Some More Kernels + Deep Learning

CPSC 532S: Modern Statistical Learning Theory 14 March 2022 <u>cs.ubc.ca/~dsuth/532S/22/</u>

Admin

- A3 is up; due next Friday the 25th
- This week only, Thursday office hours are instead Wednesday 10-11am both in X563 and on the class Zoom
- Project proposals are due by Wednesday night
 - An informal paragraph on Piazza: just tell me what you want to do
 - Again, the scope of these is meant to be small
 - A lit survey doesn't require fully understanding the proofs or anything An "extension" could be just reading the proofs and talking about when the assumptions hold, etc

Last time $\kappa(x, \cdot) \in \mathcal{H}_{\kappa}$ $\langle f, \kappa(x, \cdot) \rangle_{\mathcal{H}_{\kappa}} = f(x)$ AF JE: X-M

- 下、ステンシア • We defined the RKHS for a given kernel k
 - Representer theorem:
 - We can kernelize any algorithm that only depends on $x_i \cdot x_i$
 - - Dependence on ||x|| becomes
 - Dependence on ||w|| becomes $||f||_{\mathscr{H}}$

c increasing function $\operatorname{argmin}_{f \in \mathcal{H}} L(f(x_1), \dots, f(x_n)) + R(\|f\|_{\mathcal{H}}) \in \operatorname{span}(\{k(x_i, \cdot)\}_{i=1}^n)$

Applied previous bounds on generalization gap / suboptimality to kernels

$$\sqrt{k(x,x)}$$

• What about that $L_S(f)$ or $\inf L_{\mathscr{D}}(f)$ term?

T" Kernel": RKHS kercer (Mercerkernel) K symmetric, pod funcion T Kernel density estimation: K(x,x*) ER symmetry pod x need $Sk(x,x')dx \ge 1$ often $k(x,x') \ge 0$ Ker (A) - null space X Convolutional remels 7 Bayesians' remel of a pdf FCUDA Rennel & Linnx Kernel Fpopcorn Kernel



• What about that $L_{S}(f)$ or $\inf_{\|f\|_{\mathscr{H}_{k}}} L_{\mathscr{D}}(f)$ term?

• A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ eq. $\mathcal{T} = \mathcal{F} \times \mathcal{ER}^d : \|X\| \leq \mathbb{R}^3$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$

Configuous L'functions X-3R



• What about that $L_S(f)$ or $\inf_{\|f\|_{\mathscr{H}_k}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13) lacksquare

• What about that $L_{S}(f)$ or $\inf_{\|f\|_{\mathscr{H}_{L}}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13) \bullet
 - Separates compact sets: if $X_1 \cap X_2 = \emptyset$ are compact subsets of \mathcal{X} , there's an $f \in \mathscr{H}_k$ with f(x) > 0 for $x \in X_1$, f(x) < 0 for $x \in X_2$ (so VCdim = ∞)

• What about that $L_S(f)$ or $\inf_{\|f\|_{\mathscr{H}_L}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13)
 - Separates compact sets: if $X_1 \cap X_2 = \emptyset$ are compact subsets of \mathscr{X} ,

there's an $f \in \mathcal{H}_k$ with f(x) > 0 for $x \in X_1$, f(x) < 0 for $x \in X_2$ (so VCdim = ∞) Implies that as $B \to \infty$, get $\inf L_S(f) \to 0$, $\inf L_{\mathscr{D}}(f) \to Bayes$ error if \mathscr{D} has compact support $\mathcal{H}_{k,B} = \mathcal{H}_{k,B} = \mathcal{H$



