

# MATHEMATICAL TOPICS IN MACHINE LEARNING

## (LECTURE 4 – THE BIAS/VARIANCE DECOMPOSITION)

Professor Gavin Brown

## THE NEXT FEW WEEKS

**Week 4.** Bias-Variance decompositions

**Week 5.** Ensemble learning methods

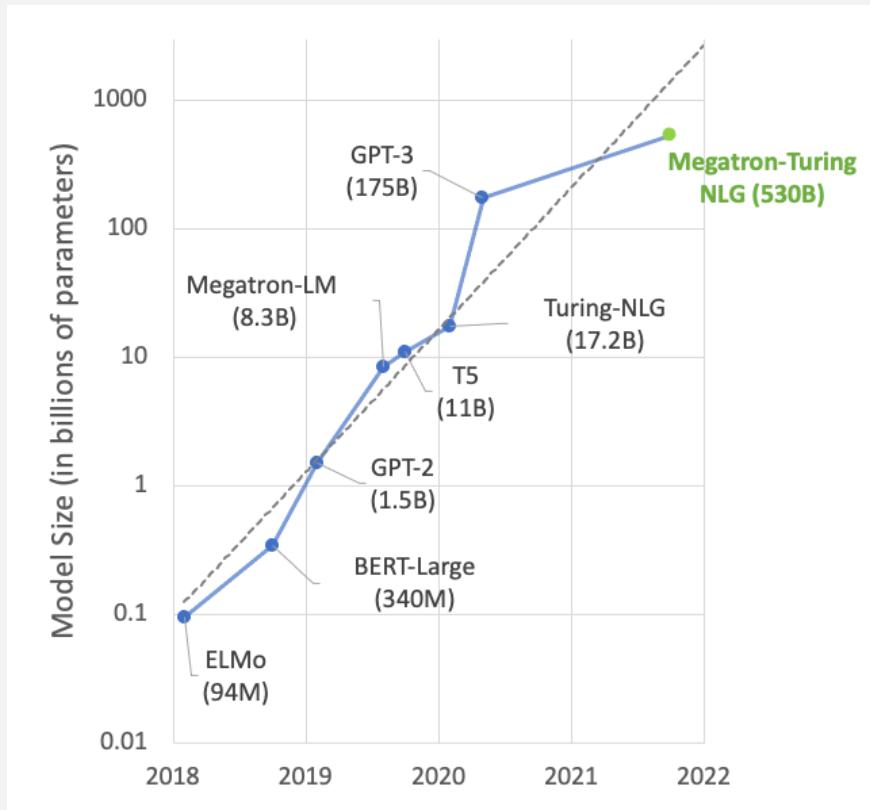
**Week 6.** Ensemble learning theory

These are my research area.  
I'm currently looking for new PhD students.

<https://profgavinbrown.github.io>

# OUR QUESTION

“Are bigger models always better models?”

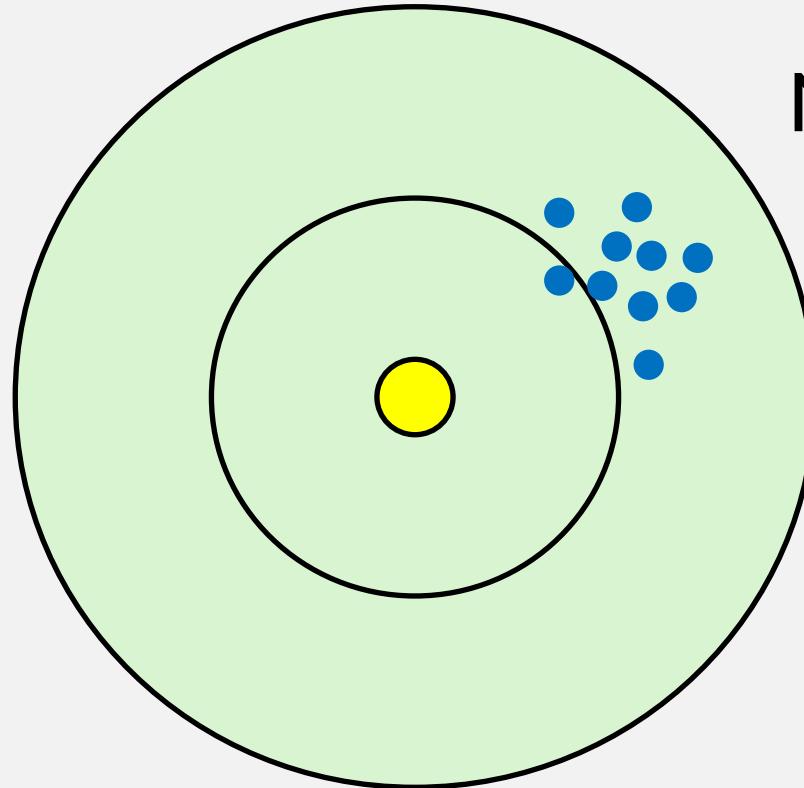


**TODAY:**

We split our expected risk in 2 parts.  
One is for the “**size**” of the model.  
One is for the “**sensitivity**” to training.

# BIAS AND VARIANCE... THE INTUITION

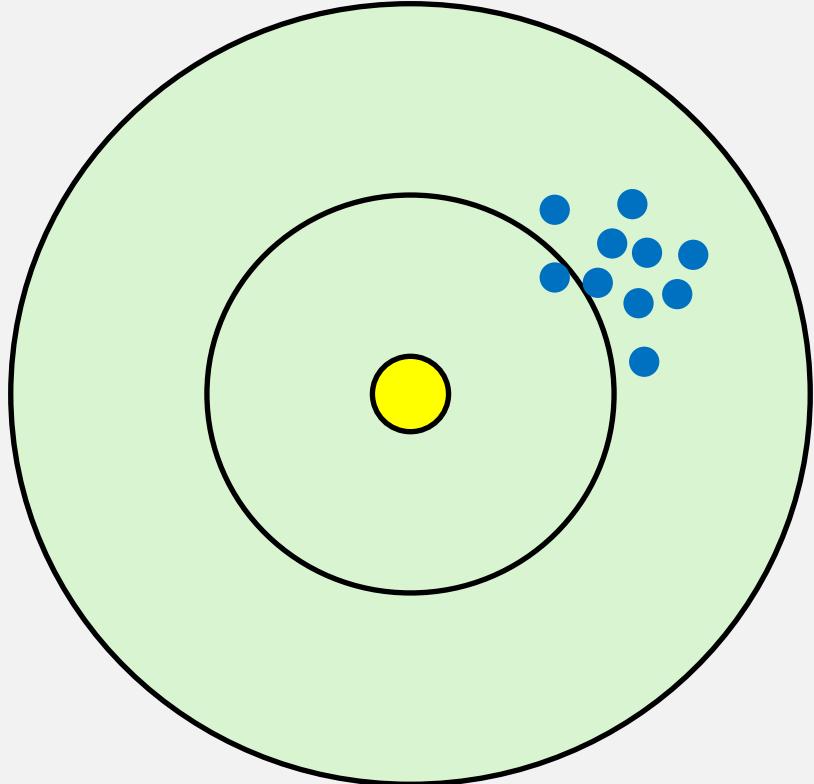
Train many copies of a model to hit a dartboard....



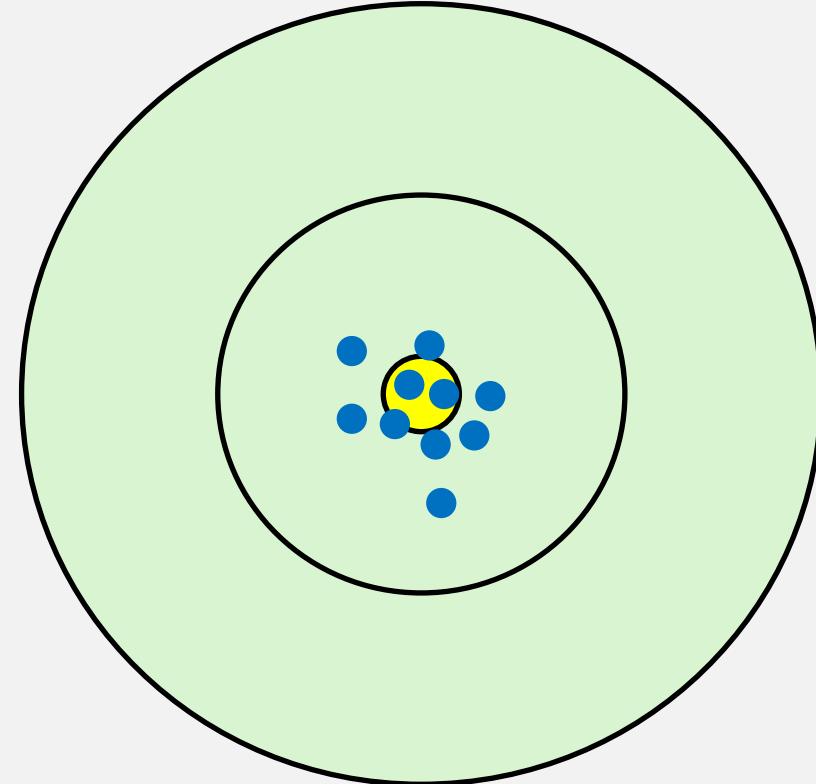
For each copy, use a  
NEW randomly sampled  
training set

This is a model with high “**bias**”, as the predictions are always a bit off target.

