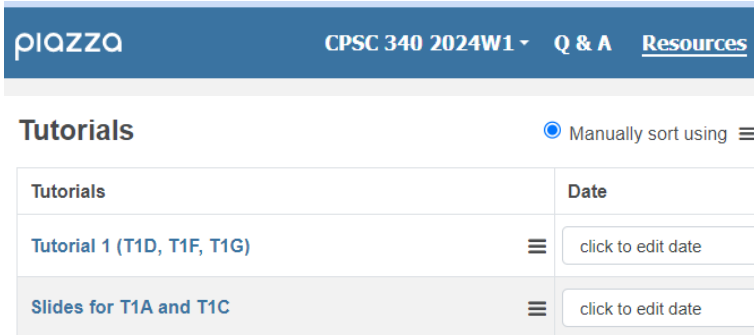# CPSC 340/540 Tutorial 2

## Winter 2024 Term 1

T1A: Tuesday 16:00-17:00;
T1C: Thursday 10:00-11:00;
Office Hour: Wednesday 15:00-16:00

<span style="color:red">Slides can be found at Piazza and my personal page after T1C.</span>

**Piazza** — CPSC 340 2024W1 ▾ **Q & A** **Resources**

**Tutorials**                    ● Manually sort using ☰

| Tutorials | Date |
|---|---|
| Tutorial 1 (T1D, T1F, T1G) ☰ | click to edit date |
| Slides for T1A and T1C ☰ | click to edit date |

## Yi (Joshua) Ren

https://joshua-ren.github.io/
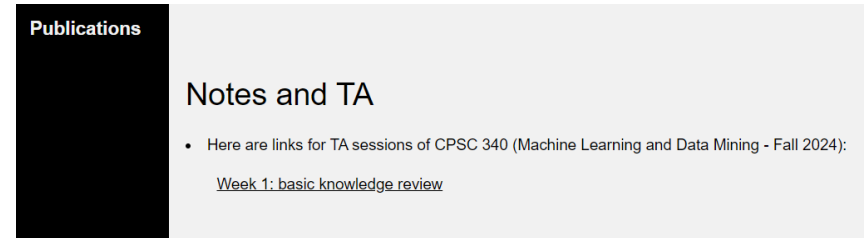renyi.joshua@gmail.com

PhD with Danica

Machine Learning:
Learning dynamics, LLM, Compositional Generalization

**Publications**

Notes and TA

- Here are links for TA sessions of CPSC 340 (Machine Learning and Data Mining - Fall 2024):
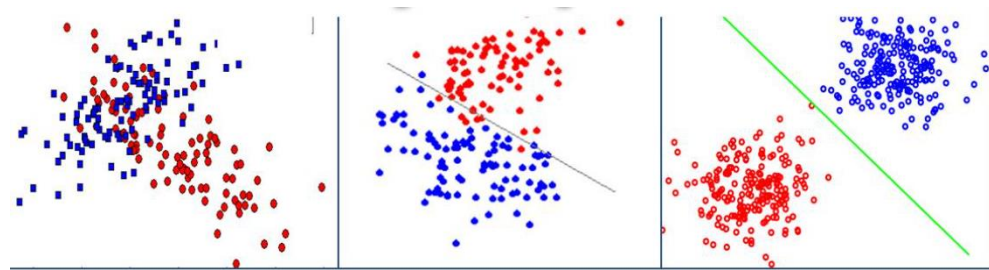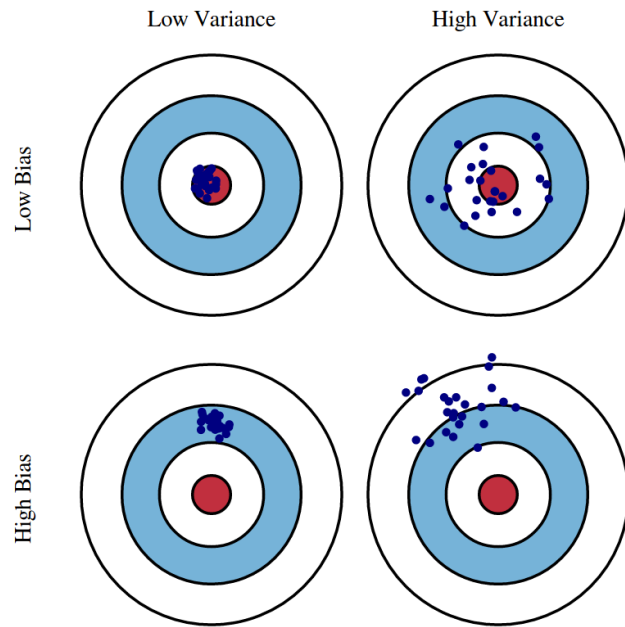    Week 1: basic knowledge review

**More helpful on theory**

**Less helpful on coding**

- Variance-bias trade-off
- KNN
- Naive Bayes

# Variance-bias trade-off (traditional discussion)

Low Variance     High Variance

Low Bias

High Bias



The "noise" is becoming smaller.

