# Lecture 2 : The curse of dimensionality; Neural Networks and approximation.

[Bach '21] [Telgarski, '20]
(DLT 'Fall20).

## Basic (Supervised) Learning Setup

**Goal.** Given data $\{(x_i, y_i)\}_{i=1\ldots n}$ with $x_i \in \mathcal{X}$ (high-dimension)
input  label.  $y_i \in \mathcal{Y}$ labels.

estimate a mapping $f : \mathcal{X} \longrightarrow \mathcal{Y}$ that

? generalises to unseen data

glorified Interpolation

→ IID Assumption : data is drawn iid from a distribution $\nu$ on $\mathcal{X} \times \mathcal{Y}$.

$y \in \{\pm 1\}$

eg: $\ell(y, y') = \log(1 + \exp(-yy'))$

→ Point-wise loss $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$

$\ell(y, y') = \frac{1}{2} | y - y'|^2$

$\hookrightarrow \mathcal{Y} = \mathbb{R}$

Given any $f : \mathcal{X} \longrightarrow \mathcal{Y}$, this defines

population risk $\quad R(f) = \mathbb{E}_\nu \left[ \ell( f(x), y) \right]$

$\left( \mathbb{E} \, \hat{R}(f) = R(f) \right)$

empirical risk $\quad \hat{R}(f) = \frac{1}{n} \sum_{i=1}^{n} \ell( f(x_i), y_i) \quad \leftarrow$ unbiased estimator of $R$.
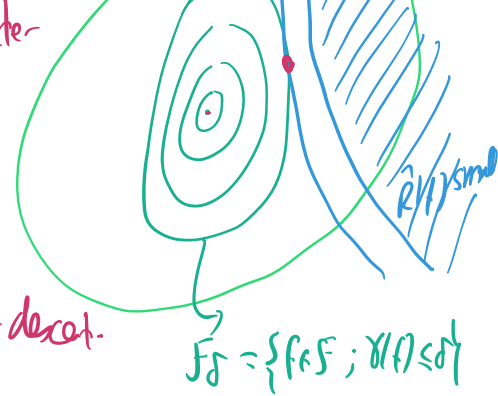
→ Hypothesis space $\mathcal{F} = \{ f : \mathcal{X} \longrightarrow \mathcal{Y} \}$

Assume that $\mathcal{F}$ comes with a norm : we can assign to each $f \in \mathcal{F}$ a complexity measure $\gamma(f)$

Examples : → norms over the NN parameter weights.

→ number of neurons

→ number of gradient-steps for models $f$ using gradient-descent.

$$\mathcal{F}_\delta = \{f \in \mathcal{F} ; \gamma(f) \le \delta\}$$

→ Empirical Risk Minimization (ERM).

↳ we search for small empirical risk using small complexity.

(C) $\min \hat{R}(f)$      (P) $\min \hat{R}(f) + \lambda \cdot \gamma(f)$
constrained $\gamma(f) \le \delta$     penalized $f \in \mathcal{F}$     ↳ Lagrange Multiplier

(I) $\min \gamma(f)$    (situations with no noise)    $f(x_i) = y_i$   $i = 1,\dots n$.
interpolant.    f s.t $\hat{R}(f) = 0$

→ <u>Basic decomposition of error/risk</u>

Let $\hat{f} \in \mathcal{F}_\delta$ produced by an arbitrary algorithm → (only access $\hat{R}$)

$\mathcal{E}_{app} = $ approximation error.

