

last week: least-squares

$$L(y, \beta(\gamma)) = \|y - \varphi(\gamma)^T \alpha\|^2 \quad \varphi(\gamma) \in \mathbb{R}^d$$

two results : excess expected risk $\frac{\sigma^2 d}{n}$ (OLS)

$$\text{model } y = \varphi(\gamma)^T \alpha + \varepsilon$$

$$\begin{aligned}\mathbb{E}\varepsilon &= 0 \\ \mathbb{E}\varepsilon^2 &= \sigma^2\end{aligned}$$

ridge regression

$$\frac{D}{\sqrt{n}}$$

generic set-up : Expected risk

$$R(\ell) = \mathbb{E}(e(y, \beta(\gamma))).$$

find bands on excess(expected) risk : $R(\ell) - R^*$

Binary classification: $\mathcal{Y} = \{-1, 1\} \subset \mathbb{R}$

Another convention $\{0, 1\}$ or $\{1, 2\}$
 \Rightarrow see chapter 13

f prediction function: $f: X \rightarrow \mathcal{Y} = \{-1, 1\}$

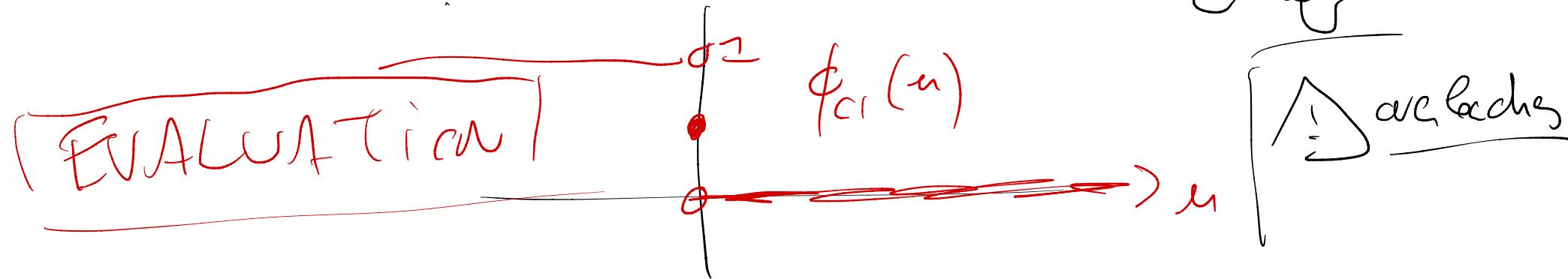
(Vapnik - Chervonenkis dimension)

Learn $g: X \rightarrow \mathbb{R}$ and define $f(x) = \text{sign}(g(x))$

$$\text{sign}(a) = \begin{cases} +1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \\ ? & \text{if } a = 0 \end{cases}$$

