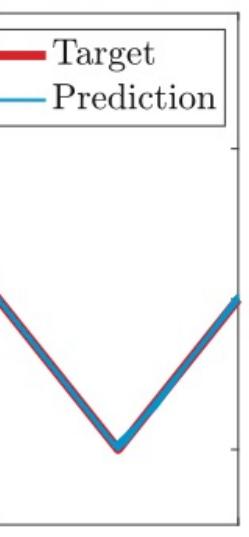
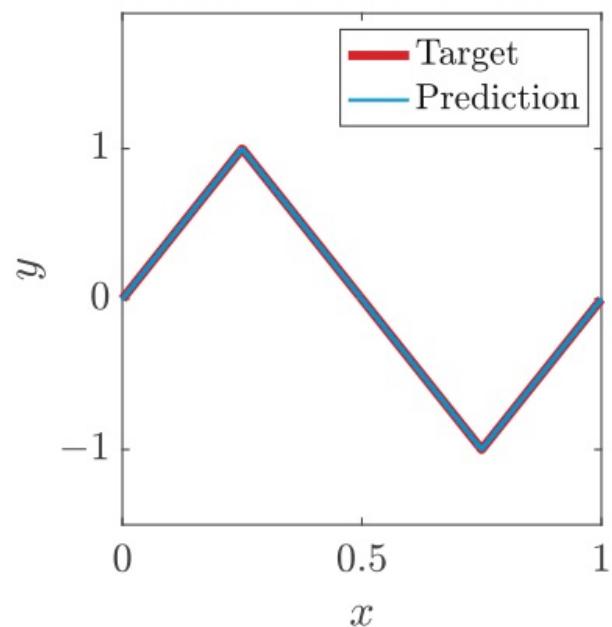


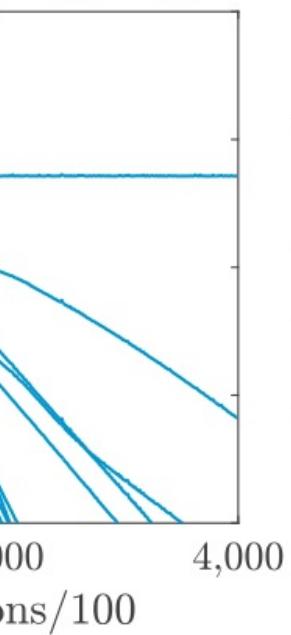
tions - $m = 20$



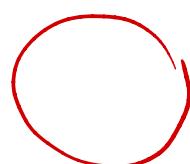
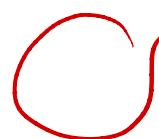
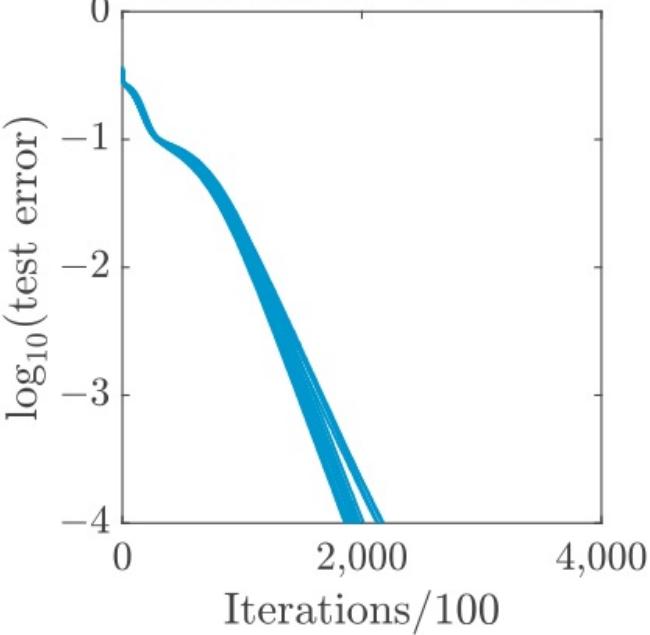
Prediction functions - $m = 100$



s - $m = 20$



Test errors - $m = 100$



Neural Networks

Empirical risk minimization:

$$\min_{f \in F} \hat{R}(f) = \sum_{i=1}^m \ell(y_i, f(x_i)) + R(f)$$

losses: ① $\ell(g, \bar{g}) = \frac{1}{2}(g - \bar{g})^2$

② $\ell(g, \bar{g}) = \frac{1}{2}|g - \bar{g}|$
for classif.