Then

$$R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f) = R(\hat{f}) - \inf_{f \in \mathcal{F}_\delta} R(f) + \boxed{\inf_{f \in \mathcal{F}_\delta} R(f) - \inf_{\mathcal{F}} R(f)}$$

$$= \boxed{R(\hat{f}) - \hat{R}(\hat{f})} + \boxed{\hat{R}(\hat{f}) - \inf_{f \in \mathcal{F}_\delta} \hat{R}(f)}_{\mathcal{E}_{opt}} +$$

$$+ \boxed{\inf_{f \in \mathcal{F}_\delta} \hat{R}(f) - \inf_{\mathcal{F}_\delta} R(f)} + \mathcal{E}_{app}$$

$$\hat{R}(\bar{f}) - R(\bar{f}) \quad \bar{f} \in \arg\min_{\mathcal{F}_\delta} R(f)$$

$$\le 2 \sup_{f \in \mathcal{F}_\delta} |\hat{R}(f) - R(f)| + \mathcal{E}_{appr} + \mathcal{E}_{opt}.$$

$\downarrow \varepsilon$ statistical error.

$$\int \to \varepsilon_{app} \; ; \quad \text{in regression} \quad R(f) = \|f - f^*\|_2^2$$

3 sources of error

$\{(x_i, y_i)\}$, with $y_i = f^*(x_i)$

$$\inf_{f \in \mathcal{F}_\delta} \|f - f^*\|_\nu^2$$

$\hookrightarrow$ As $\delta$ increases, $\varepsilon_{app}$ decreases.

$\to$ Statistical error

$$\boxed{\sup_{f \in \mathcal{F}_\delta}} \; |R(f) - \hat{R}(f)| \; \leftarrow \; \begin{array}{l} \text{a random} \\ \text{quantity.} \end{array}$$

measures <u>uniform</u> fluctuations over the ball $\mathcal{F}_\delta$.

Two main quantities drive this error:

$$\begin{cases} n = \text{the number of datapoint} \\ \delta : \text{size of } \mathcal{F}_\delta \text{ ball.} \end{cases}$$

What is the expected behavior?

Fix $f$ first.                                 $(x_i, y_i)$ iid.

$$\hat{R}(f) - R(f) = \frac{1}{n} \sum_{i=1}^{n} \left[ \ell(f(x_i), y_i) - \mathbb{E}\, \ell(f(x), y) \right]$$
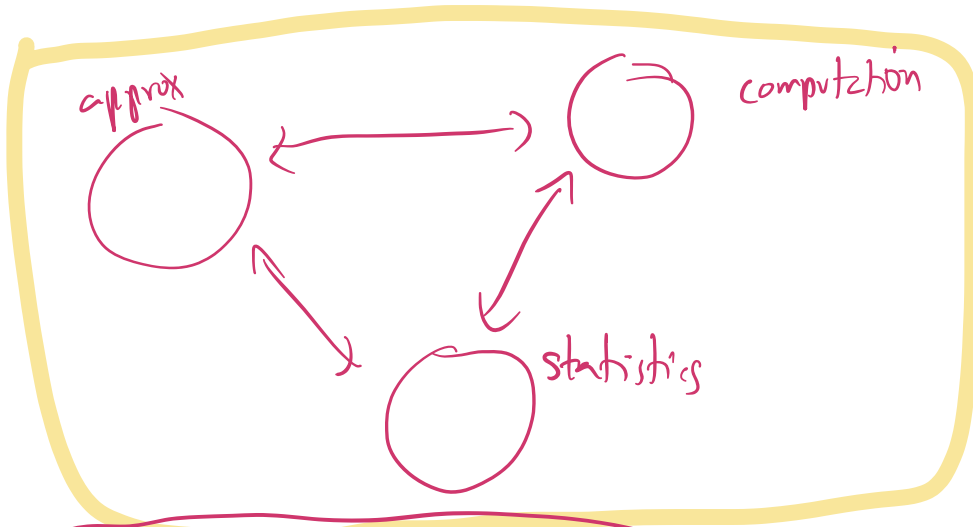
$$\Rightarrow |\hat{R}(f) - R(f)| \sim \frac{\text{std}(\ell(f(x), y))}{\sqrt{n}} \quad \begin{array}{l} \text{(non-asymptotic in)} \\ n \\ \text{(quantified in} \\ \text{high-prob using} \\ \text{classic tail bounds)} \end{array}$$

$\hookrightarrow$ From point-wise control to uniform control, requires some tools from concentration / empirical processes (Rademacher averages, etc)

$$\varepsilon_{stat} \sim \frac{h(\delta, \mathcal{F})}{\sqrt{n}} \leftarrow \quad \text{need to fight each other!}$$

$\downarrow$ Optim error        $\hat{R}(\hat{f}) - \inf_{f \in \mathcal{F}_\delta} \hat{R}(f)$        measures our ability to solve ERM.