# BIAS AND VARIANCE... THE INTUITION

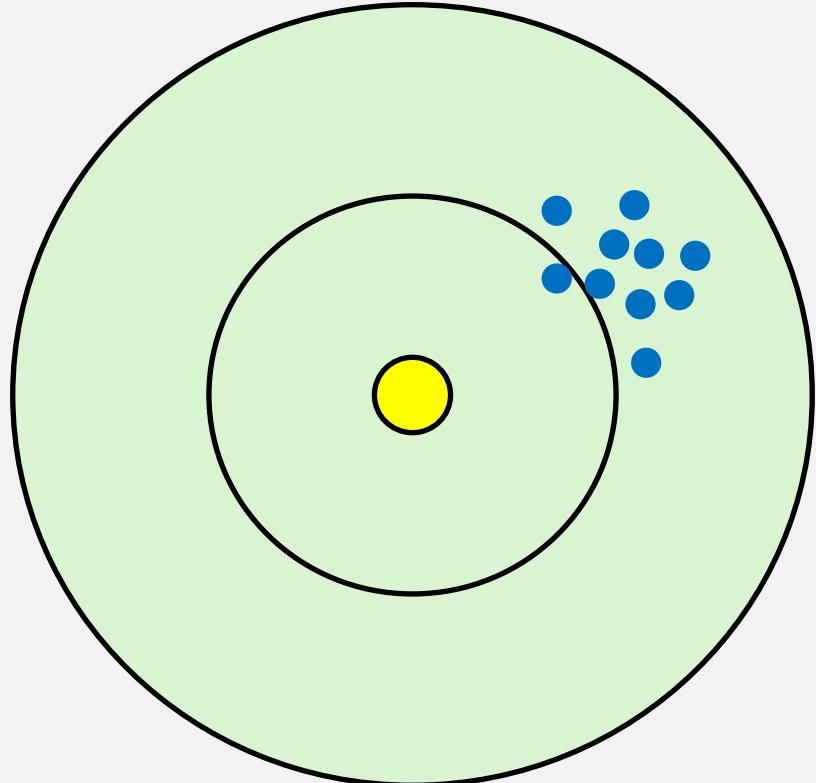


High bias

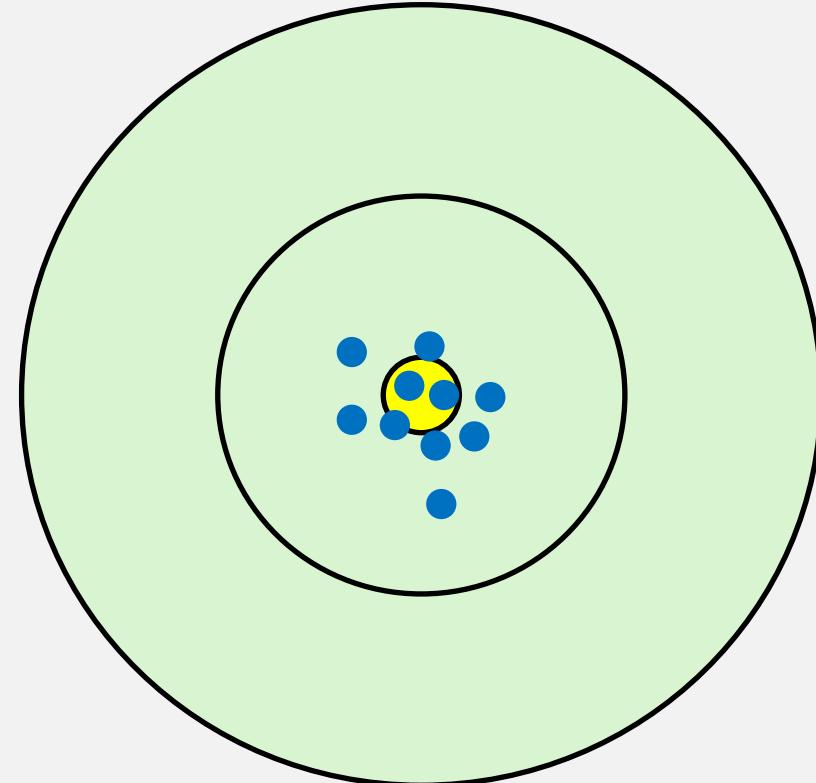


Low bias

# BIAS AND VARIANCE... THE INTUITION

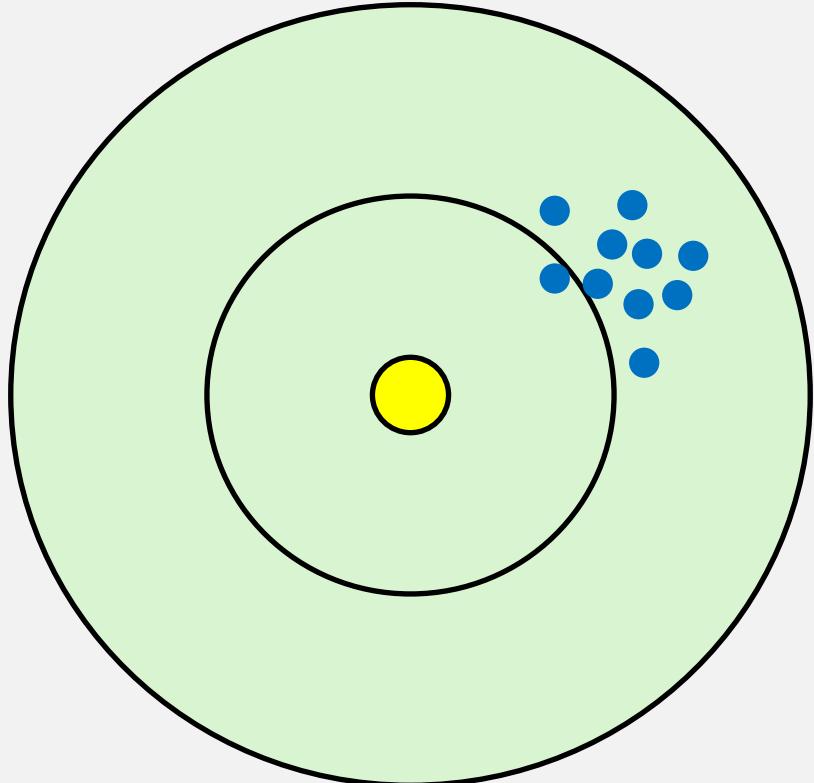


High bias  
Low variance

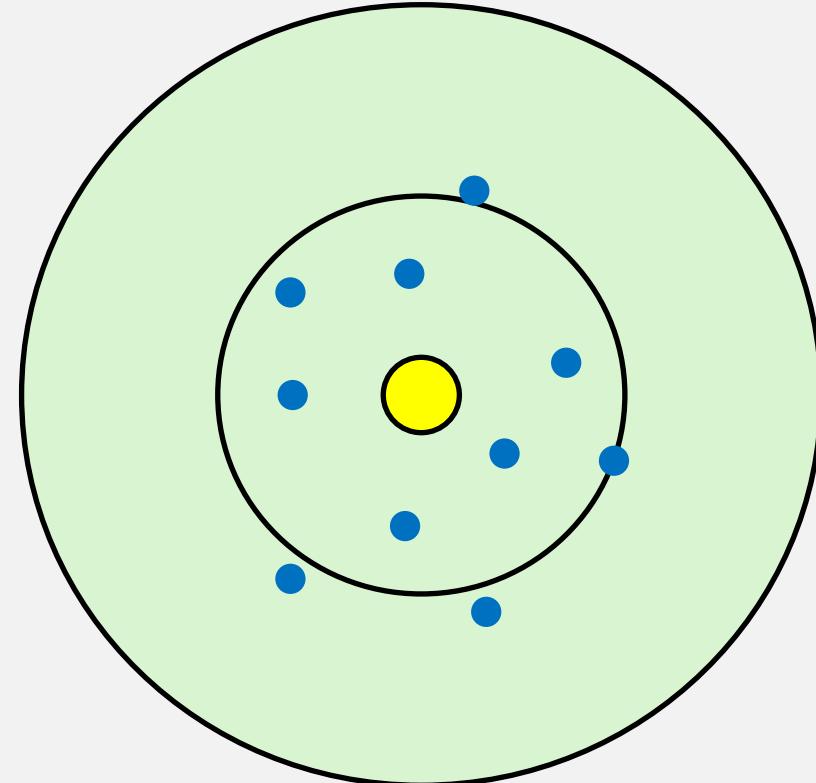


Low bias  
Low variance

# BIAS AND VARIANCE... THE INTUITION



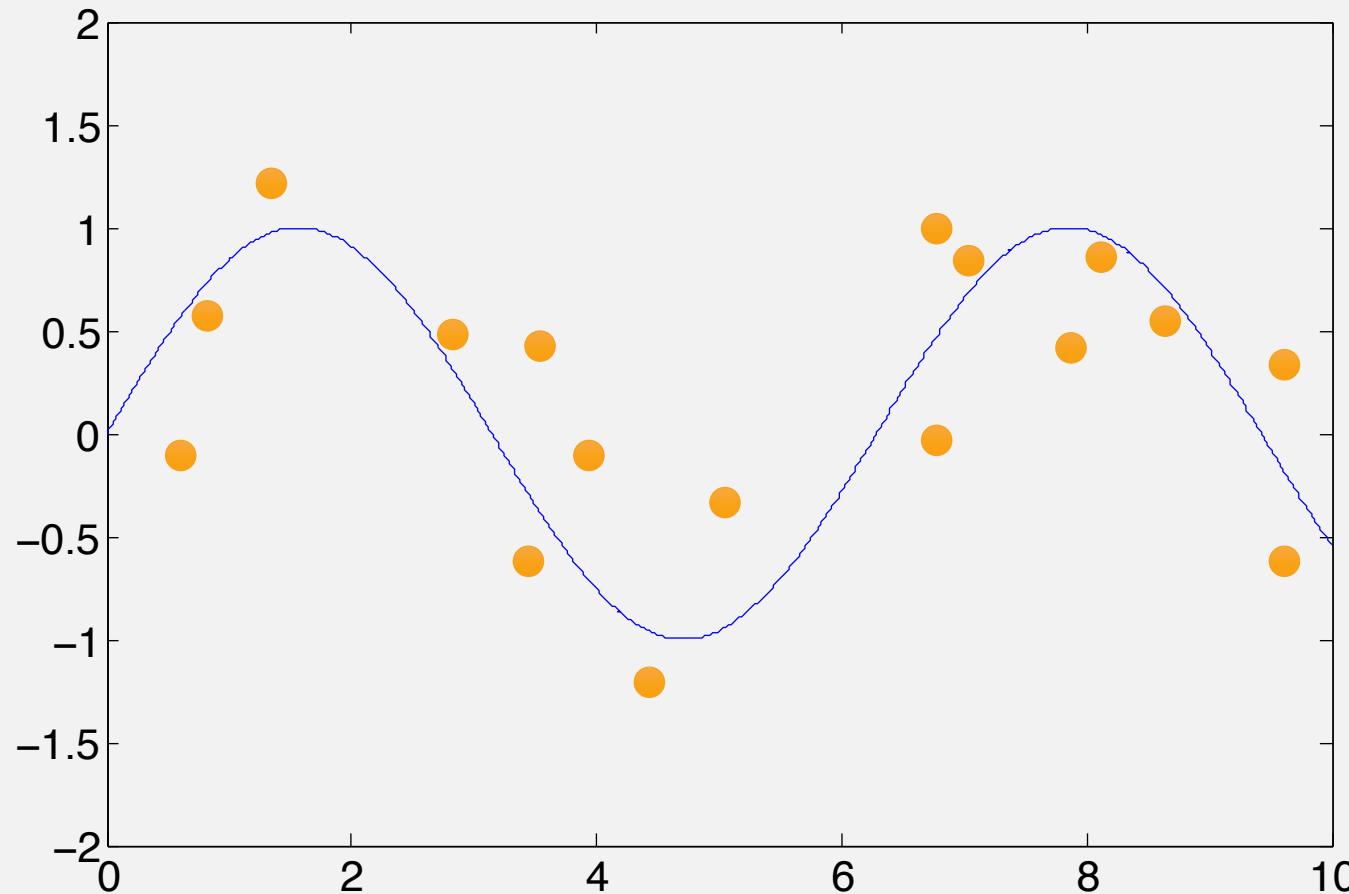
High bias  
Low variance



Low bias  
High variance

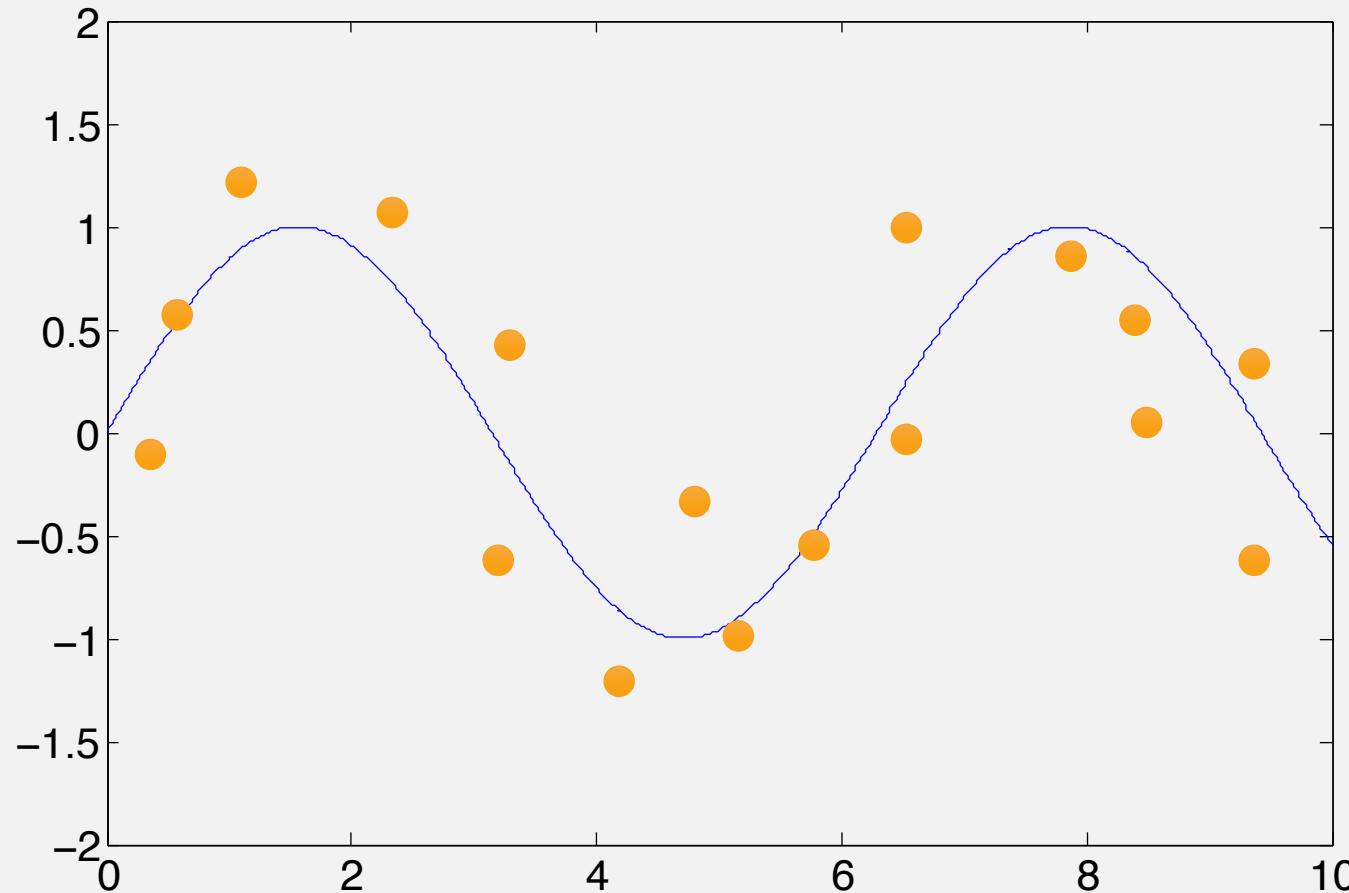
# ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE

Training set is a random draw from the true function + Gaussian noise.