convex surrogates: logistic loss

$$\ell(g, f(x_i)) = \log(1 + e^{-y f(x_i)}) \quad y \in \{-1, 1\}$$

function spaces: ① linear functions, $f(x) = \Theta^\top \varphi(x) \rightarrow$ feature vector
↓ parameters

ℓ_2 -penalty for dimension independent bands

ℓ_1 -penalty for model select-a

② kernel methods: $f(x) = \langle \Theta, \varphi(x) \rangle \quad \varphi(x) \in \text{the feature space}$

Made possible using $\langle \varphi(x), \varphi(z) \rangle = k(x, z)$.

($\varphi(x)$ fixed)

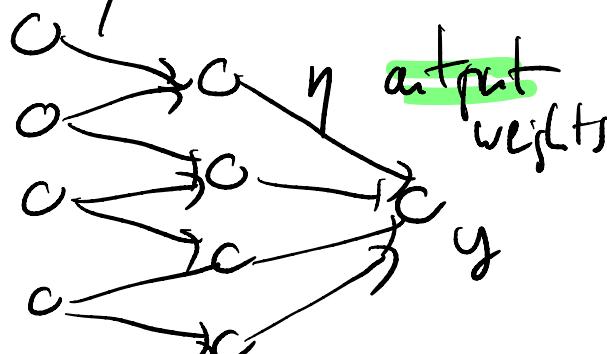
Single hidden layer neural network "Shallow"

$$x \in \mathbb{R}^d$$

$$f(x) = \sum_{j=1}^m b_j \sigma(w_j^T x + b_j)$$

↳ activation function

w, b input weights



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

sigmoid

$$\sigma(x) = \max\{0, x\}$$



↳ hidden neurons $\sigma(w_j^T x + b_j)$

x activation function at the output level.

- No activation function

$$f(x) = q^T \varphi(x) \text{ where } \varphi(u) = \sigma(w_j^T u + b_j)$$

- cross-entropy loss (\Rightarrow maximum likelihood)

$$g(x) = \sigma(f(x)) = \frac{1}{(1+e^{-f(x)})} \Rightarrow P(y=1|x) = \frac{\log g(x)}{\log(1-g(x))}$$

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

$$\begin{aligned} -\log g(x) - \log(1-g(x)) \\ = -\log \left(\frac{1}{1+e^{-g(x)}} \right) \\ = \log(1+e^{-g(x)}) \end{aligned}$$

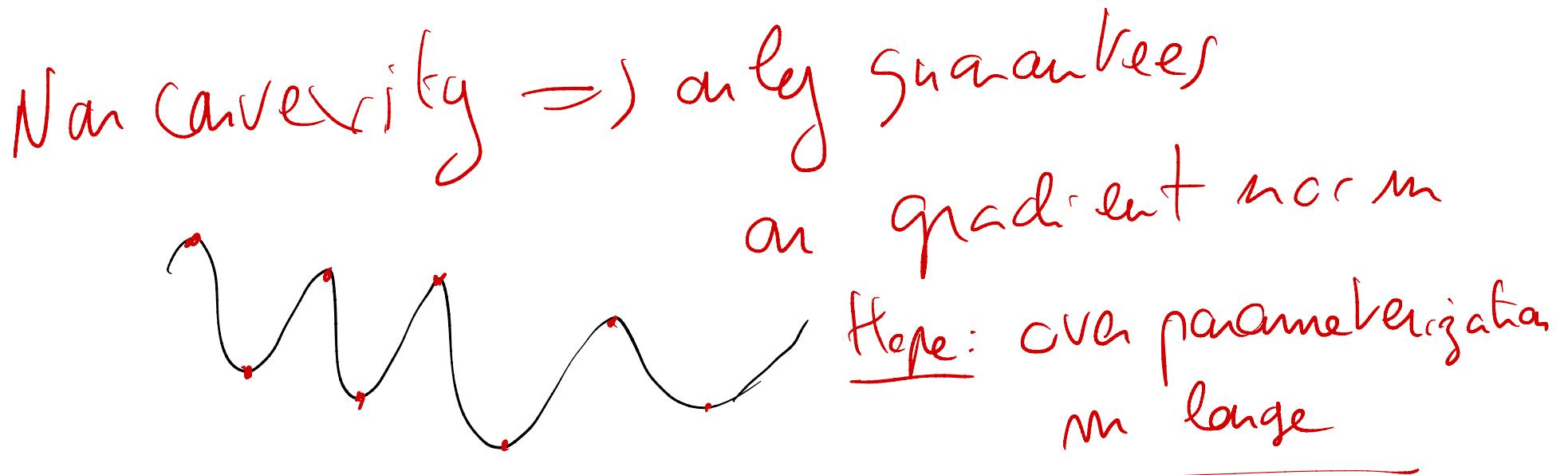
Optimization: goal: $\min_Q \mathbb{E} \ell(y_i f_Q(x_i)) = R(f_Q)$

