

Lec 24-25

Concentration "generalization" "crash course"

Why?

Concentration:

well-behaved functions  
of r.v.'s behave  
as some simple statistics.

best reference:  
lec notes  
Pamou van Handel  
APCS50

$R(\hat{f})$   
↑  
alg

E.g.

Lipschitz function, gaussian r.v.'s

→ mean

Lipschitz function  
+ convex ← false without.

bounded r.v.'s

→ mean

→ max (Gaussian matrix) concentrates

Question:

\* Classic  
(always false)

$\hat{R}(\hat{f})$

7.7K

h course

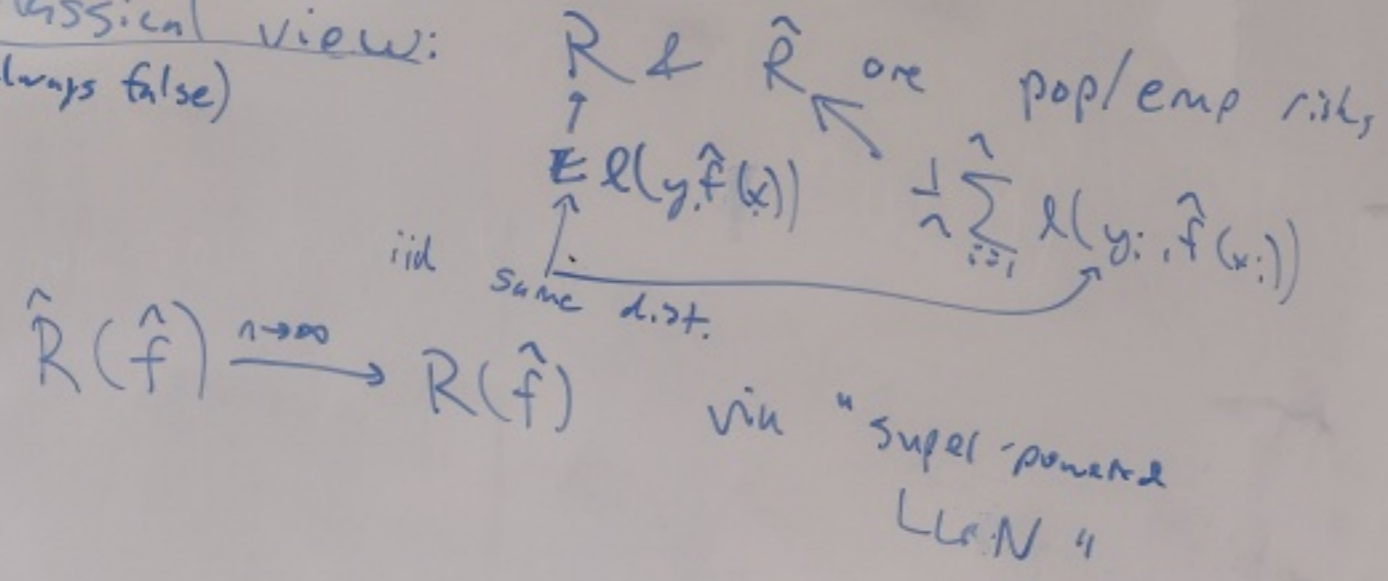
Why?

$$R(\hat{f}) - R(g) = \boxed{R(\hat{f}) - \hat{R}(\hat{f})} + \hat{R}(\hat{f}) - \hat{R}(f) \checkmark \text{opt} + \hat{R}(f) - R(f) \text{reference in } \mathcal{F} + R(f) - R(g) \checkmark \text{apx}$$

↑ alg

Question: why  $R(\hat{f}) - \hat{R}(\hat{f})$  small?

\* Classical view: (always false)



\* Modern perspective

\* Interpolation

motivation: DL e

[all prominent s unclear if current

\* OOD (out

test dist  $\neq$  train (e.g., alphafold)

\* Classical view is

(e.g., "rethinking g

\* Transfer learning

$(\hat{f})$  ?  
 $f)$  ✓ opt  
 reference in  
 $f)$  ?  $\mathcal{F}$   
 $)$  ✓ apx

nsll?  
 pop/emp risks  
 $l(y_i; \hat{f}(w_i))$   
 "penalized  
 $N$

\* Modern perspectives:

\* Interpolation  $\hat{R}(\hat{f}) \approx 0 \ll R(\hat{f}) \approx \inf_y R(y)$   
 motivation: DL experiments (Srebro et al '14, Zhang '16)

[all prominent stat/theorists working on this;  
 unclear if currently true.]

