

# Towards more trust in overparameterized models for high-dimensional data

---

Fanny Yang, Statistical Machine Learning Group, D-INFK, ETH Zurich

ELLIS Doctoral Symposium, 28. September 2021

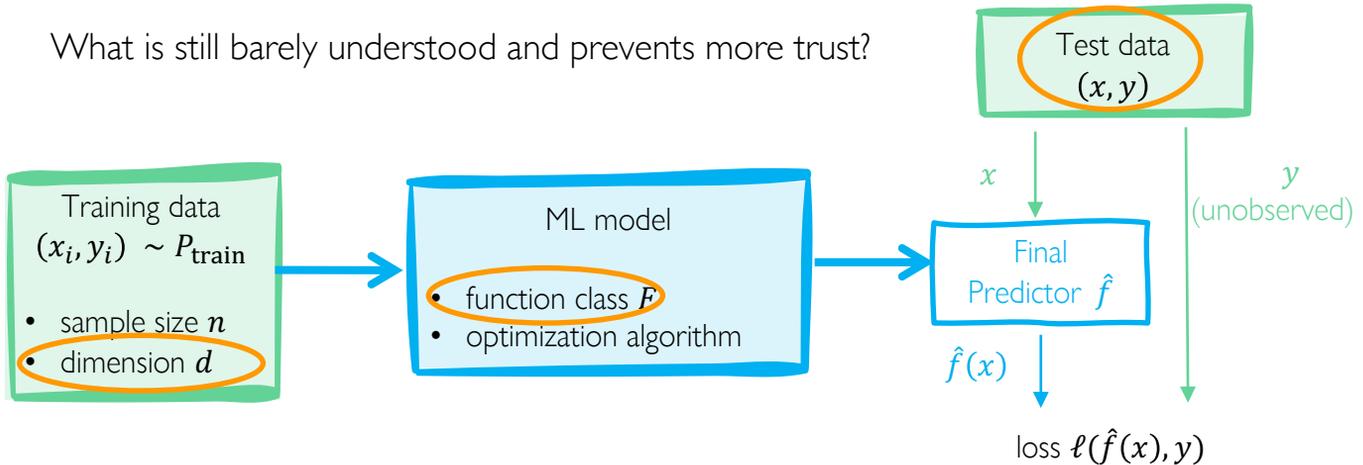
Soon AI will be replacing some humans...



...long way to go until really trustworthy 2

# Supervised learning for modern ML

What is still barely understood and prevents more trust?



## Two important settings for which understanding is missing

- 1 High-dimensionality: Data dimension  $d$  is large  
Overparameterization: large model classes  $F$  that can perfectly fit data
- 2 Reliability: how model acts when test data  $(x, y) \not\sim P_{\text{train}}$

# Plan for Part ①: high-dimensional regression

- Setting up: Regression in the modern data regime
  - Recap: regression in the classical regime
  - Regression in the modern data regime
  - Overparameterized models and estimators
- Two examples where classical intuition fails

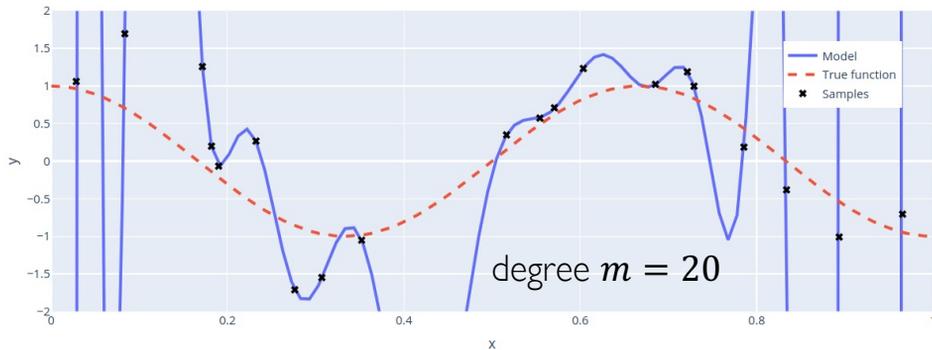
# Recap: Regression with the square loss

- **Data generation:** covariates  $x \sim P_x$ , labels  $y = f^*(x) + \epsilon$  with  $x \in R^d$  and  $y \in R$ , noise  $\epsilon$
- **Observe:** training set  $D$  with  $n$  i.i.d. data points
- **Goal:** Find  $\hat{f}_D$  that is close to  $f^*$  in some model class  $F$
- **Vanilla estimator:** minimizer of Mean Square Error

$$\hat{f}_D = \operatorname{argmin}_{f \in F} \sum_{i=1}^n (y_i - f(x_i))^2$$

# Classical setting

- sample size  $n \gg d$  dimension
- Conventional wisdom: best fit with large model  $\rightarrow$  large variance

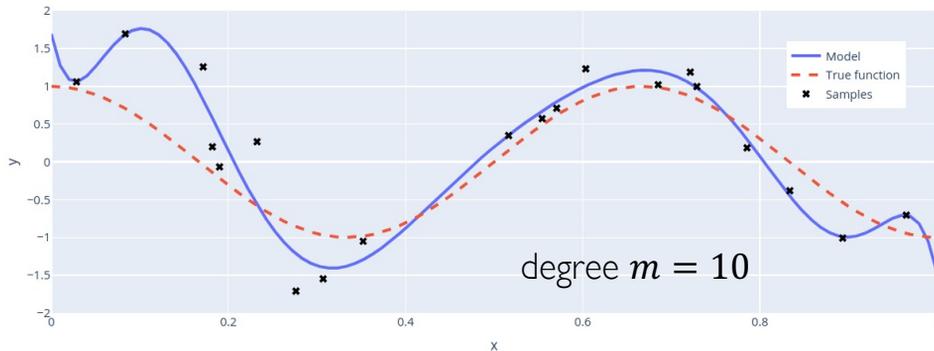


# Classical setting

- sample size  $n \gg d$  dimension
- Conventional wisdom: best fit with large model  $\rightarrow$  large variance



regularize model complexity avoiding perfect fit  $\rightarrow$  smaller variance

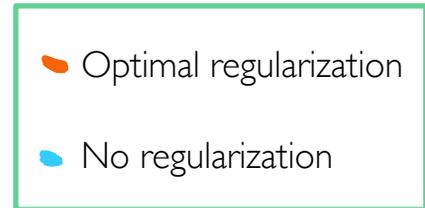
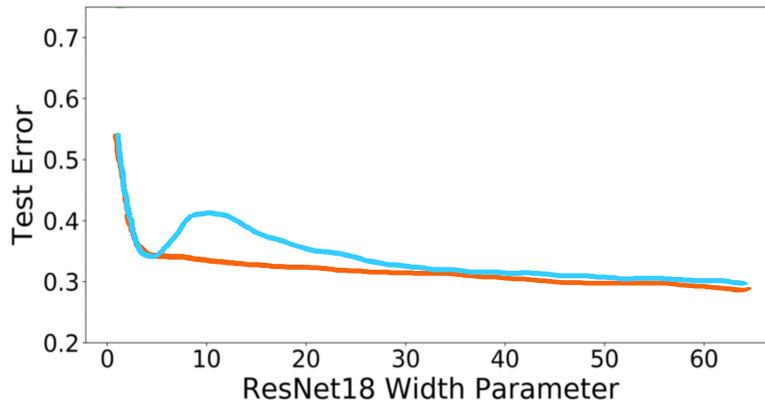


# Modern setting

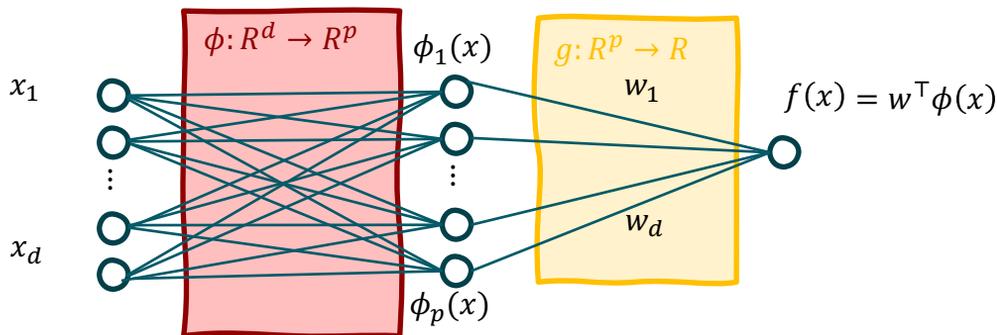
- sample size vs. dimension:  
e.g. ImageNet  $n \sim 10^6$ ,  $d > 10^4$ , gene data  $n \sim 10^1$  to  $10^4$ ,  $d > 10^4$

→ high-dimensional regime  $n \asymp d^\alpha$ , here:  $\alpha \in (0,2)$

- Modern practice: use overparameterized models (can fit data perfectly)  
**and** without regularization still seems to generalize at least as well



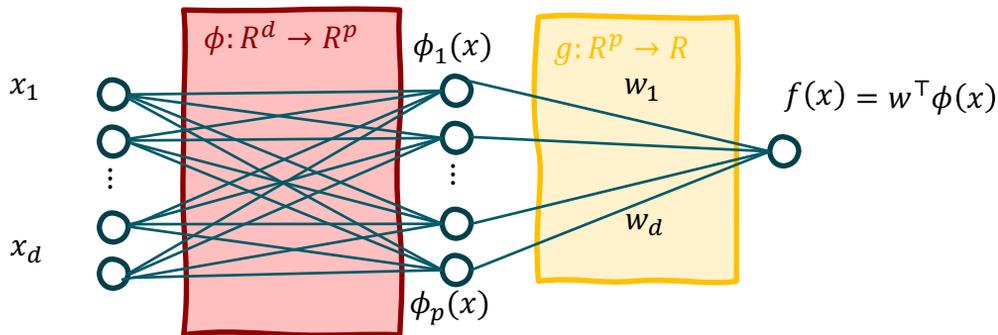
# Overparameterized models in a nutshell



Unifying diagram for overparameterized models for high dimensional data  $d \gg n$ :

- Linear models are already overparameterized for  $p = d$  for  $\phi(x) = x$
- Kernels can fit nonlinear functions using  $p > d$  fixed nonlinear features  $\phi$
- Neural networks can fit nonlinear functions using  $p$  trained features  $\phi$

# Overparameterized models in a nutshell



For overparameterized models  $F$ , many models achieve  $MSE = 0$

→ initialized at 0, first order method often yields **minimum-norm interpolator**

$$\hat{f}_D = \operatorname{argmin}_{f \in F} \|f\| \text{ s.t. } y_i = f(x_i) \text{ for all } i$$

for some norm (depends on  $F$  and algorithm)

# Evaluation of the estimators



What's the expected loss on a test sample or

$$\text{estimation error } \left\| \hat{f}_D - f^* \right\|_{L^2(P_x)}^2 = E_{X \sim P_x} \left( \hat{f}_D(X) - f^*(X) \right)^2$$

dependent on the model class  $F$ , algorithm,  $n, d$ ?

