# Lecture 5: VC bound and margin bound

# Announcements

- Homework 1 due Thursday 23:59

- Moodle finally has forums to ask questions re HW or lecture (just realized yesterday)

- Project sign-ups Monday 14:00 - find your partner on moodle If you want to present a paper not on the list, please double check with us.

 Feedback compilation

- Good: interactivity, intuition

- can be improved: handwriting, references to some results that are not explicitly noted in MW (adding some from SS), more intuition before proof but also more proof details

# About project choice

1. Identify and motivate problem - why should I / the community care? Including literature review (done-ish)

2. "Detective hat": Intuitive (not just technical level) understanding of proof, assumptions, statement in depth

3. "Reviewer hat": Which relevant questions does it shed light on and does the paper answer/shed light on it? How significant is the addition of this paper compared to existing literature? This is a key step towards Step 4.

4. "Researcher hat": What are **interesting, impactful** follow-up questions they did not answer and would be interesting and perhaps feasible to pursue?

5. Break down the identified follow-up problem into feasible chunks (e.g. lemmas, experiments) and optionally show your attempts to tackle the first few steps.

# Outline for today

- VC bound and proof

- Rademacher contraction

- Interactive: Proof using the ramp loss and contraction (students)

# Recap: Massart's lemma

Recap: Last time, we bounded the Rademacher for function classes $\mathcal{F}$ that induce a finite set $\mathcal{H}(Z^n) = \{(\ell(Z_1; f), \dots, \ell(Z_n; f)) : f \in \mathcal{F}\}$ using Massart's lemma

> **Lemma (Massart, SS Lemma 26.8)**
>
> *For n points $Z^n := \{Z_1, \dots, Z_n\}$, let all $h : \mathcal{Z} \to \{0, 1\}$ and $\mathcal{H}(Z^n) := \{(h(Z_1), \dots, h(Z_n)) : h \in \mathcal{H}\}$ with cardinality $|\mathcal{H}(Z^n)|$.*
>
> $$\tilde{\mathcal{R}}_n(\mathcal{H}(Z^n)) := \mathbb{E}_\epsilon \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \epsilon_i h(Z_i) \le \sqrt{\frac{2 \log |\mathcal{H}(Z^n)|}{n}}$$

- $|\mathcal{H}(Z^n)|$ corresponds to # labelings for $Z^n$ induced by $\mathcal{H}$
- if $|\mathcal{H}(Z^n)|$ grows exponentially $\to \tilde{\mathcal{R}}_n(\mathcal{H}(Z^n)) = O(1)$

# VC bound

We now use Massart to prove a bound for function classes of finite VC dimension (i.e. where $|\mathcal{H}(Z^n)|$ does not grow exponentially in $n$ for any $Z^n$)

Recap **definition VC dimension** for binary classification:

> **Definition (VC dimension)**
>
> Biggest $n \in \mathbb{N}$ s.t. there exists $Z^n \in \mathcal{Z}^n$ with $\mathcal{H}(Z^n) = \{0, 1\}^n$

Function classes $\mathcal{F}$ with finite VC dimension can make $\mathcal{H}$ Glivenko-Cantelli, i.e. $\mathcal{R}_n(\mathcal{H}) = o(1)$. More specifically:

> **Theorem (uniform VC bound)**
>
> *If $\mathcal{H}$ has VC dimension $d_{VC}$, w/ prob $\ge 1 - \delta$ for any estimator $f \in \mathcal{F}$*
>
> $$\mathbb{P}(yf(X) < 0) \le \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(x_i) < 0} + 4\sqrt{\frac{d_{VC} \log(n+1)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

# Proof of VC bound 1

Now we first prove a high-probability upper bound for the population 0-1 loss $\ell((x, y); f) = \mathbb{1}_{yf(x)<0}$ for finite function classes $\mathcal{F}$.

Plugging in the definition of the loss, using the uniform law, we get

$$\mathbb{P}(Yf(X) < 0) \leq \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{y_i f(x_i)<0} + 4\sqrt{\frac{\log \mathcal{R}_n(\mathcal{H})}{n}} + c\sqrt{\frac{\log(1/\delta)}{n}} \tag{1}$$

for some universal constant $c$. The proof uses the uniform law (U.L.)

$$\begin{aligned}
R(f) - R_n(f) &= \mathbb{E}\ell((x, y); f) - \frac{1}{n} \sum_{i=1}^{n} \ell((x, y); f) \\
&= \mathbb{P}(yf(x) < 0) - \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{y_i f(x_i)<0} \\
&\leq \sup_{f \in \mathcal{F}} R(f) - R_n(f) \overset{U.L.}{\leq} 2\mathcal{R}_n(\mathcal{H}) + c\sqrt{\frac{\log(1/\delta)}{n}}
\end{aligned}$$

# Proof of VC bound 2

- Note that $\mathcal{R}_n(\mathcal{H}) = \mathbb{E}_{Z^n} \tilde{\mathcal{R}}_n(\mathcal{H}) \leq \sup_{Z^n} \tilde{\mathcal{R}}_n(\mathcal{H})$ (this is crude!)