$\rightarrow$ $f_\sigma$ is non-convex in most practical situations !



$\boxed{\text{The curse of dimensionality}}$ 'How do approximation and statistical errors behave as input dimension grows ?

statistics: we observe $\{(X_i, f^\circ(X_i)\}_{i=1\cdots n}$   $X_i \sim N(0, Id)$

$f^\circ$ unknown.   $X_i \in \mathbb{R}^d$.

Q: How many samples are needed to estimate $f^\circ$ up to accuracy $\varepsilon$, ie $\mathbb{E}_X | \hat{f}(x) - f^\circ(x)|^2 \leq \varepsilon$ ?
(sample complexity).

$\rightarrow$ suppose first $f^\circ$ is linear: $f^\circ(x) = \langle x, \theta^\circ \rangle$   $\theta^\circ \in \mathbb{R}^d$.

$$\mathcal{F} = \{ f : \mathbb{R}^d \rightarrow \mathbb{R} ; f(x) = \langle x, \theta \rangle \} \cong \mathbb{R}^d$$

A: $n = d$ are sufficient for exact recovery (and necessary!)
   (solve a linear system $\langle x_i, \theta \rangle = \langle x_i, \theta^\circ \rangle$ $i=1\cdots d$.

Remark: $f^\circ(x) = \varphi(\langle x, \theta^\circ \rangle)$   $\varphi$ is even function, smooth.

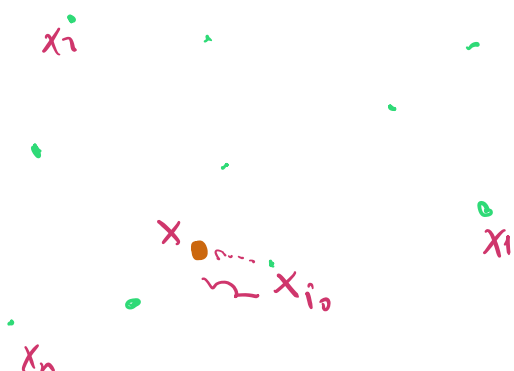A: $n = d+1$ sample are sufficient   $\varphi(t) = |t|$
   and also necessary   $= \cos t$.

↳ Suppose now that $f^\circ$ is only locally linear, ie

(β-) Lipschitz: $|f(x) - f(x')| \le \beta \cdot \|x - x'\|_2$ ←

$\mathcal{F} = \left\{ f ; \; \mathbb{R}^d \to \mathbb{R} \; ; \; f \; \begin{array}{l} \text{bounded} \\ \text{Lipschitz} \end{array} \right\}$ : a Banach space, with norm

$$\|f\|_{\mathcal{F}} = \|f\|_\infty + \text{Lip}(f).$$

$x_i$

$x_1$

$x_{i_0}$

$x_n$

We consider the smoothest interpolant:

$$\hat{f} = \underset{f}{\text{argmin}} \left\{ \text{Lip}(f) \; ; \; f(x_i) = f^\circ(x_i) \; \forall i \right\}$$