Study the average error (over data  $D$ )

using bias-variance decomposition:

$$E_D \left\| \hat{f}_D - f^* \right\|^2 = E_D \underbrace{\left\| \hat{f}_D - E_D \hat{f} \right\|^2}_{\text{variance}} + \underbrace{\left\| E_D \hat{f}_D - f^* \right\|^2}_{\text{bias}}$$

# Plan for Part ①: high-dimensional regression

- Setting up: Regression in the modern data regime
- Two examples where classical intuition fails

- Example I: Linear models where  $p = d$   
Minimum- $\ell_2$ -norm interpolation when  $d \asymp n$   
Intuition from classical setting  $d < n$ : larger  $n \rightarrow$  smaller variance

- Example II: Nonlinear models via kernels with  $p = \infty$   
Kernel estimators for  $d^\alpha \asymp n$   
Intuition for fixed  $d$ : you can learn nonlinear functions

Intuition from classical theory does not hold!

# Linear least-squares regression for $d \gg n$

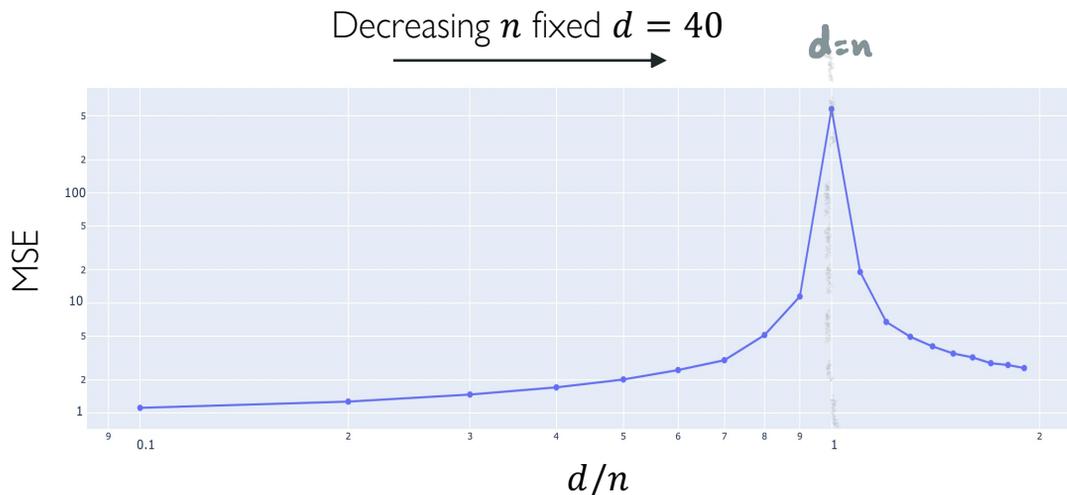
Given data matrix  $X = \begin{pmatrix} - & x_1 & - \\ & \vdots & \\ - & x_n & - \end{pmatrix}$  with rows  $x_i \sim N(0, I_d)$  and labels  $y = Xw^* + \epsilon$

For  $d < n$ : MSE minimizer

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|y - Xw\|_2$$

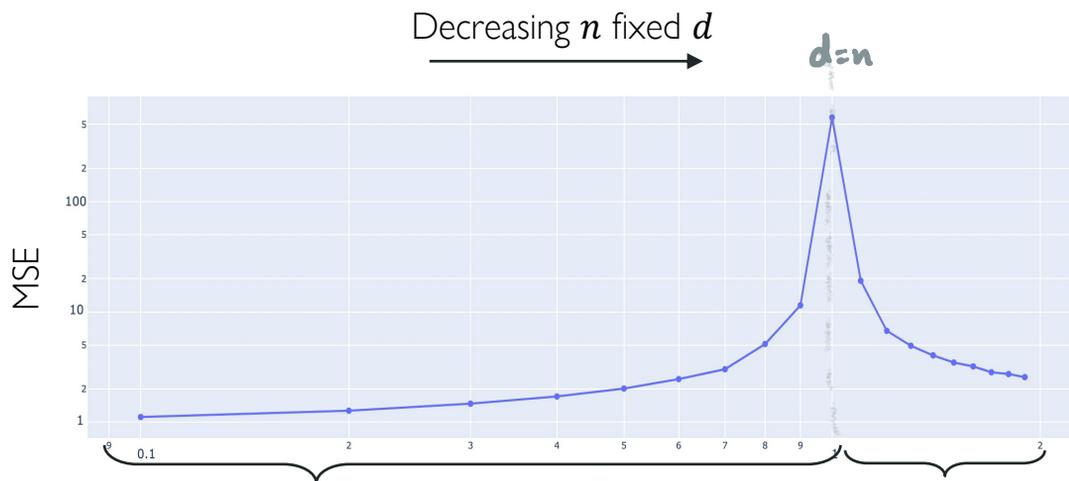
For  $d > n$ : min- $\ell_2$ -norm interpolator

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \|w\|_2 \quad \text{s.t. } y = Xw$$



# Linear least-squares regression for $d \gg n$

Goal now: Compare and build intuition for the two regimes



$n$  decreases for fixed  $d$   
 $\rightarrow$  MSE increases



?  $n$  decreases for fixed  $d$   
 $\rightarrow$  MSE decreases

# Recap: Bias and variance for linear regression

For  $d < n$ : MSE minimizer

$$\hat{w} = \operatorname{argmin}_w \|y - Xw\|_2$$

For  $d > n$ : min- $\ell_2$ -norm interpolator

$$\hat{w} = \operatorname{argmin}_w \|w\|_2 \quad \text{s.t. } y = Xw$$

Solution for both can be written as

$$\hat{w} = (X^T X)^\dagger X^T y = \underbrace{(X^T X)^\dagger X^T X w^*}_{\text{blue}} + \underbrace{(X^T X)^\dagger X^T \epsilon}_{\text{orange}}$$

$\Pi_X w^* = E_D \hat{w}$   
 $\rightarrow$  determines bias

$=: w_n$  fits noise  
 $\rightarrow$  determines variance

$$\text{Average MSE: } E_D \|\hat{w}_D - w^*\|^2 = \underbrace{\|E_D \hat{w} - w^*\|^2}_{\text{bias}} + E_D \|\hat{w}_D - E_D \hat{w}\|^2$$

= Test loss for  $X \sim N(0, I)$

$$\text{bias } \|\Pi_X w^* - w^*\|^2 \quad \text{variance } E_D \|w_n\|^2$$

# Recap: Bias and variance for linear regression

For  $d < n$ : MSE minimizer

$$\hat{w} = \operatorname{argmin}_w \|y - Xw\|_2$$

For  $d > n$ : min- $\ell_2$ -norm interpolator

$$\hat{w} = \operatorname{argmin}_w \|w\|_2 \quad \text{s.t. } y = Xw$$

Solution for both can be written as

$$\hat{w} = (X^T X)^\dagger X^T y = \underbrace{(X^T X)^\dagger X^T X w^*}_{\text{determines bias}} + \underbrace{(X^T X)^\dagger X^T \epsilon}_{\text{determines variance}}$$

$\Pi_X w^* = E_D \hat{w}$   
→ determines bias

$=: w_n$  fits noise  
→ determines variance

$$\Pi_X w^* = w^* \rightarrow \text{bias} = 0$$

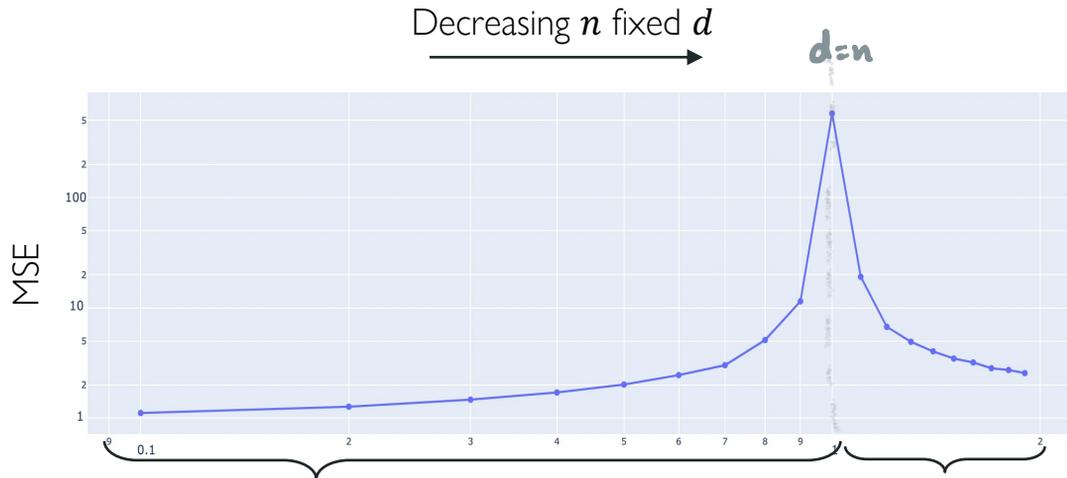
MSE minimizer for noise fit  
 $w_n = \operatorname{argmin}_w \|\epsilon - Xw\|_2$

$$\text{bias } \|\Pi_X w^* - w^*\|^2 \neq 0$$

min- $\ell_2$ -norm interpolator of noise  
 $w_n = \operatorname{argmin} \|w\|_2 \text{ s.t. } \epsilon = Xw$

# Example I: Linear LS regression for $d \gg n$

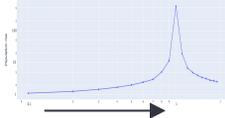
Goal now: Compare and build intuition for the two regimes



$n$  decreases for fixed  $d$   
→ variance increases



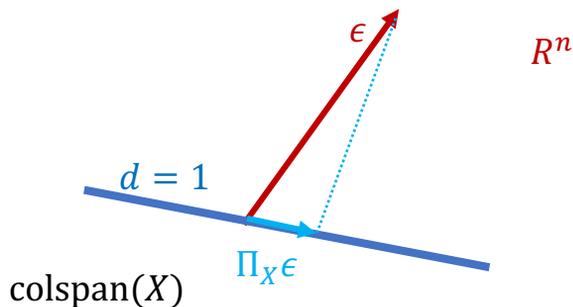
?  $n$  decreases for fixed  $d$   
→ MSE decreases

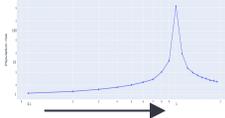


Variance for  $d < n$  increases with  $d/n$

$$w_n = \operatorname{argmin}_w \left\| \begin{matrix} n \\ \epsilon \end{matrix} - \begin{matrix} d \\ X \\ n \end{matrix} \begin{matrix} w \\ d \end{matrix} \right\| \Rightarrow Xw_n = \Pi_X \epsilon$$

projection onto  $\operatorname{colspan}(X)$





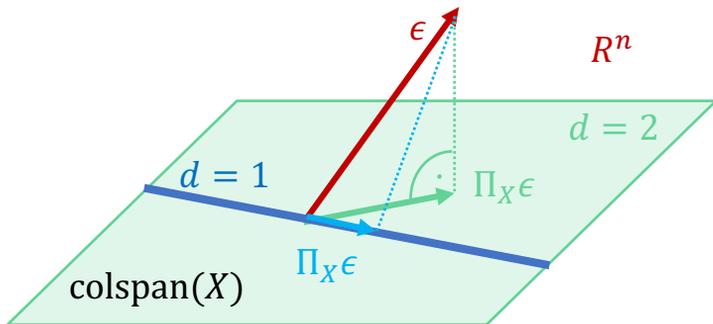
Variance for  $d < n$  increases with  $d/n$

$$w_n = \operatorname{argmin}_w \left\| \begin{matrix} n \\ \epsilon \end{matrix} - \begin{matrix} d \\ n \\ X \end{matrix} \begin{matrix} w \\ d \end{matrix} \right\|$$



$$Xw_n = \Pi_X \epsilon$$

projection onto  $\operatorname{colspan}(X)$



$\frac{d}{n}$  increases (increasing  $d$  fixed  $n$ )