• What about that $L_S(f)$ or $\inf_{\|f\|_{\mathscr{H}_{L}}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13)
 - Separates compact sets: if $X_1 \cap X_2 = \emptyset$ are compact subsets of \mathcal{X} , there's an $f \in \mathcal{H}_k$ with f(x) > 0 for $x \in X_1$, f(x) < 0 for $x \in X_2$ (so VCdim = ∞) • Implies that as $B \to \infty$, get $\inf_{\mathscr{K}_{k,B}} L_{S}(f) \to 0$, $\inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \to \text{Bayes error}$ if \mathscr{D} has compact support
 - Can show universality via Stone-Weierstrass (more later), or Fourier properties

$$K(x,y) = \Psi(x-y)$$



• What about that $L_S(f)$ or $\inf_{\|f\|_{\mathscr{H}_{L}}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13)
 - Separates compact sets: if $X_1 \cap X_2 = \emptyset$ are compact subsets of \mathcal{X} , there's an $f \in \mathcal{H}_k$ with f(x) > 0 for $x \in X_1$, f(x) < 0 for $x \in X_2$ (so VCdim = ∞) • Implies that as $B \to \infty$, get $\inf_{\mathscr{K}_{k,B}} L_{S}(f) \to 0$, $\inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \to \text{Bayes error}$ if \mathscr{D} has compact support
 - Can show universality via Stone-Weierstrass (more later), or Fourier properties

•
$$\exp(x^{\mathsf{T}}y)$$
, $\exp(-\frac{1}{2\sigma^2}||x-y||^2)$, $\exp(-\frac{1}{\sigma^2}||x-y||^2)$

 $\|x - y\|$) are universal on compact subsets of \mathbb{R}^d



• What about that $L_S(f)$ or $\inf_{\|f\|_{\mathscr{H}_{L}}} L_{\mathscr{D}}(f)$ term?

- A continuous kernel on a compact metric space \mathscr{X} is universal if \mathscr{H}_k is dense in $C(\mathscr{X})$: for every continuous $g: \mathcal{X} \to \mathbb{R}$, every $\varepsilon > 0$, there is an $f \in \mathcal{H}_k$ with $||f - g||_{\infty} = \sup |f(x) - g(x)| \le \varepsilon$ $x \in \mathcal{X}$
 - If \mathscr{X} is a topological space *not* generated by a metric, there is no universal kernel (Steinwart/Christmann exercise 4.13)
 - Separates compact sets: if $X_1 \cap X_2 = \emptyset$ are compact subsets of \mathscr{X} , there's an $f \in \mathcal{H}_k$ with f(x) > 0 for $x \in X_1$, f(x) < 0 for $x \in X_2$ (so VCdim = ∞) • Implies that as $B \to \infty$, get $\inf_{\mathscr{K}_{k,B}} L_{S}(f) \to 0$, $\inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \to \text{Bayes error}$ if \mathscr{D} has compact support
 - Can show universality via Stone-Weierstrass (more later), or Fourier properties

•
$$\exp(x^{\mathsf{T}}y)$$
, $\exp(-\frac{1}{2\sigma^2}||x-y||^2)$, $\exp(-\frac{1}{\sigma^2}||x-y||^2)$

Never true for finite-dimensional kernels

 $-\|x-y\|$) are universal on compact subsets of \mathbb{R}^d



- Know that as $B \to \infty$, get $\inf_{\mathscr{H}_{k,B}} L_S(f) \to 0$, $\inf_{\mathscr{H}_{k,B}} L_{\mathscr{D}}(f) \to Bayes$ error

for compactly supported \mathcal{D} (can use broader notion of universality in general)



- Know that as $B \to \infty$, get $\inf L_S(f)$ $\mathcal{H}_{k,B}$

 - But the rate at which this happens depends on \mathscr{D}

$$\rightarrow 0, \inf_{\mathscr{H}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$

for compactly supported \mathcal{D} (can use broader notion of universality in general)



• Know that as $B \to \infty$, get $\inf L_S(f)$ – $\mathcal{H}_{k.B}$

for compactly supported \mathscr{D} (can use broader notion of universality in general) • But the rate at which this happens depends on \mathscr{D}

- Usually compare to the regression function $f_{\mathcal{P}}(x) = \mathbb{E}[y \mid x]$

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$



• Know that as $B \to \infty$, get $\inf L_S(f)$ – $\mathcal{H}_{k.B}$

for compactly supported \mathscr{D} (can use broader notion of universality in general) - But the rate at which this happens depends on ${\mathscr D}$

- Usually compare to the regression function $f_{\mathcal{P}}(x) = \mathbb{E}[y \mid x]$
 - If $f_{\mathcal{D}} \in \mathcal{H}_k$, called well-specified:

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$



• Know that as $B \to \infty$, get $\inf L_S(f)$ – $\mathcal{H}_{k,B}$

for compactly supported \mathcal{D} (can use broader notion of universality in general) - But the rate at which this happens depends on ${\mathscr D}$

- Usually compare to the regression function $f_{\mathcal{D}}(x) = \mathbb{E}[y \mid x]$
 - If $f_{\mathcal{D}} \in \mathcal{H}_k$, called well-specified: • Stability for $B = \|f_{\mathcal{D}}\|_{\mathcal{H}_k}$: $\inf_{\|f\|_{\mathcal{H}_k} \leq B} L_{\mathcal{D}}(f)$ = Bayes error, excess error $\leq \mathcal{O}(1/\sqrt{n})$

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$



. Know that as $B \to \infty$, get $\inf_{\mathscr{K}_{k,B}} L_S(f) = \mathscr{K}_{k,B}$

for compactly supported \mathscr{D} (can use broader notion of universality in general) • But the rate at which this happens depends on \mathscr{D}

- Usually compare to the regression function $f_{\mathcal{D}}(x) = \mathbb{E}[y \mid x]$
 - If $f_{\mathcal{D}} \in \mathcal{H}_k$, called well-specified: • Stability for $B = \|f_{\mathcal{D}}\|_{\mathcal{H}_k}$: $\inf_{\|f\|_{\mathcal{H}_k} \leq B} L_{\mathcal{D}}(f)$ = Bayes error, excess error $\leq \mathcal{O}(1/\sqrt{n})$

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$

• Better rates (minimax-optimal) with "range-space condition" if $f_{\mathcal{D}}$ is "nice" in \mathscr{H}_k





. Know that as $B \to \infty$, get $\inf_{\mathscr{K}_{k,B}} L_S(f) = \mathscr{K}_{k,B}$

for compactly supported \mathscr{D} (can use broader notion of universality in general) - But the rate at which this happens depends on ${\mathscr D}$

- Usually compare to the regression function $f_{\mathcal{P}}(x) = \mathbb{E}[y \mid x]$
 - If $f_{\mathcal{D}} \in \mathcal{H}_k$, called well-specified: • Stability for $B = \|f_{\mathcal{D}}\|_{\mathcal{H}_k}$: $\inf_{\|f\|_{\mathcal{H}_k} \leq B} L_{\mathcal{D}}(f)$ = Bayes error, excess error $\leq \mathcal{O}(1/\sqrt{n})$

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$

- Better rates (minimax-optimal) with "range-space condition" if $f_{\mathcal{D}}$ is "nice" in \mathscr{H}_k • Pretty different style of analysis, based on $\|\hat{f} - f_{\mathcal{D}}\|_{\mathcal{H}_k}$







• Know that as $B \to \infty$, get $\inf L_S(f)$ – $\mathcal{H}_{k,B}$

for compactly supported \mathscr{D} (can use broader notion of universality in general) - But the rate at which this happens depends on ${\mathscr D}$

- Usually compare to the regression function $f_{\mathcal{D}}(x) = \mathbb{E}[y \mid x]$
 - If $f_{\mathcal{D}} \in \mathcal{H}_k$, called well-specified: • Stability for $B = \|f_{\mathscr{D}}\|_{\mathscr{H}_k}$: $\inf_{\|f\|_{\mathscr{H}_k} \leq B} L_{\mathscr{D}}(f)$ = Bayes error, excess error $\leq \mathcal{O}(1/\sqrt{n})$
 - - Pretty different style of analysis, based on $\|\hat{f} f_{\mathcal{D}}\|_{\mathcal{H}_{k}}$

Approximation error

$$\rightarrow 0, \inf_{\mathscr{K}_{k,B}} L_{\mathscr{D}}(f) \rightarrow \text{Bayes error}$$

- Better rates (minimax-optimal) with "range-space condition" if $f_{\mathcal{D}}$ is "nice" in \mathscr{H}_k

Misspecified case: more complicated analyses based on "approximation spaces"



Gaussian processes

• $f \sim GP(m, k)$ is a random function $f \colon \mathscr{X} \to \mathbb{R}$ s.t., for any x_1, \ldots, x_n , $\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(x_1) \\ \vdots \\ m(x_n) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} \right)$