$\left[ \text{ERM in interpolant form} \right]$

For any $x$,

$$|\hat{f}(x) - f^\circ(x)| \le |\hat{f}(x) - \hat{f}(x_{i_0})| + |\hat{f}(x_{i_0}) - f^\circ(x_{i_0})| \overset{0.}{\diagup}$$

$$+ |f^\circ(x_{i_0}) - f^\circ(x)|$$

$$\underset{\beta\|x - x_{i_0}\|}{\wedge}$$

$$\le 2\beta \|x - x_{i_0}\|$$

$$\Rightarrow \underset{\underset{\mu}{\underset{|}{x \sim N(0, I_d)}}}{\mathbb{E}} |\hat{f}(x) - f^\circ(x)|^2 \le 4\beta^2 \underset{x, \, \exists x_i}{\mathbb{E}} \|x - x_{i_0}\|^2$$

$$\sim W_2^2(\mu, \hat{\mu}_n)$$

training set $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

$$\sim n^{-1/d} \quad \text{[Dudley '68]}$$

As a consequence, if $\varepsilon \sim n^{-1/d}$ $\quad$ [BLG '14].

$$\Rightarrow \quad n \sim \varepsilon^{-d} \quad \text{we can ensure that we}$$

learn with accuracy $\varepsilon$. $\quad \to$ | cursed by dimension |

ⓑ Is this sample complexity also necessary?

. Consider the box $B = [-1/2, 1/2]^d$
and the function $\Psi: B \to \mathbb{R}$
$$\Psi(x) = \text{dist}(x; \partial B)$$
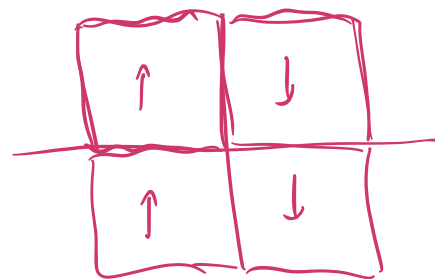Exercise: $\Psi$ is 1-Lipschitz.

· Consider for each $z \in \{-1/2, 1/2\}^d$
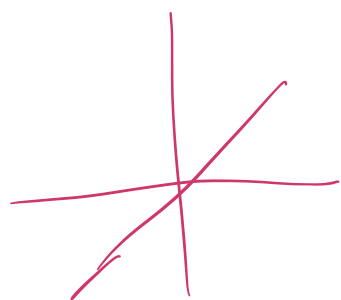an arbitrary sign $g(z) = \pm 1$

$$f^\circ(x) = \sum_{z \in \mathcal{H}_d} g(z) \, \Psi(x - z)$$

$f^\circ$ is 1-Lipschitz and supported
in $(-1, 1]^d$.

→ **Claim:** if $\underline{n \ll 2^d}$, then $\underline{\text{any}}$ estimator will
incur in a relative error

$$\frac{\mathbb{E}_x \, |\hat{f}(x) - f^\circ(x)|^2}{\mathbb{E}_x |f^\circ(x)|^2} = \Theta(1)$$

ⓑ To summarise:

→ linear functions (or generalised linear functions)
$$n \sim d \quad \text{(easy)}$$

→ Lipschitz functions → $n \sim \varepsilon^{-d}$ (impossible)

In between, Sobolev class:

$$W^{(s)} = \left\{ f : \mathbb{R}^d \to \mathbb{R} \; ; \; f \text{ has } s\text{-derivatives bounded on } L_p \right\}$$

Sample complexity $n \sim \varepsilon^{-d+2s/s}$ [Tsibakov].

$$\tilde{d} = d/s \qquad \sim \varepsilon^{-\tilde{d}}$$

$\to$ unless $s = \Theta(d)$, no real change.

Conclusion: we will need to search for alternatives.

## Shallow Neural Networks and Approximation

$\to$ Consider $\sigma : \mathbb{R} \to \mathbb{R}$ (Lipschitz) activation function

$$f_k(x \, ; \, \theta) = \sum_{k=1}^{k} \underbrace{c_k \, \sigma(\langle x, a_k \rangle + b_k)}_{\text{ridge function.}} \qquad \begin{array}{l} x \in \mathbb{R}^d \\ a_k \in \mathbb{R}^d \\ b_k, c_k \in \mathbb{R}. \end{array}$$

$\theta = \{ a_k, b_k, c_k \}_k$.