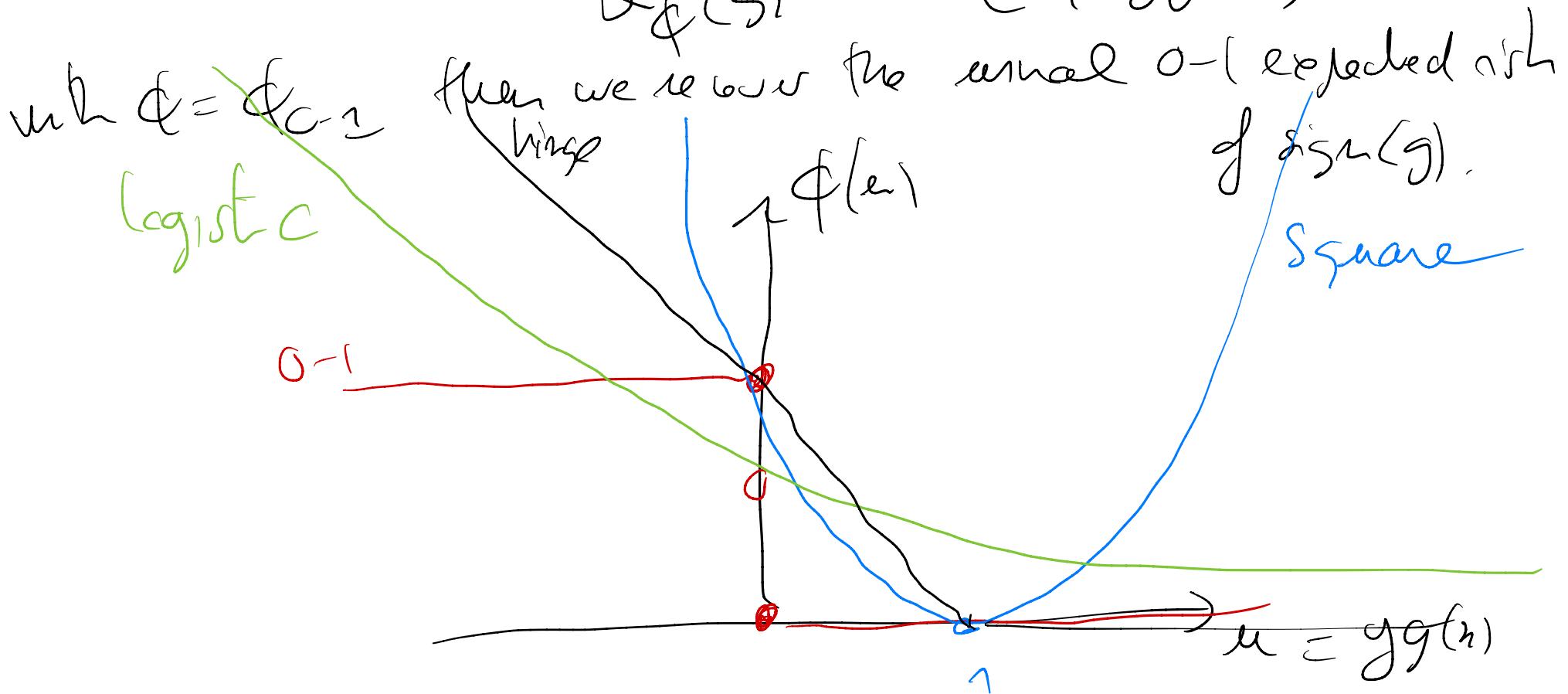
Question: What is the risk of f as a function of g ?

$$\begin{aligned}
 L(g) &= R(\text{sign} \circ g) = \mathbb{E}\left[\frac{1}{\text{sign}(g(x)) - 1}\right] \\
 u \in \mathbb{R} &\quad \left[\begin{array}{l} \text{values in } \{-1, 1\} \\ 0-1 loss \end{array} \right] \\
 &= P(\text{sign}(g(x)) \neq y) \\
 &= \mathbb{E}\left[\frac{1}{g(x) \neq y} \frac{1}{g(x) y < 0}\right] + \frac{1}{2} g(x) = \text{margin} \\
 &\quad \left[\begin{array}{l} \text{margin} \\ \text{margin} \end{array} \right] \\
 &= \mathbb{E}\left[\frac{1}{g(x) y < 0}\right] + \mathbb{E}\left[\frac{1}{g(x) = 0}\right] \\
 &= \Phi_{0-1}(y g(x)) \quad \text{with } \Phi_{0-1}(u) = \begin{cases} 1/2 & u < 0 \\ 1/2 & u = 0 \\ 0 & u > 0 \end{cases}
 \end{aligned}$$



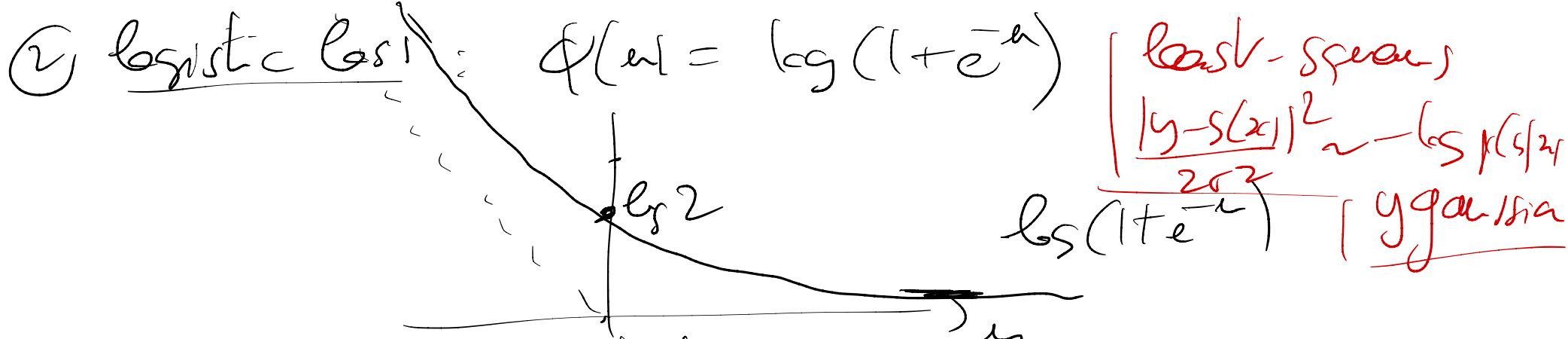
Convex surrogates: ϕ -risk.