- Mean function $m: \mathcal{X} \to \mathbb{R}$ can be any function; usually use 0
 - will see that we can just shift everything by m so that this is WLOG



- Mean function $m: \mathcal{X} \to \mathbb{R}$ can be any function; usually use 0
 - will see that we can just shift everything by m so that this is WLOG

• Covariance function $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ can be any psd function, i.e. any kernel



- Mean function $m: \mathcal{X} \to \mathbb{R}$ can be any function; usually use 0
 - will see that we can just shift everything by m so that this is WLOG
- Note: samples f are almost surely not in \mathcal{H}_k , for infinite-dim \mathcal{H}_k

• Covariance function $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ can be any psd function, i.e. any kernel



- Mean function $m: \mathcal{X} \to \mathbb{R}$ can be any function; usually use 0
 - will see that we can just shift everything by m so that this is WLOG
- Note: samples f are almost surely not in \mathcal{H}_k , for infinite-dim \mathcal{H}_k
 - but they are almost surely in a "slightly bigger" RKHS
 - see e.g. Section 4 of Kanagawa et al. (2018)

• Covariance function $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ can be any psd function, i.e. any kernel

• Assume a **prior** $f \sim GP(m, k)$

- Assume a **prior** $f \sim GP(m, k)$
- Assume likelihood of observations by $y_i \sim \mathcal{N}(f(x_i), \sigma^2)$

- Assume a **prior** $f \sim GP(m, k)$
- Assume likelihood of observations by $y_i \sim \mathcal{N}(f(x_i), \sigma^2)$
 - $\mathbb{E}[y_i] = \mathbb{E}[f(x_i)], \quad \operatorname{Cov}(y_i, y_i) = \operatorname{Cov}(f(x_i), f(x_i)) + \sigma^2 \delta_{ii}$

- Assume a prior $f \sim GP(m, k)$
- $\mathbb{E}[y_i] = \mathbb{E}[f(x_i)], \quad \operatorname{Cov}(y_i, y_i) = \operatorname{Cov}(f(x_i), f(x_i)) + \sigma^2 \delta_{ii}$ $f \mid S \sim GP\left(\left[x \mapsto y^{\top}(K_{S} + \sigma^{2}I)^{-1}k_{S}(x)\right], \left[(x, x') \mapsto k(x, x') - k_{S}(x)^{\top}(K_{S} + \sigma^{2}I)^{-1}k_{S}(x')\right]\right)$
- Assume likelihood of observations by $y_i \sim \mathcal{N}(f(x_i), \sigma^2)$ • The **posterior** works out to be (via Kolmogorov Extension Theorem)



- Assume a prior $f \sim GP(m, k)$
- Assume likelihood of observations
 - $\mathbb{E}[y_i] = \mathbb{E}[f(x_i)], \quad \operatorname{Cov}(y_i, y_j) =$
- The **posterior** works out to be (via K $f \mid S \sim \text{GP}\left(\left[x \mapsto y^{\top}(K_S + \sigma^2 I)^{-1}k_S(x)\right]\right)$



by
$$y_i \sim \mathcal{N}(f(x_i), \sigma^2)$$

 $\operatorname{Cov}(f(x_i), f(x_j)) + \sigma^2 \delta_{ij}$
Kolmogorov Extension Theorem)
 $p_i^{-1}, [(x, x') \mapsto k(x, x') - k_S(x)^{\top}(K_S + \sigma^2 I)^{-1}k_S(x))$



More Gaussian Processes

- GP regression: can get posterior contraction rates
- Understanding posterior variance can be very useful!
 - e.g. Bayesian optimization / active learning / bandits / …