$$\min_Q \sum_i^m \ell(y_i, f_Q(x_i)) = \hat{R}(f_Q)$$

with $f_Q(x_i) = \sum_{j=1}^m q_j b(w_j^T x_i + b_j)$.

$$Q = (q_j, w_j, b_j)_{m(d+2)}$$

- 3 algorithms:
- GD on $\hat{R}(f_Q)$
 - SGD on $\hat{R}(f_Q)$ with multiple passes $\frac{1}{M}$ on train set
 - SGD on $R(f_Q)$ with a single pass $\frac{1}{\sqrt{m}}$ on test set



Estimator error: ReLU Rectified linear Unit
 $\sigma(u) = u_+ = \max\{u, 0\}$

positively-homogeneous: $\sigma(\lambda u) = \lambda \sigma(u)$ if $\lambda > 0$



$$\eta_j (\underbrace{w_j^\top n + s_j}_d)_+ = d_j q_j \left(\frac{w_j^\top n + s_j}{d_j} \right)_+ \quad \text{if } d_j > c$$

$$d_j q_j \frac{w_j^\top n + s_j}{d_j}$$

consider ℓ_2 -penalties on weights

"weight decay"
explicit

implicit

$$\sum_j \eta_j^2 + \|w_j\|_2^2 + \frac{s_j^2}{R^2} \quad \left| \begin{array}{l} \text{Assumption: inputs } \|n\|_2 \leq R \text{ almost surely} \\ \text{or} \\ \sum_j \eta_j^2 + \|w_j\|_2^2 + \frac{s_j^2}{R^2} \end{array} \right.$$

penalty after rescaling = $\sum_j d_j^2 \eta_j^2 + \frac{1}{d_j^2} \left(\|w_j\|_2^2 + \frac{s_j^2}{R^2} \right)$ ADD CONSTANT: $\frac{\|w_j\|_2^2 + \frac{s_j^2}{R^2}}{\eta_j^2} = 1, k_j$

$$\left| \frac{\|w_j\|_2^2 + \frac{s_j^2}{R^2}}{\eta_j^2} = 1, k_j \right.$$

pencely a weights

$$\sum_{j=1}^m |u_j| \sqrt{\|w_j\|_2^2 + \frac{s_j^2}{D^2}}$$

check of normalization : $\sqrt{\|w_j\|_2^2 + \frac{s_j^2}{D^2}} = 1$] input weights

\Rightarrow pencely $\Rightarrow \sum_{j=1}^m |u_j| = \|u\|_1$] output weights

Goal = $F = \left\{ f(x) = \sum_{j=1}^m u_j (w_j^T x + s_j)_+, \quad \|w_j\|_2^2 + \frac{s_j^2}{D^2} \leq 1, \quad \|u\|_1 \leq D \right\}$

Recall: linear models

$$F = \left\{ f(x) = \phi^T \varphi(x), \quad \|\phi\|_2 \leq D \right\}$$

Total Rademacher complexity

$$f = \left\{ f_{\theta}(x) = \sum_{j=1}^m u_j (\omega_j^T x + s_j)_+, \quad \|u_j\|_2^2 + \frac{s_j^2}{B^2} \leq 1, \forall j \right\}$$

$\|u_j\|_2 \leq D$

Rademacher average

$$R_n = \mathbb{E}_{\text{data}} \left[\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \ell(y_i, f_{\theta}(x_i)) \right]$$

ε_i : Rademacher r.v. $P(\varepsilon_i = \pm 1) = \frac{1}{2}$, iid

$$\mathbb{E} \left[R(f_{\hat{\theta}}) - \inf_{\theta \in \Theta} R(f_{\theta}) \right] \leq 4R_n$$