→  $\operatorname{colspan} X$  larger

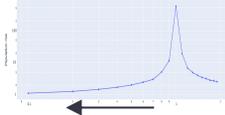
→ projection  $\Pi_X \epsilon$  closer to  $\epsilon$

→ norm of  $\Pi_X \epsilon$  increases

For simplicity: orthogonal  $X^T X = nI_d$  such that  $\|w_n\|_2 = \frac{1}{\sqrt{n}} \|Xw_n\|_2 = \frac{1}{\sqrt{n}} \|\Pi_X \epsilon\|_2$

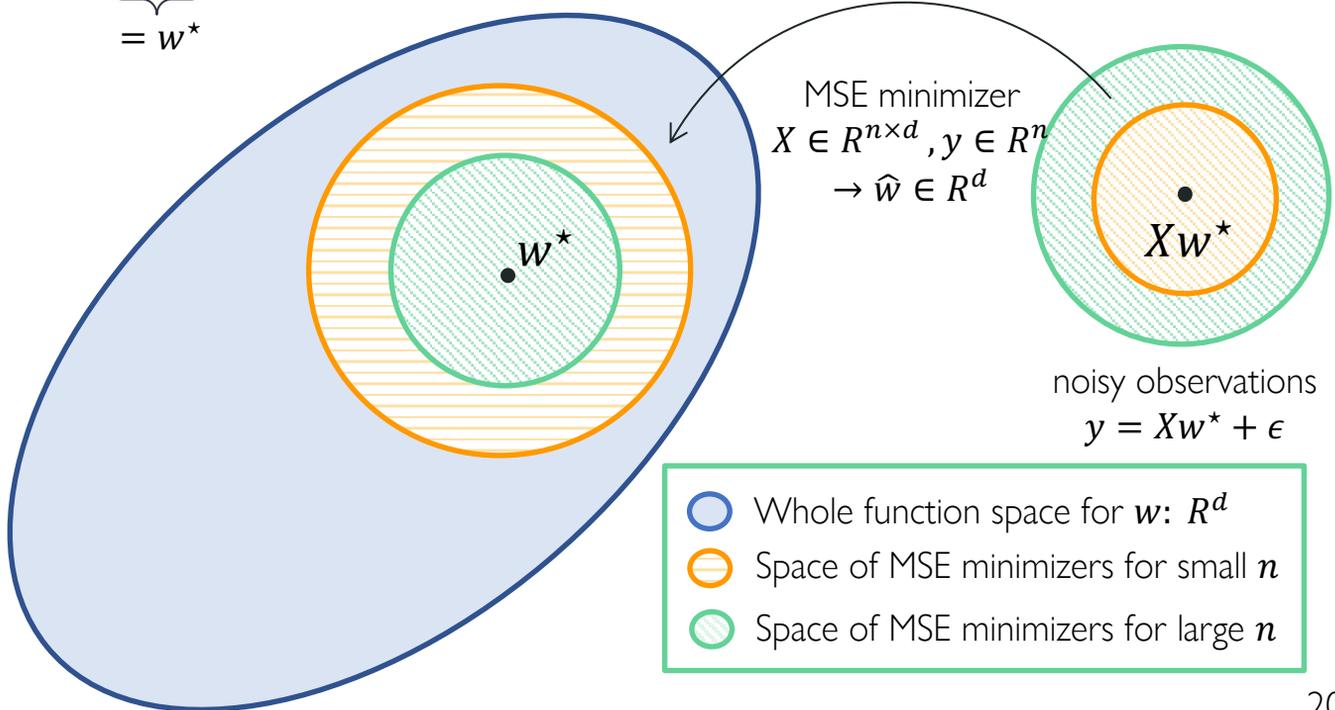


variance  $E_D \|w_n\|^2$  increases with  $\frac{d}{n}$



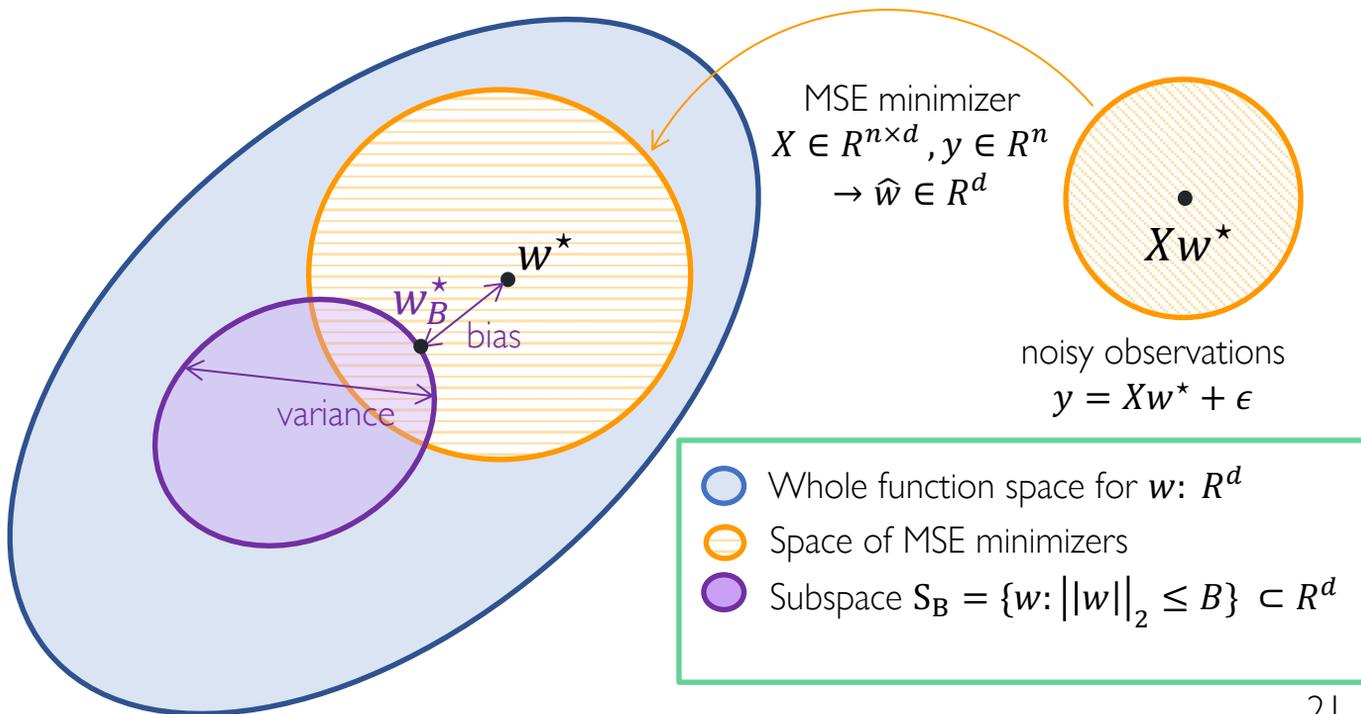
# Variance for $d < n$ decreases with $n$

$\hat{w} = \underbrace{\Pi_X w^*}_{= w^*} + w_n$  with bias  $\|\Pi_X w^* - w^*\|^2 = 0$  and variance  $E_D \|w_n\|^2 = E_D \frac{\|\Pi_X \epsilon\|_2^2}{n}$



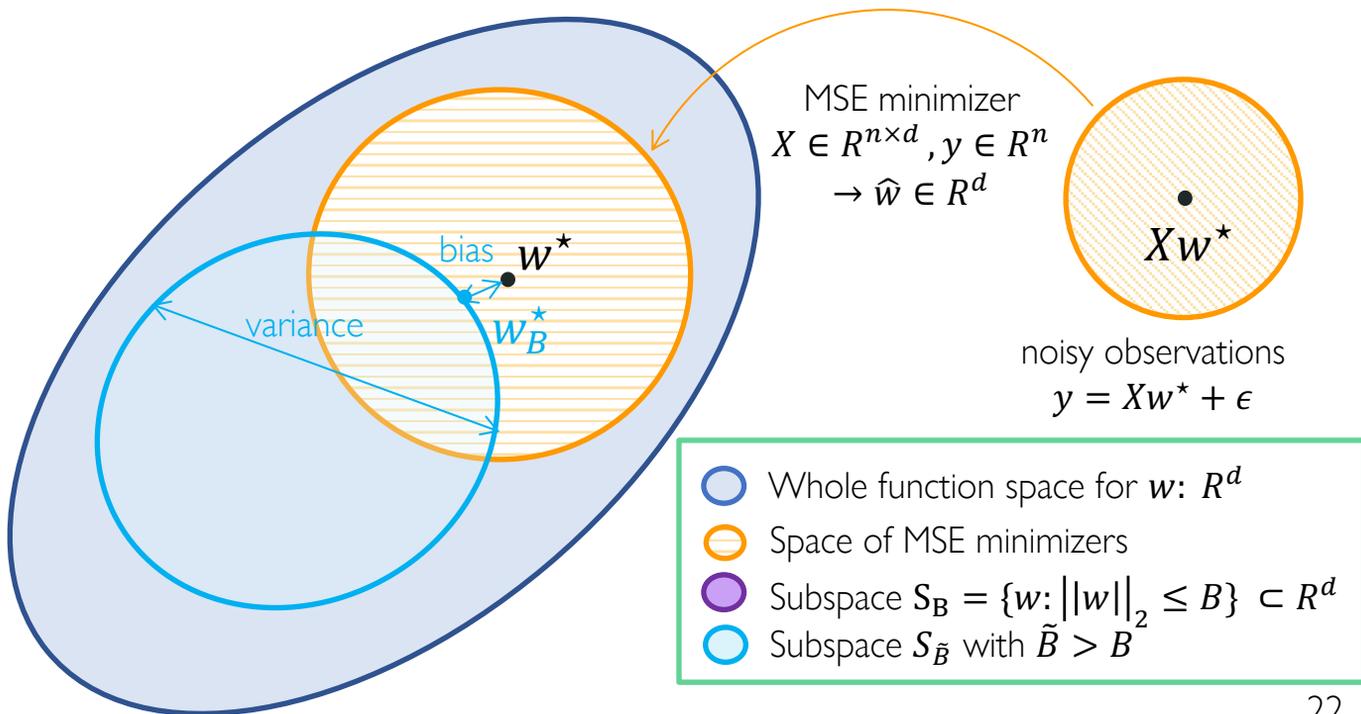
# Recap: Bias variance trade-off for $d < n$

How did we achieve the bias-variance tradeoff? → explicit regularization!



