

Neural networks: Empirical risk minimization

predictors: $f_{\theta}: X \rightarrow \mathbb{R}$

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i))$$

Linear models: $f_{\theta}(x) = \theta^T \varphi(x)$ with $\varphi(x) \in \mathbb{R}^d$
= $\langle \theta, \varphi(x) \rangle$ with $\varphi(x) \in \mathbb{H}$ Hilbert space

+': convex optimization

- generalization guarantees

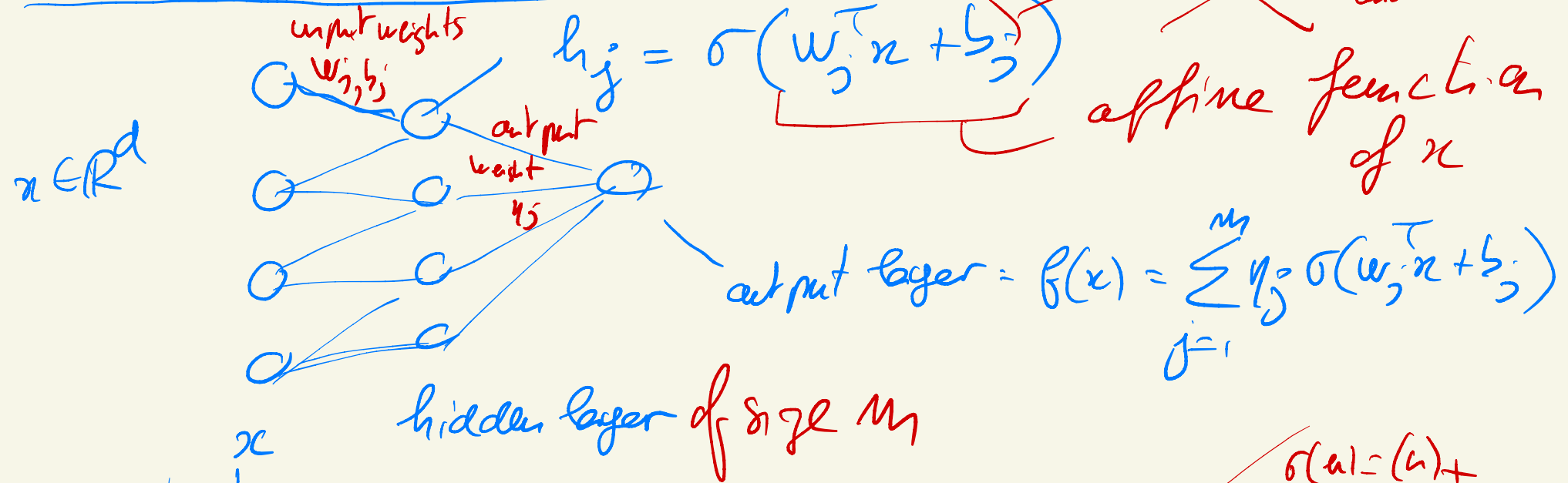
- kernels: escape curse of dimensionality if f^* smooth

-': still exponential in dimension if f^* non-smooth
excess risk = $O(n^{-\frac{1}{d+2}})$

Neural networks: "feature learning" \Rightarrow learn $\varphi(x)$

Today: 3 types of errors \rightarrow optimization
 \rightarrow estimation
 \rightarrow approximation