$\to$ Space of functions represented with shallow

NNs is $\quad H_\sigma = \left\{ f_k(\cdot \, ; \, \theta) \, ; \, \sigma \, ; \, k \in \mathbb{N} \right\}$

Q: How expressive is this set?

## Universal Approximation

For a given metric $d$ defined over continuous functions,

$\forall f \in CC(\mathbb{R}^d)$ and $\forall \varepsilon > 0$, there exist $\hat{f} \in H_\sigma$

such that $d(f, f) \leq \varepsilon$.

$\rightarrow$ To get intuition, let us first consider a 3-layer approximation.

**Theorem** Let $g: \mathbb{R}^d \rightarrow \mathbb{R}$ continuous, $\varepsilon > 0$, $\Omega = [0,1)^d$, and let $\delta > 0$ so that $\|x - x'\|_\infty < \delta$, $x, x' \in \Omega$, then $|g(x) - g(x')| \leq \varepsilon$. Then there exists a 3-layer ReLu network $f$ with $\bigoplus (\delta^{-d})$ and $\|f - g\|_1 \leq 2\varepsilon$.

$$\overset{\|\ }{L^1(\Omega)}$$

**Proof**: [Telgarski '20]

**Lemma**: Let $U \subseteq \mathbb{R}^d$, with a partition $P$ of $U$ into rectangles ; $P = (R_1 \cdots R_N)$, each with lengths $\leq \delta$. Then $\exists$ scalars $\alpha_1 \cdots \alpha_N$ with

$\boxed{|P| \sim \delta^{-d}}$ $\displaystyle\sup_{x \in U} |g(x) - h(x)| \leq \varepsilon$, with

$$h(x) = \sum_i \alpha_i \mathbb{1}_{[R_i]}$$

(in other words, piece-wise constant approximation has small error provided pieces are sufficiently small.

$\underline{\text{let}}$ $P$ the partition from above of $U = [0, 1+c)^d$ into rectangles $R_i = \prod_j [a_j, b_j)$ with $b_j - a_j \leq \delta$.

$$h(x) = \sum_i \alpha_i \mathbb{1}_{R_i} , \quad \text{so} \quad \|g - h\|_1 \leq \varepsilon.$$

. We construct a network of the form

$$f(x) = \sum \alpha_i g_i(x) \quad \text{where each } g_i \text{ is}$$

a Relu net that approximates $\mathbb{1}_{R_i}$.

$$\|f-g\|_1 \leq \overbrace{\| f- h\|_1} + \overbrace{\|h - g\|_1}^{\leq \varepsilon} \text{ from Lemma.}$$

$$= \left\| \sum_i \alpha_i \left( \mathbb{1}_{R_i} - g_i \right) \right\|_1 + \varepsilon$$

$$\leq \sum_i |\alpha_i| \; \left\| \mathbb{1}_{R_i} - g_i \right\|_1 + \varepsilon \qquad \leq 2\varepsilon$$

If we can build $g_i$ so that

$$\rightarrow \quad \left\| \mathbb{1}_{R_i} - g_i \right\|_1 \leq \frac{\varepsilon}{\sum_i |\alpha_i|} \quad , \quad \text{then we are done.}$$

- Fix $i$ , $R_i = \otimes [a_j, b_j)$ . For $\gamma > 0$ , and for each coordinate $j \in \{1 \ldots d\}$, let

$$g_{\gamma, j}(z) = \sigma\left( \frac{z - (a_j - \delta)}{\gamma} \right) - \sigma\left( \frac{z - a_j}{\gamma} \right)$$

$$- \sigma\left( \frac{z - b_j}{\gamma} \right) + \sigma\left( \frac{z - (b_j + \delta)}{\gamma} \right)$$

$$\hookrightarrow \sigma(t) = \max(0, t)$$



$a_j - \gamma \quad a_j \qquad b_j \quad b_j + \gamma$

$$g_\gamma(x) = \sigma\left( \left[ \sum_j g_{\gamma, j}(x) \right] - (d-1) \right)$$