- Expected squared test error in this setting is

$$E\left[(\tilde{y}_i - \hat{y}_i)^2\right] = E\left[(\hat{y}_i - \bar{y}_i)\right]^2 + \left(E[\hat{y}_i^2] - E[\hat{y}_i]^2\right) + \sigma^2$$

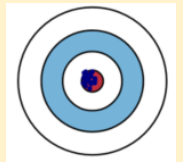"test squared error"          "bias"          "variance"          "noise"

  – Where expectations are taken over possible training sets of 'n' examples.
  – Bias is expected error due to having wrong model.
  – Variance is expected error due to sensitivity to the training set.
  – Noise (irreducible error) is the best can hope for given the noise ($E_{best}$).

- Some learning theory results use $E_{best}$ to further decompose $E_{test}$:

$$E_{test} = \left(E_{test} - E_{train}\right) + \left(E_{train} - E_{best}\right) + E_{best}$$

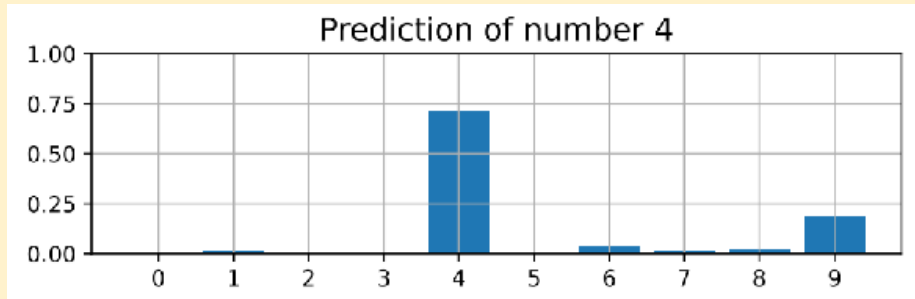$E_{gap}$          $E_{model}$          "noise"

- $E_{gap}$ measures *how sensitive we are to training data.*
- $E_{model}$ measures *if our model is complicated enough to fit data.*
- $E_{best}$ measures how low can **any** model make test error.
  – $E_{best}$ does not depend on what model you choose.

# Variance-bias trade-off (but when data is non-seperable)
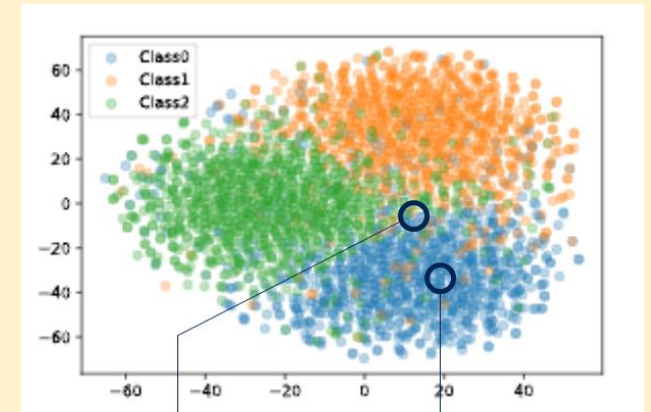
**Not always the best!**

- Network's **calibration**: low variance (high confidence) is not always good

  ➢ Fact: most of the time, our model gives a probabilistic prediction, e.g., spam-filter, MNIST, …

  ➢ Different samples with the same label can be different.

  ➢ Then, we want **confidence** aligns well with **facts**.
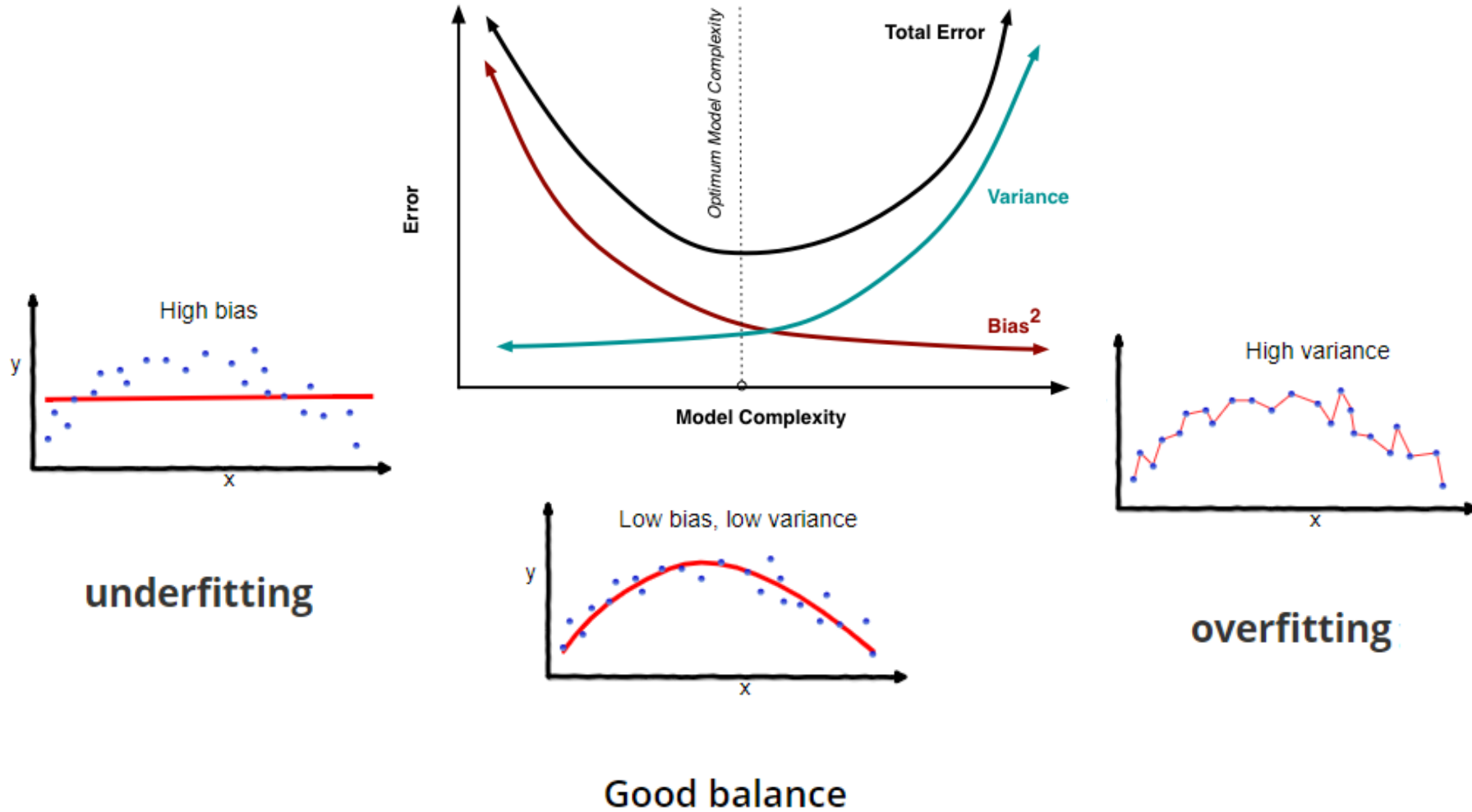
$p^*(y|x) \approx [1,0,0]$