$$\mathcal{R}_\phi(s) = \mathbb{E}(\phi(yg(z)))$$



① quadratic loss: $\phi(u) = (u-1)^2$

$$\phi(yg(z)) = (yg(z)-1)^2 = (g(z)-y)^2$$



probabilistic interpretation:

define a model $p(y=1|x) = \text{sigmoid}(g(x))$

$\text{Sigmoid} = \sigma$

$$= \frac{1}{1 + e^{-g(x)}}$$

$$\sigma(-u) = 1 - \sigma(u)$$

$$\frac{1}{1 + e^u} = 1 - \frac{1}{1 + e^{-u}}$$

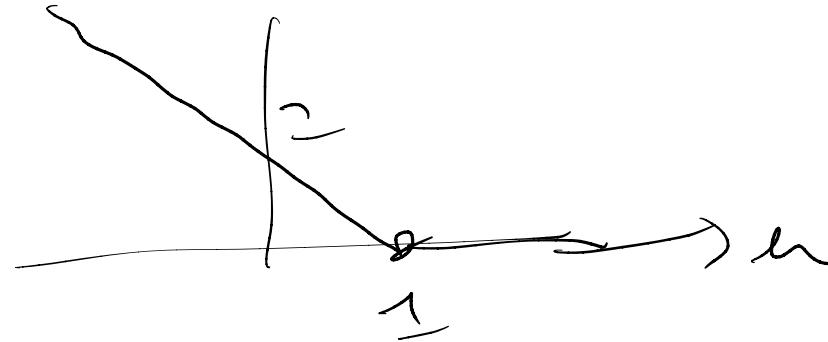
$$p(y|x) = \frac{1}{1 + e^{-yg(x)}}$$

$$-\log p(y|x) = \log(1 + e^{-yg(x)}) = \text{logistic loss}$$

MAXIMUM LIKELIHOOD / CROSS-ENTROPY LOSS

③ hinge loss (Support vector machine (SVM))

$$\varphi(u) = \max\{0, 1-u\} = ((-u) +$$



Convex
not
differentiable

if \hat{s} is the minimizer of $R_\varphi(s) = \mathbb{E}[\varphi(y s)]$

if $\text{sign } s^* = \hat{g}^*$?

Ex: least squares = $R_\varphi(s) = \mathbb{E}((y - g(x))^2)$

$$f(x) = \underbrace{\mathbb{E}_y P(y=1|x)}_{= -\mathbb{E}_y P(y=1|x)} \quad \text{if } P(y=1|x) > \frac{1}{2}$$

$$g(x) = \mathbb{E}[y|x]$$

$$\text{sign}(g^*(x)) = \text{sign}(\mathbb{E}[y|x])$$

$$\mathbb{E}[y|x] = \frac{-\mathbb{P}(y=1|x)}{-\mathbb{P}(y=1|x)} = \frac{2\mathbb{P}(y=1|x)}{-1} = g^*$$

All bands will be of the form - $R_{\Phi}(S) - R_{\Phi}^* \leq \frac{2}{\sqrt{n}}$

$$R_{\Phi_{0-1}}(S) - R_{\Phi_{0-1}}^* \leq \text{f} \left(R_{\Phi}(S) - R_{\Phi}^* \right)$$

Calibration function

$f(u) = \sqrt{u}$ for quadratic

$f(u) = u$ for logistic

$$R(f) = E(\ell(y, f(x))) = E(y - f(x))^2 \text{ regression}$$
$$= E \phi(y - f(x)) \text{ classifier}$$

$$f: X \rightarrow \mathbb{R}$$



Training data: $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ iid from distribution
 Empirical risk $\hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$ $i=1, \dots, n$ $p(\beta)$
 Expected risk $R(\beta) = \mathbb{E}_{\mathcal{D}}[\ell(y, f(x))]$. $\hat{\beta}$ \approx arg min $\hat{R}(\beta)$ approximate
 Empirical risk minimization: $\hat{\beta} \in \arg \min_{\beta} \hat{R}(\beta)$
 fEF function
 dev