$\hookrightarrow$ we verify that

(i) $g_\gamma(x) = \begin{cases} 1 & \text{if } x \in R_i \\ 0 & \text{if } x \notin \otimes [a_j - \delta, b_j + \delta) \\ [0,1] & \text{otherwise.} \end{cases}$

(ii) $\|g_\gamma - \mathbb{1}_{R_i}\|_1 \leq \mathcal{O}(\gamma)$ $\square$

Remark: Are the two hidden layers necessary?
question.

A: Hell no! Recall a classic result in polynomial approximation

Theorem (Stone-Weierstrass) Let $\Omega = [0,1]^d$. Let $\mathcal{F}$
be a function class satisfying:
(i) Each $f \in \mathcal{F}$ is continuous)
(ii) $\forall x \in \Omega, \exists f \in \mathcal{F}$ such that $f(x) \neq 0$.
(iii) $\forall x \neq x' \in \Omega, \exists f \in \mathcal{F}$ $f(x) \neq f(x')$.
(iv) $\mathcal{F}$ is an __algebra__ closed under multiplication and
vector space operations.

Then $\mathcal{F}$ enjoys universal approx: any continuous $g: \mathbb{R}^d \to \mathbb{R}$
$\forall c > 0, \exists \underline{f \in \mathcal{F}}$ with $\sup_x |f(x) - g(x)| \leq \varepsilon$.

↳ This theorem can be used to establish UAT for $H_\sigma$
with general choices of $\sigma$:

↳ $\sigma$ sigmoidal $\lim_{t \to -\infty} \sigma(t) = 0$ , $\lim_{t \to +\infty} \sigma(t) = 1$
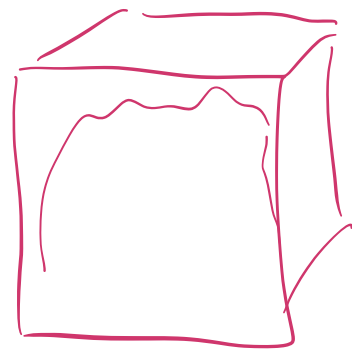
[Hornik et al 89]

↳ $\sigma \neq$ polynomial (Leshno)

↳ Is this surprising?
↳ Are the approximation rates cursed using $H_\sigma$?
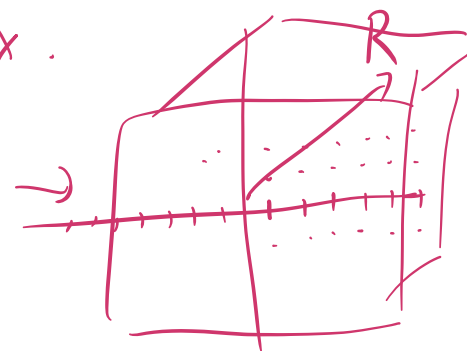
# The fourier perspective

→ Let $f \in C(\mathbb{R}^d)$ and consider its restriction to a compact $\Omega$ (eg $\Omega = [0,1]^d$).

For $\zeta \in \mathbb{Z}^d$, consider its fourier decomposition

$$\hat{f}(\zeta) = \left\langle f, e^{i\langle x, \zeta\rangle}\right\rangle_{L^2(\Omega)} = \int_\Omega f(x) e^{-i\langle x, \zeta\rangle} dx.$$

↳ Fourier Inversion lemma:

$$\boxed{f_M(x) = \sum_{\|\zeta\| \le R} \hat{f}(\zeta) e^{i\langle x, \zeta\rangle}}$$

$\downarrow$

$f$ in $L^2(\Omega)$ as $M, R \to \infty$

$M$ is the number of frequencies inside $\{\|\zeta\| < R\}$.

$$e^{i\langle x, \zeta\rangle} = \sigma(\langle x, \zeta\rangle) \quad \text{with} \quad \sigma(t) = e^{it} = \cos t + i\sin t.$$

$f_M$ is a shallow NN with M periodic neurons.