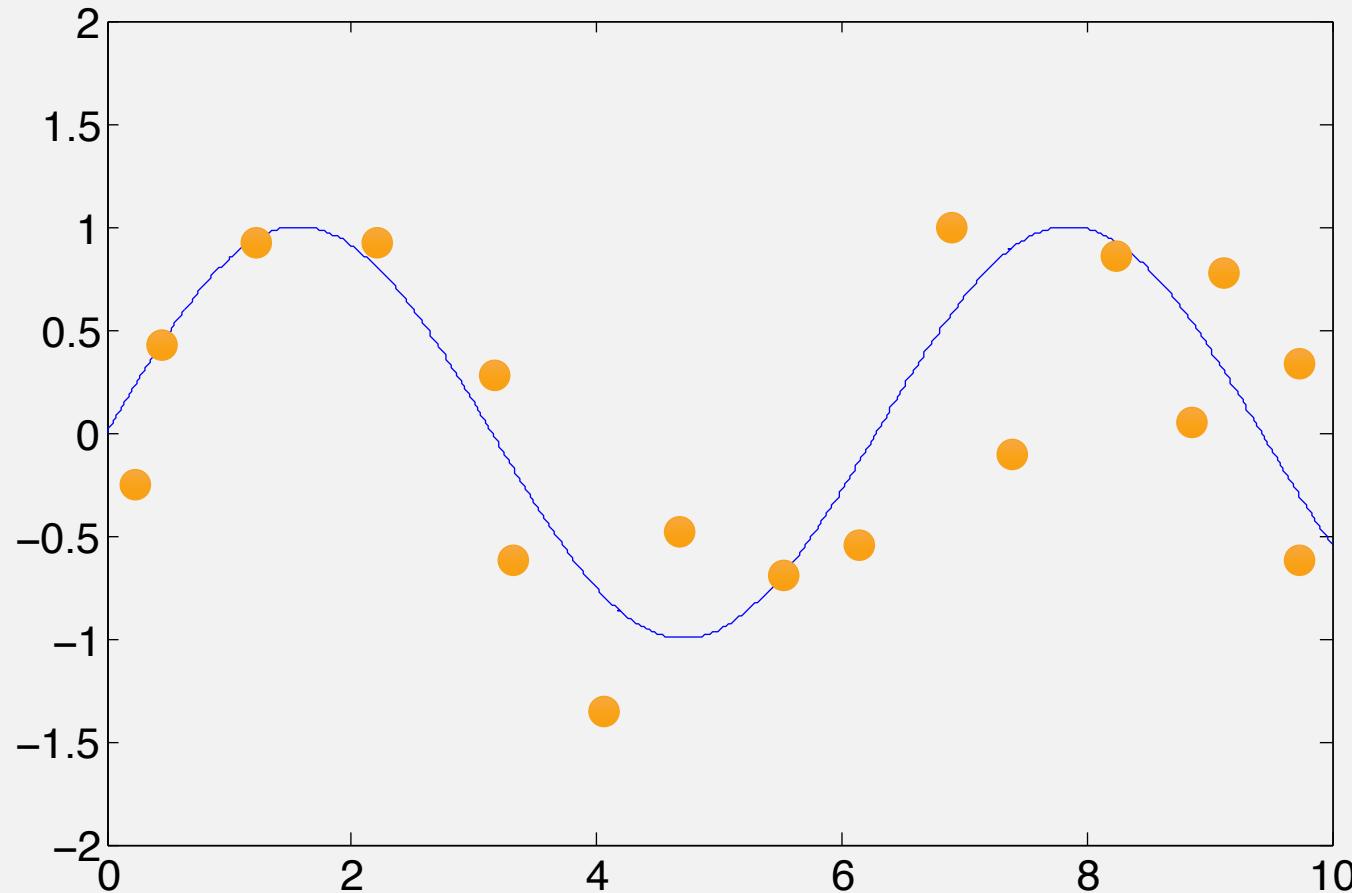
# ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE

Training set is a random draw from the true function + Gaussian noise.



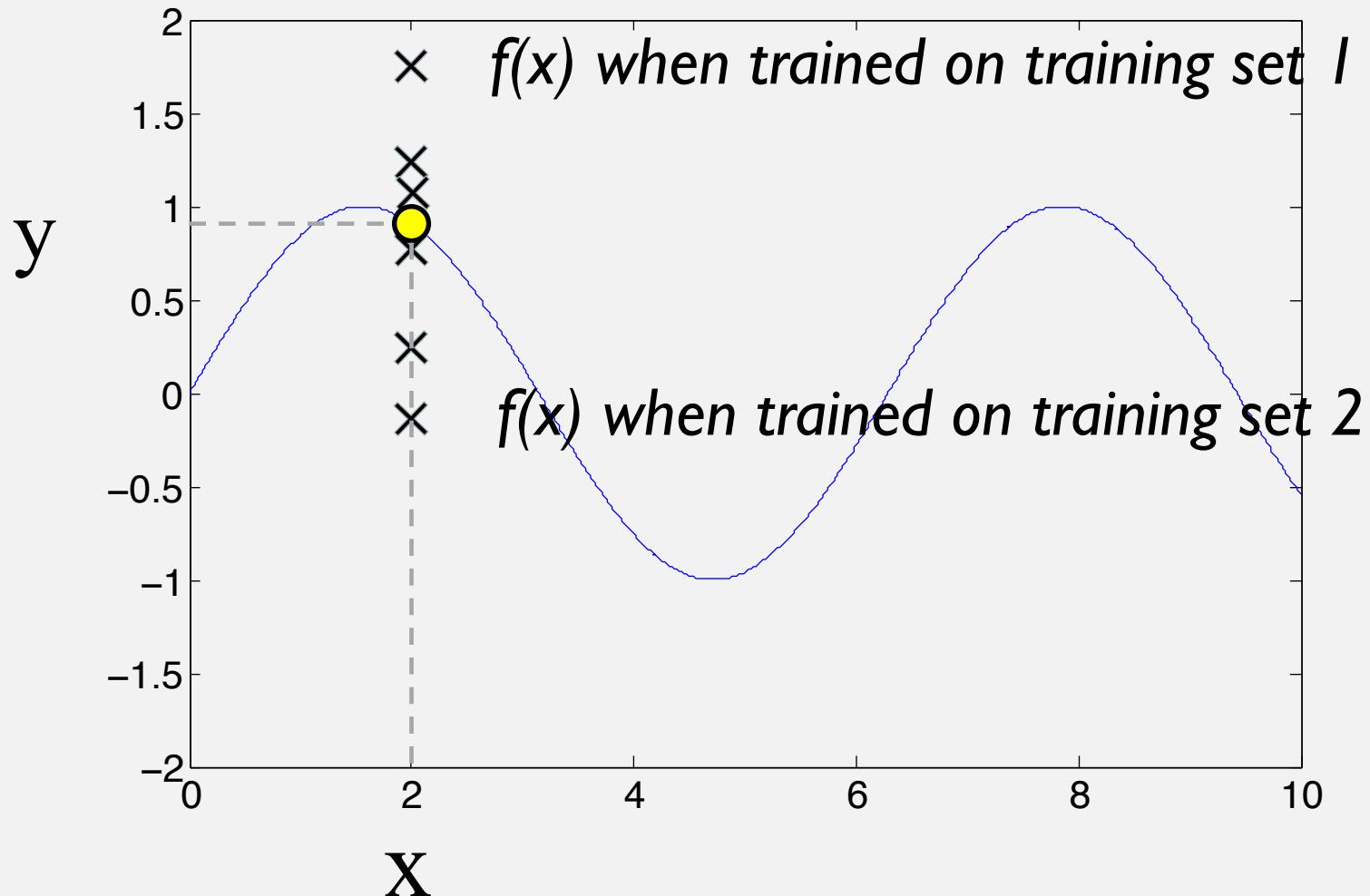
# ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE

Training set is a random draw from the true function + Gaussian noise.

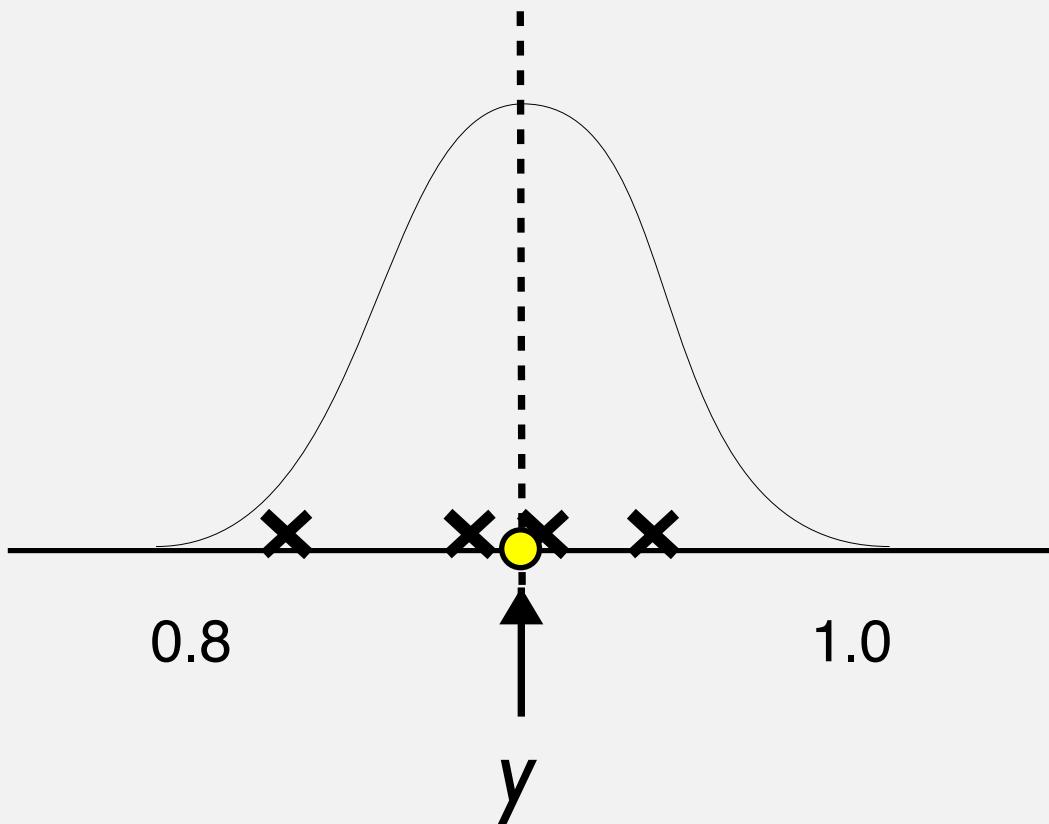


## ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE

Test point: predict  $\sin(2)$ . True value is about  $y = 0.9092$ .