$p^*(y|x) \overset{?}{\approx} [1,0,0]$
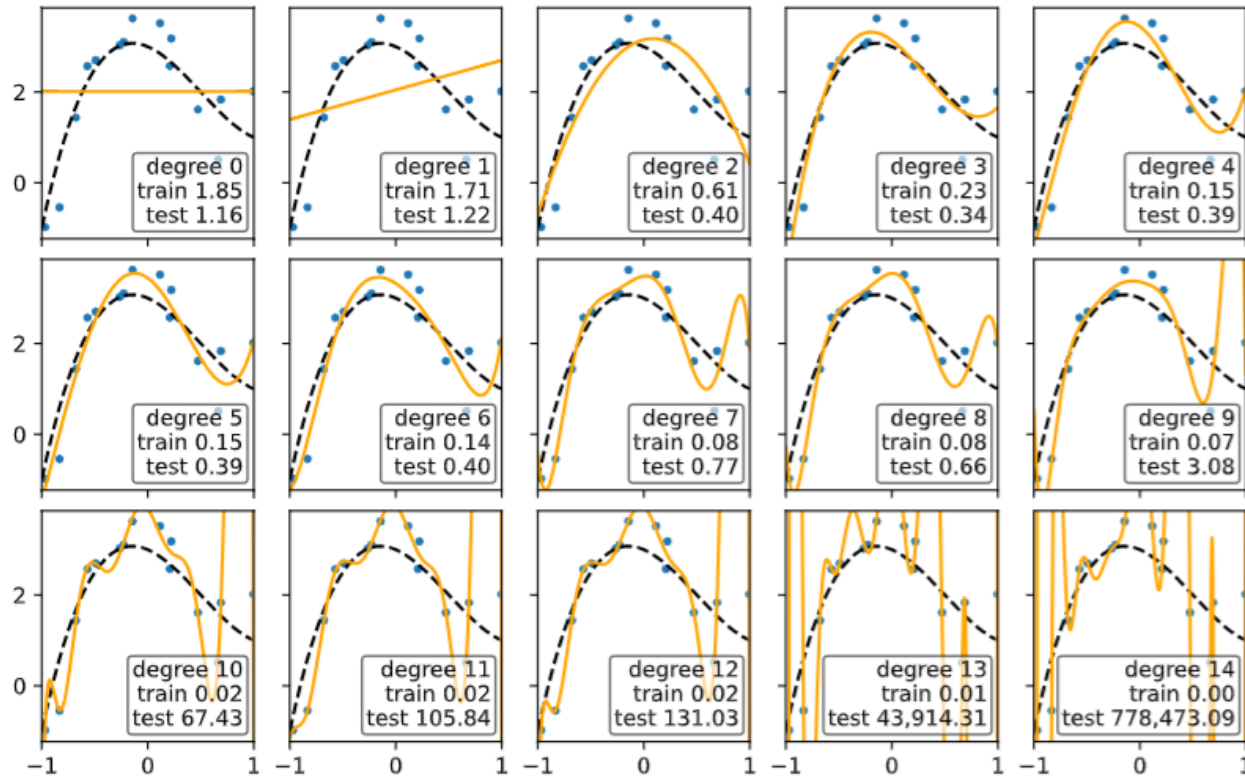
$p^*(y|x) = [0.98, 0.01, 0.01]$
$e_{y_n}^{\mathrm{T}} = [1, \quad 0, \quad 0]$