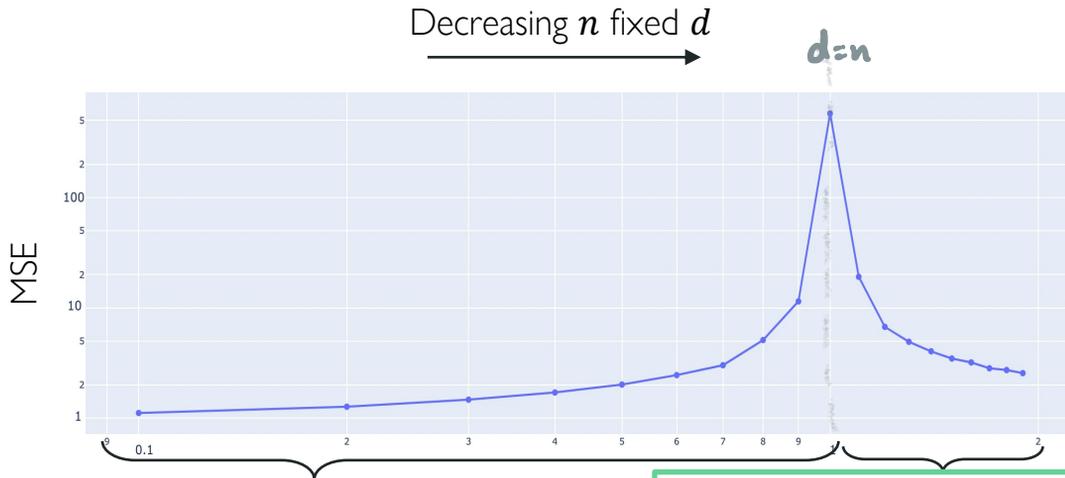
# Recap: Bias variance trade-off for $d < n$

How did we achieve the bias-variance tradeoff? → explicit regularization!



# Example I: Linear LS regression for $d \gg n$

Goal now: Compare and build intuition for the two regimes



$n$  decreases for fixed  $d$   
→ variance increases



$n$  decreases for fixed  $d$   
→ MSE decreases



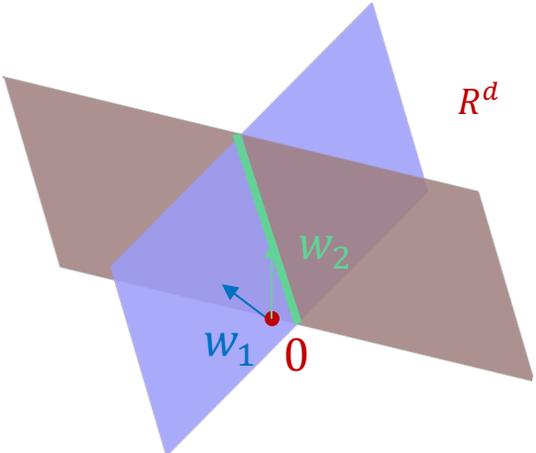
# Bias and variance for $d > n$

Bias  $\|\Pi_X w^* - w^*\|^2$  increases with  $\frac{d}{n}$  since  $\Pi_X w^*$  projects on smaller space  $\rightarrow$  closer to  $w^*$

Variance  $E_D \|w_n\|_2^2$  with  $w_n = \operatorname{argmin}_w \|w\|_2$  s.t.  $n \begin{bmatrix} \epsilon \end{bmatrix} = n \begin{matrix} d \\ X \end{matrix} \begin{bmatrix} w \end{bmatrix} d$

- is a point in the intersection of  $n$  hyperplanes  $x_i^T w = \epsilon_i$

- the one that has **minimum distance**  $\|w_n\|_2$  to origin



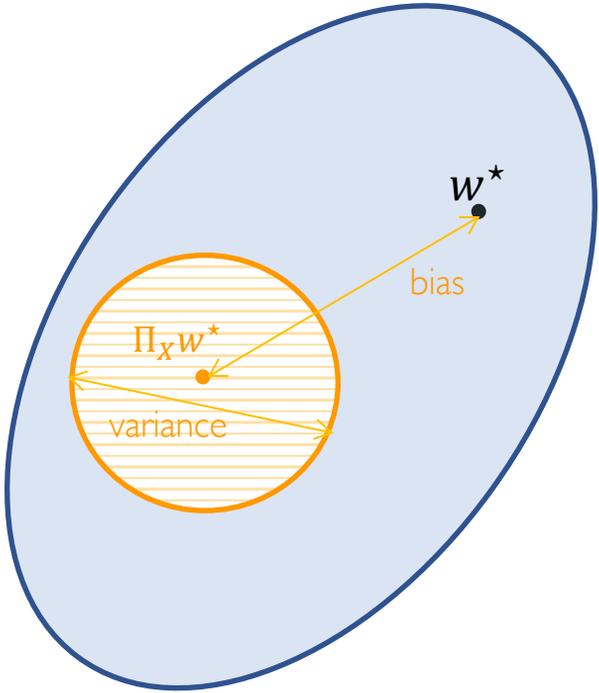
$\frac{d}{n}$  decreases (increasing  $n$  fixed  $d$ )  
 $\rightarrow$  intersection smaller  
 $\rightarrow$  minimum distance =  $\|w_n\|_2$  larger

$\Rightarrow$  variance decreases with  $\frac{d}{n}$  !

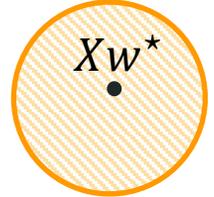


# Bias variance trade-off for $d > n$

$\hat{w} = \Pi_X w^* + w_n$  has bias  $\|\Pi_X w^* - w^*\|^2$  and variance  $E_D \|w_n\|^2$



min- $\ell_2$ -norm interpolators:  
 $X \in R^{n \times d}, y \in R^n \rightarrow \hat{w} \in R^d$



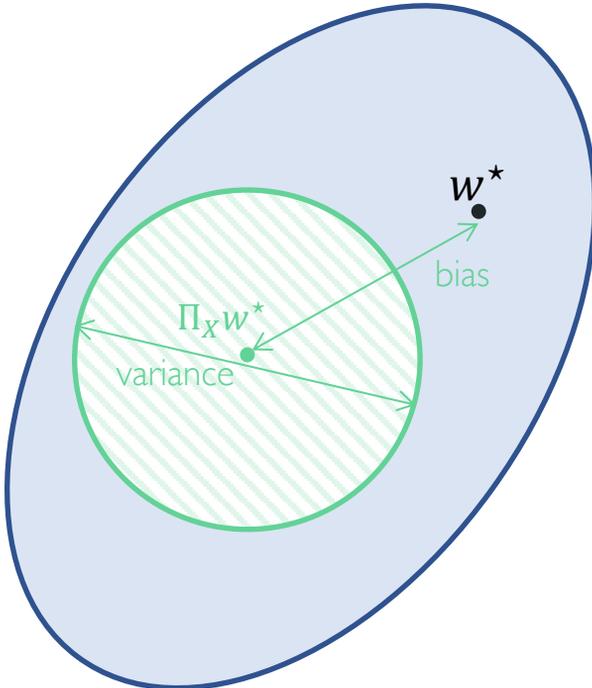
noisy observations  
 $y = Xw^* + \epsilon$

- Whole function space  $R^d$
- Effectively reached for small  $n$

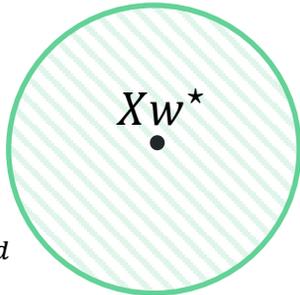


# Bias variance trade-off for $d > n$

$\hat{w} = \Pi_X w^* + w_n$  has bias  $\|\Pi_X w^* - w^*\|^2$  and variance  $E_D \|w_n\|^2$



min- $\ell_2$ -norm interpolators:  
 $X \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n \rightarrow \hat{w} \in \mathbb{R}^d$



noisy observations  
 $y = Xw^* + \epsilon$

- Whole function space  $\mathbb{R}^d$
- Effectively reached for small  $n$
- Effectively reached for large  $n$

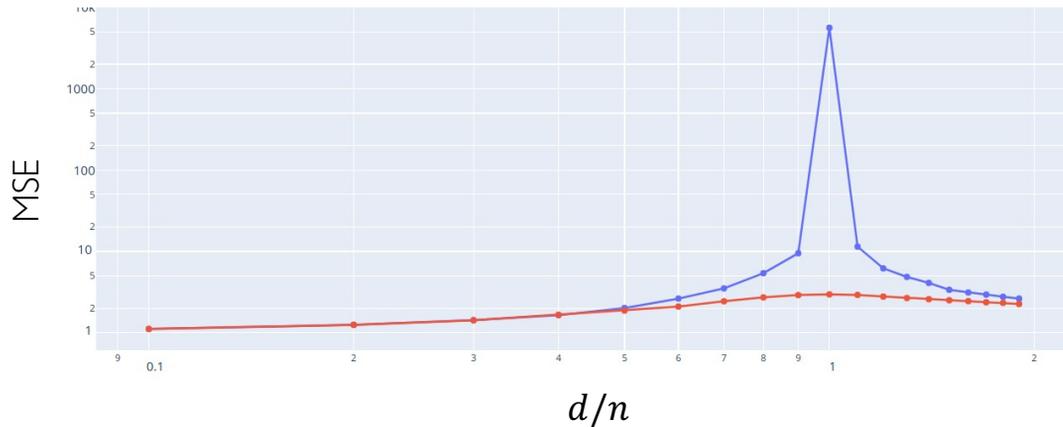
⇒ Bias variance trade-off via  $\frac{d!}{n}$

# High-dimensional (non)-asymptotic rates

- for min- $\ell_2$ -norm and different covariances  $\Sigma_d$ 
  - $d \asymp n$ : Hastie et al. '19
  - $d \asymp n \log n$ : Bartlett et al. '20; Muthukumar et al. '20; Koehler et al. '21
- for min- $\ell_1$ -norm for sparse ground truth  $\Sigma_d = I_d$ 
  - $d \asymp n \log n$ : Chinot et al. '20, WDY '21 (in preparation)
- Related work on logistic regression: Deng et al. '19, Chinot et al. '21
- For adversarial robustness: Javanmard et al. '20, DTAHY '21

# Addendum: What to do in practice?

... just regularize!



— ridge regression with optimal regularization parameter

— minimum  $\ell_2$ -norm interpolator

# Plan for Part ①: high-dimensional regression

- Regression in the modern data regime
- Two examples where classical intuition fails

- Example I: Linear models where  $p = d$

Minimum- $\ell_2$ -norm interpolation when  $d \asymp n$

Intuition in the modern regime: larger  $n \rightarrow$  larger variance



- Example II: Nonlinear models via kernels with  $p = \infty$

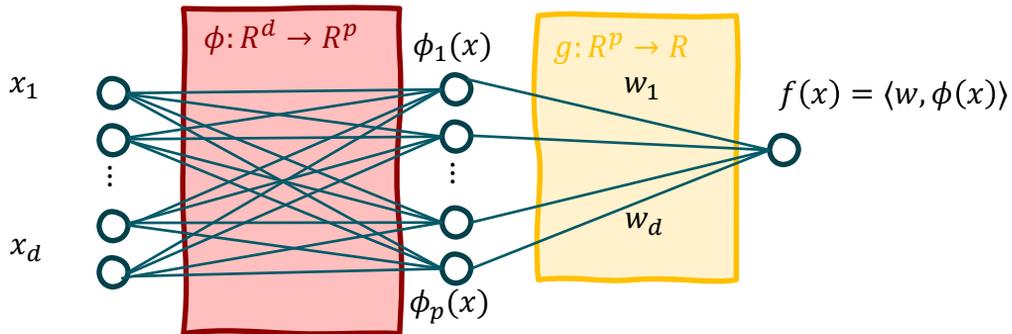
Kernel estimators for  $d^\alpha \asymp n$  for large  $d$

Intuition for fixed  $d$ : you can learn nonlinear functions

Intuition from classical theory does not hold!