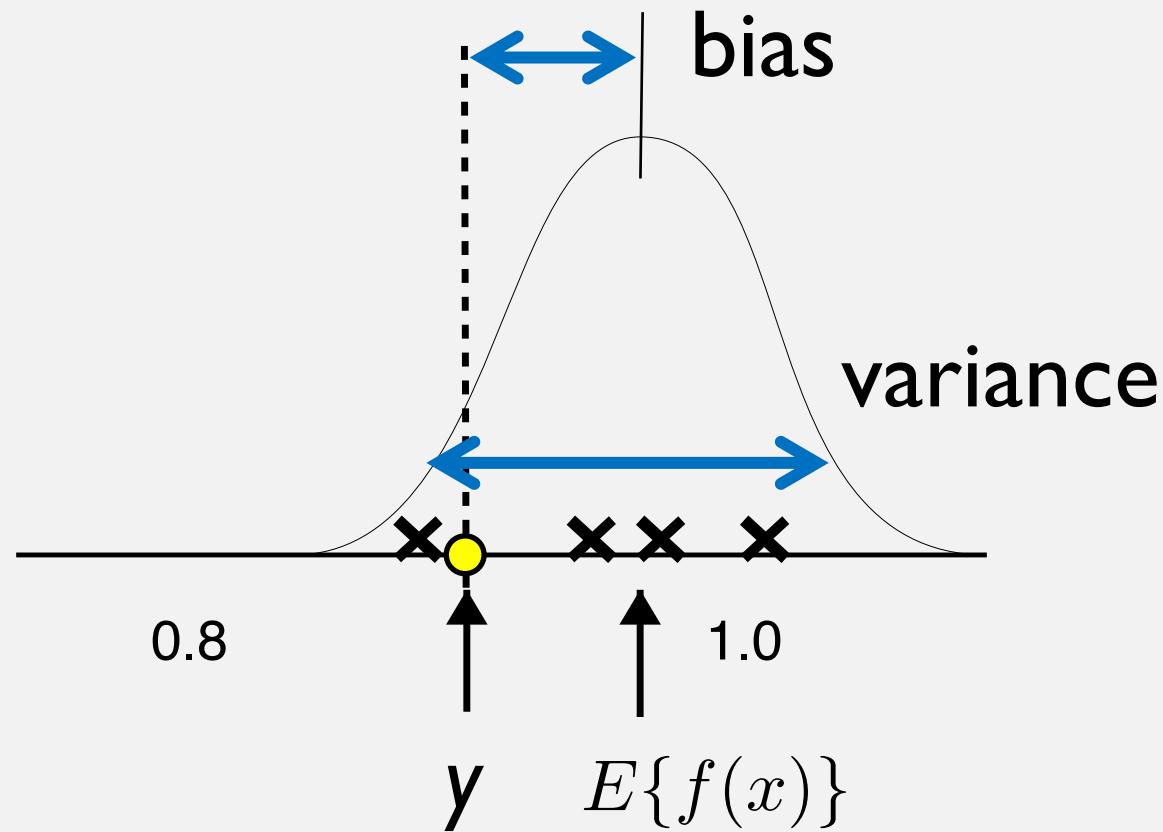
## ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE



Unbiased predictor.

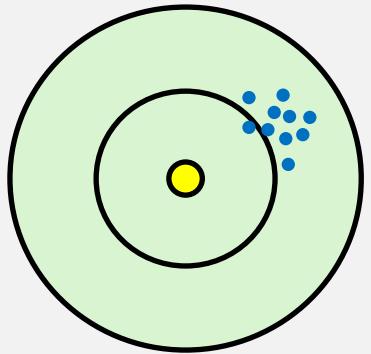
$$\sin(2) = 0.9092$$

## ANOTHER EXAMPLE: LEARNING A NOISY SINE WAVE

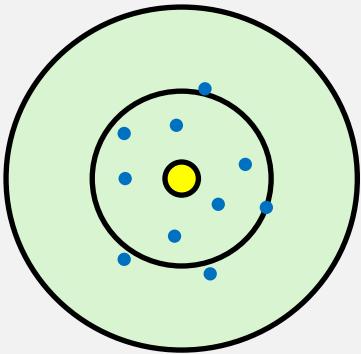


Biased predictor.

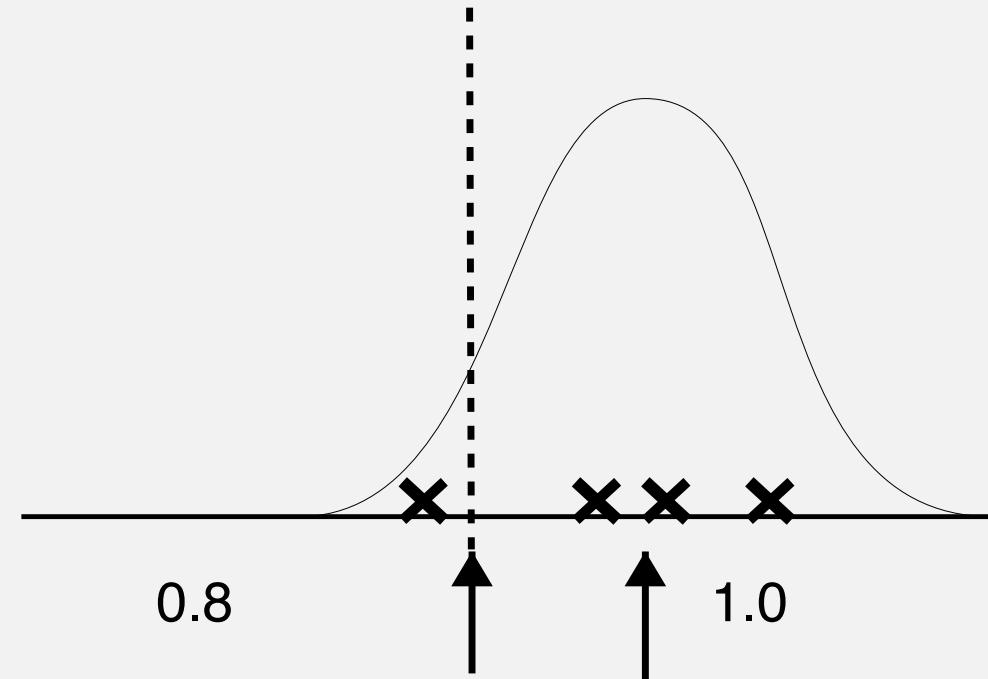
# THIS IS THE INTUITIVE IDEA.



High bias  
Low  
variance



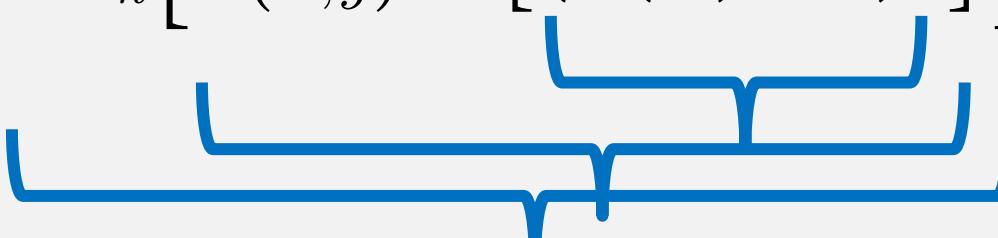
Low bias  
High  
variance



BUT, AS ALWAYS THERE IS SOME MATHS BEHIND IT.

The “expected” risk :

$$\mathbb{E}_{\mathcal{S}_n} [R(f)] = \mathbb{E}_{\mathcal{S}_n} \left[ \mathbb{E}_{(x,y) \sim \mathcal{D}} [(f(x) - y)^2] \right]$$



Squared loss of a **single** trained model,  
at a **single** test point.  
averaged over all **possible** test points.  
averaged over **all possible** training sets,  
averaged over **all possible** test points.

# DECOMPOSING THE EXPECTED RISK

$$\mathbb{E}_{\mathcal{S}_n} [R(f)] = \mathbb{E}_{\mathbf{x}} \left[ \underbrace{\mathbb{E}_{y|\mathbf{x}} [(y - \mathbb{E}_{y|\mathbf{x}}[y])^2]}_{\text{noise}} + \left( \mathbb{E}_{\mathcal{S}_n} [f(\mathbf{x})] - \mathbb{E}_{y|\mathbf{x}}[y] \right)^2 + \mathbb{E}_{\mathcal{S}_n} \left[ (f(\mathbf{x}) - \mathbb{E}_{\mathcal{S}_n} [f(\mathbf{x})])^2 \right] \right].$$

The noise is a constant that measures the squared difference between the trained model, Notice that the model is not mentioned, and the “expected” model, in response and the Bayes model getting different training sets.

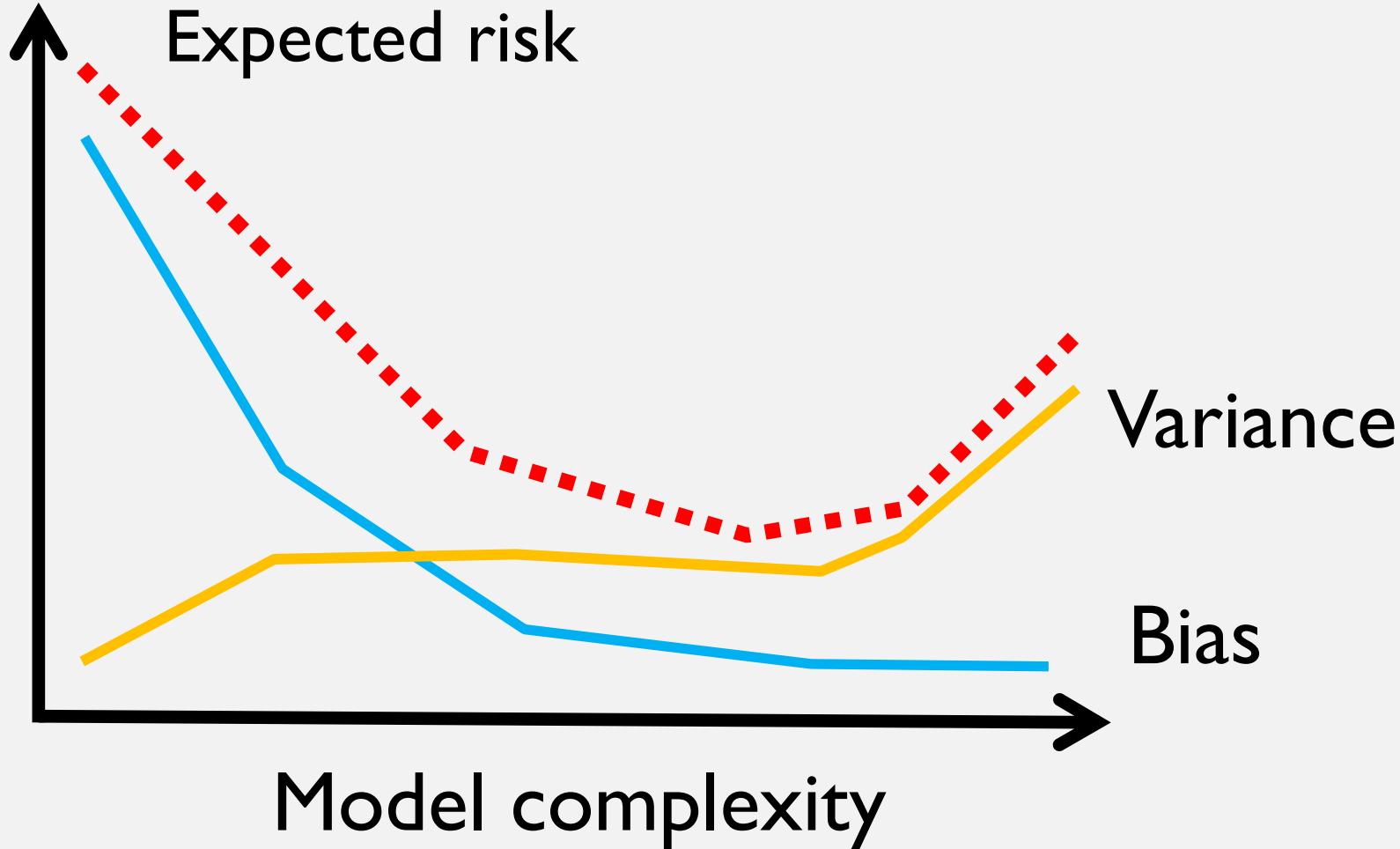
It is the population risk of the Bayes model, i.e.,  $R(y^*)$ .

