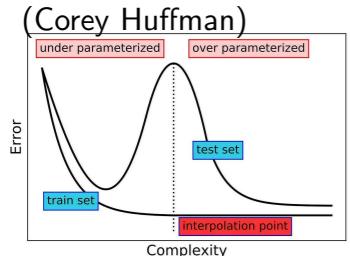


## Commentary on Double Descent



## Connections Between Differentially Private Learning and Online Learning (Elias Lindgren)

**Definition** (Private Learning Algorithm) A randomized learning algorithm

$$A : (X \times \{\pm 1\})^m \rightarrow \{\pm 1\}^X$$

is  $\epsilon, \delta$ -differentially private if for every two samples

$S, S' \in (X \times \{\pm 1\})^m$  that disagree on a single example, the output distributions  $A(S)$  and  $A(S')$  are  $(\epsilon, \delta)$ -indistinguishable.

**Littlestone Dimension** of  $\mathcal{H}$  is largest  $d$  s.t.  $\mathcal{H}$  shatters all depth  $d$  trees

Finite Littlestone Dimension  $\equiv$  online learnable

## Understanding a Randomized Extension of Littlestone Dimension

(Mary Monroe)

**Lemma 1.** For any hypothesis class  $\mathcal{H}$ ,  $M^*(\mathcal{H}) \geq RL(\mathcal{H})$ .

*Proof.* We construct an adversary who forces all randomized algorithms to make at least  $RL(\mathcal{H})$  mistakes. Pick any finite tree  $T$  shattered by  $\mathcal{H}$ . Then the adversary constructs a random walk down

## On the cross-validation bias due to unsupervised preprocessing

(Bart Chen)

**Theorem 2.** Let each element of the matrix  $X \in \mathbb{R}^{(n+m) \times p}$  and each coefficient  $\beta_j$  of  $\beta \in \mathbb{R}^p$  is drawn i.i.d. from a standard normal distribution, and let  $Y = X\beta \in \mathbb{R}^{n+m}$ . If  $j, k \in [p]$  are such that  $\|X_{1:n+m,j}\| \geq \|X_{1:n+m,k}\|$  for all  $i \in [p] \setminus \{j, k\}$ , and an OLS model with no intercept is trained on  $\tilde{X} \in \mathbb{R}^{n \times 2}$  and  $Y_{1:n}$  where  $\tilde{X}_{1:n,1} = X_{1:n,j} = \mathbf{p}_1$  and  $\tilde{X}_{1:n,2} = \mathbf{p}_2$ , then the bias of the validation error with respect to the generalization error of the model is:

$$\text{bias} = \mathbb{E}[(p-2) \left( (X_{n+1,j}^2 - 1)(Z_{j,0}^{(1)})^2 + (X_{n+1,k}^2 - 1)(Z_{j,0}^{(2)})^2 \right) - X_{n+1,j}^2 - X_{n+1,k}^2 + 2]$$

Where  $j_0 \in [p] \setminus \{j, k\}$  is arbitrary,  $Z_j^{(1)} = \frac{(\|\mathbf{p}_1\|^2 \mathbf{p}_1^T - \mathbf{p}_2^T \mathbf{p}_1 \mathbf{p}_1^T) \mathbf{p}_1}{\det(X^T X)}$ , and  $Z_j^{(2)} = \frac{(\|\mathbf{p}_1\|^2 \mathbf{p}_1^T - \mathbf{p}_2^T \mathbf{p}_1 \mathbf{p}_1^T) \mathbf{p}_2}{\det(X^T X)}$  for all  $j \in [p] \setminus \{j, k\}$

## Neural Tangent Kernel

(Madi Yerlanov)

$$\nabla_x f_\theta(x) = -\frac{1}{m} \sum_{i=1}^m \nabla_\theta f_\theta(x)^T \nabla_\theta f_\theta(x_i) \nabla_x \ell(f_\theta(x_i), y_i)$$

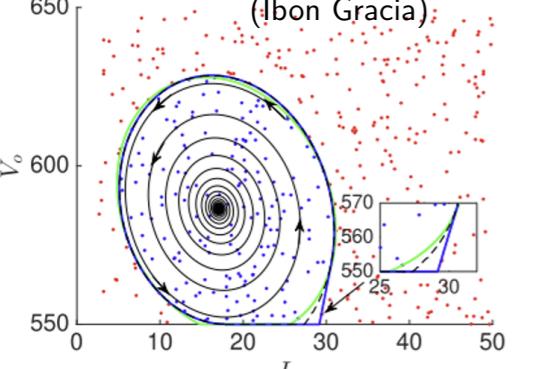
the Neural Tangent Kernel  $\Theta : \mathbb{R}^{n_0} \times \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_0 \times n_0}$  is defined as:

$$\Theta(x, x') = \nabla_\theta f_\theta(x)^T \nabla_\theta f_\theta(x'),$$

$$\Theta_{i,j}^{(l)}(x, x') = \sum_{p=1}^P \frac{\partial f_l(x)}{\partial \theta_p} \frac{\partial f_l(x')}{\partial \theta_p},$$