(From Chap. 4)

Estimator
empirical
risk minimizer

$$R_n = \mathbb{E}_{\text{data}} \left[\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \ell(y_i, g_\theta(x_i)) \right]$$

ε : random error r.v. $P(\varepsilon = \pm 1) = \frac{1}{2}$, i.i.d

Extra-assumption: loss is G -Lipschitz cont. means $|\ell(y, u) - \ell(y, v)| \leq G|u - v|$

$$R_n \leq G \mathbb{E}_{\text{data}} \left[\sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i g_\theta(x_i) \right]$$

contract a principle

Lemma: $\sup_{\|u\|_1 \leq D} \bar{y}^\top u$
 $= D \cdot \|u\|_2$

$$\begin{aligned} & \|u\|_1 \leq D \\ & \|w\|_2^2 + \frac{s^2}{D^2} = 1 \end{aligned}$$

$$\sum_{j=1}^m (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_+$$

$$\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_+ \right|$$

uniformly in
 $j \in \{1, \dots, m\}$

$$= GD \mathbb{E} \left[\sup_{\substack{\text{mp} \\ \|w_j\|_2^2 + \frac{s_j^2}{D^2} = 1}} \sum_{i=1}^n \varepsilon_i (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_+ \right]$$

$\forall j$

$$\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_+ \right|$$

contract a principle
 $\leq 2GD$

$$\mathbb{E} \left[\sup_{\substack{\text{mp} \\ \|w_j\|_2^2 + \frac{s_j^2}{D^2} = 1}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_+ \right| \right]$$

$$\left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (\bar{w}_j^\top \bar{x}_i + \bar{s}_j)_- \right|$$

$$R_m \leq 2GD \mathbb{E}_{\text{data}} \left[\frac{\eta_{ij}}{\|w\|_2^2 + \frac{s^2}{R^2}} \right]$$

lemma: $\frac{\eta_{ij}}{\|w\|_2^2 + \frac{s^2}{R^2}} = \frac{\varepsilon_i(x_i)}{\|w\|_2^2 + \frac{s^2}{R^2}}$
 $\|w\|_2^2 = \|w\|_2$
 $\|(\varepsilon_i(x_i))\|_2 = 1$

$$\begin{aligned}
 &= 2GD \mathbb{E}_{\text{data}} \left\| \frac{1}{n} \sum_i \varepsilon_i \left(\frac{x_i}{R} \right) \right\|_2 \\
 &\leq 2GD \mathbb{E} \left[\left\| \frac{1}{n} \sum_i \varepsilon_i x_i \right\|_2 + \left\| \frac{1}{n} \sum_i \varepsilon_i R \right\|_2 \right]
 \end{aligned}$$

Finsler's $\leq 2GD$

$$\begin{aligned}
 &\leq \sqrt{\mathbb{E} \left\| \frac{1}{n} \sum_i \varepsilon_i x_i \right\|_2^2} + R \sqrt{\mathbb{E} \left\| \frac{1}{n} \sum_i \varepsilon_i \right\|^2} \\
 &= \frac{1}{n^2} \sum_i \mathbb{E} \|x_i\|_2^2 \leq \frac{R^2}{n}
 \end{aligned}$$

$$\leq 2GD \left[\frac{R}{\sqrt{n}} + \frac{R}{\sqrt{n}} \right] = \frac{4GDn}{\sqrt{n}}$$

Summary

$$f = \left\{ f(x) = \sum_{j=1}^m u_j (\omega_j^T x + s_j)_+, \quad \|u_j\|_2^2 + \frac{s_j^2}{D^2} \leq 1, \quad \|u_j\|_1 \leq D \right\}$$

Estimation error for the ERM a \hat{F}

$$\leq 4$$

$$\frac{4 \text{ GRD}}{\sqrt{n}} = 16 \frac{\text{ GRD}}{\sqrt{n}}$$

comes from
dist between estimation
error & True measure



Approximation error:

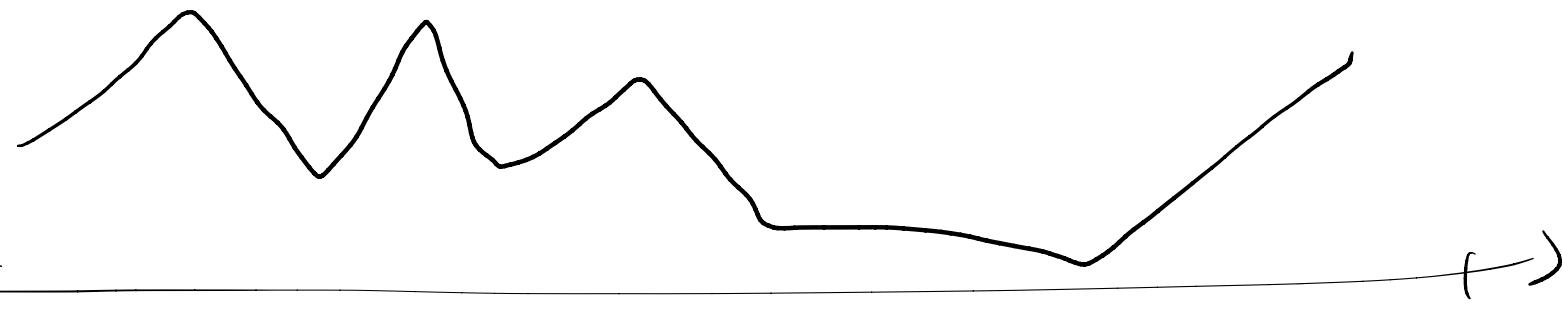
- ① universality
- ② bands as approximation error

$$m = +\infty$$

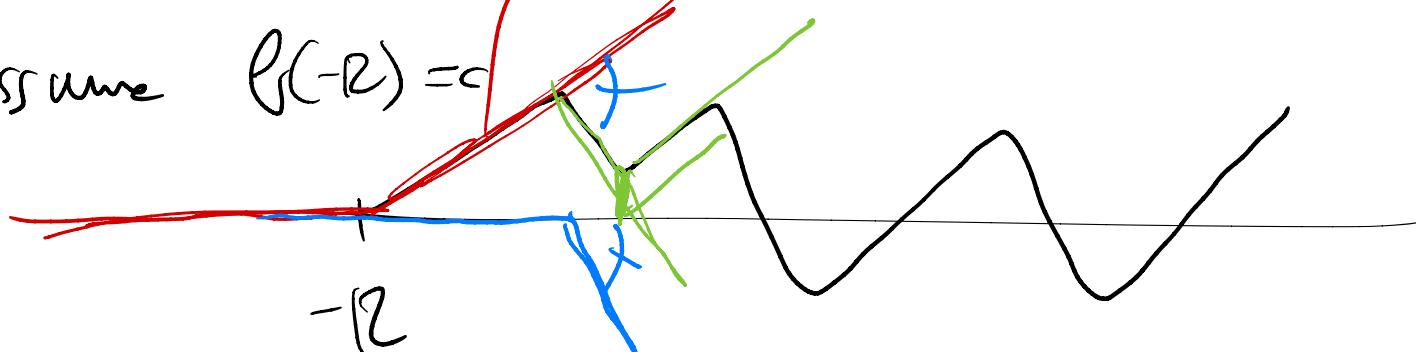
- ③ what can be achieved with m finite?

Universality in \mathbb{C} | Given $f_{\infty}: \mathbb{R}^d \rightarrow \mathbb{R}$
Can we approach it with measures $m \rightarrow \mu$