\* OOD (out-of-domain/distribution) } Risteski, Zhao  
 test dist  $\neq$  train dist  
 (e.g., alphafold)

\* Classical view is broken  
 (e.g., "rethinking generalization" (Zhang '16))

\* Transfer learning } Raghunathan

Classical view  
 General form:

$R(f)$

Rem. \* We'll do  
 \*  $\frac{1}{n}$  no

\* "Uniform deviation"

$R(g)$

Classical view

General form:

w/  $pr \geq 1 - \delta$ ,

$\forall f \in \mathcal{F}$

"uniform term"  
modern/popular view:  
this is an error

$$R(f) \leq \hat{R}(f) + \sqrt{\frac{1-\delta}{n}} + \sqrt{\frac{\text{comp}(\mathcal{F})}{n}}$$

Rem. \* We'll define  $\text{comp}(\mathcal{F})$  later

\*  $\frac{1}{n}$  not  $\frac{1}{\sqrt{n}}$  is possible, not really relevant to us.

\* "Uniform deviation" because "uniform if  $\mathcal{F}$ " (we'll explain later).

Theorem (Hoeffding)

Suppose  $(Z_1, \dots, Z_n)$  ind,  $Z_i \in [a_i, b_i]$  a.s.

Then

$$P\left[\frac{1}{n} \sum (Z_i - \mathbb{E}Z_i) \geq \varepsilon\right] \leq \exp\left(-\frac{2n^2 \varepsilon^2}{\sum (b_i - a_i)^2}\right) (\delta)$$

alternatively, w/  $p \geq 1 - \delta$ ,

$$\frac{1}{n} \sum \mathbb{E}Z_i \leq \frac{1}{n} \sum Z_i + \sqrt{\frac{\sum (a_i - b_i)^2}{2n^2} \ln \frac{1}{\delta}}$$

Plan: (a) use it for  $\hat{R}(f) - R(f)$   
(b) sketch proof

$R(\hat{F}) - R(F)$  ①

Given fixed  $f$ ,

data  $(x_i, y_i)_{i=1}^n$   
iid

define

$$Z_i := \ell(y_i, f(x_i)) \underbrace{e^{[a, b]}}_{\text{assume}}$$

$\Rightarrow$  w/  $p \geq 1 - \delta$ ,

$$R(f) = \frac{1}{n} \sum \mathbb{E}Z_i \leq \frac{1}{n} \sum Z_i + \sqrt{\frac{n(b-a)^2}{2n^2} \ln \frac{1}{\delta}}$$

$$= \hat{R}(f) + (b-a) \sqrt{\frac{\ln \frac{1}{\delta}}{2n}}$$

$$\exp\left(-\frac{2n^2 \epsilon}{\sum_i (b_i - a_i)^2}\right) (s)$$

$$Z_i + \sqrt{\frac{\sum_i (a_i - b_i)^2}{2n^2}} \ln \frac{1}{\delta}$$

$(x_i, y_i)_{i=1}^n$  assume  $x_i \in [a, b]$

$$\sum_i Z_i + \sqrt{\frac{n(b-a)^2}{2n^2}} \ln \frac{1}{\delta}$$

$$a) \sqrt{\frac{\ln 1/\delta}{2n}}$$

② Given  $\mathcal{F}$  with  $|\mathcal{F}| < \infty$ ,  
define  $\delta_f := \frac{\delta}{|\mathcal{F}|}$

for any test  $w$  / pr  $\geq 1 - \delta_f$ ,  $R(f) \leq \hat{R}(f) + (b-a) \sqrt{\frac{\ln |\mathcal{F}| + \ln 1/\delta_f}{2n}}$

$\Rightarrow$  via union bound

$$w/ \text{pr } 1 - \delta = 1 - \sum_{f \in \mathcal{F}} \delta_f, \forall f \in \mathcal{F}$$

$$R(f) \leq \hat{R}(f) + (b-a) \sqrt{\frac{\ln |\mathcal{F}| + \ln 1/\delta}{2n}}$$