We have tight control of $\Big|$ how regularity of $f$

$\to$ $\uparrow$

decay of $\hat{f}$ as $\|\zeta\|$ grows.

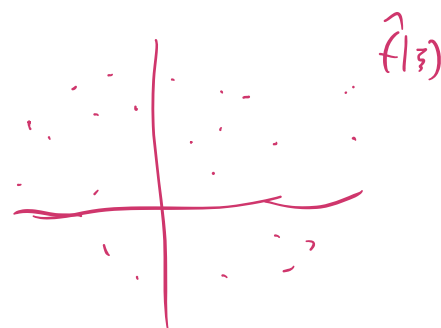$\boxed{f \in W^s}$ subolev. $\Rightarrow \|f - f_M\| = O(M^{-s/d})$ [Tsybakov] [de Vore]

curse of dim in approximation!

ↄ Overcome this curse?

Fourier Representation, in $L_1$

$$\boxed{\hat{f}(\zeta)} = \int f(x) \, e^{-i\langle x, \zeta\rangle} \, dx$$

$$\boxed{f(x) = \int \boxed{\hat{f}(\zeta)} \, e^{+i\langle x, \zeta\rangle} \, d\zeta}$$


$\hat{f}(\zeta)$

If $\|\hat{f}\|_1$ is finite, $\|\hat{f}\|_1 = \int |\hat{f}(\zeta)| \, d\zeta < +\infty$.

Idea: $f(x) = \int \sigma(\langle x, w\rangle) \, g(w) \, dw$ with $g$ integrable
$\|g\|_1 < +\infty$.

$$g(w) = sign(g(w)) \cdot \|g\|_1 \cdot \underbrace{\boxed{\frac{|g(w)|}{\|g\|_1}}}_{q(w)}$$

$q \geq 0$

$\int q(w) \, dw = 1$

$f(x) = \int \sigma(\langle x, w\rangle) \cdot sign(g(w)) \cdot \|g\|_1 \cdot q(w) \, dw$

$\Rightarrow q$ is a proba density.

$$= \mathbb{E}_q \left[ \underbrace{\sigma(\langle x, w\rangle) \cdot sign(g(w))}_{\phi(x,w)} \right] \cdot \|g\|_1$$

$|\phi| \leq 1$

ↄ $\hat{f}_M(x) = \frac{1}{M} \sum_{m=1}^{M} \phi(x, w_i)$, $w_i \sim q$ iid. $\boxed{\begin{array}{c}\text{Monte Carlo}\\ \text{estimator}\end{array}}$

$$\mathbb{E} \| f - \hat{f}_M \|^2 \leq \|g\|_1^2 \frac{\mathbb{E}_q \left( \mathbb{E}_x \, \phi(x, w)^2 \right)}{M} \leq \frac{\boxed{\|g\|_1^2} \, \sup_w \left( \mathbb{E}_x \, \phi^2 \right)}{M}$$

Curse of dim is avoided, provided $\|g\|_1 < +\infty$.

Theorem: (Barron '93) Suppose $C = \int \|\widehat{\nabla f}(\xi)\| \, d\xi < +\infty$

with $f, \hat{f} \in L_1$. Then we can approximate $f$ ~~with~~

to accuracy $\varepsilon$ with Relu/ Sigmoid units with $\sim \dfrac{C}{\varepsilon^2}$

(No curse!).

$\hookrightarrow$ Main q: When/ Why is $\underline{C}$ small.

Conclusions / Take- home:

(✶) Learning in high-dim efficiently requires new function

spaces, adapted to the physical world ( images,

sounds, etc.)