Look like KRR analysis for the mean, plus posterior variance decreasing

GP classifiers: usual choice corresponds to kernel logistic regression

More kernel resources

- Foundations: Berlinet and Thomas-Agnan, <u>RKHSes in Probability and Stats</u> (2004) Including more hardcore details: Steinwart and Christmann, <u>SVMs</u> (2008)
- Ridge regression analyses:
 - <u>Smale and Zhou (2007)</u> fairly readable
 - <u>Caponnetto and de Vito (2007)</u> minimax rate for "mostly"-well-specified, including regression with a kernel output as well
 - <u>Steinwart et al. (2009)</u> minimax in Sobolev spaces with Matérn kernels (hard)
- Rasmussen and Williams, <u>Gaussian Processes for Machine Learning</u> (2006)
- Connections between kernels and GPs: <u>Kanagawa et al. (2018)</u>
- Mean embeddings: <u>Muandet et al.</u> (2016)
 - v. related to a lot of my research; there are slides in L16, but won't get to them



 $K(x,x) = J_A(x) J_A(x)$ Q: X -> IRP

 $K(x, x') = \mathcal{X}(\mathcal{U}(x), \mathcal{U}(x'))$ $\int \mathcal{T}$ (pause) $\mathcal{U}: \mathcal{X} \to \mathbb{R}^{D}$ Graussian or stn.



10

Deep learning • Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"):

 $\begin{bmatrix} x \\ w_i \\ x + b_i \\ \frac{2i}{i} \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ \frac{2}{i} \\ \frac{2}{i} \end{bmatrix} = \begin{bmatrix} w_2 \\ \frac{2}{i} \\ \frac{2}{i} \\ \frac{2}{i} \end{bmatrix}$

Deep learning

Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell}f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$



- $f^{(0)}(x) = x$
 - $W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$ $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\mathcal{P}}: \mathbb{R}^{d'_{\mathcal{C}}} \to \mathbb{R}^{d_{\mathcal{C}}} \text{ (usually } d'_{\mathcal{P}} = d_{\mathcal{C}})$



$f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ • $f^{(0)}(x) = x$

•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers



•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers



• $f^{(0)}(x) = x$ $f^{(\ell)}(x) = \sigma_{\mathcal{P}}(W_{\mathcal{P}} f^{(\ell-1)}(x) + b_{\mathcal{P}})$

•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers
- Common choices for activations σ :

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers



•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers
- Common choices for activations σ:

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers

• Componentwise: $\operatorname{ReLU}(z) = \max\{z, 0\}, \operatorname{sigmoid}(z) = 1/(1 + \exp(-z))$



•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers
- Common choices for activations σ:
 - softmax(z)_i = exp(z_i)/ $\sum \exp(z_j)$, max pooling, attention, ...

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers

• Componentwise: $\operatorname{ReLU}(z) = \max\{z, 0\}, \operatorname{sigmoid}(z) = 1/(1 + \exp(-z))$



•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers
- Common choices for activations σ:
 - softmax(z)_i = exp(z_i)/ $\sum \exp(z_j)$, max pooling, attention, ...

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers

• Componentwise: $\operatorname{ReLU}(z) = \max\{z, 0\}, \operatorname{sigmoid}(z) = 1/(1 + \exp(-z))$

Usually train via SGD, but it's non-convex: in general, possibility of local minima





•
$$W_{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$$
 $b_{\ell} \in \mathbb{R}^{d'_{\ell}}$

- Usually $\sigma_I(x) = x$; intermediate layers called hidden layers
- Common choices for activations σ:
 - softmax(z)_i = exp(z_i)/ $\sum \exp(z_j)$, max pooling, attention, ...

Deep learning

 Mostly assuming fully-connected, feedforward nets ("multilayer perceptrons"): $f^{(\ell)}(x) = \sigma_{\ell}(W_{\ell} f^{(\ell-1)}(x) + b_{\ell})$ $f(x) = f^{(L)}(x)$

 $\sigma_{\ell} : \mathbb{R}^{d'_{\ell}} \to \mathbb{R}^{d_{\ell}}$ (usually $d'_{\ell} = d_{\ell}$)

Can think of this as a directed, acyclic computation graph, organized in layers

• Componentwise: $\operatorname{ReLU}(z) = \max\{z, 0\}, \operatorname{sigmoid}(z) = 1/(1 + \exp(-z))$

 Usually train via SGD, but it's non-convex: in general, possibility of local minima • ERM is NP-hard, even with 1 ReLU, even for square loss (Goel et al. ITCS 2021) 11









$$(p_i) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_i \mathbb{I}(x_i \ge b_i)$





|g(x) - f(x)|

$$(p_i) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_i \mathbb{I}(x_i \ge b_i)$