The “sensitivity” to its training data.

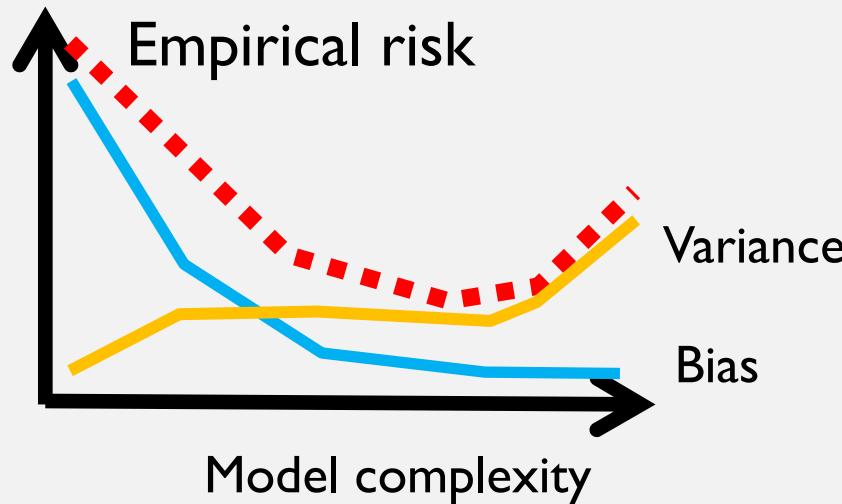
$$\mathbb{E}_{\mathcal{S}_n} [f(\mathbf{x})]$$

“Expected” model = the response that the model would give if we could average over all possible training sets.

# WHY DOES THIS MATTER? OVERFITTING, THAT'S WHY.



# WHY DOES THIS MATTER? OVERFITTING, THAT'S WHY.

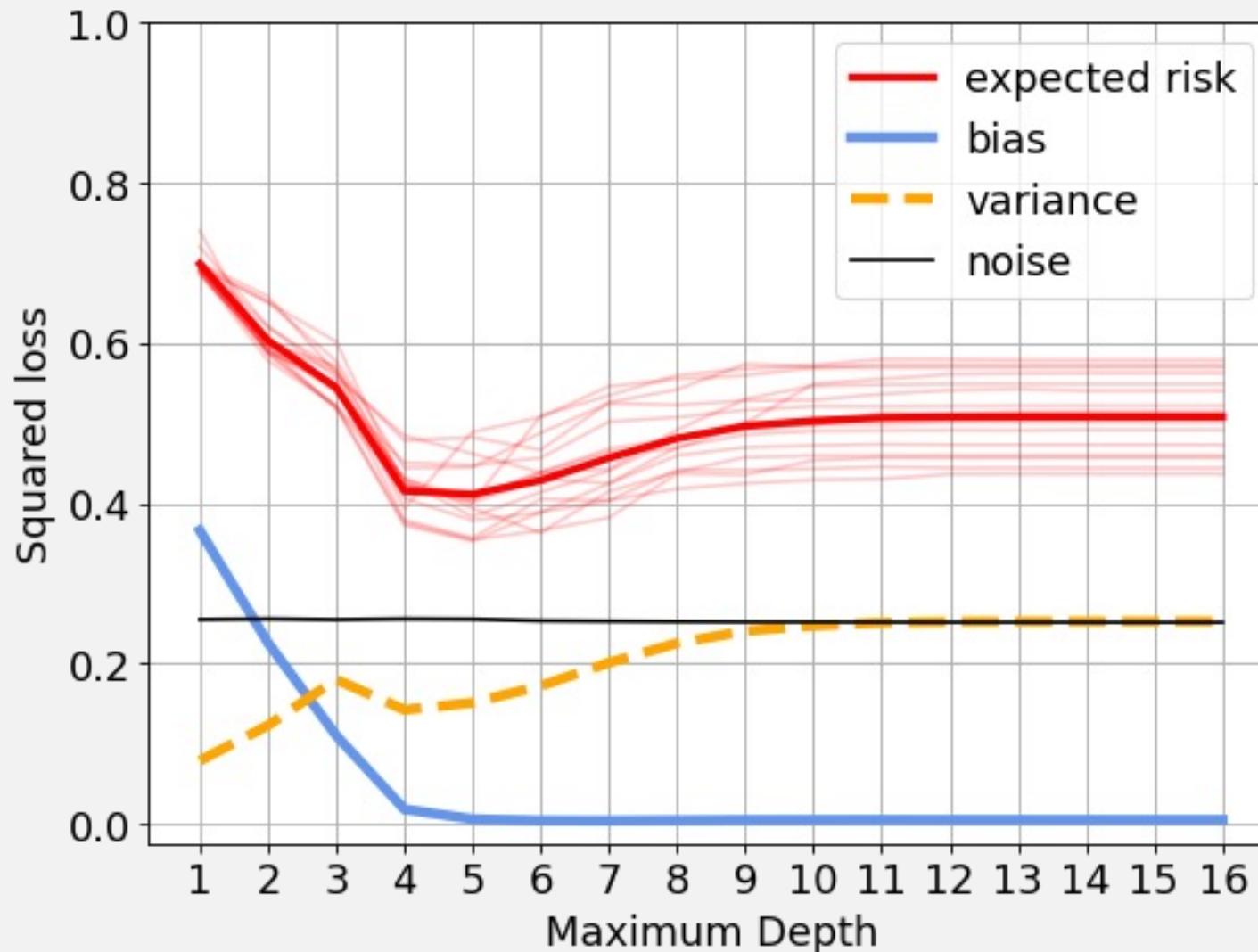


Big and complex models tend to have **LOW BIAS**, i.e. strong..

But, they also overfit to their training data.

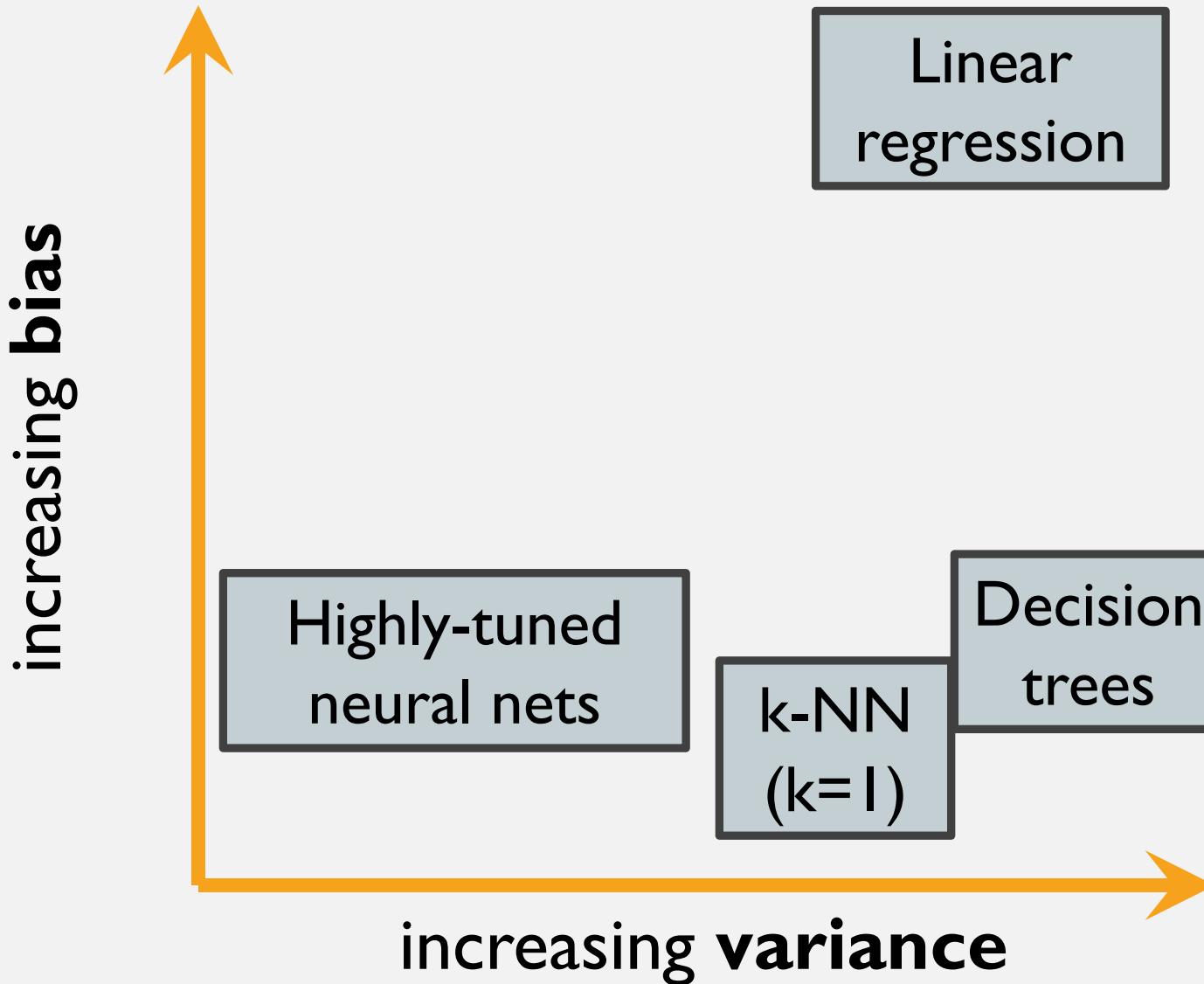
Hence, are usually **HIGH VARIANCE**.

# BUILDING DECISION TREES OF INCREASING DEPTH

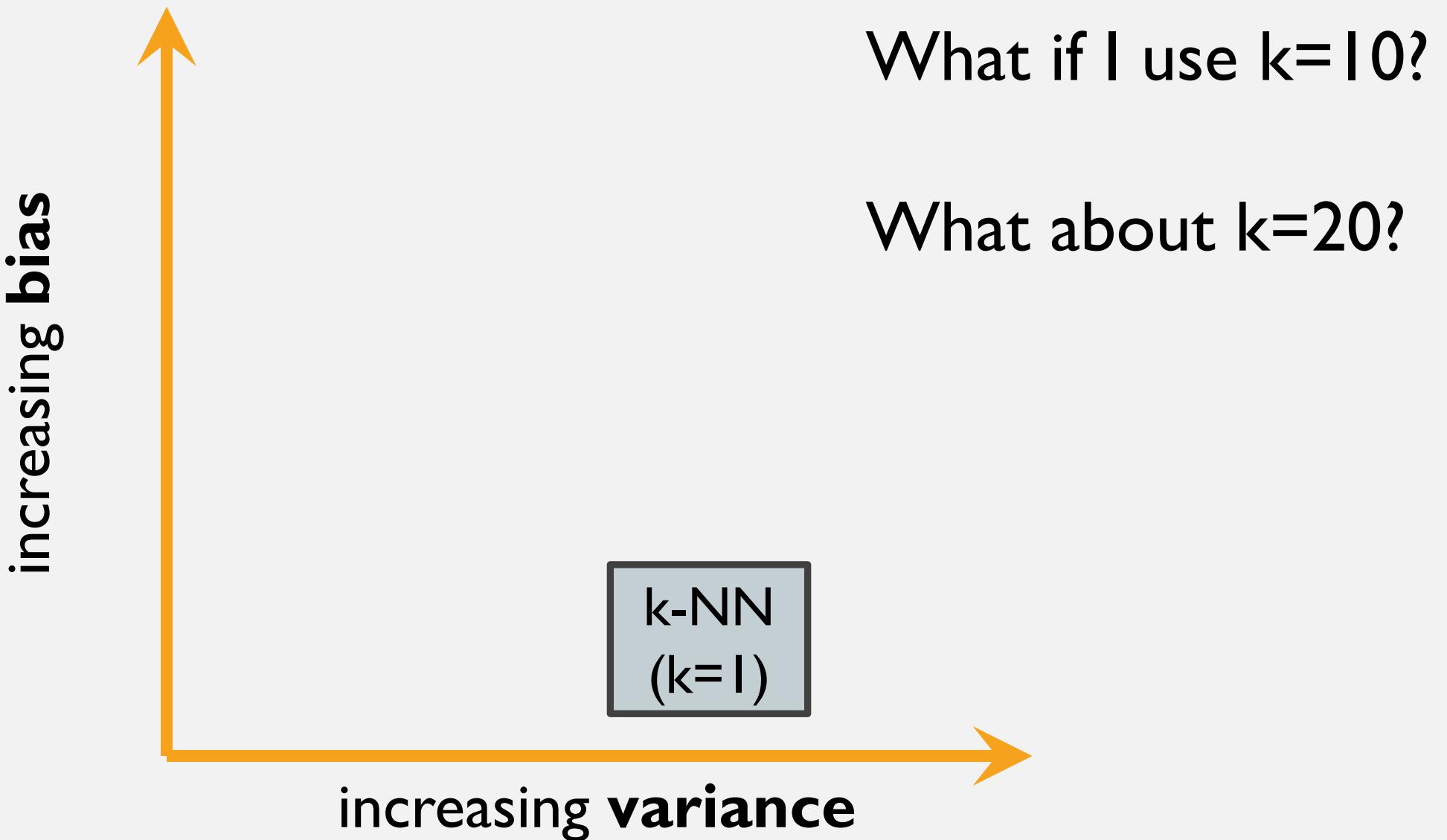


Code provided.