$p^*(y|x) = [0.50, 0.30, 0.20]$
$e_{y_n}^{\mathrm{T}} = [1, \quad 0, \quad 0]$

Guo, Chuan, et al. "On calibration of modern neural networks." *ICML*, 2017. (6k+ citations)
Ren, Yi, Shangmin Guo, and Danica J. Sutherland. "Better supervisory signals by observing learning paths." ICLR 2022

# Variance-bias trade-off (traditional)
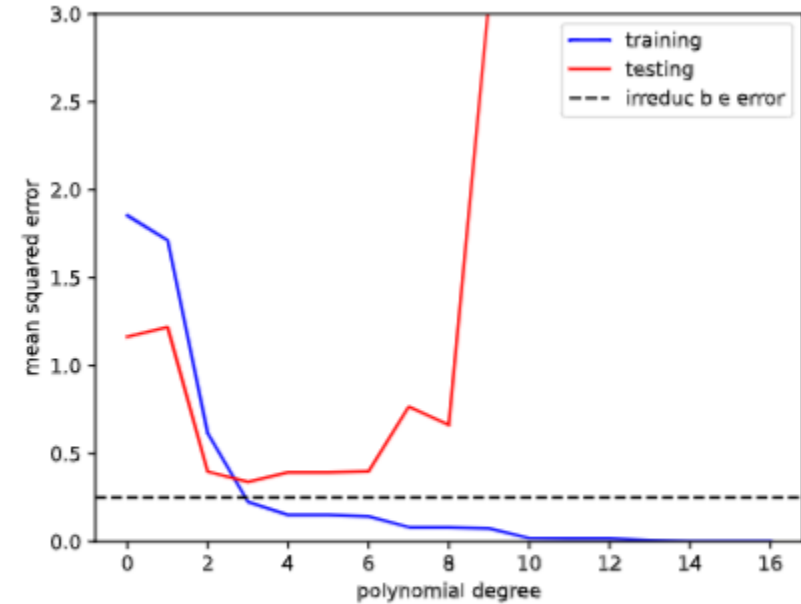
# Variance-bias trade-off (double descent, benign overfitting)