# Recap: Kernel regression

- Data generation:  $x \sim P_x$ ,  $y = f^*(x) + \epsilon$  with  $x \in R^d$  and  $y \in R$
- Observe:  $n$  i.i.d. data points in training set  $D$
- Goal: Find  $\hat{f}_D$  that is close to  $f^*$  in kernel class  $F$  induced by a kernel  $K(x, x') = \langle \phi(x), \phi(x') \rangle$  with  $\phi(x) \in R^p$



# Recap: Kernel regression

- Data generation:  $x \sim P_x$ ,  $y = f^*(x) + \epsilon$  with  $x \in R^d$  and  $y \in R$
- Observe:  $n$  i.i.d. data points in training set  $D$
- Goal: Find  $\hat{f}_D$  that is close to  $f^*$  in kernel class  $F$  induced by a kernel  $K(x, x') = \langle \phi(x), \phi(x') \rangle$  with  $\phi(x) \in R^p$

Consider universal kernel estimators for  $p \rightarrow \infty$ . Implementations yield

- Avoiding perfect fitting: kernel ridge regression

$$\hat{f}_D = \operatorname{argmin}_{f \in F} \sum_i^n (y_i - f(x_i))^2 + \lambda \|f\|_F^2$$

- Perfect fitting: minimum-Hilbert-norm interpolator

$$\hat{f}_D = \operatorname{argmin}_{f \in F} \|f\|_F \text{ s.t. } y_i = f(x_i) \text{ for all } i$$

(the solution of gradient descent on square loss upon convergence)

# Kernels and neural networks— previous work

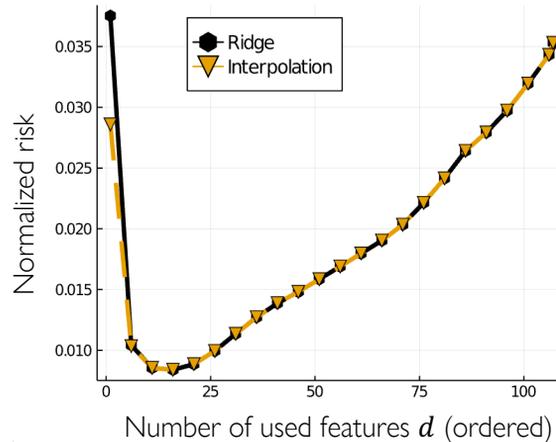
- **Practice:** neural networks can learn highly nonlinear functions very well  
But don't know which interpolating solution it finds!
- **Theory:** kernel estimators can learn arbitrary nonlinear functions  
solution of simple convex problem  $\rightarrow$  analyzable (a lot of previous work)  
where  $\|\cdot\|_K$ -norm induces structure dependent on  $K$  (e.g. smoothness)
- **Recent trend:** infinite-width neural networks (NN) behave like certain kernels (NTK)  
 $\rightarrow$  use kernel learning to understand why NN work well\*

We show: Futile effort when considering vanilla fully-connected NN!

# ... vanilla kernels fail in high dimensions

Vanilla fully connected-NTK behave “similarly” to Laplace kernels<sup>1</sup> - a good thing?

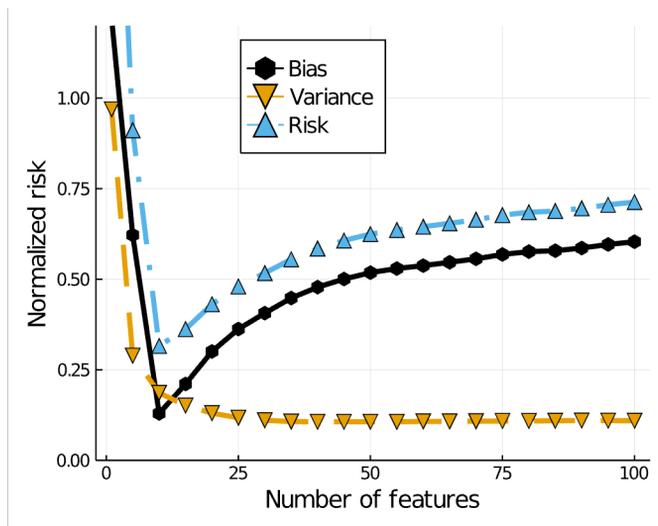
- on CIFAR10: FC-NTK : ~52%<sup>2</sup>, Laplace kernel: ~52%<sup>2</sup>
- Laplace kernel for fitting  $f^*(x)$ : true housing price,  $n = 371$  (on basically noiseless data)



<sup>1</sup>[Geifman et al. '20, Bietti et al. '20] <sup>2</sup>[Lee et al. '20, Belkin et al. '18]

# Bias variance trade-off for kernels

Laplace kernel for  $f^*(x) = 0.5 \sum_{i=1}^4 x^2_{(2i+1)} - \sum_{i=1}^4 x_{(2i)}$  for fixed  $n = 500$



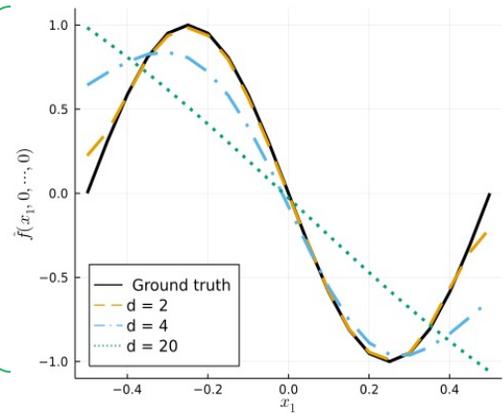
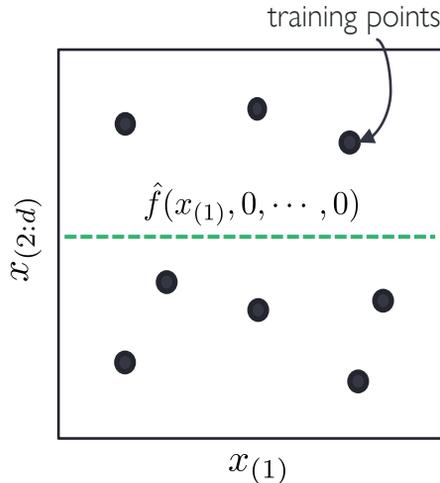
Similar to linear as  $\frac{d}{n}$  increases, but bias increase dominates variance decrease!

Goal now: Characterize the bias as a function of  $d$  vs.  $n$ !

# Kernels learn low-degree polynomials

## Setting

- $x \sim U([-0.5, 0.5]^d)$
- $y = \sin(x_1)$
- $n = 100$  i.i.d.
- Laplace kernel



As dimension grows, the estimator degenerates to a low degree polynomial

# Main result: Polynomial approx. barrier



Theorem (DWY '21, ICML) - *simplified, informal version*

Assume simplest setting  $\mathbf{x} \sim N(\mathbf{0}, I_d)$ , then as  $d, n \rightarrow \infty, \frac{d^\alpha}{n} \rightarrow c > 0$

$$\|\hat{f} - f^*\| \geq \inf_{p \in P_{\leq 2\alpha}} \|f^* - p\| \quad \text{almost surely}$$

where  $P_{\leq 2\alpha}$  is the set of polynomials of degree at most  $2\alpha$ , any  $\alpha > 0$

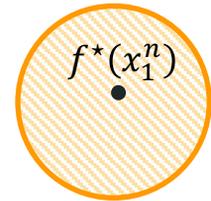
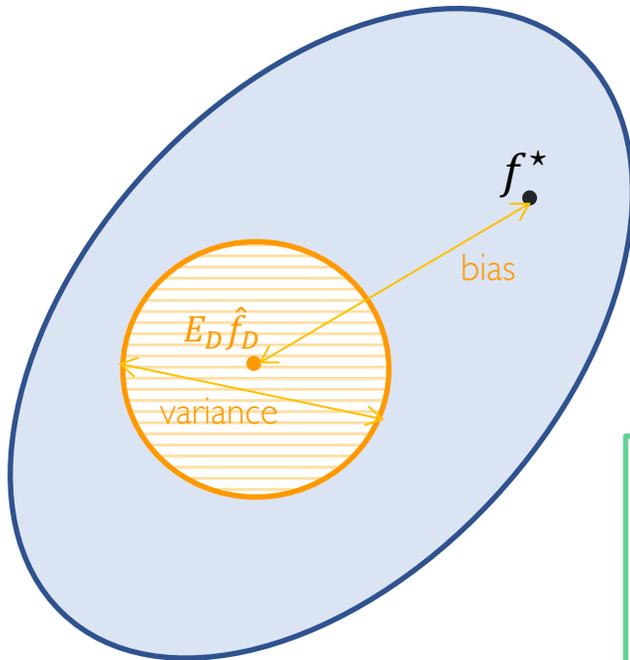
- more generally can assume  $\mathbf{x} \sim N(\mathbf{0}, \Sigma_d)$  and replace  $d$  by  $\text{tr}(\Sigma_d)$
- for **rotationally invariant** kernels  $k_\tau(\mathbf{x}, \mathbf{y}) = h\left(\frac{\|\mathbf{x}\|^2}{\tau}, \frac{\|\mathbf{y}\|^2}{\tau}, \frac{\mathbf{x}^\top \mathbf{y}}{\tau}\right)$ 
  - different functions  $h$  such as RBF (Laplacian, Gaussian), inner product, fully connected NTK of any depth
  - for different scalings  $\tau$

\*poly. barrier restricted to particular distr., scaling  $\tau = d$ :  
[Ghorbani et al. '19, '20]

# Bias variance trade-off for kernel estimators

Same story as for linear models...

kernel estimators:  
 $y \in R^n \rightarrow \hat{f}_D \in F$

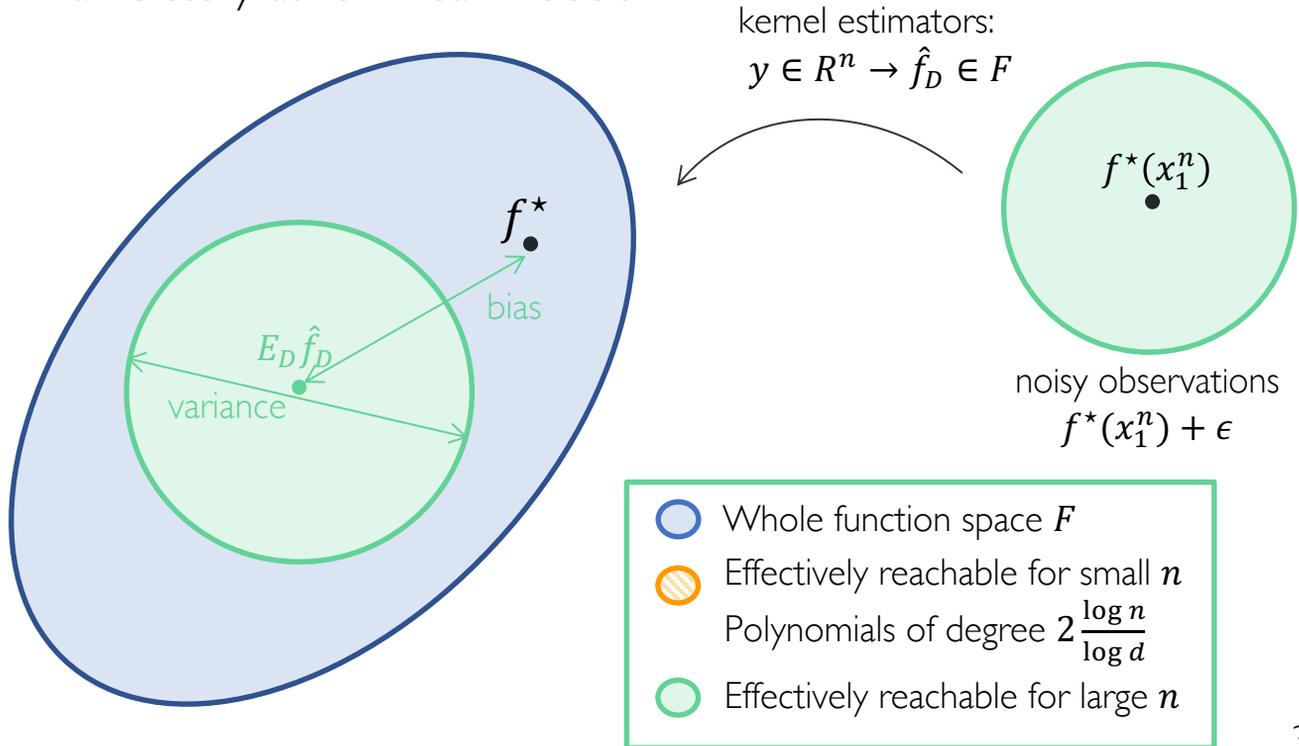


noisy observations  
 $f^*(x_1^n) + \epsilon$

- Whole function space  $F$
- Effectively reachable for small  $n$   
Polynomials of degree  $2 \frac{\log n}{\log d}$

# Bias variance trade-off for kernel estimators

Same story as for linear models...