Remark \*  $\ln |\mathcal{F}|$  was our  $\text{comp}(\mathcal{F})$

\* Why did need union bound?  
(Why not apply Hoeffding to output alg?)  
Alg output depends data  $\Rightarrow (Z_1, \dots, Z_n)$  might be dependent

$$\hat{f}(x) = \begin{cases} y_i & \text{if } x = x_i \\ 0 & \text{o.w.} \end{cases}$$

Can construct joint pr on  $X \times Y$   
s.t.  $\hat{R}(f) \xrightarrow{n \rightarrow \infty} R(f)$ , but also  $\hat{R}(f) = 0 \neq f = R(f)$ .

\* Fixes:

- (a) use cross-validation
- (b) alg-specific anal
- (c) Brute force / uniform control all poss.

How to prove Hoeffding

Definition. rv  $Z$  is  $\sigma^+$ -  
(sub-Gaussian with var  $\sigma^2$ )

if  $\forall t \in \mathbb{R}, \mathbb{E} \exp(\dots)$

Rem. \*  $\exists 4$  equivalent definitions  
(vershynin ndp)

\*  $\exists$  other examples of " (5)

$$\frac{\ln \frac{1}{\delta}}{2n}$$

\* Fixes

- (a) use cross-validation error
- (b) algo-specific analysis (typically one pass)
- (c) Brute force / uniform deviation: control all possible outputs of alg.

How to prove Hoeffding

Definition. rv  $Z$  is  $\sigma^2$ -sub-Gaussian  
 (sub-Gaussian with variance proxy  $\sigma^2$ )  
 if  $\forall t \in \mathbb{R}$ ,  $\mathbb{E} \exp(t(X - \mathbb{E}X)) \leq \exp\left(\frac{t^2 \sigma^2}{2}\right)$

Rem. \*  $\exists 4$  equivalent definitions  
 (Vershynin hdp book)

\*  $\exists$  other examples of "well concentrated r.v.s"  
 (sub-exponential)

$$\frac{\ln \frac{1}{\delta} + \ln \frac{1}{\delta}}{2n}$$

at be dependent  
 joint pr on  $X \times Y$   
 output  $\hat{R}(f) = 0 \neq R(f)$ .

Examples.

- \*  $\mathcal{N}(\mu, \sigma^2)$
- \*  $X \in [a, b]$  as r.v.  
 One subcase:  $\mathbb{E} \exp(tX) = \frac{1}{2} \exp(at) + \frac{1}{2} \exp(bt)$

\* If  $(X_1, \dots, X_n)$  i.i.d.  
 $\Rightarrow \frac{1}{n} \sum X_i$  is  $\frac{\sigma^2}{n}$ -sub-Gaussian

## Examples.

\*  $\mathcal{N}(\mu, \sigma^2)$  is  $\sigma^2$ -sub-Gaussian (complete the square)

\*  $X \in [a, b]$  a.s., then  $X$  is  $\left(\frac{b-a}{2}\right)^2$ -sub-Gaussian.  
One subcase: if  $X \sim \text{Unif}([-1, 1])$

$$\begin{aligned} \mathbb{E} \exp(tX) &= \frac{1}{2} [\exp(-t) + \exp(t)] \\ &= \frac{1}{2} \left[ \sum_{i \geq 0} \frac{(-t)^i}{i!} + \sum_{i \geq 0} \frac{(t)^i}{i!} \right] = \frac{1}{2} \cdot 2 \cdot \sum_{i \geq 0} \frac{t^{2i}}{(2i)!} \\ &= \sum_{i \geq 0} \frac{(t)^{2i}}{2^i i!} = \sum_{i \geq 0} \frac{(t^2/2)^i}{i!} \end{aligned}$$

\* If  $(X_1, \dots, X_n)$  are  $\sigma_i^2$ -sub-Gaussian & independent

$$\Rightarrow \frac{1}{n} \sum X_i \text{ is } \boxed{\frac{\sum \sigma_i^2}{n^2}} \text{-sub-Gaussian}$$

behaves like  $\frac{1}{n}$ ;

"Concentration phenomenon"

$$= \exp(t^2/2).$$

$$\exp\left(\frac{t^2 \sigma^2}{2}\right)$$

abd rivis"