 $k = \max\{k : b_k \le x\}$ $|g(x) - f(x)| \le |g(x) - g(b_k)| + |g(b_k)|| \le |g(x) - g(b_k)|| \le |g($

$$(y_i) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_i \mathbb{I}(x_i \ge b_i)$

$$(p_k) - f(b_k) | + | f(b_k) - f(x)$$





 $k = \max\{k : b_k \le x\}$ $|g(x) - f(x)| \le |g(x) - g(b_k)| + |g(b_k)|$ $\leq \rho |x - b_k|$

$$(y_i) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_i \mathbb{I}(x_i \ge b_i)$

$$(p_k) - f(b_k) | + |f(b_k) - f(x)|$$





 $k = \max\{k : b_k \le x\}$ $|g(x) - f(x)| \le |g(x) - g(b_k)| + |g(b_k)|$ $\leq \rho |x - b_k|$ $\leq \rho \frac{\varepsilon}{\rho} = \varepsilon$

$$(y_i) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_i \mathbb{I}(x_i \ge b_i)$

$$(p_k) - f(b_k) | + |f(b_k) - f(x)|$$





 $k = \max\{k : b_k \le x\}$ $|g(x) - f(x)| \le |g(x) - g(b_k)| + |g(b_k)|$ $\leq \rho |x - b_k|$ $\leq \rho \frac{\varepsilon}{\rho} = \varepsilon$ Can do better by depending on *total variation* of g

$$f_{i}(x_{i}) - g(b_{i-1})$$
 $f(x) = \sum_{i=0}^{m-1} a_{i} \mathbb{I}(x_{i} \ge b_{i})$

$$(p_k) - f(b_k) | + | f(b_k) - f(x) |$$





Universal approximation in \mathbb{R}^d **Theorem:** Let $g : \mathbb{R}^d \to \mathbb{R}$ be continuous. For any $\varepsilon > 0$, choose $\delta > 0$ so that $||x - x'||_{\infty} \le \delta$ implies $|g(x) - g(x')| \le \varepsilon$. Then there is a three-layer ReLU network f with $\Omega\left(\frac{1}{\delta^d}\right)$ nodes satisfying $\int_{[0,1]^d} |f(x) - g(x)| dx \le 2\varepsilon$.



Universal app
Theorem: Let
$$g : \mathbb{R}^{d} \to \mathbb{R}$$
 be continuous
 $\|x - x'\|_{\infty} \leq \delta$ implies $\|g(x) - g(x')\| \leq \delta$
with $\Omega\left(\frac{1}{\delta^{d}}\right)$ nodes satisfying $\int_{[0,1]^{d}} |f(x)|^{2}$

ALL

Proof approximates continuous g by piecewise-constant h, then uses a two-layer ReLU net to check if x is in each piece, roughly like in 1d. (Telgarsky's Theorem 2.1.)

roximation in \mathbb{R}^d

us. For any $\varepsilon > 0$, choose $\delta > 0$ so that

 ε . Then there is a three-layer ReLU network f

 $(x) - g(x) | \mathrm{d}x \leq 2\varepsilon.$



Stone-Weierstrass Theorem: Let \mathcal{F} be a set of functions such that

- 1. Each $f \in \mathcal{F}$ is continuous.
- 2. For each x, there is at least one $f \in \mathcal{F}$
- 4. \mathscr{F} is an algebra: for $f, g \in \mathscr{F}$, $\alpha f + g \in \mathscr{F}$ and $fg = (x \mapsto f(x)g(x)) \in \mathscr{F}$. Then R is dense in C(R) w.r.t. [[·[[00.

$$\mathcal{F}$$
 with $f(x) \neq 0$.