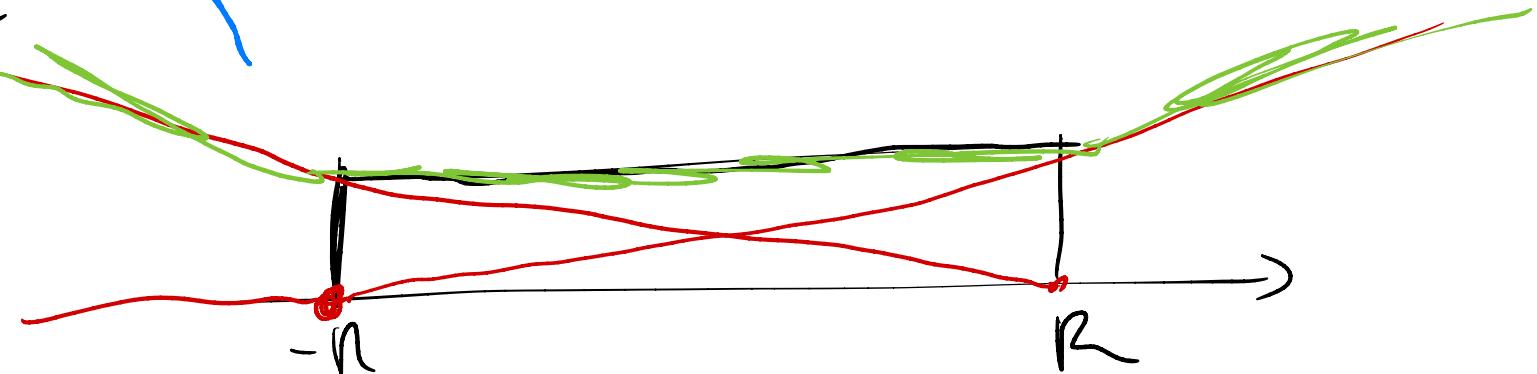
① True for piecewise affine functions

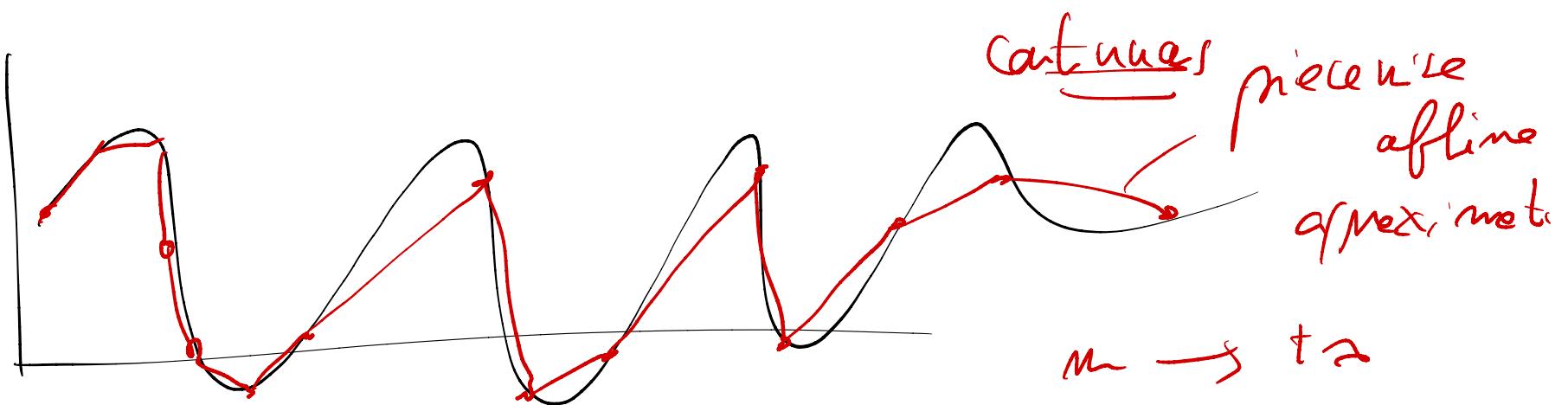


④ Assume $f(-R) = c$



⑤ If $f(-R) \neq c$



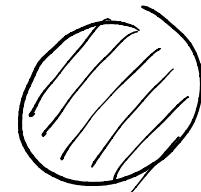


Definition of \mathcal{C} function space

$$f(x) = \sum_{j=1}^m a_j (w_j^T x + s_j)_+$$

$$= \int_K (w_i^T x + s)_+ d\nu(w, s)$$

$(w_j, s_j) \in K$ compact set $\subset \mathbb{R}^{d+n}$



$$\nu(w, s) = \sum_{j=1}^m q_j \delta_{(w_j, s_j)}$$

$$\|q\|_2 = \sum_{j=1}^m |q_j| = \text{total variation of } \nu$$

Dirac measure at (w_j, s_j)

Consider = $f(x) = \int (w^T n + b)_+ d\nu(u, s)$
 penalty $\int_K |\nu(u, s)|$ total variation
ex: if ν has density $\mu(u, s)$ with respect to Lebesgue

$$\int_K |\mu(u, s)| du ds$$

$\gamma_2(f) = \inf \int |\nu(w, s)| \text{ s.t. } f(x) = \int_n (w^T n + b)_+ d\nu(u, s)$