Fully connected single hidden layer NN



~~"bias"~~ offset constant term

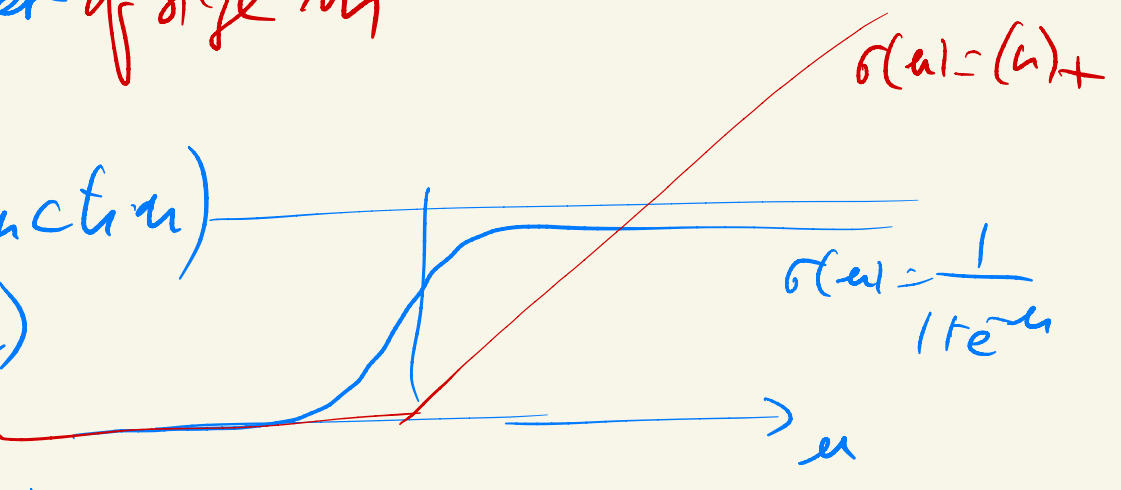
input layer

Choice of σ (activation function)

- ① Sigmoid (20th century)
- ② ReLU (21th —)

Rectified linear unit
 $\sigma(u) = \max(0, u) = (u)_+$

$\Rightarrow \sigma$ can be anything except a polynomial



Optimization = $R(\text{local}) = \frac{1}{n} \sum_{i=1}^n \ell(g_i, \sum_{j=1}^m \eta_j \sigma(w_j^T x + b_j))$

NON-CONVEX

\Rightarrow (S)GD may get trapped in local minimum

\perp non linearity in param.



① No guarantees

② Power of overparametrization
 m is large

"Proposition" (see blog post)

if $m \rightarrow \infty$, and weights are initialized randomly then GD converges to the global minimum

m large $\left\{ \begin{array}{l} \text{slow and sad for the planet} \\ \text{overfitting} \end{array} \right.$

Estimation error = $R(\hat{G}) = \mathbb{E} \ell(y, f(\hat{G}))$ expected risk

$\hat{R}(\hat{G}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{G}(x_i))$ empirical risk

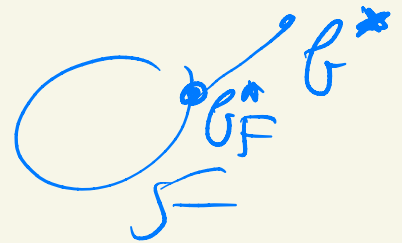
\mathcal{F} space of functions

$$\underbrace{R(\hat{G}) - \inf_{f \text{ measurable}} R(f)}_{\text{excess risk}} = R(\hat{G}) - R(G) + R(G) - \inf_{f \text{ measurable}} R(f) \quad \forall G \in \mathcal{F}$$

$$= \underbrace{R(\hat{G}) - R(G_F^*)}_{\text{estimation error}} + \underbrace{R(G_F^*) - \inf_{f \text{ measurable}} R(f)}_{\text{approximation error}} \quad \begin{matrix} f^* = \arg \min_{\mathcal{F}} R(f) \\ \mathcal{F} \subseteq \mathcal{F} \end{matrix}$$

$$\underbrace{R(\hat{G}) - \hat{R}(\hat{G})}_{\leq C} + \underbrace{\hat{R}(\hat{G}) - \hat{R}(G_F^*)}_{\leq C} + \underbrace{\hat{R}(G_F^*) - R(G_F^*)}_{\leq 0 \text{ if } \hat{R} = \mathbb{E} R \text{ minimizer}}$$

$$\leq \underbrace{\sup_{\mathcal{F}} R(f) - \hat{R}(f)}_{\text{estimation error}} + \sup_{\mathcal{F}} \hat{R}(f) - R(f)$$



Rademacher complexity

$$R_n(F) = \mathbb{E}_{\text{data}} \sum_{\epsilon \in \{-1, 1\}^n}$$

↳ "Rademacher variables"

$$\sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \epsilon_i \ell(y_i, f(x_i))$$

$$\mathbb{E} \sup_{f \in F} (R(f) - \hat{R}(f)) \leq 2 R_n(F) \Rightarrow \text{see class 3}$$

Contraction principle: if $\varphi_i: \mathbb{R} \rightarrow \mathbb{R}$ are 1-Lipschitz continuous

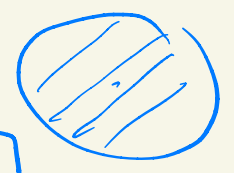
$$\begin{aligned} & \mathbb{E} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \epsilon_i \varphi_i(f(x_i)) \\ & \leq \mathbb{E} \sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \end{aligned}$$

$$|\varphi_i(u) - \varphi_i(v)| \leq |u - v|$$

For neural nets:

Neural networks: Assumptions: $\|w\|_2 \leq R$ and most sure key

$$\mathcal{F} = \left\{ \sum_{j=1}^m \eta_j \sigma(w_j^T x + b_j), \|w\|_1 \leq D_\eta, \Omega \left(\begin{matrix} w_j \\ b_j \end{matrix} \right) \leq D_{w,b} \right\}$$



\mathcal{F}_+ :

\hookrightarrow norm to be defined later

$\mathcal{F} \subset \mathcal{F}_+ - \{-\mathcal{F}_+\}$ because

$\eta \geq 0$

$$R_n(\mathcal{F}) \leq \underbrace{R_n(\mathcal{F}_+)} + R_n(-\mathcal{F}_+).$$

$$\left[\begin{array}{l} \eta_j = \overbrace{(\eta_j)_+}^{\geq 0} - \overbrace{(-\eta_j)_+}^{\geq 0} \\ |\eta_j| = |(\eta_j)_+| + |(-\eta_j)_+| \end{array} \right]$$

$$\mathcal{F} = \left\{ \sum_{j=1}^m \eta_j \sigma(w_j^T x + b_j), \|w\|_2 \leq D_w, \Omega(w_j, b_j) \leq D_{w,b} \right\} = \max \left\{ \|w_j\|_2 \frac{b_j}{R} \right\}$$

$$P_n(\mathcal{F}) \leq \mathbb{E}_{\xi, \text{data}} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i)$$

$$= \mathbb{E}_{\xi, \text{data}} \sup_{\substack{\eta_j, w_j, b_j \\ j \in \{1, \dots, m\}}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \sum_{j=1}^m \eta_j \sigma(w_j^T x_i + b_j)$$

$$= \mathbb{E}_{\xi, \text{data}} \sup_{\hat{\sigma}} \sup_{w_j, b_j} \frac{1}{n} \sum_{i=1}^n \epsilon_i \sigma(w_j^T x_i + b_j) \times D_w$$

if σ is 1-Lipschitz continuous $\leq \mathbb{E}$

$$\sup_i \sup_{w_j, b_j} \frac{1}{n} \sum_{i=1}^n \epsilon_i \left[w_j^T x_i + b_j \right] \times D_w$$

$$\|w_0\|_2 \leq D_w, b$$

$$|b_j| \leq R D_w, b$$

$$\sup_{x \in \mathcal{D}} = \sup_{x \in \mathcal{D}} \sup_{y \in \mathcal{Y}}$$

\leq

$$\leq D_y \mathbb{E} \sup_i \sup_{w_i, b_j} \frac{1}{n} \sum_{i=1}^n \epsilon_i [w_i^T x_i + b_j] = b_j \left[\frac{1}{n} \sum \epsilon_i \right] + w_j^T \left[\frac{1}{n} \sum \epsilon_i x_i \right]$$

$\|w\|_2 \leq D_{w,b}$
 $|b_j| \leq R D_{w,b}$

$$\leq D_y \mathbb{E} \sup_i R D_{w,b} \left| \frac{1}{n} \sum \epsilon_i \right| + D_{w,b} \left\| \frac{1}{n} \sum \epsilon_i x_i \right\|_2$$

$$\leq D_y D_{w,b} \left[\sqrt{\mathbb{E} \left(\frac{1}{n} \sum \epsilon_i \right)^2} + \sqrt{\mathbb{E} \left\| \frac{1}{n} \sum \epsilon_i x_i \right\|_2^2} \right]$$

by Jensen's inequality

$$\boxed{C^h R \rightarrow C^h C}$$

ϵ_i are independent and zero mean

$$= \frac{1}{n} \text{var}(\epsilon_i) = \frac{1}{n}$$

$$\begin{aligned} &= \frac{1}{n} \mathbb{E} \epsilon_i^2 \|x_i\|^2 \\ &\leq \frac{R^2}{n} \mathbb{1} \leq R^2 \end{aligned}$$