Decomposition

$$\hat{R}(\hat{\beta}) - R^* = \hat{R}(\hat{\beta}) - \inf_{g \in \mathcal{G}} R(g) + \inf_{g \in \mathcal{G}} R(g) - R^*$$
 Estimata eric

Approximate error
 decreasing in
 model size

Approximation error: $\bar{F} = \int f(\theta, \theta \in \Theta)$



ex: $f_\theta(x) = \phi(x)^T \theta$ linear model
neural network

$$\inf_{\theta \in \Theta} R(f_\theta) - R^* = \inf_{\theta \in \Theta} R(f_\theta) - \inf_{\theta \in \Theta} R(\theta_\theta) + \inf_{\theta \in \Theta} R(\theta_\theta) - R^*$$

Assumpta: $\inf_{\theta \in \Theta} R(\theta_\theta) = R(f_{\theta_\theta})$

$$\boxed{f \neq f^*}$$

"incomparable"
sym. env

$$\begin{aligned} R(f_\theta) - R(\theta_\theta) &= E \left(\ell(y, f_\theta(z)) - \ell(y, \theta_\theta(z)) \right) \\ &\leq E \left(G[\theta_\theta(z) - \theta_\theta(x)] \right) \end{aligned}$$

Assumptions: loss is G -Lipschitz continuous wrt to second variable

$$\inf_{\theta \in \Theta} R(b_0) - R^* = \inf_{\theta \in \Theta} R(b_0) - \inf_{\theta' \in \mathbb{R}^d} R(\theta_0) + \inf_{\theta' \in \mathbb{R}^d} R(\theta_0)$$

Assumpta: $\inf_{\theta' \in \mathbb{R}^d} R(\theta_0) = R(\theta_{0*})$

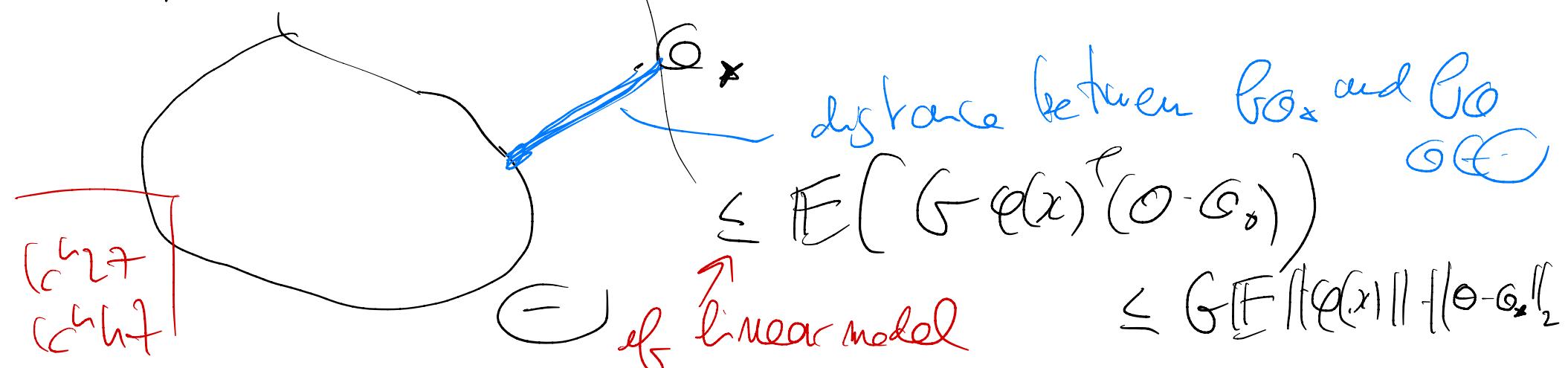
$$\boxed{f_{\theta^*}}$$

"unbiased & symmetric"

$$\inf_{\theta \in \Theta} R(b_0) - R(b_{0*}) = \mathbb{E}(e(g, b_0(x)) - e(g, b_{0*}(x)))$$

$$\leq \mathbb{E}(G(b_0(x) - b_{0*}(x)))$$

Assumptas: loss is G -Lipschitz continuous wrt to second variable



Estimator error:

$$R(\hat{f}) - \inf_{\text{fEF}} R_G = R(\hat{G}) - R(G_F) \quad \text{if optimal risk obtained}$$

$$\begin{aligned} &= R(\hat{f}) - \hat{R}(\hat{G}) + \hat{R}(G) - \hat{R}(G_F) + \underline{\hat{R}(G_F) - R(G_F)} \\ &\quad / \end{aligned}$$