## Data-driven Invariant Sets of Black-Box Systems

(Ibon Gracia)



## Continuum-armed Bandits: Incentivization and Non-Stationarity (Amit Kiran Rege)

**Theorem 1.** The expected regret for Algorithm (i) on the metric space  $(A, \Phi)$ , with  $\psi$ -covering  $A_0$  of size  $|A_0| \leq \lambda_d/\psi^d$  where  $d$  is the covering dimension of  $A$ , for some constant  $\lambda_0 > 0$ , is

$$\mathbb{E}[R_T^A] \leq \lambda_e^c \cdot T^{(d+2)/(d+3)} \cdot L^{d/(d+3)} \cdot (\log T)^{1/(d+3)}$$

where  $\lambda_e^c = (8\lambda_0^3 \lambda_d)^{1/(d+3)}$ . The expected compensation for the same setting, for some constant  $\lambda > 0$  is

$$\mathbb{E}[C_T^A] \leq \lambda_e^c \cdot T^{(d+1)/(d+3)} \cdot L^{2d/(d+3)} \cdot (\log T)^{2/(d+3)}$$

where  $\lambda_e^c = (4\sqrt{2})(8)^{(1-d)/(2d+6)}(\lambda_d)^{2/(d+3)}$

Algorithm: **CAL(n)**

1. While  $t < n$  and  $m < 2^n$
2.  $m \leftarrow m + 1$
3. If  $X_m \in \text{DIS}(V)$
4. Request label  $Y_m$ ; let  $V \leftarrow \{h \in V : h(X_m) = Y_m\}$ ,  $t \leftarrow t + 1$
5. Return any  $\hat{h} \in V$

### B. Rademacher Processes and Data Dependent Bounds on Excess Risk

It is important to be able to upper bound the excess risk with bounds that are data-dependent and not distribution dependent since the joint distribution of  $(X, Y)$ ,  $\mathcal{D}$  is unknown. The main idea of [3] is to characterize the complexity of class  $\mathcal{F}$  by using Rademacher complexity. In general,

## Active Learning

(K. Aditi)

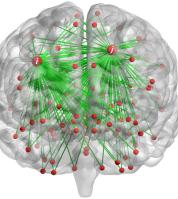
Proved:

- $\mathcal{L}_m$  is a bounded, closed, and convex subset in  $\mathbb{R}^{m^2}$  of dimension  $m(m-1)/2$
- $m_G(x) = P_{\mathcal{L}_m}(E[s_G(x)L])$
- $\hat{m}_G(x) = P_{\mathcal{L}_m} \left( \frac{1}{n} \sum_{k=1}^n s_{kG}(x)L_k \right)$

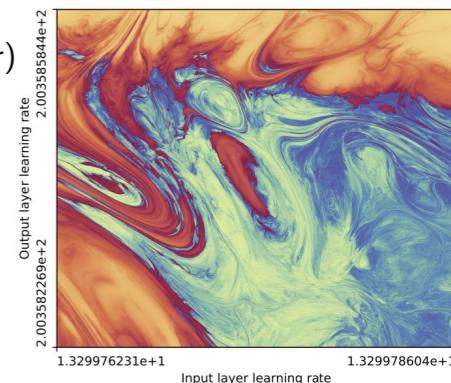
Thus, the population and sample global regressions exist and are unique for all  $x$ . (And it is a convex learning problem!)

## Network-Valued Regression

(Addie McCurdy)

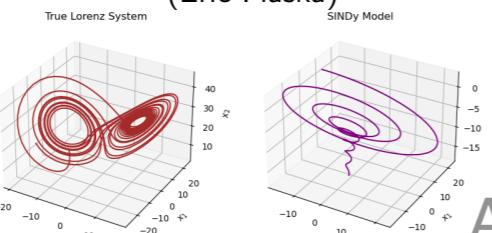


tanh full batch (fractal dim 1.66)



## Explorations in SINDy

(Eric Flaska)



## Learning Finite Automata

(Jasdeep Singh)

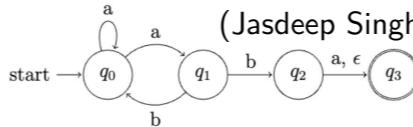


Figure 2: Nondeterministic Finite Automaton (NFA)

## Neural Networks as Gaussian Processes

(Walter Virany)

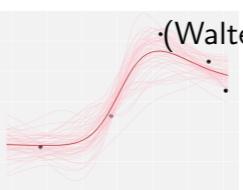


Figure: Ensemble of neural networks converges to a Gaussian process. Created using the Neural Tangents package ([Novak et al., 2019]).