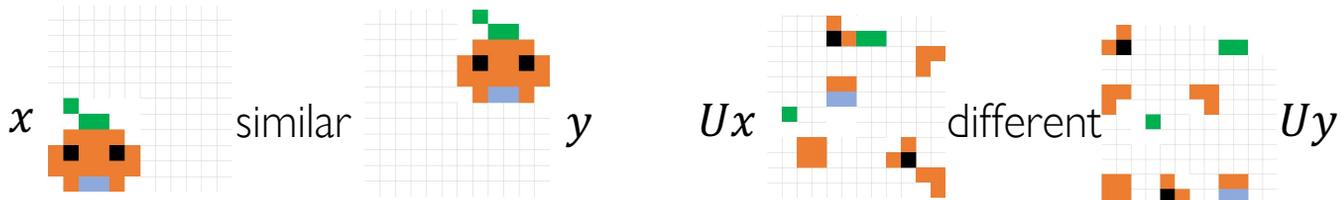


# Why rotation invariance increases bias with $d$

- rotationally invariant kernels should satisfy  $k(x, y) = k(Ux, Uy)$  with  $U$  orthonormal matrix (think of  $U$  permutation matrix)
- kernel method  $\hat{f}$  predicts similar values for similar samples, similarity defined by  $k(x, y)$

**Example I:** ground truth are human labelers, e.g. images

→ good classifier  $\hat{f}$  should find similar what we find similar, i.e.  $k(x, y)$  should reflect that



not using spatial structure, pays attention to every pixel independently

# Why rotation invariance is bad for large $d$

- rotationally invariant kernels should satisfy  $k(\mathbf{x}, \mathbf{y}) = k(U\mathbf{x}, U\mathbf{y})$   
 $U$ : orthonormal matrix (e.g. permutation matrix)
- kernel method predicts similar values for similar samples, similarity defined by  $k(\mathbf{x}, \mathbf{y})$
- further can often write  $k(\mathbf{x}, \mathbf{y}) \sim \sum_{j=0}^{\infty} g_j \left(\frac{\mathbf{x}^\top \mathbf{y}}{d}\right)^j$  in high dimensions for some  $g_j$

**Example II:** ground truth depends only on first variable

→ good predictor should map any two  $\mathbf{x}, \mathbf{y}$  with  $x_1 = y_1$  to same value,  $k(\mathbf{x}, \mathbf{y})$  high

However, assume  $\mathbf{x}_{2:d}, \mathbf{y}_{2:d}$  random and large  $d \rightarrow \frac{\mathbf{x}^\top \mathbf{y}}{d} \approx \frac{x_1 y_1}{d}$  small →  $k(\mathbf{x}, \mathbf{y})$  small

generally: low-dimensional features dominated by irrelevant features

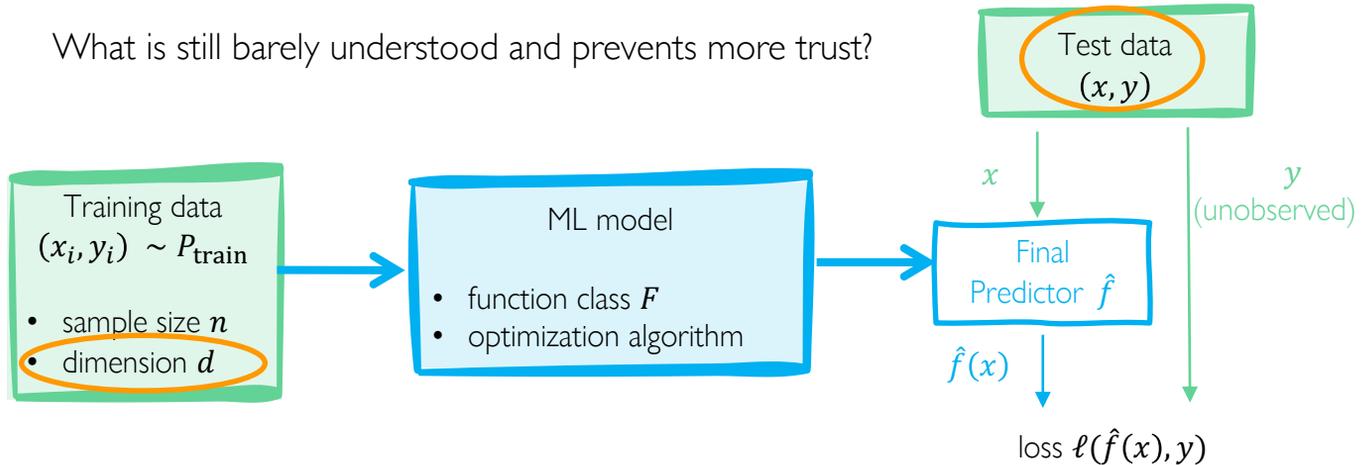
# Is there any hope for kernels in large $d$ ?

- For images: structured kernels actually work well in practice...?
  - convolutional kernels can achieve ~90% on CIFAR10<sup>1</sup>
  - can we analyze the asymptotic limits of such kernels?
- For functions depending on few variables:
  - nonlinear feature selection before kernel
  - or A(utomatic) R(elevance) D(etermination) kernels<sup>2</sup>

<sup>1</sup>e.g. [Shankar et al. '20], <sup>2</sup>e.g. [MacKay '94]

# Supervised learning for modern ML

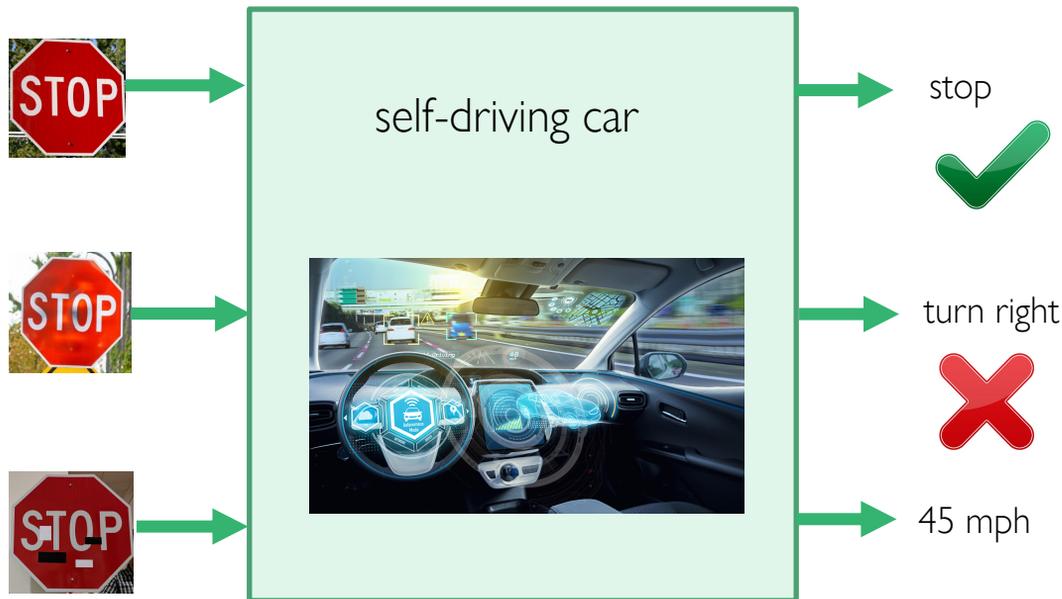
What is still barely understood and prevents more trust?



Two important settings for which understanding is missing

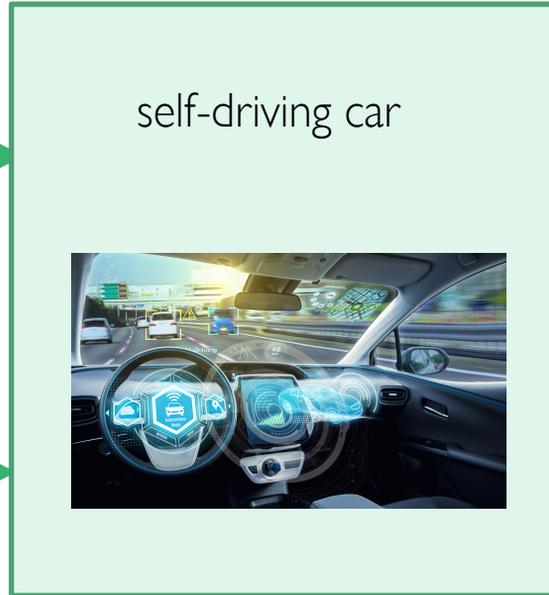
- 1 High-dimensionality: how small average loss can be when  $d$  large
- 2 Reliability: how model acts when test data  $(x, y) \neq P_{\text{train}}$