$$\sum_i \mathbb{E} \ell(g_i, \hat{g}(x_i))$$

$O(\frac{1}{\sqrt{n}})$

$$\leq \sup_{\text{fEF}} R(f) - \hat{R}(f) + \underbrace{\hat{R}(\hat{G}) - \inf_{\text{SEF}} \hat{R}(s)}_{\text{optimal error}} + \sup_{\text{fFP}} \hat{R}(d) - R(d)$$

optimal error

\leftarrow optimization error

$$\sup_G R(G) - \hat{R}(G) + \sup_f \hat{R}(f) - R(f)$$

uniform deviation

Goal: $\mathbb{E} D \leq \alpha$, or $D \leq \alpha(\delta)$ with proba. 1 - δ

Concentrate inequalities : Hoeffding & McDiarmid

Hoeffding's inequality : $Z_i \in [c, c']$ independent

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n Z_i - \frac{1}{n} \sum_{i=1}^n E(Z_i) \\ = \frac{1}{n} \sum_{i=1}^n (Z_i - E(Z_i)) \end{aligned}$$

$$\sim \frac{1}{\sqrt{n}}$$

Central limit theorem

$$P\left(\frac{1}{n} \sum_{i=1}^n Z_i - \frac{1}{n} \sum_{i=1}^n E(Z_i) \geq t\right) \leq e^{-2nt^2} = \delta$$

$$t = \frac{1}{2n} \log \frac{1}{\delta}$$

(\Leftarrow) with probability greater than $1-\delta$,

$$\frac{1}{n} \sum_{i=1}^n Z_i - \frac{1}{n} \sum_{i=1}^n E(Z_i) \leq \sqrt{\frac{1}{2n} (\delta \frac{1}{\delta})}$$

uniform deviates : $\sup_{\text{GEF}} \hat{R}(e) - R(e)$, f finite

$$\Pr_{\substack{\text{S}_{\text{GF}} \\ \text{GEF}}} (\hat{R}(e) - R(e) > \epsilon) = \Pr_{\substack{\text{U} \\ \text{GEF}}} (\sum_{\text{GEF}} \Pr(\hat{R}(e) - R(e) > \epsilon))$$

$\leq \sum_{\text{GEF}} \Pr(\hat{R}(e) - R(e) > \epsilon)$

union bound \rightarrow

- $2n\epsilon^2/\delta_x^2$

$$e(q, p_m) \in [c, \epsilon_2]$$

$\leq \sum_{\text{GEF}} e$

flipping \rightarrow

= $|F| e^{-2n\epsilon^2/\delta_x^2}$

This implies that with prob $1 - \delta$

$$\Pr_{\substack{\text{S}_{\text{GF}} \\ \text{GEF}}} \hat{R}(e) - R(e) \leq \frac{\delta_x}{\sqrt{2n}} \sqrt{\log \frac{1}{\delta} + \log |F|}$$

+ expected : $\sqrt{\frac{\log |F|}{n}}$

$$\epsilon_{\text{eff}} = \begin{cases} \delta_x^2 & \delta_x \leq \delta \\ \frac{\delta_x^2}{2n} \log |F| & \delta_x > \delta \end{cases}$$

Rademacher average complexity:

$$f = \{ (x, s) \mapsto \epsilon_x f(x) \}$$

F set of functions from $X \rightarrow \mathbb{R}$
less function Lipschitz continuous
(wrt to 2nd arg.)

Goal: find $\sup_{f \in F} R(f) - \tilde{R}(f) = \sup_{f \in F} E\left(\epsilon_x f(x)\right) - \frac{1}{m} \sum_{i=1}^m \epsilon_i f(x_i)$

$= \sup_{h \in \mathcal{H}} E(h(s)) - \frac{1}{m} \sum_{i=1}^m h(s_i)$

Def: Rademacher average

$$R_m(f) = E_{\delta_1, \delta_2, \dots, \delta_m} \frac{1}{m} \sum_{i=1}^m \epsilon_i h(s_i)$$

data $\leftarrow h \in \mathcal{H}$

$$\delta = (x, s)$$

where $\epsilon_i \in \{-1, 1\}$ are independent
and zero mean, $P(\epsilon_i = 1) = P(\epsilon_i = -1) = \frac{1}{2}$
"Rademacher" random variables

$$\underline{\text{Score}} \quad \overline{H}(f(\text{uniform data})) \leq 2G \underset{\geq}{\underline{\text{R}_n(f)}}$$