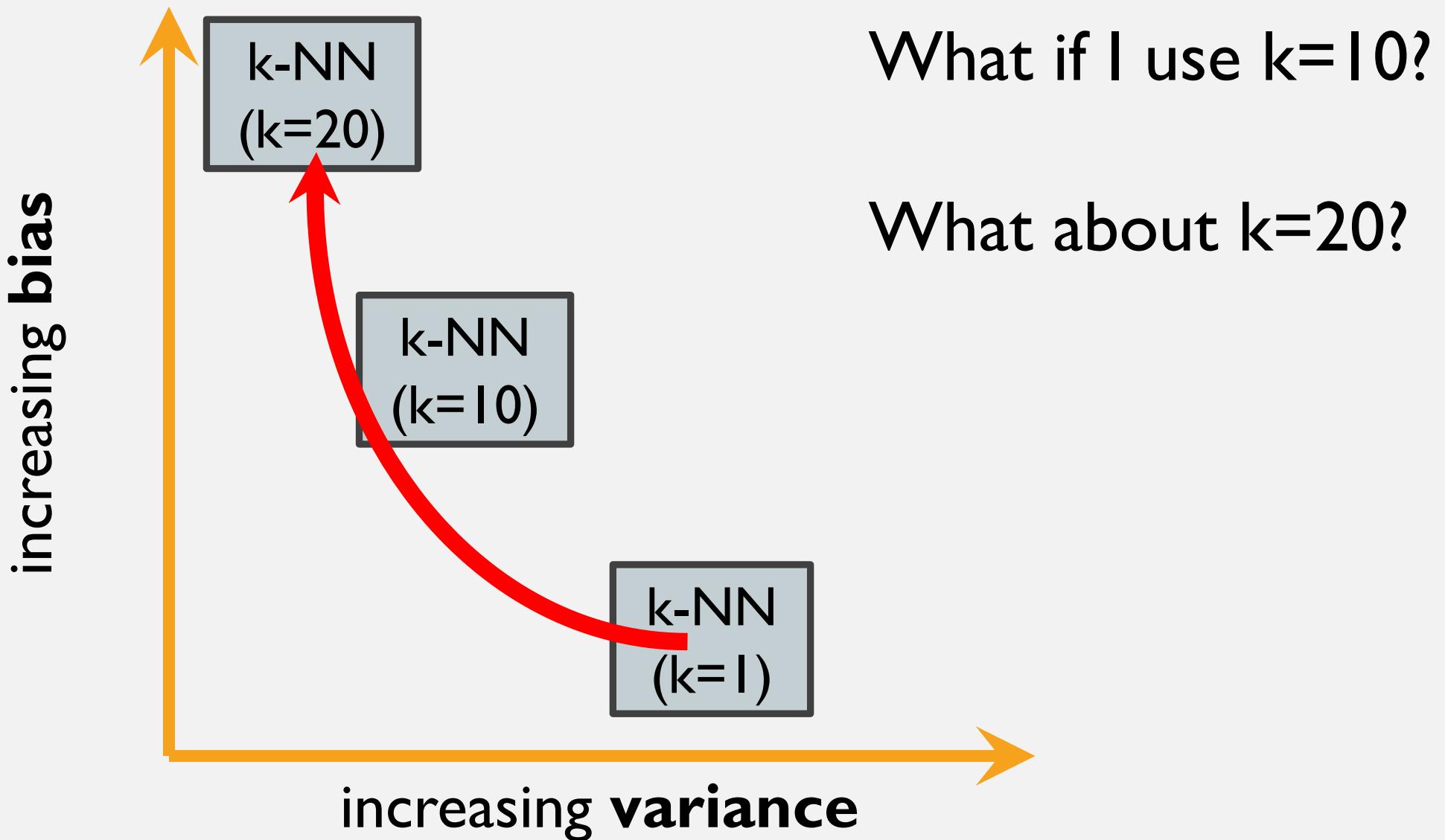
# BIAS AND VARIANCE IN REAL MODELS



# BIAS AND VARIANCE IN REAL MODELS



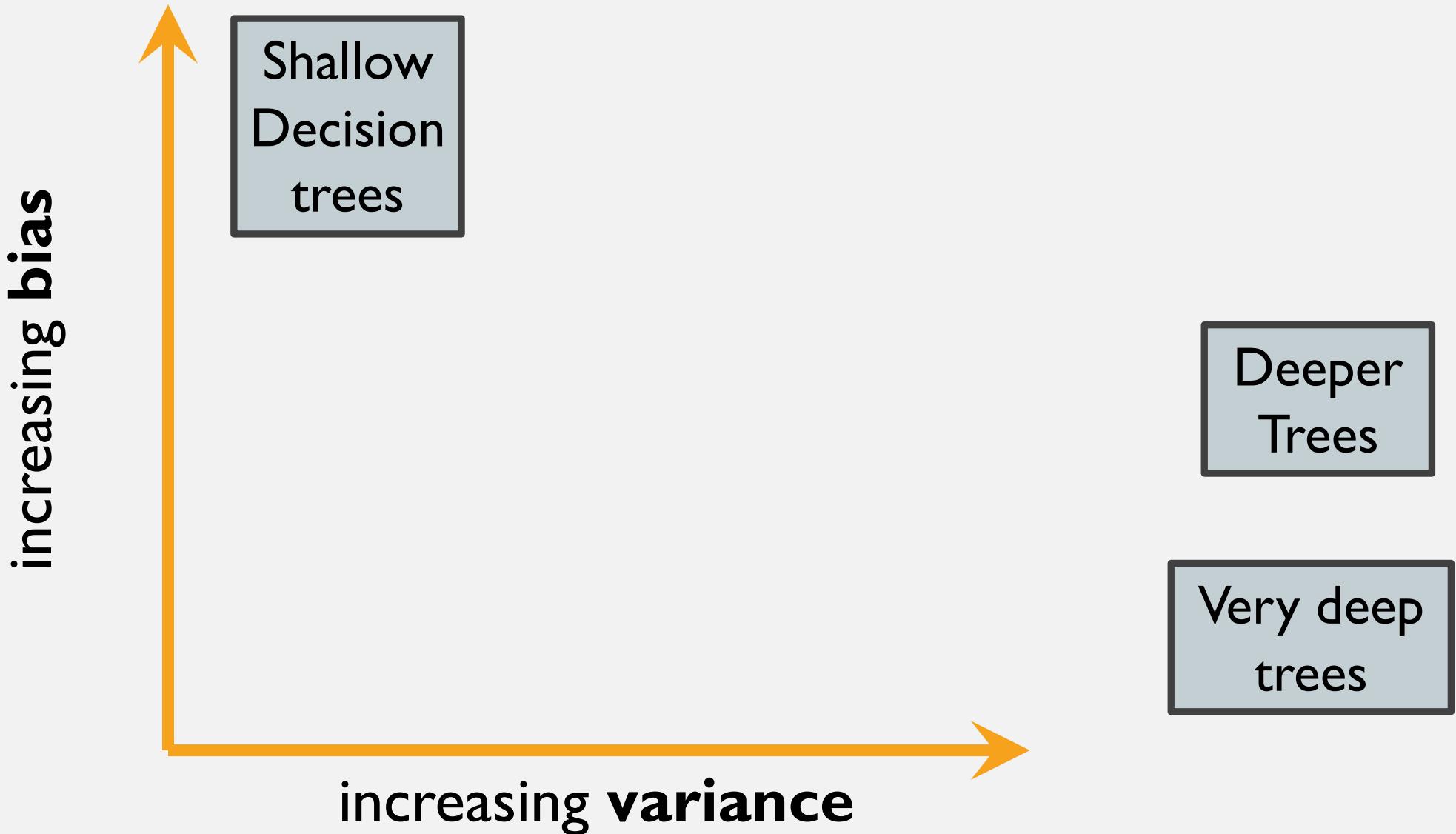
# BIAS AND VARIANCE IN REAL MODELS



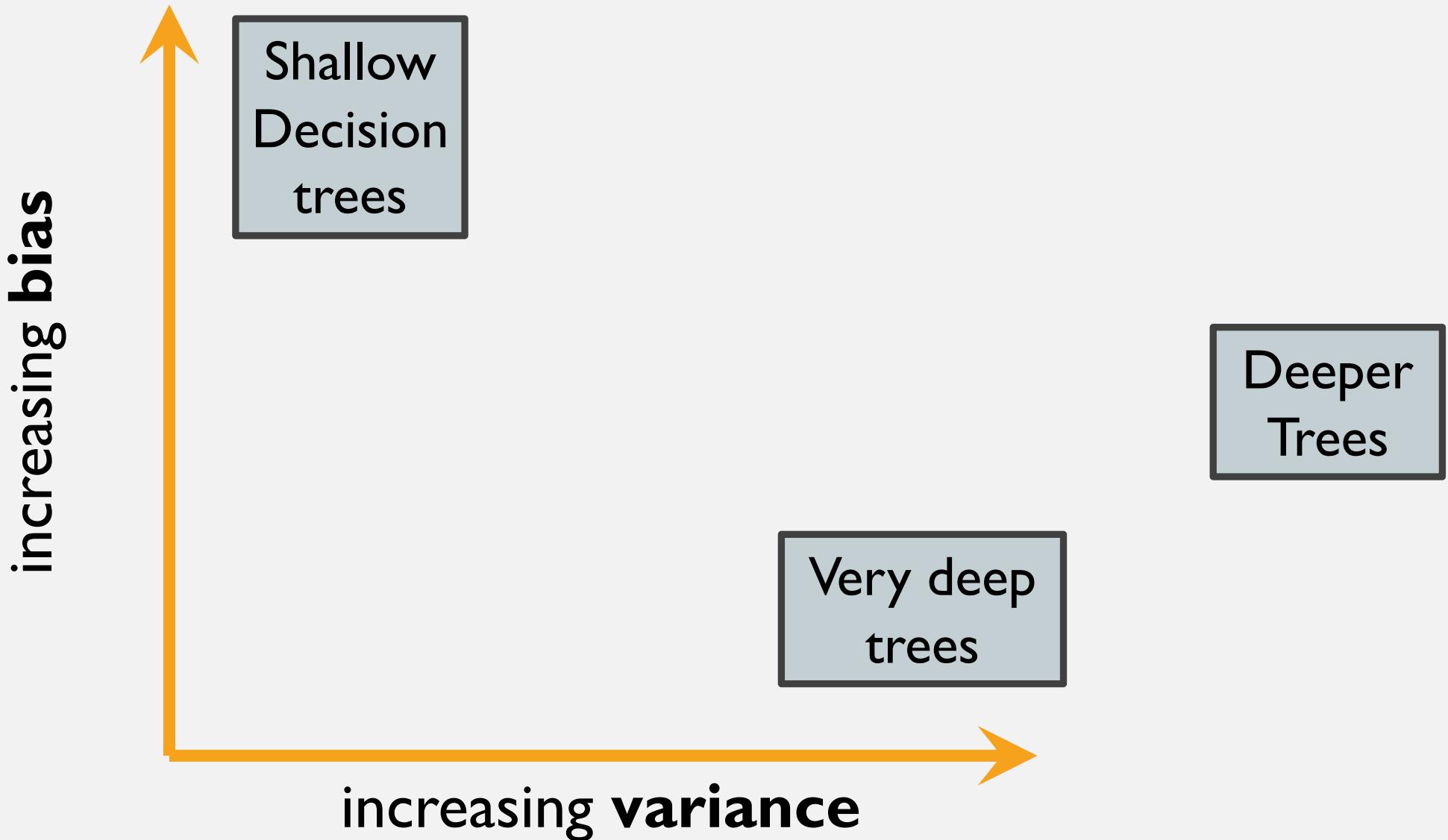
# WHERE DO THEY GO?