- Further by Massart, $\sup_{Z^n} \tilde{\mathcal{R}}_n(\mathcal{H}) \leq \sup_{Z^n} \sqrt{\frac{2 \log |\mathcal{H}(Z^n)|}{n}}$ yielding

$$\mathcal{R}_n(\mathcal{H}) \leq \sqrt{\frac{2 \log \sup_{Z^n} |\mathcal{H}(Z^n)|}{n}} \tag{2}$$

(loose since distribution independent!)

Furthermore, we have the following upper bound on the size of $\mathcal{H}(Z^n)$

> **Lemma (Sauer-Shelah, MW Prop 4.18.)**
>
> *If $\mathcal{F}$ has VC dimension $d_{VC}$, then for any $Z_1, \ldots, Z_n$ we have growth function $N_{\mathcal{H}}(n) := \sup_{Z^n \in \mathcal{Z}^n} |\mathcal{H}(Z^n)| \leq (n+1)^{d_{VC}}$ for all $n \geq d_{VC}$.*

Plugging Sauer-Shelah into eq. 2, and that into eq. 1 in the uniform law to yield result

# Rademacher contraction

A useful property of Rademacher complexities (and Gaussian!) holds for losses $\ell : \mathbb{R}^n \to \mathbb{R}^n$ with $\ell(\theta) = (\ell_1(\theta_1), \ldots, \ell_n(\theta_n))$ with $L-$Lipschitz $\ell_j : \mathbb{R} \to \mathbb{R}$, i.e.

$$|\ell_j(a) - \ell_j(b)| \leq L|a - b| \text{ for all } a, b \in \mathbb{R}.$$

> **Lemma (Rademacher contraction, SS Lemma 26.9)**
>
> *For any $\mathbb{T} \subset \mathbb{R}^n$ and $\ell : \mathbb{R}^n \to \mathbb{R}^n$ with univariate L-Lipschitz functions it holds that*
> $$\tilde{\mathcal{R}}_n(\ell \circ \mathbb{T}) \leq L\tilde{\mathcal{R}}_n(\mathbb{T})$$

We will need Rademacher contraction to prove the margin bound theorem that we will now collectively prove.

# Skipped during lecture: Proof ingredients

Let $\epsilon$ be the vector of $n$ i.i.d. Rademacher r.v. and define the shorthand $\epsilon_{2:n} = (\epsilon_2, \ldots, \epsilon_n)$ and same for $\theta$.

The following holds for all $n$

- Key 1: de-symmetrize using the tower property: For any $g$ we have $\mathbb{E}_\epsilon g(\epsilon) = \mathbb{E}_{\epsilon_1}[\mathbb{E}[g(\epsilon)|\epsilon_1]] = \frac{1}{2}\mathbb{E}[g(\epsilon)|\epsilon_1 = 1] + \frac{1}{2}\mathbb{E}g(\epsilon)|\epsilon = -1]$

- Key 2: Lipschitz property $\ell_i(\theta_i) - \ell_i(\tilde{\theta}_i) \leq L|\theta_i - \tilde{\theta}_i|$ for all $i$

- Key 3: For each $\epsilon$ we can define $h(\theta_{2:n}) = \sum_{i=2}^n \epsilon_i \ell_i(\theta_i)$. One can prove via contradiction that

$$\sup_{\theta, \tilde{\theta} \in \mathbb{T}} |\theta_1 - \tilde{\theta}_1| + h(\theta_{2:n}) + h(\tilde{\theta}_{2:n}) = \sup_{\substack{\theta, \tilde{\theta} \in \mathbb{T} \\ \theta_1 \geq \tilde{\theta}_1}} \theta_1 - \tilde{\theta}_1 + h(\theta_{2:n}) + h(\tilde{\theta}_{2:n})$$

## Skipped during lecture: R.C. contraction proof

$$n\tilde{\mathcal{R}}_n(\ell \circ \mathbb{T}) = \mathbb{E}_\epsilon \sup_{\theta \in \mathbb{T}} \sum_{i=1}^n \epsilon_i \ell_i(\theta_i)$$

$$\stackrel{1.}{=} \frac{1}{2}\left[\mathbb{E}_{\epsilon_{2:n}} \sup_{\theta \in \mathbb{T}} \ell_1(\theta_1) + \sum_{i=2}^n \epsilon_i \ell_i(\theta_i) + \sup_{\tilde{\theta} \in \mathbb{T}} -\ell_1(\tilde{\theta}_1) + \sum_{i=2}^n \epsilon_i \ell_i(\tilde{\theta}_i)\right]$$

$$= \frac{1}{2}\left[\mathbb{E}_{\epsilon_{2:n}} \sup_{\theta, \tilde{\theta} \in \mathbb{T}} \ell_1(\theta_1) - \ell_1(\tilde{\theta}_1) + \sum_{i=2}^n \epsilon_i \ell_i(\theta_i) + \sum_{i=2}^n \epsilon_i \ell_i(\tilde{\theta}_i)\right]$$