## Review of “Towards Size-Independent Generalization

### Bounds for Deep Operator Nets

(Kal Parvanov)

#### Theorem (2.2)

Considering the same class of DeepONets as in Theorem 2.1 and using the Huber loss  $\ell_{H,\delta}(x) := \begin{cases} \frac{1}{2}x^2 & \text{for } |x| \leq \delta \\ \delta \cdot (|x| - \frac{1}{2}\delta) & \text{for } |x| > \delta \end{cases}$  with  $\delta = (\frac{1}{2})^{n-1}$  as the loss function.

The expectation over data of the supremum of the generalization error over the above class of DeepONets can then be bounded by,

$$\mathcal{O} \left( \frac{C_{n,n-1}}{\sqrt{m}} \left( \prod_{j=2}^{n-1} C_{j,-j} \right) M_{j,B} M_{X,T} \right)$$

where  $C_{n,n-1}$  and  $C_{j,-j}$  are the constants defined in Theorem 2.1.

### Stability based Generalization

### Bounds for Adversarial Training

(Karan Muvvala)

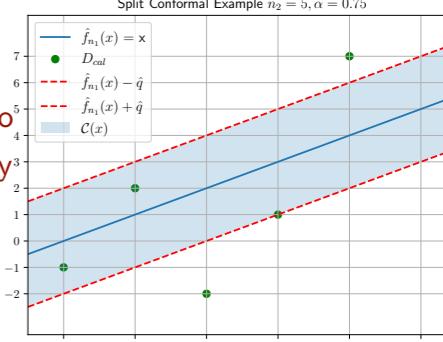
Table: Comparison of the upper bounds of  $\epsilon_{gen}$ .

Paper	Assumptions	Gen. Gap
Farnia et al. [4]	Convex-strongly concave	$O(T)$
This work [1]	Convex - nonconcave	$O(\epsilon T + T)$
Alternate bound	Convex-nonconcave	$O(\epsilon T^2)$

Takeaways:

- This project: Bounds are  $\epsilon$ -dependent. As  $\epsilon \rightarrow 0$  bounds reduce standard training bounds
- As  $n \rightarrow \infty$  [1] suggests  $O(\epsilon T) \rightarrow$  robust test accuracy should decrease
- Alternate bound: as  $n \rightarrow \infty$ ,  $\epsilon_{gen} \rightarrow 0$  - robust test accuracy should increase (contrasting results)

Split Conformal Example  $n_2 = 5, \alpha = 0.75$



## Probably Approximately Correct

### Constrained Learning

(Yiting Chen)

#### Constrained Statistical Learning Problem

$$P^* = \min_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{D}_0} [\ell_0(h, (x, y))]$$

subject to  $\mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell_i(h, (x, y))] \leq c_i, i = 1, \dots, k,$

## Effective Theory of Deep Linear

### Networks at Initialization

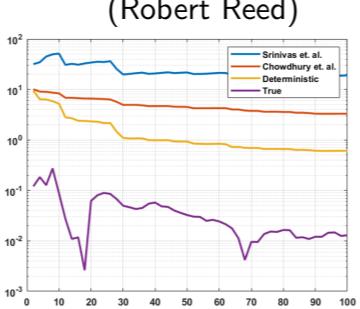
(Noah Francis)

$$\begin{aligned} \mathbb{E} \left[ z_{i_1:i_2}^{(1)} z_{j_1:j_2}^{(1)} \right] &= \sum_{j_1, j_2=1}^{n_0} \mathbb{E} \left[ W_{i_1 j_1}^{(1)} x_{j_1:i_2} W_{i_2 j_2}^{(1)} x_{j_2:i_2} \right] \\ &= \sum_{j_1, j_2=1}^{n_0} \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{C_W}{n_0} x_{j_1:i_1} x_{j_2:i_2} \\ &= \delta_{i_1 i_2} \frac{C_W}{n_0} \sum_{j_1=1}^{n_0} x_{j_1:i_1} x_{j_1:i_2} \\ &=: \delta_{i_1 i_2} W G_{i_1 i_2}^{(0)}. \end{aligned}$$

## A PAC framework for regression...

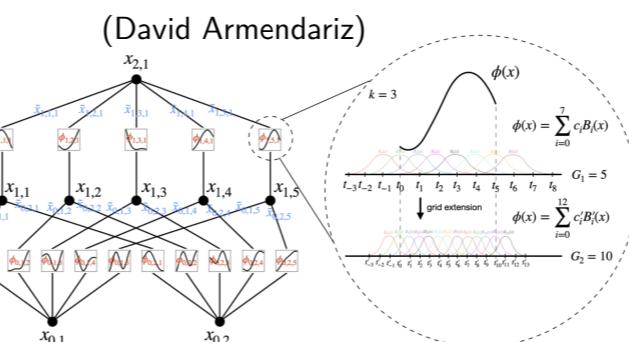
### via kernel methods

(Robert Reed)



## Kolmogorov-Arnold Networks

(David Armendariz)



## Conformal Prediction:

### Background and Comparison to

### Gaussian Processes Uncertainty

### Quantification

(Tyler Jensen)