3. Separates points: for each $x \neq x'$, there is at least one $f \in \mathscr{F}$ with $f(x) \neq f(x')$.



Stone-Weierstrass Theorem: Let \mathcal{F} be a set of functions such that

- 1. Each $f \in \mathcal{F}$ is continuous.
- 2. For each x, there is at least one $f \in \mathcal{C}$
- 4. \mathscr{F} is an algebra: for $f, g \in \mathscr{F}$, $\alpha f + g \in \mathscr{F}$ and $fg = (x \mapsto f(x)g(x)) \in \mathscr{F}$. Conditions hold for $\sigma_1 = \exp, \sigma_2 = \operatorname{Id}$, so that $\mathscr{F}_{\exp} = \{x \mapsto \sum a_i \exp(w_i^T x)\}$ *i*=1

$$\mathcal{F}$$
 with $f(x) \neq 0$.

3. Separates points: for each $x \neq x'$, there is at least one $f \in \mathscr{F}$ with $f(x) \neq f(x')$.



Stone-Weierstrass Theorem: Let \mathcal{F} be a set of functions such that

- 1. Each $f \in \mathcal{F}$ is continuous.
- 2. For each x, there is at least one $f \in \mathcal{F}$
- 4. \mathscr{F} is an algebra: for $f, g \in \mathscr{F}$, $\alpha f + g \in \mathscr{F}$ and $fg = (x \mapsto f(x)g(x)) \in \mathscr{F}$. Conditions hold for $\sigma_1 = \exp, \sigma_2 = \operatorname{Id}$, so that $\mathscr{F}_{\exp} = \{x \mapsto \sum a_i \exp(w_i^{\mathsf{T}} x)\}$ If $\sigma : \mathbb{R} \to \mathbb{R}$ is continuous, $\lim \sigma(z) = 0$, $\lim \sigma(z) = 1$, works too:

 $z \rightarrow -\infty$ $z \rightarrow -\infty$ Approximate g by $h \in \mathscr{F}_{exp}$ with $\frac{\varepsilon}{2}$ error, and replace each exp with a 1d σ -based net

$$\mathcal{F}$$
 with $f(x) \neq 0$.

3. Separates points: for each $x \neq x'$, there is at least one $f \in \mathcal{F}$ with $f(x) \neq f(x')$.





Stone-Weierstrass Theorem: Let \mathcal{F} be a set of functions such that

- 1. Each $f \in \mathcal{F}$ is continuous.
- 2. For each x, there is at least one $f \in \mathcal{F}$
- 4. \mathscr{F} is an algebra: for $f, g \in \mathscr{F}$, $\alpha f + g \in \mathscr{F}$ and $fg = (x \mapsto f(x)g(x)) \in \mathscr{F}$. Conditions hold for $\sigma_1 = \exp, \sigma_2 = \operatorname{Id}$, so that $\mathscr{F}_{\exp} = \{x \mapsto \sum a_i \exp(w_i^{\mathsf{T}} x)\}$ If $\sigma : \mathbb{R} \to \mathbb{R}$ is continuous, $\lim \sigma(z) = 0$, $\lim \sigma(z) = 1$, works too:

 $z \rightarrow -\infty$ Approximate g by $h \in \mathscr{F}_{exp}$ with $\frac{\varepsilon}{2}$ error, and replace each exp with a 1d σ -based net

Generally: universal approximator iff σ is **not** a polynomial

$$\mathcal{F}$$
 with $f(x) \neq 0$.

3. Separates points: for each $x \neq x'$, there is at least one $f \in \mathcal{F}$ with $f(x) \neq f(x')$.

$$z \rightarrow -\infty$$





SSBD chapter 20:

• 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$

SSBD chapter 20:

- - (remember that computers always represent things as $\{0,1\}^d$...)

• 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$

SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$ • (remember that computers always represent things as $\{0,1\}^d$...) • ...but, it takes exponential width to do that

SSBD chapter 20:

- 2 layer nets with sign activations can represent all functions $\{\pm 1\}^d \rightarrow \{\pm 1\}$
 - (remember that computers always represent things as $\{0,1\}^d...$)
- ...but, it takes exponential width to do that
- ...but, there's a network of size $\mathcal{O}(T^2)$ that can implement all boolean functions that can be computed in maximum runtime T

Limits of universal approximation

- Curse of dimensionality: usually requires # of units exponential in dimension Also usually requires exponential norm of weights
- Doesn't say anything about whether ERM finds a good network, just that one exists lacksquare• Let alone anything about whether (S)GD finds it