$$\text{Symmetrization lemma: } \mathbb{E} \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(x_i) - \mathbb{E}[h] \right) \leq 2R_u(\mathcal{H})$$

$$\begin{aligned} \mathbb{E} \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(x_i) - \mathbb{E}[h] \right) &= \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \left(h(x_i) - \mathbb{E}[h(x_i)] \right) \\ &= \mathbb{E} \sup_{h \in \mathcal{H}} \cdot \mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n \left(h(x_i) - \mathbb{E}[h(x_i)] \right) \middle| \mathcal{D} \right). \\ &\leq \mathbb{E} \left(\mathbb{E} \sup_{h \in \mathcal{H}} \left(\frac{1}{n} \sum_{i=1}^n \left(h(x_i) - \mathbb{E}[h(x_i)] \right) \middle| \mathcal{D} \right) \right) \\ &\quad \downarrow \\ \text{Assume an independent copy} \\ \text{of data } \mathcal{D}' = (x'_1, \dots, x'_n) \\ \text{with same distribution as } \mathcal{D} \\ \mathbb{E}(h(x)) = \mathbb{E}(h(x_i)) = \mathbb{E}(h(x'_i) | \mathcal{D}) \\ \mathbb{E}(h(x)) = \mathbb{E}(h(x_n)) = \mathbb{E}(h(x'_n) | \mathcal{D}) \end{aligned}$$

$$= \mathbb{E}_{\mathcal{D}, \mathcal{D}'} \left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \left(h(x_i) - \mathbb{E}[h(x_i)] \right) \right]$$

by symmetry

$$\begin{aligned} &\leq \mathbb{E}_{\mathcal{D}, \mathcal{D}', \varepsilon} \sup_{h \in \mathcal{H}} \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i h(x_i) \right) + \mathbb{E}_{\mathcal{D}, \mathcal{D}', \varepsilon} \sup_{h \in \mathcal{H}} \left(\frac{1}{n} \sum_{i=1}^n \varepsilon'_i h(x'_i) \right) = 2R_u(\mathcal{H}) \end{aligned}$$

$$R_m(F) = \mathbb{E}_{\substack{\text{data} \in D, \epsilon \\ \text{left}}} \sum_{i=1}^m \varepsilon_i h(\beta_i) = \mathbb{E}_{D, \epsilon} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^m \varepsilon_i \ell(y_i, f(x))$$

G is Lipschitz
constant of f(x)

$$\leq G \cdot \mathbb{E}_{D, \epsilon} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^m \varepsilon_i f(x)$$

CONTRACTIVE PRINCIPLE

$$\Rightarrow \mathbb{E}_{\substack{\text{uniform deviation} \\ f \in F}} \leq 2G R_m(F)$$

(Ledoux + Talagrand)

$\mathbb{E}_{\mathcal{D}, \mathcal{E}} \sum_{i=1}^n \varepsilon_i \varphi(x)$ for linear prediction
 $F = \{u \mapsto \varphi(u), \|u\|_2 \leq 1\}$

$\|\varphi(x)\|_2 \leq R$ almost
 $\varphi(x) \in \mathbb{R}^d$ surely

R , R , R
 r, r
 Redemecker and on the data

$$\begin{aligned}
 &= \mathbb{E}_{\mathcal{D}, \mathcal{E}} \max_{\|\varphi\|_2 \leq D} \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i \varphi(x_i) \right)^T \mathbb{E}_{\mathcal{D}, \mathcal{E}} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \varphi(x_i) \right\|_2 \\
 &\leq D \sqrt{\mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \varphi(x_i) \right\|_2^2} = D \sqrt{\frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \varepsilon_i^2 \|\varphi(x_i)\|^2_2} \leq R^2
 \end{aligned}$$

$$\begin{aligned}
 &\leq DR / \sqrt{n} \\
 &\text{Sengal's inequality } \mathbb{E} \|u\|_2 \leq \sqrt{\mathbb{E} \|u\|_2^2}
 \end{aligned}$$

Sengal's inequality $\mathbb{E} \|u\|_2 \leq \sqrt{\mathbb{E} \|u\|_2^2}$

Put everything together
 $\hat{\theta} = \text{Minimizer of } \hat{R}(f_{\theta}) \text{ an } \| \theta \|_2 \leq D$
(linear model)

$$E \hat{L}(f_{\hat{\theta}}) - \inf_{\| \theta \|_2 \leq D} R(f_{\theta}) \leq \frac{4GRD}{\sqrt{n}}$$

$$\text{if } \| \varphi(z) \|_2 \leq R$$