mathcal{R}}_S(F) \equiv \mathbb{E}_\sigma \left[\sup_f \frac{1}{m} \sum_i \sigma_i f(z_i) \right]$, so we apply McDiarmid's inequality again (because $f \in [0, 1]$):

$$\mathcal{R}_m(F) \leq \hat{\mathcal{R}}_S(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (22)$$

Putting the Pieces Together

With probability $\geq 1 - \delta$:

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (21)$$

Recall that $\hat{\mathcal{R}}_S(F) \equiv \mathbb{E}_\sigma \left[\sup_f \frac{1}{m} \sum_i \sigma_i f(z_i) \right]$, so we apply McDiarmid's inequality again (because $f \in [0, 1]$):

$$\mathcal{R}_m(F) \leq \hat{\mathcal{R}}_S(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (22)$$

Putting the two together:

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_S(F) + \mathcal{O} \left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) \quad (23)$$

What about hypothesis classes?

Define:

$$Z \equiv X \times \{-1, +1\} \quad (24)$$

$$f_h(x, y) \equiv \mathbb{1}[h(x) \neq y] \quad (25)$$

$$F_H \equiv \{f_h : h \in H\} \quad (26)$$

What about hypothesis classes?

Define:

$$Z \equiv X \times \{-1, +1\} \quad (24)$$

$$f_h(x, y) \equiv \mathbb{1} [h(x) \neq y] \quad (25)$$

$$F_H \equiv \{f_h : h \in H\} \quad (26)$$

We can use this to create expressions for generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1} [h(x) \neq y]] = \mathbb{E} [f_h] \quad (27)$$

$$\hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1} [h(x_i) \neq y] = \hat{\mathbb{E}}_S [f_h] \quad (28)$$

What about hypothesis classes?

Define:

$$Z \equiv X \times \{-1, +1\} \quad (24)$$

$$f_h(x, y) \equiv \mathbb{1} [h(x) \neq y] \quad (25)$$

$$F_H \equiv \{f_h : h \in H\} \quad (26)$$

We can use this to create expressions for generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1} [h(x) \neq y]] = \mathbb{E} [f_h] \quad (27)$$

$$\hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1} [h(x_i) \neq y] = \hat{\mathbb{E}}_S [f_h] \quad (28)$$

We can plug this into our theorem!

Generalization bounds

- We started with expectations

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\hat{\mathcal{R}}_S(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (29)$$

- We also had our definition of the generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1}[h(x) \neq y]] = \mathbb{E}[f_h] \quad \hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1}[h(x_i) \neq y] = \hat{\mathbb{E}}_S[f_h]$$

Generalization bounds

$$\hat{\mathcal{R}}_S(F_H) = \frac{1}{2} \hat{\mathcal{R}}_S(H) \quad (30)$$

$$\hat{\mathfrak{R}}_S(\mathcal{G}) = \frac{1}{2} \hat{\mathfrak{R}}_{S_x}(\mathcal{H}). \quad (3.16)$$

Proof: For any sample $S = ((x_1, y_1), \dots, (x_m, y_m))$ of elements in $\mathcal{X} \times \{-1, +1\}$, by definition, the empirical Rademacher complexity of \mathcal{G} can be written as:

$$\begin{aligned} \hat{\mathfrak{R}}_S(\mathcal{G}) &= \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i 1_{h(x_i) \neq y_i} \right] \\ &= \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \frac{1 - y_i h(x_i)}{2} \right] \\ &= \frac{1}{2} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m -\sigma_i y_i h(x_i) \right] \\ &= \frac{1}{2} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] = \frac{1}{2} \hat{\mathfrak{R}}_{S_x}(\mathcal{H}), \end{aligned}$$

where we used the fact that $1_{h(x_i) \neq y_i} = (1 - y_i h(x_i))/2$ and the fact that for a fixed $y_i \in \{-1, +1\}$, σ_i and $-y_i \sigma_i$ are distributed in the same way. \square

Generalization bounds

- We started with expectations

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\hat{\mathcal{R}}_S(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (29)$$

- We also had our definition of the generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\mathbb{1}[h(x) \neq y]] = \mathbb{E}[f_h] \quad \hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1}[h(x_i) \neq y] = \hat{\mathbb{E}}_S[f_h]$$

- Combined with the previous result:

$$\hat{\mathcal{R}}_S(F_H) = \frac{1}{2} \hat{\mathcal{R}}_S(H) \quad (30)$$

- All together:

$$R(h) \leq \hat{R}(h) + \hat{\mathcal{R}}_S(H) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{m}}\right) \quad (31)$$

Wrapup

- Interaction of data, complexity, and accuracy
- Still very theoretical
- Next up: How to evaluate generalizability of specific hypothesis classes

Recap

- Rademacher complexity provides nice guarantees

$$R(h) \leq \hat{R}(h) + \hat{\mathcal{R}}_S(H) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{2m}}\right) \quad (32)$$

- But in practice hard to compute for real hypothesis classes
- Is there a relationship with simpler combinatorial measures?