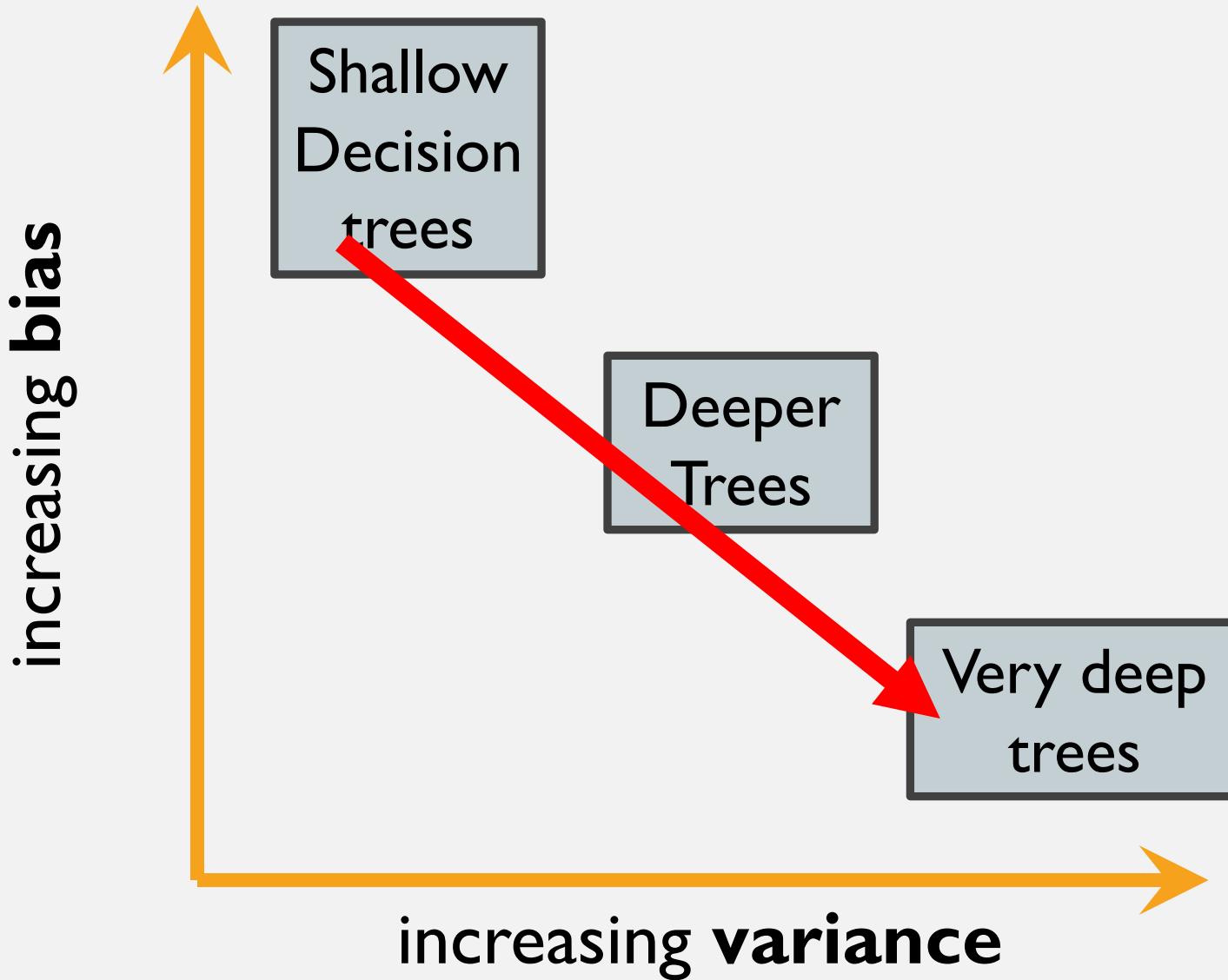
# WHERE DO THEY GO?



# WHERE DO THEY GO?



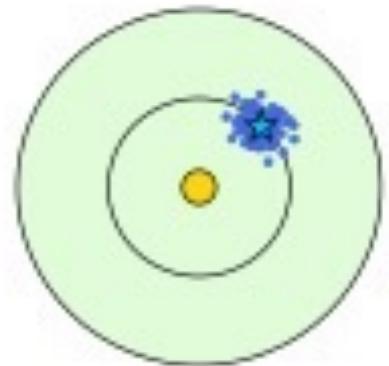
# WHERE DO THEY GO?



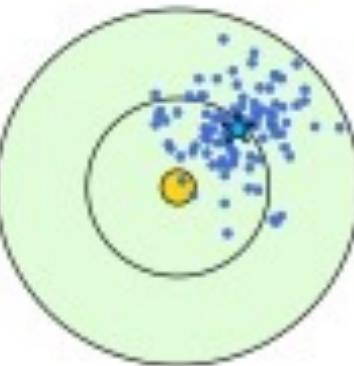
# SO, NOW WE KNOW WHAT BIAS/VARIANCE ARE...

1. We talked about squared loss. **What about other losses?**
2. What can we do about **high bias**? Or **high variance**?
3. What happens with **really BIG models**?
4. What is the relation to the **approximation-estimation** decomposition?

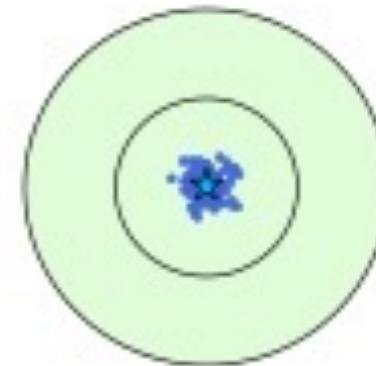
# IT HOLDS FOR OTHER LOSSES – E.G. ... CROSS-ENTROPY



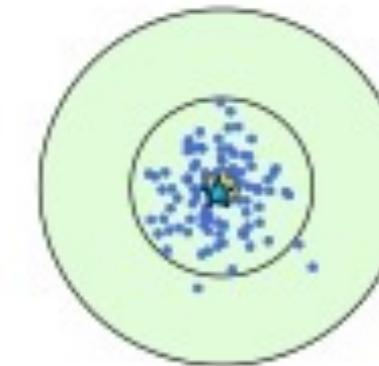
high bias,  
low variance



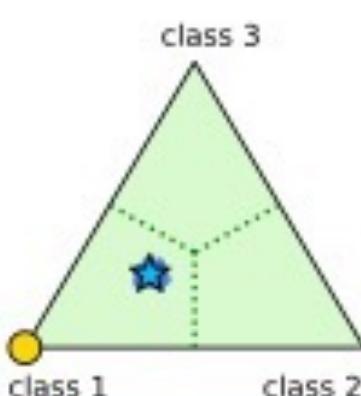
high bias,  
high variance



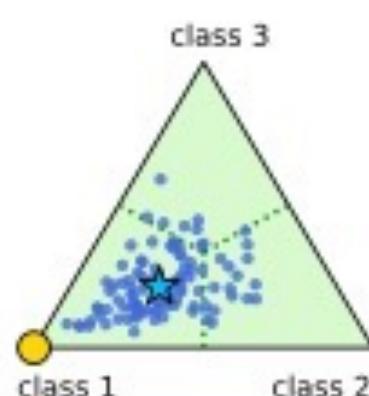
low bias,  
low variance



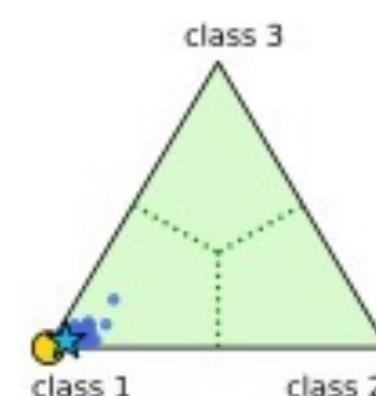
low bias,  
high variance



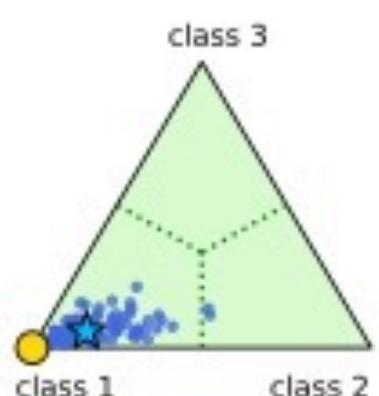
high bias,  
low variance



high bias,  
high variance



low bias,  
low variance



low bias,  
high variance

# IT HOLDS FOR OTHER LOSSES ...

$$\underbrace{\mathbb{E}_{\mathcal{S}_n}[R(f)]}_{\text{expected squared risk}} = \mathbb{E}_{\mathbf{x}} \left[ \underbrace{\mathbb{E}_{y|\mathbf{x}}[(y - \mathbb{E}_{y|\mathbf{x}}[y])^2]}_{\text{noise}} + \underbrace{\left(\mathbb{E}_{\mathcal{S}_n}[f(\mathbf{x})] - \mathbb{E}_{y|\mathbf{x}}[y]\right)^2}_{\text{bias}} + \underbrace{\mathbb{E}_{\mathcal{S}_n}[(f(\mathbf{x}) - \mathbb{E}_{\mathcal{S}_n}[f(\mathbf{x})])^2]}_{\text{variance}} \right].$$

With squared loss, we used the ... ARITHMETIC mean  $\mathbb{E}_{\mathcal{S}_n}[f(\mathbf{x})]$

$$\underbrace{\mathbb{E}_{\mathcal{S}_n}[\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\ell(\mathbf{y}, f(\mathbf{x}))]]}_{\text{expected cross-entropy risk}} = \mathbb{E}_{\mathbf{x}} \left[ \underbrace{\mathbb{E}_{\mathbf{y}|\mathbf{x}}[K(\mathbf{y} \parallel \mathbf{y}^*)]}_{\text{noise}} + K(\mathbf{y}^* \parallel \overset{\circ}{f}(\mathbf{x})) + \underbrace{\mathbb{E}_{\mathcal{S}_n}[K(\overset{\circ}{f}(\mathbf{x}) \parallel f(\mathbf{x}))]}_{\text{variance}} \right],$$

With cross-entropy, we have instead... GEOMETRIC mean  $\overset{\circ}{f}(\mathbf{x})$

Other losses... other types of mean... e.g. harmonic mean.

# SO WHAT CAN WE DO IN PRACTICE?

If you think you have....