$$\stackrel{2.}{\le} \frac{1}{2}\left[\mathbb{E}_{\epsilon_{2:n}} \sup_{\theta, \tilde{\theta} \in \mathbb{T}} L|\theta_1 - \tilde{\theta}_1| + \sum_{i=2}^n \epsilon_i \ell_i(\theta_i) + \sum_{i=2}^n \epsilon_i \ell_i(\tilde{\theta}_i)\right]$$

$$\stackrel{3.}{=} \frac{1}{2}\left[\mathbb{E}_{\epsilon_{2:n}} \sup_{\theta \in \mathbb{T}} L\theta_1 + \sum_{i=2}^n \epsilon_i \ell_i(\theta_i) + \sup_{\tilde{\theta} \in \mathbb{T}} (-L\tilde{\theta}_1) + \sum_{i=2}^n \epsilon_i \ell_i(\tilde{\theta}_i)\right]$$

$$\stackrel{1.}{=} \mathbb{E}_\epsilon \sup_{\theta \in \mathbb{T}} L\epsilon_1 \theta_1 + \sum_{i=2}^n \epsilon_i \ell_i(\theta_i)$$

Use the same argument for the RHS inductively on each coordinate. $\square$

## Mimicking proof-based research in collaboration

- Learning objectives: Both for actual guarantees and presentation, collaboration

  1. Get intuition why a problem / conjecture should be true

  2. Break down a proof to parts

  3. Prove individual parts

- Matching questions in the interactive session today

  1. Intuitively why should enforcing a large margin yield better generalization? Show graphically (no right or wrong)

  2. Given contraction inequality, ramp loss and Rademacher complexity for linear functions, prove the margin bound

  3. Prove Rademacher complexity for linear function class

# Instructions

- Groups:
  - We will divide the class into three groups of $\approx 4$ people each.
  - Each group will solve one of the three questions jointly.
  - Once you know your group, choose a representative to present later

- Group work:
  - 15 minutes of discussion to solve the question - if done early, feel free to solve another groups' question
  - Another 5 minutes to prepare the representative's blackboard presentation

- Final presentation
  - 30 minutes of 3 short presentations (7 min presentation, 3 min Q&A)
  - Introduce yourself and group members by names
  - Present your results.

# Primer on margins for linear classifiers

- Class of linear classifiers $\mathcal{F} = \{f : f(x) = w^\top x \; w \in \mathbb{R}^d\}$

- Intuition in introductory lectures for linearly separable data: large minimum distance to the boundary is good that can be computed as

$$d_{\min} = \min_i y_i \frac{w^\top x_i}{\|w\|_2}$$

where $\min_i y_i \langle w, x_i \rangle$ is called the margin

- Can obtain set of maximizing directions by solving

$$\max_{\gamma, w} \gamma \text{ s.t. } y_i \langle \frac{w}{\|w\|_2}, x_i \rangle \geq \gamma$$

which for bounded $\|w\|_2 \leq B$ is the same as solving

$$\max_{\gamma', \|w\|_2 \leq B} \gamma' \text{ s.t. } y_i \langle w, x_i \rangle \geq \gamma'$$

- We will look the generalization performance of feasible $w$ with $\|w\|_2 \leq B$ which achieve a margin of at least some $\gamma$

# Margin bound for binary classification

Key ingredient of proof (in interactive session)

> **Definition (ramp loss)**
>
> The ramp loss $\ell_\gamma$ is defined as
> $$\ell_\gamma(u) = \begin{cases} 1 & u \in (-\infty, 0) \\ 1 - \frac{u}{\gamma} & u \in [0, \gamma] \\ 0 & u \in (\gamma, \infty) \end{cases}$$
> and $\frac{1}{\gamma}$-Lipschitz.

# Margin bound for linear classifiers

Definitions

- Set of linear functions $\mathcal{F}_B = \{f(x) = \langle w, x \rangle : \|w\|_2 \leq B\}$
- Define the risk $R_n^\gamma(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(x_i) \leq \gamma}$ and
  $R^\gamma(f) = \mathbb{E}_{X,Y} \mathbb{1}_{Yf(X) \leq \gamma}$

Assumption (A): Boundedness of covariates $\mathbb{P}(\|x\|_2 \leq D) = 1$

> **Theorem (margin bound)**
>
> *If the assumptions are valid for any fixed $\gamma$, w/ prob. at least $1 - \delta$, for any $f \in \mathcal{F}_B$ we have*
> $$R^0(f) = \mathbb{P}[y \neq sign(f(x))] \leq R_n^\gamma(f) + \frac{2DB}{\gamma\sqrt{n}} + c\sqrt{\frac{\log(1/\delta)}{n}}$$
> *for some constant $c > 0$.*

# References

- Massart, Rademacher for classification: Shalev-Schwartz Chapter 26