Growth Function

Define the **growth function** $\Pi_H : \mathbb{N} \rightarrow \mathbb{N}$ for a hypothesis set H as:

$$\forall m \in \mathbb{N}, \Pi_H(m) \equiv \max_{\{x_1, \dots, x_m\} \in X} |\{(h(x_1), \dots, h(x_m)) : h \in H\}| \quad (33)$$

Growth Function

Define the **growth function** $\Pi_H : \mathbb{N} \rightarrow \mathbb{N}$ for a hypothesis set H as:

$$\forall m \in \mathbb{N}, \Pi_H(m) \equiv \max_{\{x_1, \dots, x_m\} \in X} \left| \{ (h(x_1), \dots, h(x_m)) : h \in H \} \right| \quad (33)$$

i.e., the number of ways m points can be classified using H .

Rademacher Complexity vs. Growth Function

If G is a function taking values in $\{-1, +1\}$, then

$$\mathcal{R}_m(G) \leq \sqrt{\frac{2 \ln \Pi_G(m)}{m}} \quad (34)$$

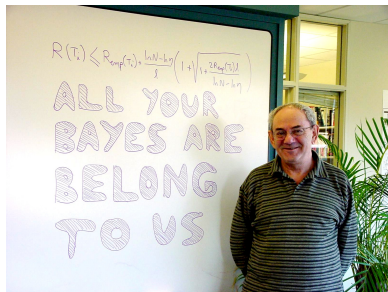
Uses Masart's lemma (Theorem 3.7)

Corollary 3.9 (Growth function generalization bound) *Let \mathcal{H} be a family of functions taking values in $\{-1, +1\}$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, for any $h \in \mathcal{H}$,*

$$R(h) \leq \widehat{R}_S(h) + \sqrt{\frac{2 \log \Pi_{\mathcal{H}}(m)}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}. \quad (3.22)$$

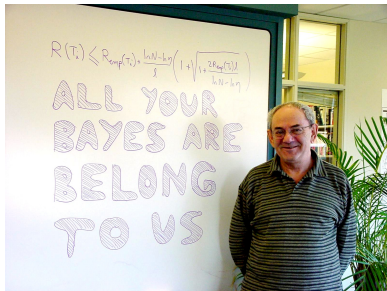
Not very convenient since it requires computing $\Pi_H(m), \forall m$

Vapnik-Chervonenkis Dimension



$$VC(H) \equiv \max \{m : \Pi_H(m) = 2^m\} \quad (35)$$

Vapnik-Chervonenkis Dimension



$$VC(H) \equiv \max \{m : \Pi_H(m) = 2^m\} \quad (35)$$

The size of the largest set that can be fully shattered ($\theta\rho\upsilon\mu\mu\alpha\tau\iota\sigma\tau\epsilon\iota$) by H .

Entropy Properties of a Decision Rule Class with ML abilities - Alexey Chervonenkis lecture

VC Dimension for Hypotheses

- Need upper and lower bounds
- Lower bound: example
- Upper bound: Prove that no set of $d + 1$ points can be shattered by H (harder)

Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

Intervals

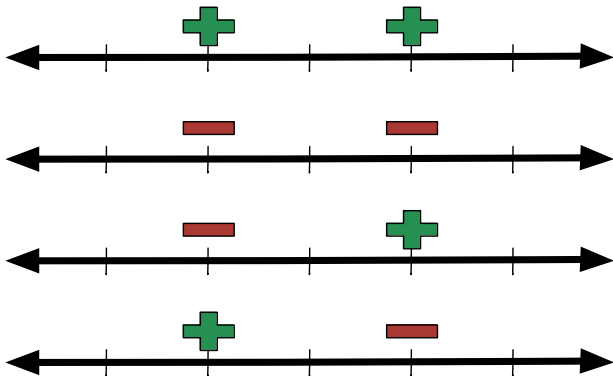
What is the VC dimension of $[a, b]$ intervals on the real line.

- What about two points?

Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

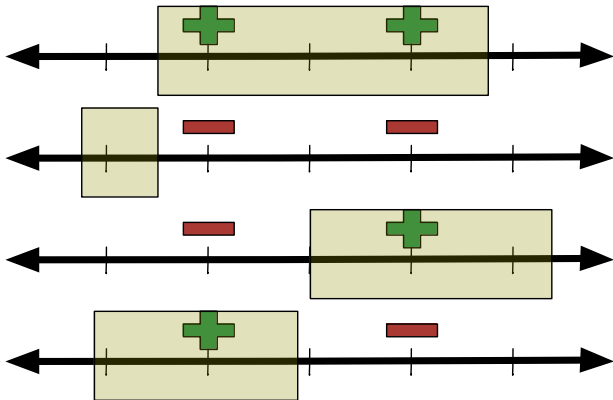
- What about two points?



Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

- What about two points?



Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

- Two points can be perfectly classified, so VC dimension ≥ 2

Intervals

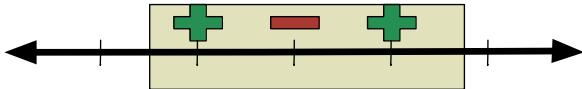
What is the VC dimension of $[a, b]$ intervals on the real line.

- Two points can be perfectly classified, so VC dimension ≥ 2
- What about three points?

Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

- Two points can be perfectly classified, so VC dimension ≥ 2
- What about three points?



Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

- Two points can be perfectly classified, so VC dimension ≥ 2
- What about three points?
- **No set** of three points can be shattered

Intervals

What is the VC dimension of $[a, b]$ intervals on the real line.

- Two points can be perfectly classified, so VC dimension ≥ 2
- What about three points?
- **No set** of three points can be shattered
- Thus, VC dimension of intervals is 2

Hyperplanes

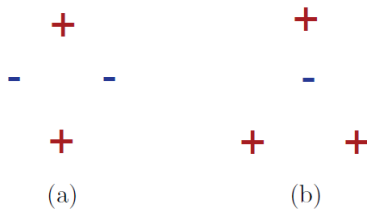


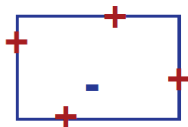
Figure 3.2

Unrealizable dichotomies for four points using hyperplanes in \mathbb{R}^2 . (a) All four points lie on the convex hull. (b) Three points lie on the convex hull while the remaining point is interior.

Axis-aligned-rectangles



(a)



(b)

Figure 3.3

VC-dimension of axis-aligned rectangles. (a) Examples of realizable dichotomies for four points in a diamond pattern. (b) No sample of five points can be realized if the interior point and the remaining points have opposite labels.

Sine Functions

- Consider hypothesis that classifies points on a line as either being above or below a sine wave

$$\{t \rightarrow \sin(\omega x) : \omega \in \mathbb{R}\} \quad (36)$$

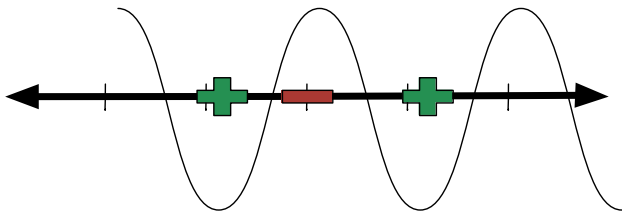
- Can you shatter three points?

Sine Functions

- Consider hypothesis that classifies points on a line as either being above or below a sine wave

$$\{t \rightarrow \sin(\omega x) : \omega \in \mathbb{R}\} \quad (36)$$

- Can you shatter three points?



Sine Functions

- Consider hypothesis that classifies points on a line as either being above or below a sine wave

$$\{t \rightarrow \sin(\omega x) : \omega \in \mathbb{R}\} \quad (36)$$

- Can you shatter four points?

Sine Functions

- Consider hypothesis that classifies points on a line as either being above or below a sine wave

$$\{t \rightarrow \sin(\omega x) : \omega \in \mathbb{R}\} \quad (36)$$

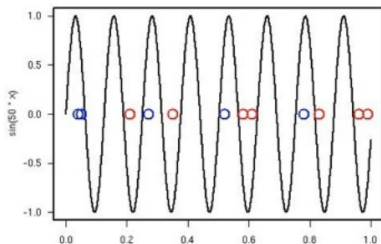
- How many points can you shatter?

Sine Functions

- Consider hypothesis that classifies points on a line as either being above or below a sine wave

$$\{t \rightarrow \sin(\omega x) : \omega \in \mathbb{R}\} \quad (36)$$

- Thus, VC dim of sine on line is ∞



Connecting VC with growth function

VC dimension obviously encodes the complexity of a hypothesis class, but we want to connect that to Rademacher complexity and the growth function so we can prove generalization bounds.

Connecting VC with growth function

VC dimension obviously encodes the complexity of a hypothesis class, but we want to connect that to Rademacher complexity and the growth function so we can prove generalization bounds.

Theorem

Sauer's Lemma *Let H be a hypothesis set with VC dimension d . Then*

$\forall m \in \mathbb{N}$

$$\Pi_H(m) \leq \sum_{i=0}^d \binom{m}{i} \equiv \Phi_d(m) \quad (37)$$

Connecting VC with growth function

VC dimension obviously encodes the complexity of a hypothesis class, but we want to connect that to Rademacher complexity and the growth function so we can prove generalization bounds.