# What is this traffic sign?



Fails for small perturbations that don't change the class

# Which object is this?



truck

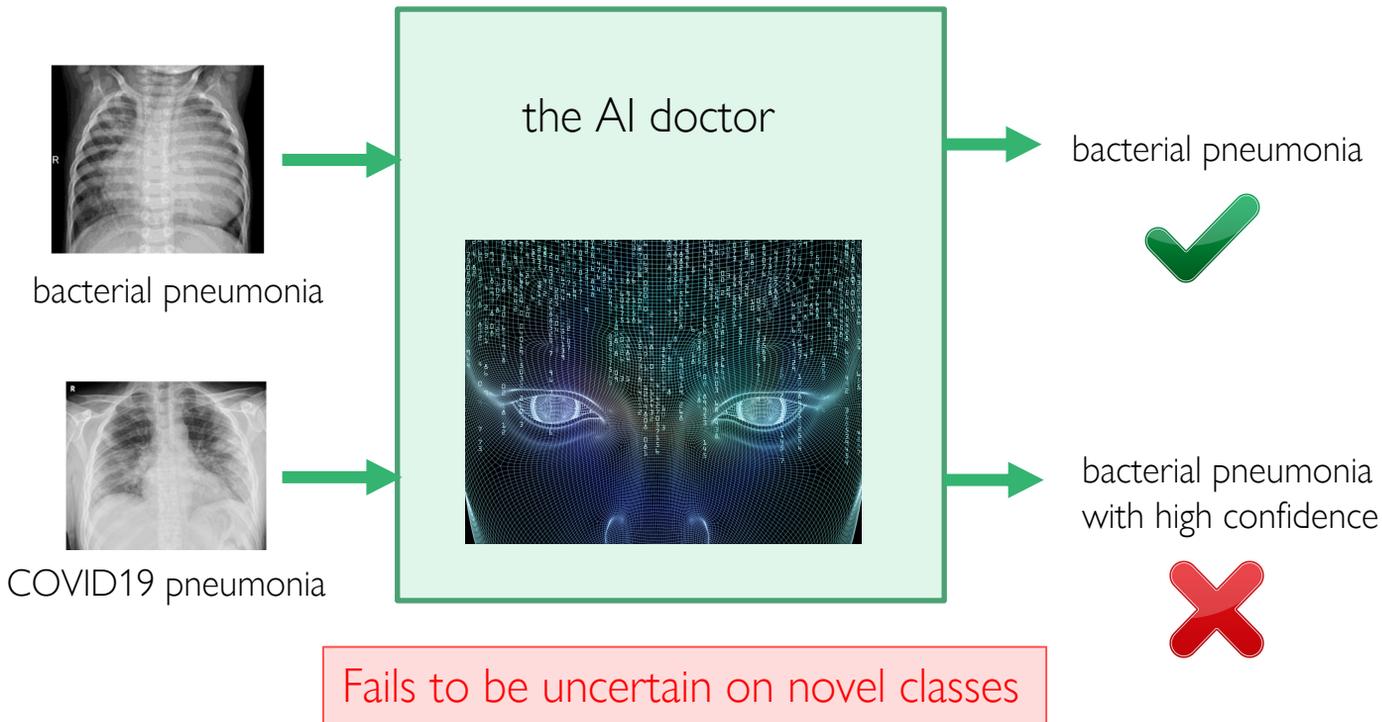


sky



Fails on unseen scenarios

# Which disease does this patient have?



# Reliability when $(x, y) \not\sim P_{train}$

Adversarially transformed  
 $T(x)$  with  $x \sim P_{train,x}$



Same  $P_{y|x}$  and classes  
 $x \sim P_{test,x} \neq P_{train,x}$



Same  $P_{y|x}$  but new classes  
 $y \notin \text{supp}(P_{train,y})$



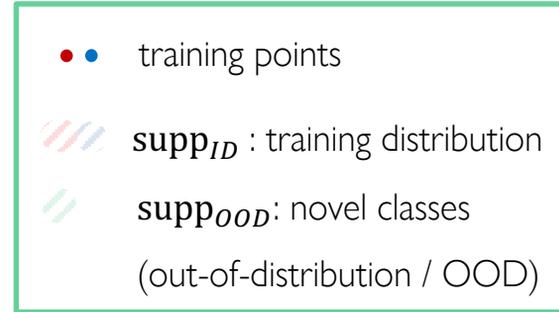
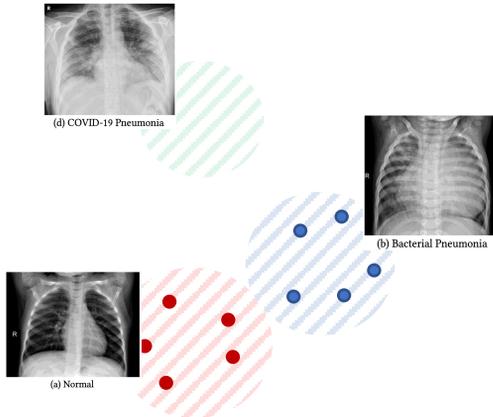
Goals:

- robust generalization: have low expected error on  $x$
- calibration: predictive uncertainty should be accurate
- detect novel class samples

# Plan for Part ②: novel class detection

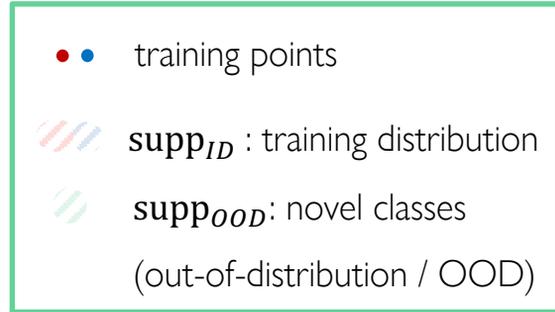
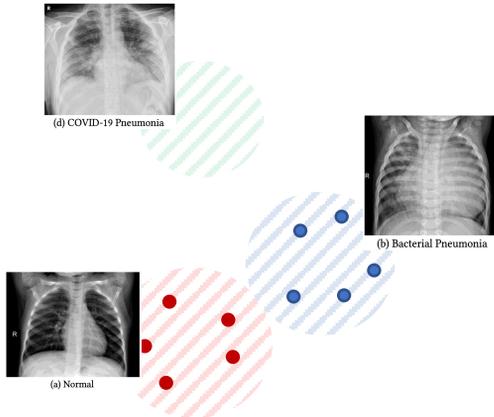
- Easy vs. hard novel class (OOD) detection
- OOD detection using ensembles
- Unknown OOD setting
- Idea: regularized diversity with unlabeled data

# Novel classes in the test set



- Given labeled training points from  $\text{supp}_{ID}$ , test point  $x \in \text{supp}_{test} = \text{supp}_{ID} \cup \text{supp}_{OOD}$
- Goal: Flag if  $x \in \text{supp}_{OOD}$ , predict if  $x \in \text{supp}_{ID}$   
(also known as anomaly detection, open set recognition, one-class classification)

# Two types of test statistics



Can view it as classifying between OOD and ID without OOD labels

- view as density estimation problem  $\rightarrow$  flag if probability of  $x$  too low
- by-product of predictive uncertainty problem  $\rightarrow$  flag if uncertainty too high

works better  
with neural networks

# Easy vs. hard novel image classes

*supp<sub>ID</sub>*: CIFAR10 classes 1-5



easy OOD: different dataset

- *supp<sub>OOD</sub>*: SVHN classes 6-10
- ID far from OOD data



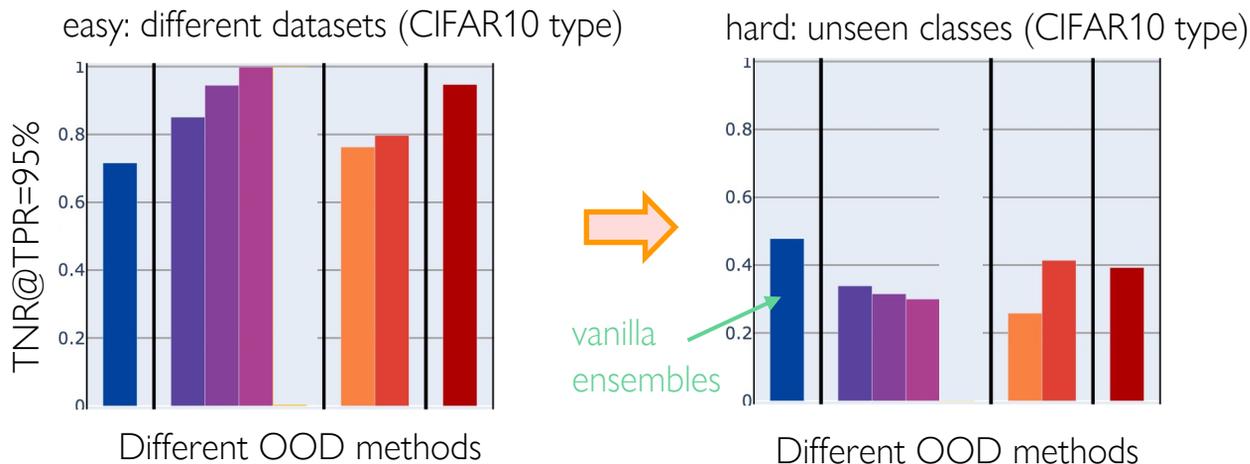
hard OOD: unseen classes

- *supp<sub>OOD</sub>* : CIFAR10 classes 6-10
- ID close to OOD data



# SOTA methods fail for novel classes

- True positive rate (TPR): percentage of truly novel classes marked as OOD
- True negative rate (TNR): percentage of seen classes marked as ID



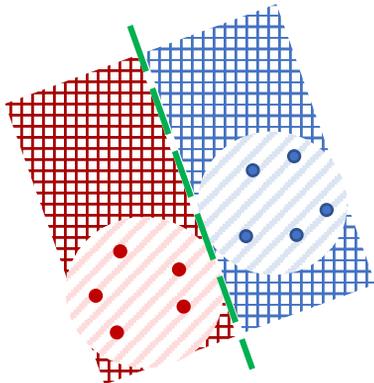
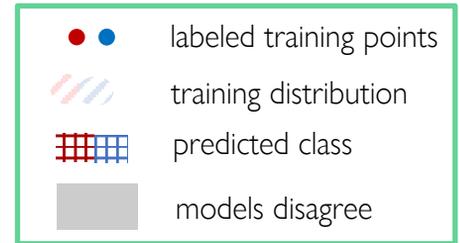
large drop for unseen (novel) class from the same dataset!

# Novel class detection using ensembles

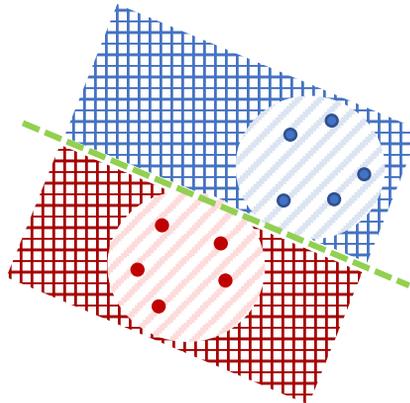
For a binary classification problem consider two models that

- have good validation accuracy on old classes
- are different outside of training distribution

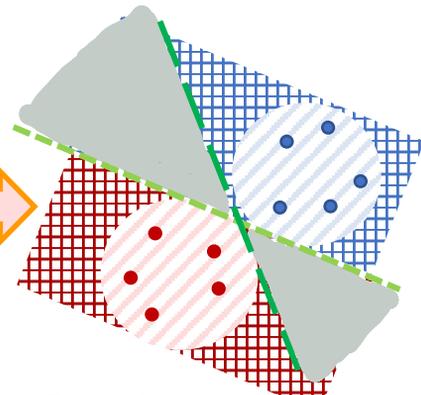
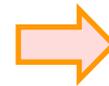
→ can mark samples with disagreement as novel



Hypoth. model 1 predictions



Hypoth. model 2 predictions

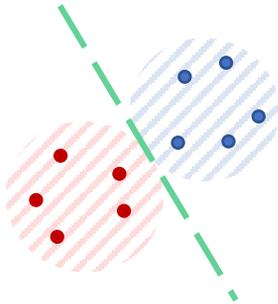


Ensemble predictions

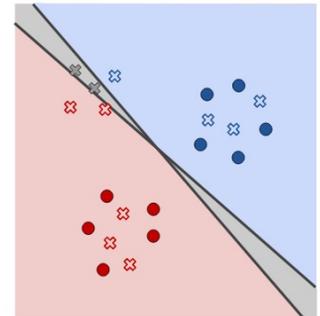
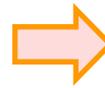
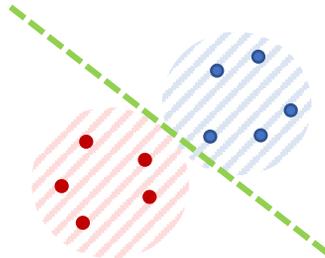
# Vanilla ensembles are not diverse enough

SOTA neural network ensembles tend to agree where they can

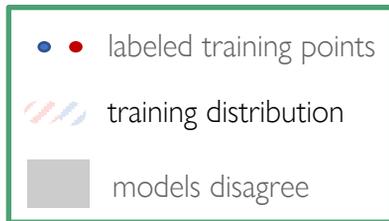
Model 1: random init 1



Model 2: random init 2



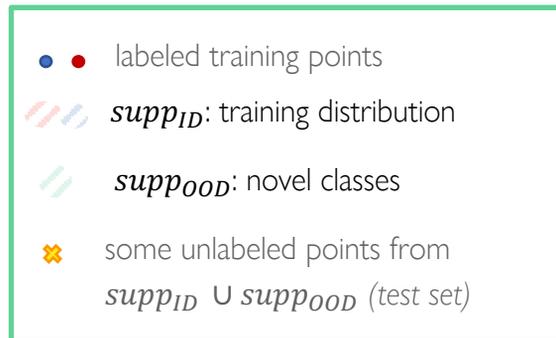
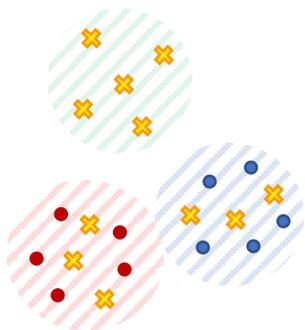
vanilla ensemble



Hard OOD requires disagreement on larger region  
→ more diverse ensembles!