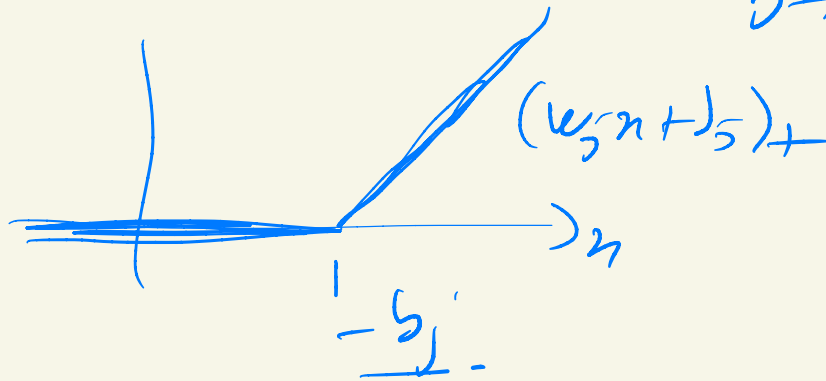
$$\boxed{\leq D_y D_{w,b} \frac{2R}{\sqrt{n}}}$$

$$\rightarrow O\left(\frac{1}{\sqrt{n}}\right)$$

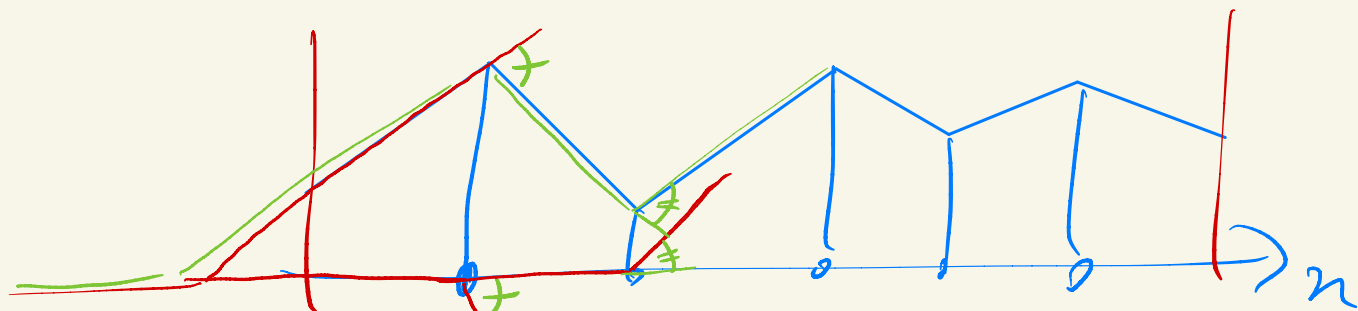
independent of n (# of neurons)

Approximating a function with neural networks

Universality: \mathbb{R}^D : approximate a function $f(x) \approx \sum_{j=1}^m \alpha_j (w_j x + b_j)_+ + C$



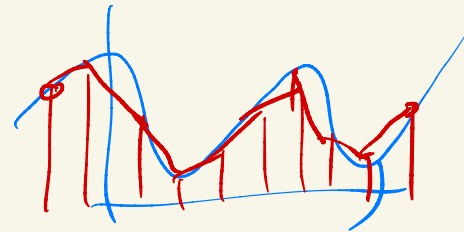
α_m is piecewise affine with m "kinks"



Proposition: any piecewise affine function in a finite interval and $m+1$ pieces can be written as a NN with $m+1$ neurons

ReLU

(2) All continuous functions can be approximated by piecewise affine functions



③ Universality in d dimensions

$$f(x) = \frac{1}{(2\pi)^d} \int \hat{f}(w) e^{i w^T x} dw$$

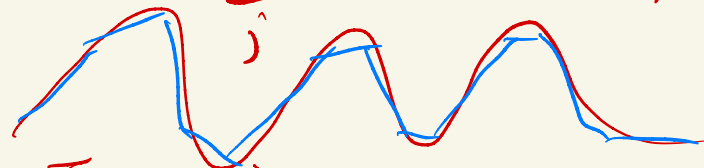
if f sufficiently regular

$$\text{where } \hat{f}(w) = \int f(z) e^{-i w^T z} dz$$

↳ Fourier transform

$$\varphi(w^T x)$$

cosine \approx piecewise affine function
sine $= \sum \eta_j (a_j u + b_j)_+ = \cos u$