(a) Polynomial regression, $h(x) = w_0 + w_1 x + w_2 x^2 + \cdots + w_k x^k$, for increasing $k$, to data points shown in blue. ERM fits are in orange; dashed black lines show $\mathbb{E}[y \mid x]$, a cubic function. Text gives mean squared error for training and testing sets.
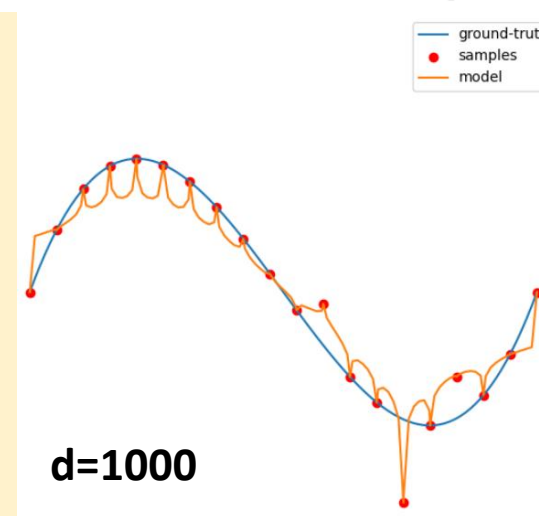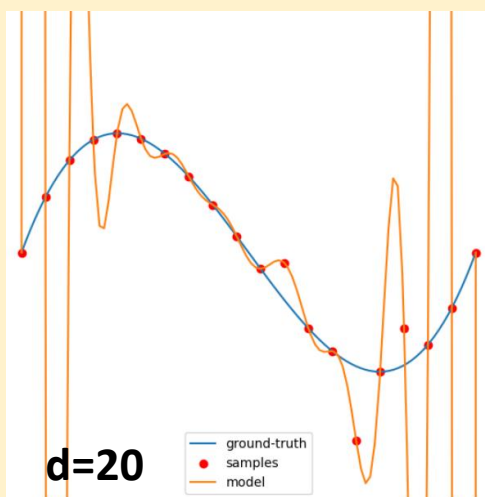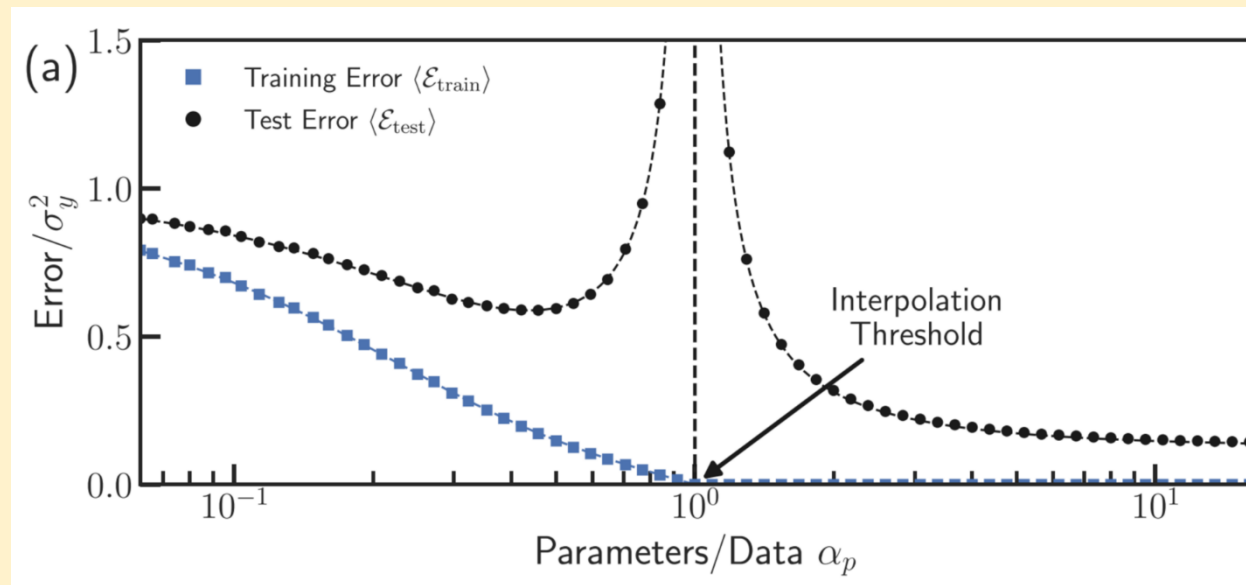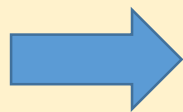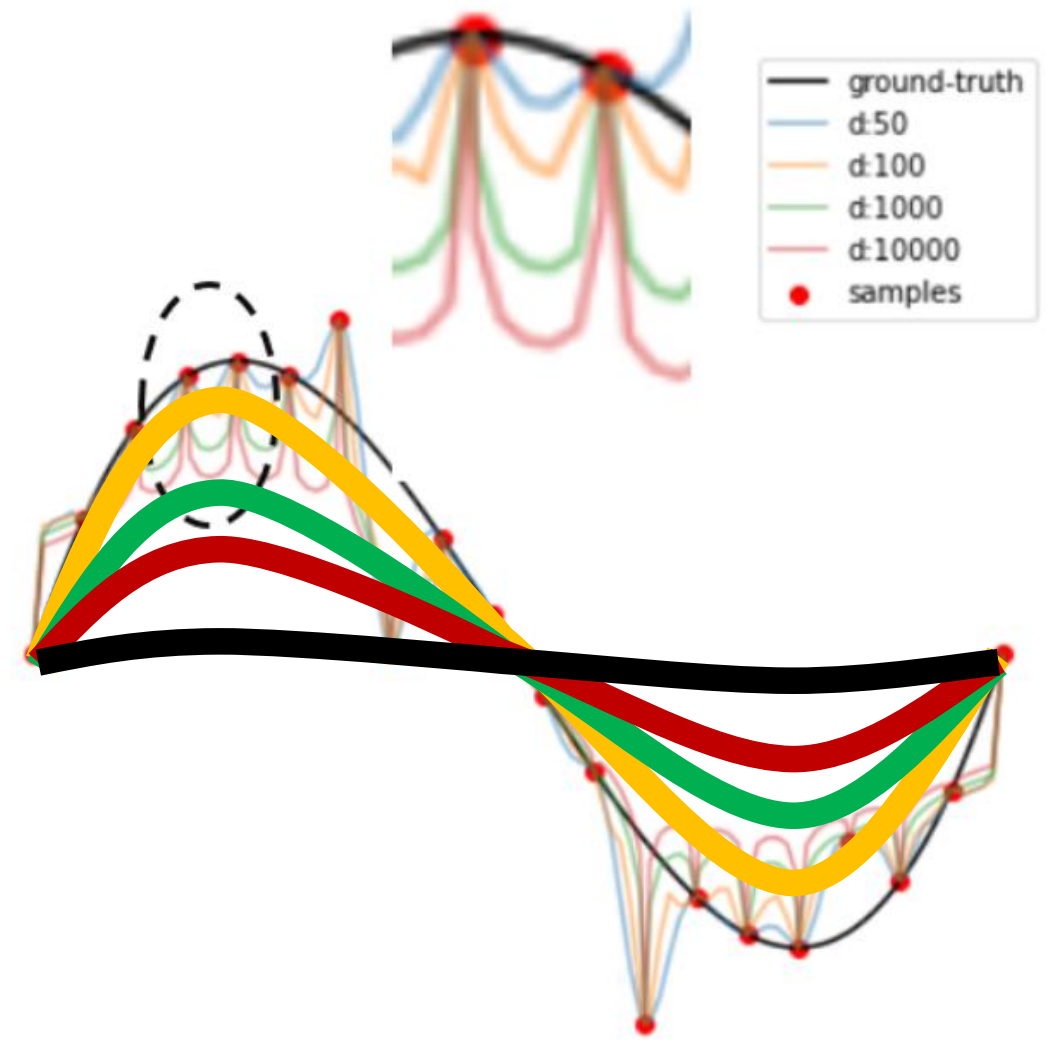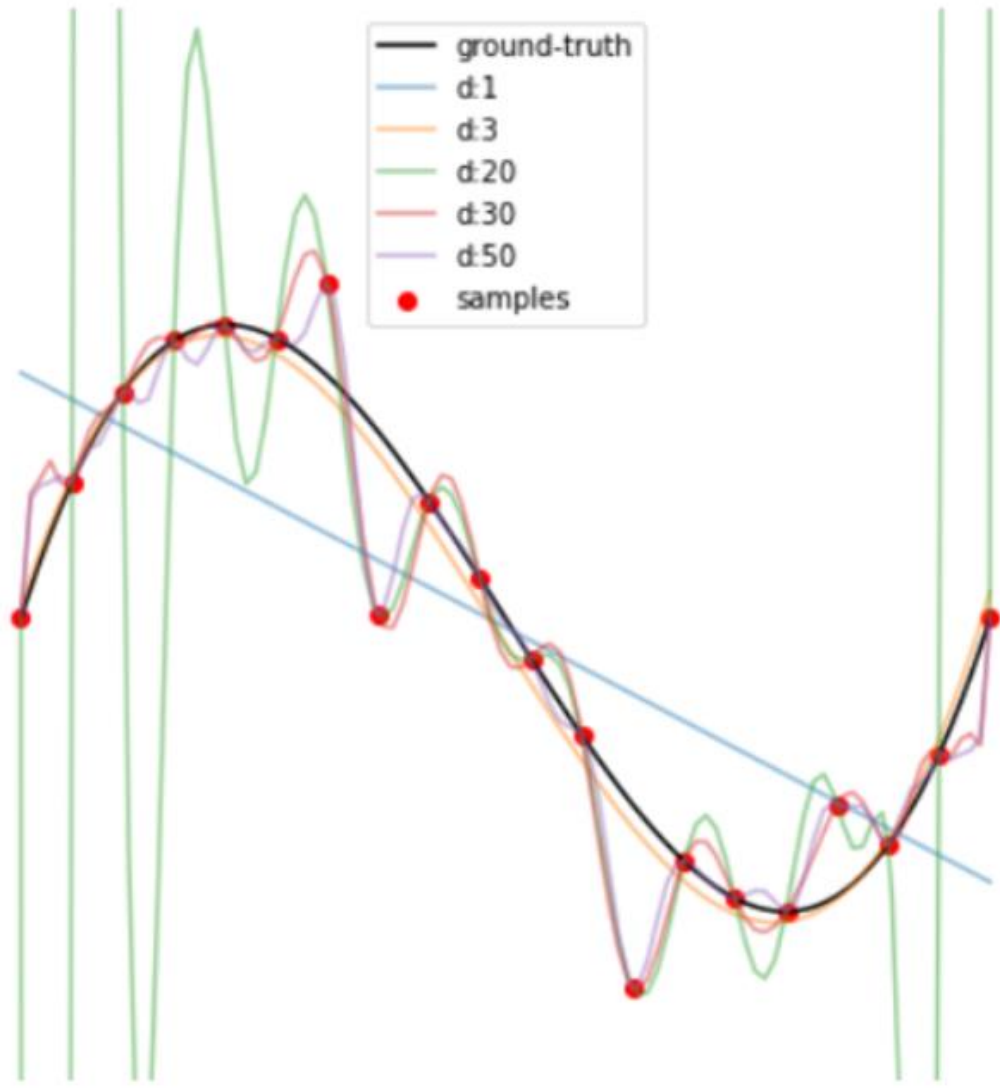
(b) Training and test errors from Figure 1.1a.