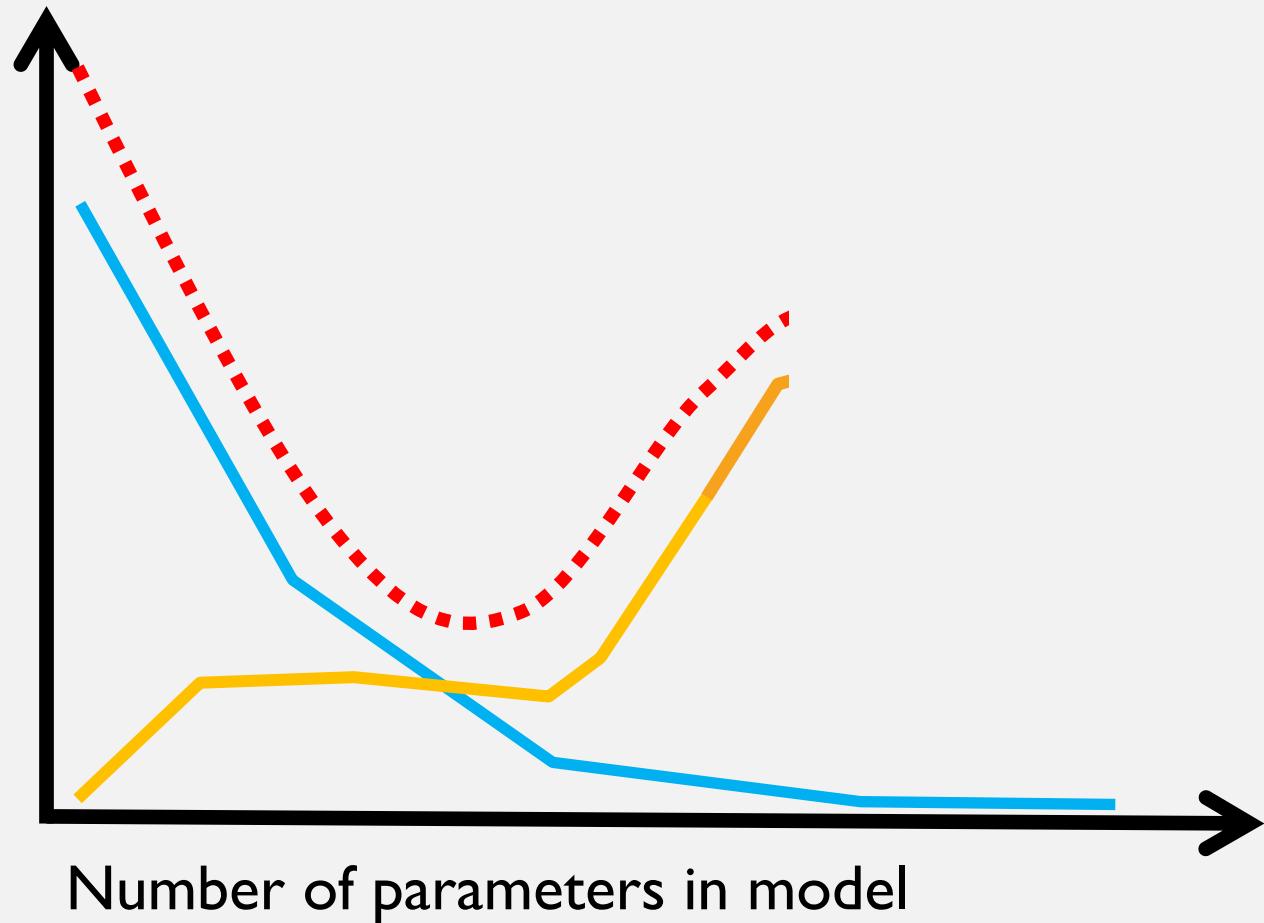
**High bias...**

1. Make the model more complex.
2. Get more labelled training data.
3. Use “ensemble” method: *Boosting* (coming next week)

**High variance...**

1. Make the model simpler.
2. Use regularization (e.g. weight decay)
3. Use “ensemble” method: *Bagging* (coming next week)

# IT GETS WEIRD WITH REALLY HUGE MODELS



Imagine you have more parameters than training data points...

= “**over-parameterized**” model

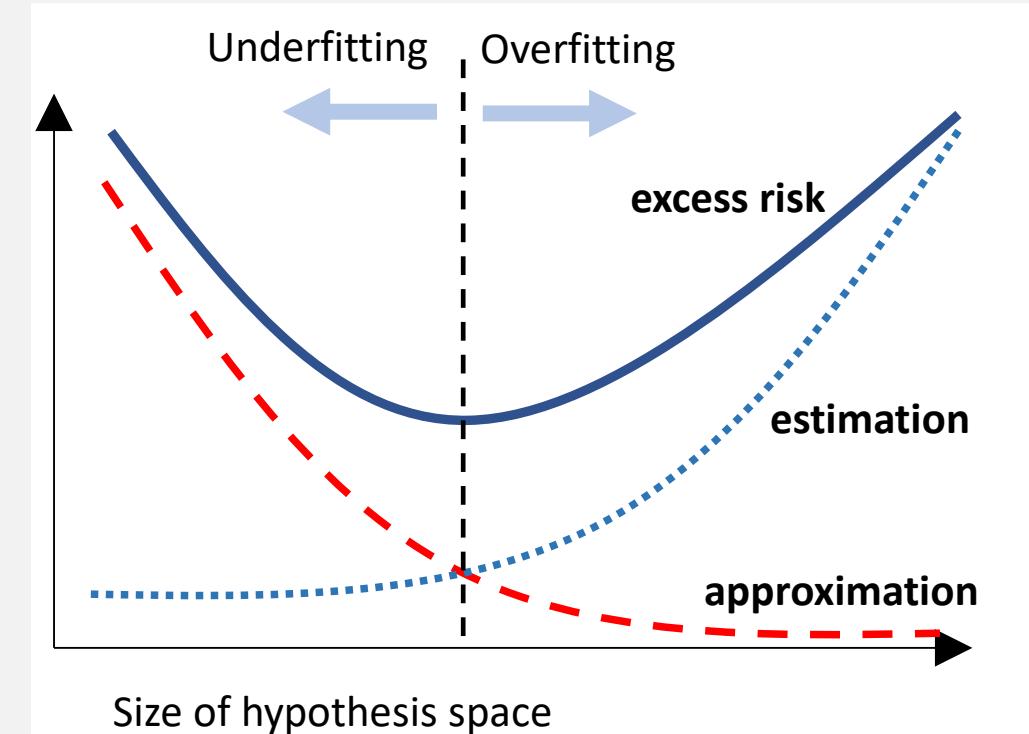
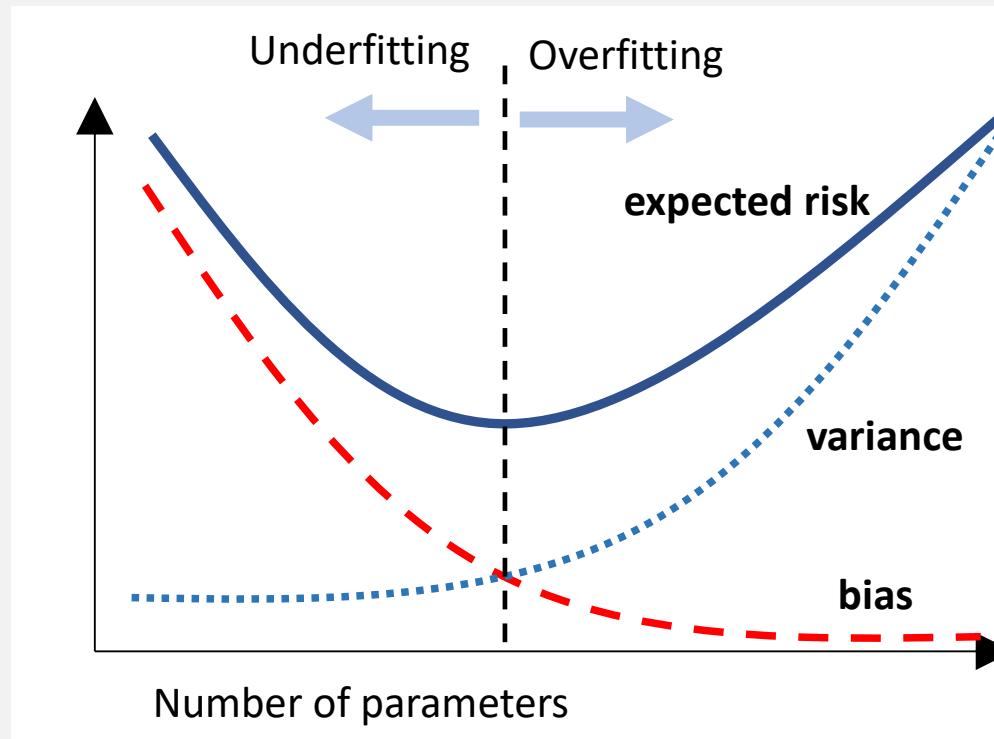
The “**double descent**” phenomenon.

Risk goes down ...then up ...then down again!

Somewhat explained by “peaking” variance.

**Many open research questions.**

# BIAS/VARIANCE == APPROXIMATION/ESTIMATION ???



Hmmmm...

# BIAS/VARIANCE IS NOT THE SAME AS APPROXIMATION/ESTIMATION

## OPTIONAL READING

To be published any day now. This shows just how close you are to cutting-edge research.

### Bias/Variance is not the same as Approximation/Estimation

Gavin Brown  
Riccardo Ali

[gavin.brown@manchester.ac.uk](mailto:gavin.brown@manchester.ac.uk)  
[riccardo.ali.it@gmail.com](mailto:riccardo.ali.it@gmail.com)

#### Abstract

We study the relation between two classical results: the bias-variance decomposition, and the approximation-estimation decomposition. Both are important conceptual tools in Machine Learning, helping us describe the nature of model fitting. It is commonly stated that they are “closely related”, or “similar in spirit”. However, sometimes it is said they are equivalent. In fact they are different, but have subtle connections cutting across learning theory, classical statistics, and information geometry, that (very surprisingly) have not been previously observed. We present several results for losses expressible as a Bregman divergence: a broad family with a known bias-variance decomposition. Discussion and future directions are presented for more general losses, including the 0/1 classification loss.

#### 1 Introduction

Geman et al. (1992) introduced the bias-variance decomposition to the Machine Learning community, and Vapnik & Chervonenkis (1974) introduced the approximation-estimation decomposition, founding the field of statistical learning theory. Both decompositions help us understand model fitting: referring to model size, and some kind of trade-off. The terms are often used interchangeably. And yet, they are different things. Given their fundamental nature and similar purposes, it is surprising that explicit connections are not widely known—perhaps due to differing notations and conventions of their respective communities.

The main intuition behind this decomposition is that a model's performance can be broken down into three components: bias, variance, and irreducible error.

This is Riccardo.

**He was a student in this class,  
in February 2023.**



# TODAY WE DID...

The bias-variance decomposition for SQUARED LOSS

**Homework:** the same decomposition for CROSS ENTROPY LOSS

This week :

- ...Read the notes, thoroughly.
- ...Go over the proofs.
- ...Play with the code.

COMP34312: Mathematical Topics in Machine Learning

Gavin.Brown@manchester.ac.uk  
Anirbit.Mukherjee@manchester.ac.uk

**What is this module about?**

This has been an interesting new module to design. We didn't want to just teach you a bunch of fashionable advanced ML models. You can learn that yourself from other online resources. Even if we did, given how fast the field moves, they'll be out of date in a couple of years.

Instead, we decided to help you understand a fundamental open question that challenges the state-of-the-art in our field. The topic was selected to be relatively close to our research interests, meaning we can help you understand some of the very latest issues. So, what's the question? That's it! What's the question?

Modelling all is going ITC. It seems like Google or Facebook put out a press release every other week, about their latest 100 billion parameter neural learning model. Or it's now 200, or 500 billion parameters? This module will give a formal, mathematical treatment to the following question:

"Are bigger models always better models?"

**Lecture:** Tuesdays, 12pm, Manchester Cooper Building, room G20.

**Example class:** Thursdays 1pm, Kilburn Building, room 1.8.

**Assessment:** 20% open-book online MCQ (Fri 17th March & Fri 12th May).  
80% closed-book exam in late May.

The module will be delivered in weeks 1-6 by Gavin, and in weeks 7-12 by Anirbit.  
This pack of notes contains material for Gavin's part. Anirbit's will follow in due course.

THANKS