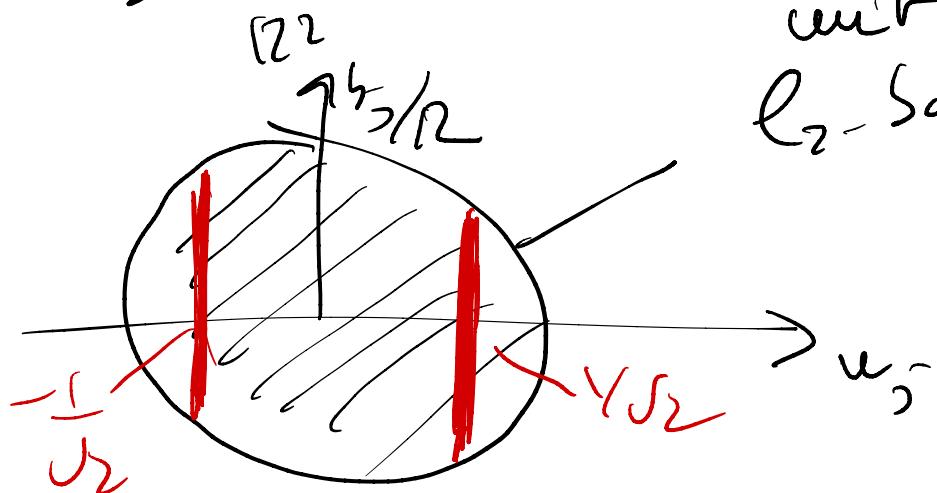
ex: if $f(x) = (w^T n + b)_+$ with $(x, s) \in K$
 $\gamma_2(f) \leq 1$

$$K = \left\{ (u, s) \in \mathbb{R}^{d+n}, \quad \|u\|_2 = \frac{1}{\sqrt{2}}, \quad |s| \leq \frac{R}{\sqrt{2}} \right\}$$

Earlier : $\bar{K} = \left\{ \|u_j\|_2^2 + \frac{s_j^2}{n} \leq 1 \right.$

with
 ℓ_2 -Ball

KCF



$$\left\{ \|u_j\|_2 = \frac{1}{\sqrt{2}}, \quad \frac{|s_j|}{n} \leq \frac{R}{\sqrt{2}} \right.$$

$$\underline{K} = \left\{ (u, s) \in \mathbb{R}^{d+n}, \quad \|u\|_2 = 1, \quad |s| \leq n \right\}$$

How to write $f(x) = \int_K (w^n + s)_{+} dv(u, s)$

KCF $\sqrt{2}$

$$\int \{dv(u, s)\}$$

1D Sect. a 9.3.3 .

Taylor's formula with integral remainder

$$f(z) = f(y) + \int_y^n f'(t) dt \quad \text{order 1 -}$$

$$f(z) = f(y) + (z-y) f'(y) + \frac{1}{2} (z-y)^2 f''(c) \quad c \in (y, z)$$

classical
remainder

$$+ \int_y^n f''(t)(z-t) dt$$

$$+ \int_y^x f''(t)(z-t) dt$$

(y, z)

$$\underbrace{y = -R : f(z) = f(-R) + (z+R) f'(-R) + \int_{-R}^R f''(t)(z-t) dt}$$

$$|f(z)| \leq \int_{-R}^R |f''(z)| + \left(f(-R), f(R) \right)$$

Sectra 9.33

Beyond 1D = $f(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \hat{f}(u) e^{iux} du$

$$\text{if } \|x\|_2 \leq R$$

inverse
Fourier transform

e^{iu} function $(-R, R) \rightarrow \mathbb{C}$.

$$e^{iu} = \int_{-R}^R \boxed{\text{FT}}_{(t)} (u-t)_+ dt \Rightarrow f(u) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \hat{g}(w) R(u)$$

ADAPTIVITY To LINEAR LATENT VARIABLE

if $f^\delta(z) = g^\delta(w_{\tau_n}^T)$, then $\gamma_\delta(f^\delta)$ "small".

|
lipschitz-continuity \Rightarrow CV rate $\frac{1}{n^{1/3}}$

lipschitz-continuity \Rightarrow CV rate $\frac{1}{n}$

Extended to. $f^\delta(z) = g^\delta(w_{\tau_n}^T)$ $n^{1/d}$

|
 $w \in \mathbb{R}^{d \times h}$

b_i -penalties on w

$$f(z) = \sum_{j=1}^m q_j (w_{\tau_n}^T + s_j)_+ = \underline{\gamma} \varphi(z)$$