What happens when d=2000?

Figures from: https://www.cs.ubc.ca/~dsuth/532D/24w1/notes/1-intro-erm.pdf

https://colab.research.google.com/drive/1UJYKXj317aJGeIgwV0qqL3MhUnHf9VJK#scrollTo=PrA4y-mEJZZW

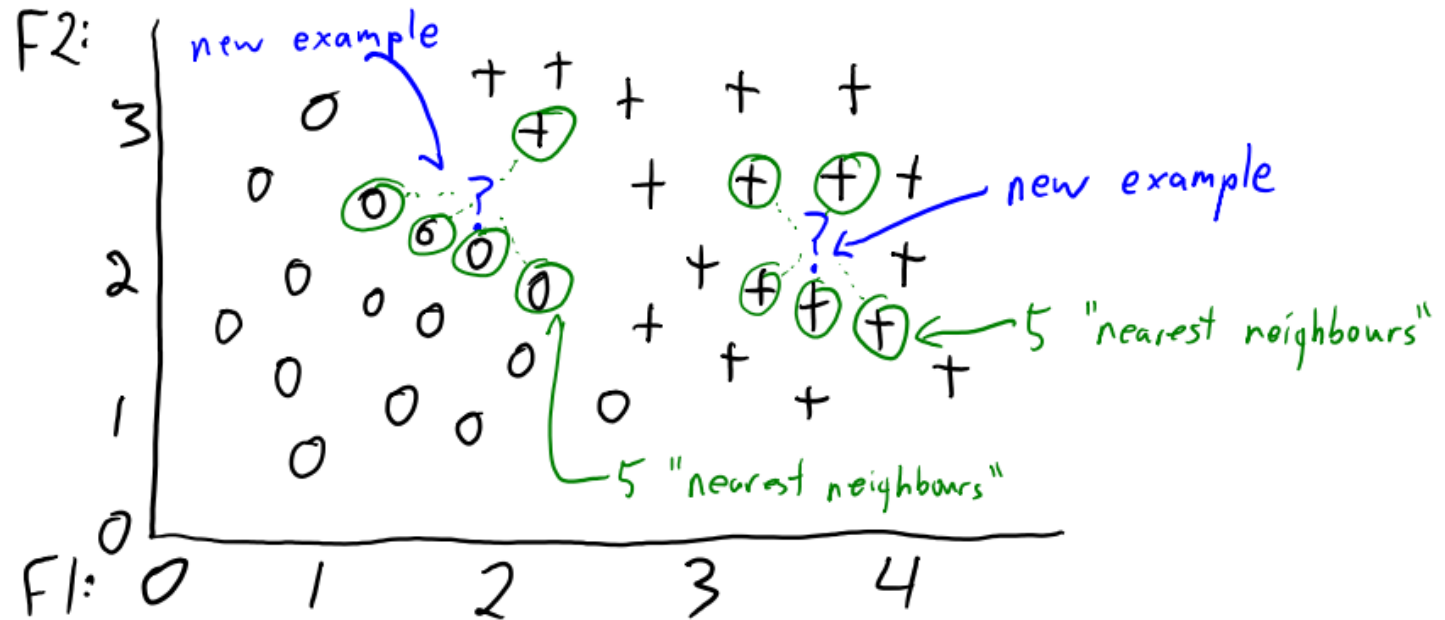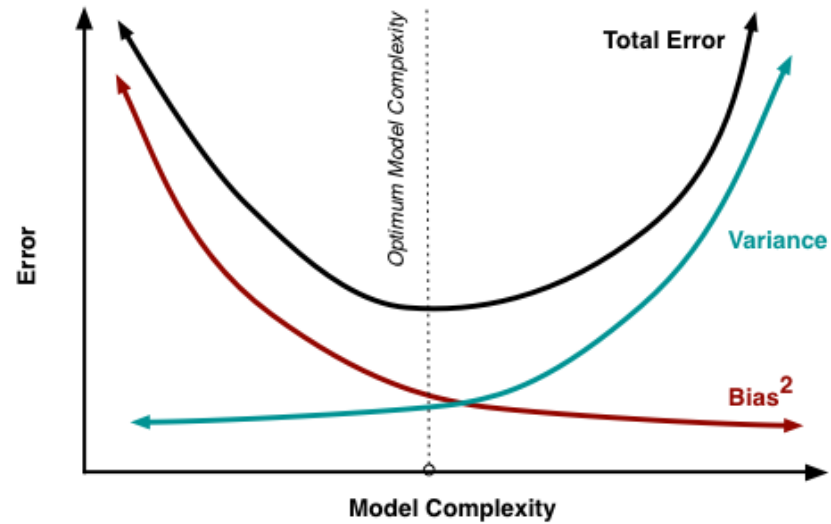# Variance-bias trade-off (double descent (DNN), benign overfitting)

# KNN (Algorithm and implementation)

- To classify an example $\tilde{x}_i$:
    0. Define distance
    1. Find the 'k' training examples $x_i$ that are "nearest" to $\tilde{x}_i$.
    2. Classify using the most common label of "nearest" training examples.
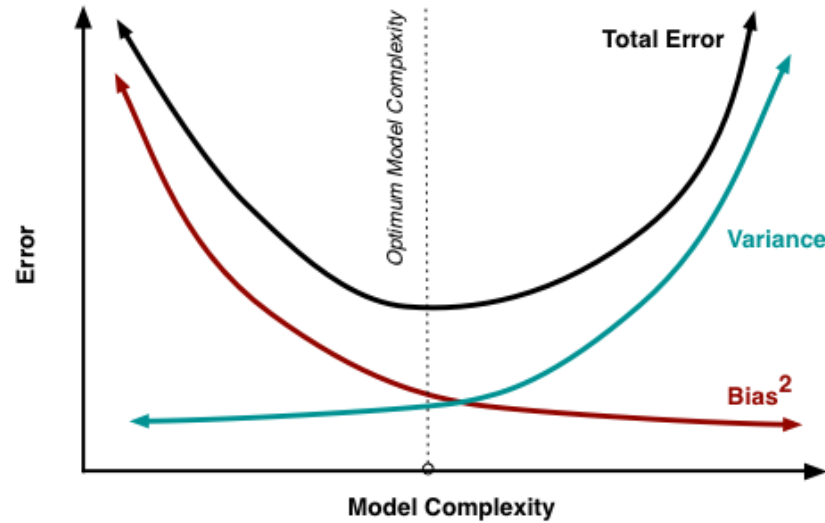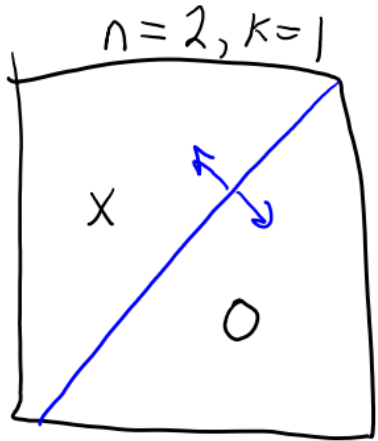
# KNN (bias-variance trade-off)
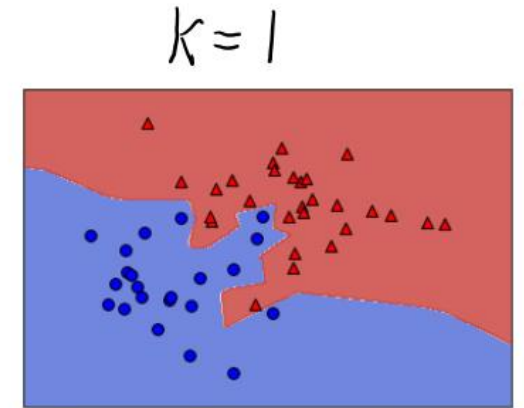