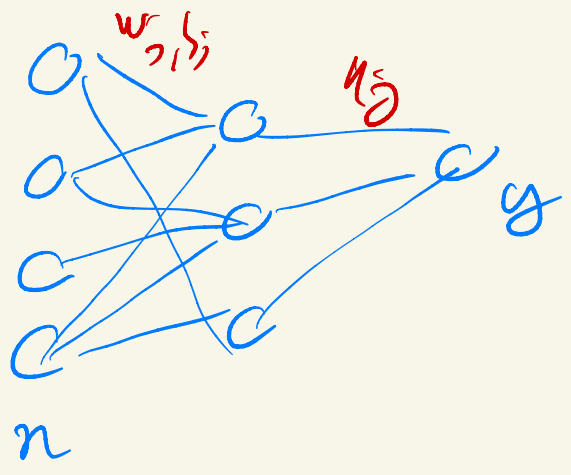
$$\cos(w^T x) \approx \sum_j \eta_j (a_j \underbrace{w^T x}_{\text{nodes}} + b_j)_+$$

Neural networks with potentially many neurons

$$B = \{ w_j^2 + b_j^2 \leq 1 \}$$

$$f(x) = \sum_{j=1}^m \eta_j \sigma(w_j^T x + b_j) = \int_B \sigma(w^T x + b) d\mu(w, b)$$

where $d\mu(w, b) = \sum_{j=1}^m \eta_j \delta_{w_j, b_j}$



$$\int g(w, b) \delta_{w_j, b_j} = g(w_j, b_j)$$

Two function spaces:

$$f(x) = \int \sigma(w^T x + b) d\mu(w, b)$$

$$\|\cdot\|_2 : \Omega_2(B) = \int |d\mu(w, b)| \text{ total variation} = \|\eta\|_1$$

$$\|\cdot\|_2 : \Omega_2(B)^2 = \int \left| \frac{d\mu(w, b)}{d\tau(w, b)} \right|^2 d\tau(w, b)$$

uniform measure on the ball B

Similar to

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) = \hat{R}(f)$$

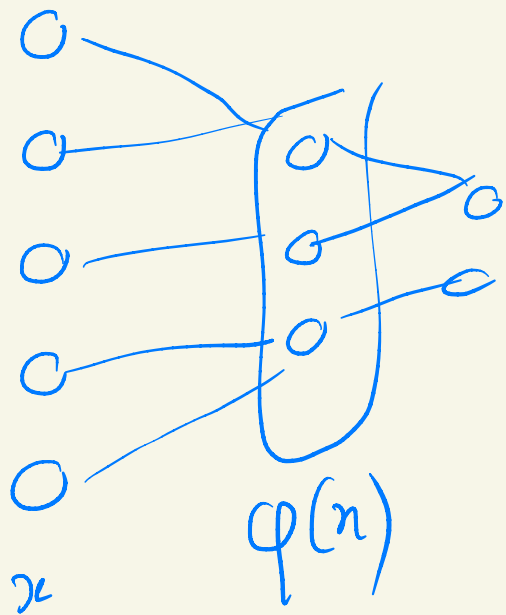
$$\int \ell(y, f(x)) d\hat{f}_n(x, y)$$

$$d\hat{f}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i, y_i} \rightarrow$$

empirical measure

fast distribution

link between neural network and kernel method



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

$$= \Theta^T \varphi(x) \text{ if } \begin{aligned} \varphi(x) &\in \mathbb{R}^m \\ \varphi(x)_j &= \frac{1}{\sqrt{m}} \sigma(w_j^T x + b_j) \\ \Theta &\in \mathbb{R}^m, \Theta_j = \eta_j \sqrt{m} \end{aligned}$$

in a kernel method, $h(x, x') = \varphi(x)^T \varphi(x')$ is fixed

in a neural net, $h(x, x')$ is parameterized (by w_j, b_j)

implies kernel learning and learned

$$h(x, x') = \varphi(x)^T \varphi(x') = \frac{1}{m} \sum_{j=1}^m \sigma(w_j^T x + b_j) \sigma(w_j^T x' + b_j)$$

converge to the expected
what if (w_j, b_j) are drawn?

$$h(x, x') = \varphi(x)^T \varphi(x') = \frac{1}{m} \sum_{j=1}^m \sigma(w_j^T x + b_j) \sigma(w_j^T x' + b_j)$$

converge to the expectation
what if (w_j, b_j) are random?

$$\text{Expectation} = \mathbb{E}_{w_j} \sigma(w_j^T x + b_j) \sigma(w_j^T x' + b_j) = h(x, x')$$

if $m \rightarrow \infty$, and input weights are not optimized
 \Rightarrow MK \Rightarrow kernel method

\Rightarrow see log posts (june/july 2020).