Theorem

Sauer's Lemma *Let H be a hypothesis set with VC dimension d . Then*

$\forall m \in \mathbb{N}$

$$\Pi_H(m) \leq \sum_{i=0}^d \binom{m}{i} \equiv \Phi_d(m) \quad (37)$$

This is good because the sum when multiplied out becomes

$\binom{m}{i} = \frac{m \cdot (m-1) \dots}{i!} = \mathcal{O}(m^d)$. When we plug this into the learning error limits:
 $\log(\Pi_H(2m)) = \log(\mathcal{O}(m^d)) = \mathcal{O}(d \log m)$.

Sauer's Lemma

Definition. Growth Function:

$$\Pi_F(n) = \max\{|F|_s : s \subseteq \mathcal{X}, |s| = n\}$$

Definition. VC dimension

$$d_{VC}(F) = \max\{|s| : s \subseteq \mathcal{X}, f \text{ shatters } s\}$$

Here, we say that a family of binary functions F shatters a set $S \in \mathcal{X}$ if $|F|_S = 2^{|S|}$.

Theorem 2.1. Sauer's Lemma: If $F \subseteq \{\pm 1\}^{\mathcal{X}}$ and $d_{VC} = d$, then $\Pi_F(n) \leq \sum_{i=0}^d \binom{n}{i}$. And for $n \geq d$, $\Pi_F(n) \leq \left(\frac{en}{d}\right)^d$

That means: if $d_{VC}(F)$ is ∞ , we always get exponential growth function; however, if $d_{VC}(F) = d$ is finite, the growth function increases exponentially up to d and polynomially for $n > d$.

PROOF. Fix $(x_1, \dots, x_n) \in \mathcal{X}$, and consider a table containing the values of functions in the class $F|_{x_1^n}$ restricted to the sample. For instance, consider the following example:

	x_1	x_2	x_3	x_4	x_5
f_1	-	+	-	+	+
f_2	+	-	-	+	+
f_3	+	+	+	-	+
f_4	-	+	+	-	-
f_5	-	-	-	+	-

Each row is one possible evaluation of the functions in F on the fixed sample, and the cardinality of $F|_{x_1^n}$ equals to the number of rows. We transform the table by "shifting" columns.

Definition. shifting column i : for each row, replace a "+" in column i with a "-" unless it would produce a row that is already in the table.

After applying the shifting operation in order from x_1 to x_5 , we get the table($F|_{x_1^n}^*$):

	x_1	x_2	x_3	x_4	x_5
f_1	-	+	-	-	-
f_2	-	-	-	+	+
f_3	-	-	-	-	+
f_4	-	-	-	-	-
f_5	-	-	-	+	-

Observations:

- (1) Size of the table unchanged, because the rows in $F|_{x_1^n}^*$ are still distinct;
- (2) The table $F|_{x_1^n}^*$ exhibits "closed below" property, i.e., for each row containing a "+", replacing that "+" with a "-" produces another row in the table.
- (3) $d_{VC}(F|_{x_1^n}^*) \leq d_{VC}(F|_{x_1^n})$. To see this, consider the application of the shifting operation to a single column, and notice that if F^* (after shifting) shatters a subset of columns, then so does F (before shifting).

Therefore,

$$(3) \text{ and } (2) \Rightarrow F^* \text{ can not have more than } d \text{ "+"'s in a row. Hence, } \#\text{row of } F^* \leq \sum_{i=0}^d \binom{n}{i};$$

$$(1) \Rightarrow |F|_{x_1^n} \leq \sum_{i=1}^d \binom{n}{i}$$

Wait a minute ...

Is this combinatorial expression really $\mathcal{O}(m^d)$?

$$\begin{aligned} \sum_{i=0}^d \binom{m}{i} &\leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \quad \frac{m}{d} \geq 1 \\ &\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \quad \text{περιόριση θετικών όρων} \\ &= \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i (1+x)^m = \sum_{i=0}^m \binom{m}{i} x^i \\ &= \left(\frac{m}{d}\right)^d \left(1 + \frac{d}{m}\right)^m \leq \left(\frac{m}{d}\right)^d e^d \quad (1+x) \leq e^{-x} \end{aligned}$$

Generalization Bounds

Combining our previous generalization results with Sauer's lemma, we have that for a hypothesis class H with VC dimension d , for any $\delta > 0$ with probability at least $1 - \delta$, for any $h \in H$,

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{2d \log \frac{em}{d}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (43)$$

Whew!

- Infinite hypothesis class is PAC-learnable iff it has finite VC dimension
- We're now going to see if we can find an algorithm that has good VC dimension
- And works well in practice . . . Support Vector Machines