# Our setting: Unknown OOD

Often unlabeled test data can be available including OOD (PU-learning)



In medical example

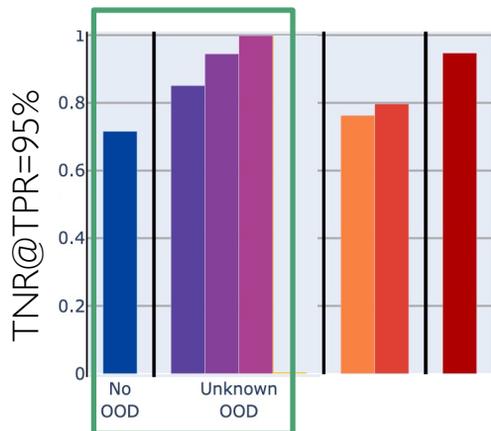
- test set: unlabeled X-rays collected during the week when new disease arrives
- even though predict using old model, valuable to detect new diseases by end of week

Is the problem now trivial?

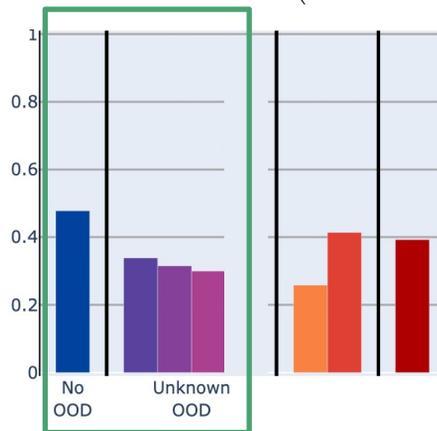
# Previous unknown OOD methods still fail

- True positive rate (TPR): percentage of truly novel classes marked as OOD
- True negative rate (TNR): percentage of seen classes marked as ID

easy: different datasets (CIFAR10 type)



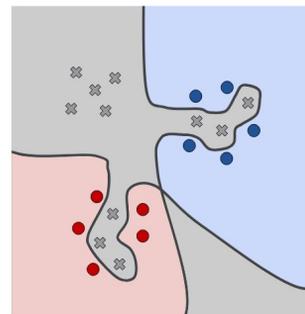
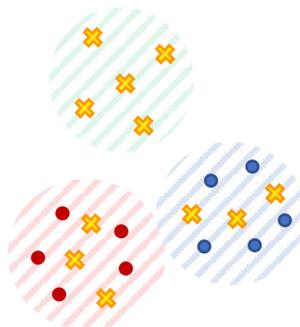
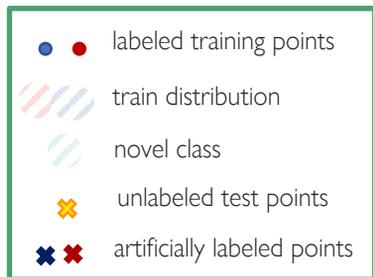
hard: unseen classes (CIFAR10 type)



- Vanilla Ensembles
- Max. Clasif. Discrep.
- Mahalanobis (unknown OOD)
- nnPU

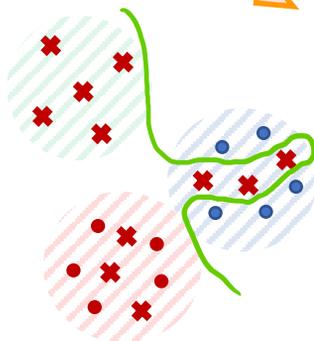
Not obvious how to leverage unknown OOD!

# Idea: regularized diversity with unlabeled data



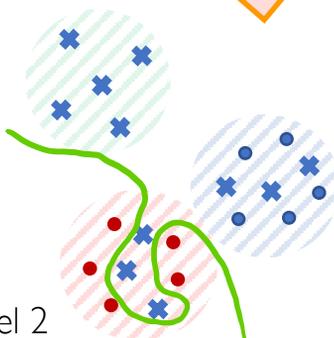
without regularization

label ✕ as ●



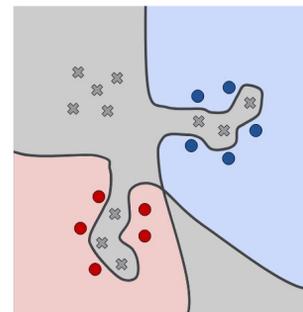
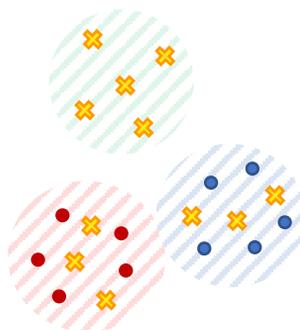
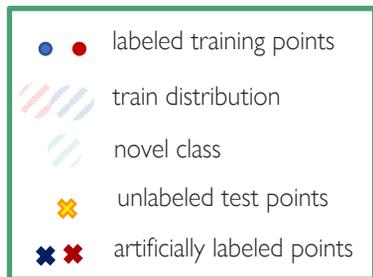
Model 1

label ✕ as ●



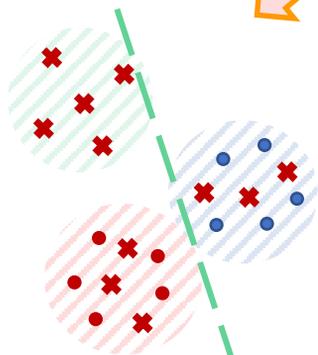
Model 2

# Idea: regularized diversity with unlabeled data



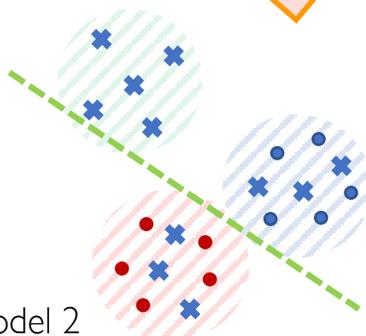
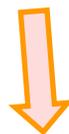
without regularization

label ✕ as ●

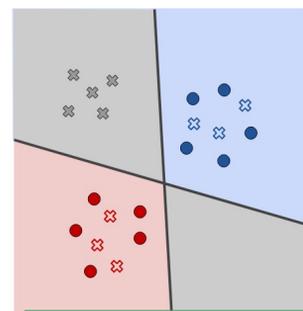


Model 1

label ✕ as ●



Model 2



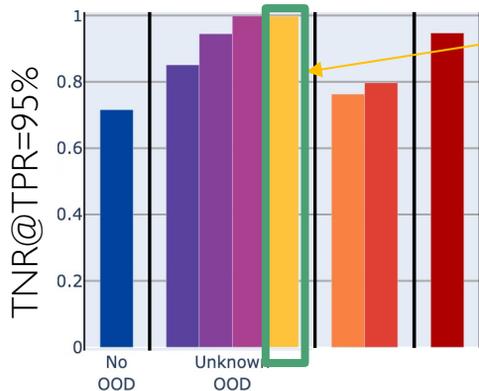
with regularization

# Learning ensembles that disagree

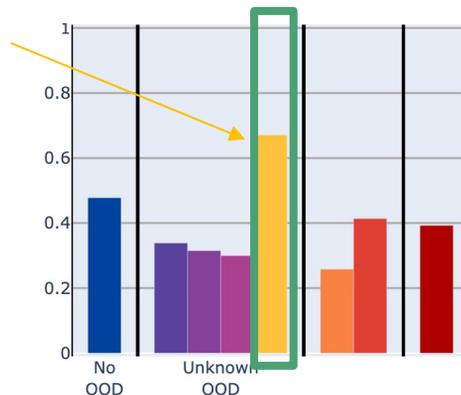


Algorithm (TSY '20) – ensembles with regularized disagreement

1. Learn model that fit different artificial labels  $c = 1$  to  $K$  on unlabeled set but are regularized to have high validation accuracy (e.g. fine-tuning + early stopping)
2. For new point  $x$ , average pairwise disagreement between classifiers  $c = 1$  to  $K$
3. flag as OOD if disagreement  $>$  some threshold



Ours



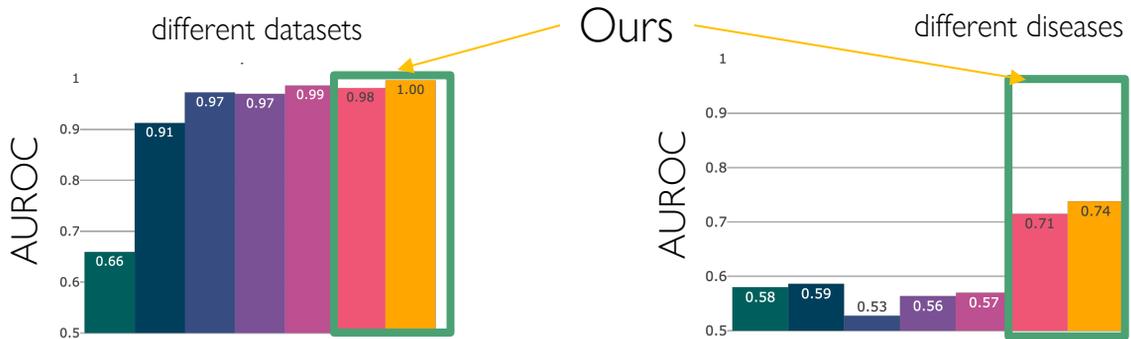
# Learning ensembles that disagree



Algorithm (TSY '20) – ensembles with regularized disagreement

1. Learn model that fit different artificial labels  $c = 1$  to  $K$  on unlabeled set but are regularized to have high validation accuracy (e.g. fine-tuning + early stopping)
2. For new point  $x$ , average pairwise disagreement between classifiers  $c = 1$  to  $K$
3. flag as OOD if disagreement  $>$  some threshold

Novel classes in chest X-ray + retinal datasets



Thanks!