Q: how to put the value of "n" and "k" in this diagram?

# KNN (bias-variance trade-off)

$$f = WX + k|W|_2^2$$

$$= W \begin{bmatrix} x \\ ... \\ x^n \end{bmatrix} + k|W|_2^2$$



$n=2, k=1$



$n=20, k=1$

A: **larger n** means higher model complexity,
   **larger k** behaves like stronger regularization

$k=10$

$k=3$

$k=1$

# KNN (Curse of Dimensionality)

- Fact 1: need exponential more examples to get reasonable good neighbours
  - Volume of space grows exponentially with dimension.
    - Circle has area $O(r^2)$, sphere has area $O(r^3)$, 4d hyper-sphere has area $O(r^4)$,...
  - Need exponentially more points to 'fill' a high-dimensional volume.
    - "Nearest" neighbours might be really far even with large 'n'.

  -- Assume $r < 0.05$ is a reasonable choice on unit ball



neighbor $= 10^{-1}$     $= 10^{-2}$     $= 10^{-3}$

1 d      2 d      3 d

- **That is why many learning methods want DENSE representations and low-rank manifold**

# KNN (Curse of Dimensionality)

- Fact 2: if samples are **uniformly generated**, most samples are on "surface"



nonsurface = 0.9     $= 0.9^2$     $= 0.9^3$

1 d       2 d       3 d

- Split the whole space into several 0.1*0.1*0.1 ... blocks.
- Random select one
- Higher prob. that block comes from the "surface"

- Support 3: since samples are **NOT uniformly generated**, they are on "**low-dim manifold**"



**Think about their distance**

# Naive Bayes (Algorithm and implementation)

- Use bag of words to create features, gets users to label them

| $ | Hi | CPSC | 340 | Vicodin | Offer | ... |
|---|----|------|-----|---------|-------|-----|
| 1 | 1  | 0    | 0   | 1       | 0     | ... |
| 0 | 0  | 0    | 0   | 1       | 1     | ... |
| 0 | 1  | 1    | 1   | 0       | 0     | ... |
| ... | ... | ... | ... | ... | ... | ... |

$x_1 = [110010]$

$x_2 = [000011]$

$x_3 = [011100]$

| Spam? |
|-------|
| 1     |
| 1     |
| 0     |
| ...   |

$y_1 = 1$

$y_2 = 1$

$y_3 = 0$

- Intuition:

$$\text{if } p(y_i = 1|x_i) > p(y_i = 0|x_i)$$

  - return "spam"
else
  - return "not spam"

# Naive Bayes (Algorithm and implementation)

- Supervise learning usually model $p(\mathbf{y}_i|\mathbf{x}_i)$ directly, but here we use Bayes to decompose that:

$$p(\mathbf{y}_i|\mathbf{x}_i) = \frac{p(\mathbf{x}_i|\mathbf{y}_i)p(\mathbf{y}_i)}{p(\mathbf{x}_i)} = \frac{p(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iB}|\mathbf{y}_i)p(\mathbf{y}_i)}{p(\mathbf{x}_i)}$$

- $p(\mathbf{x}_i)$ is usually hard to calculate, because we might not have enough data when "Bag size" $B$ is large
  <span style="color:red">(Recall curse of dimensionality)</span>

$$p(\mathbf{y}_i = 1|\mathbf{x}_i) > p(\mathbf{y}_i = 0|\mathbf{x}_i)$$

$$\frac{p(\mathbf{x}_i|\mathbf{y}_i = 1)p(\mathbf{y}_i = 1)}{p(\mathbf{x}_i)} > \frac{p(\mathbf{x}_i|\mathbf{y}_i = 0)p(\mathbf{y}_i = 0)}{p(\mathbf{x}_i)}$$

$$p(\mathbf{x}_i|\mathbf{y}_i = 1)p(\mathbf{y}_i = 1) > p(\mathbf{x}_i|\mathbf{y}_i = 0)p(\mathbf{y}_i = 0)$$

- Then, $p(\mathbf{x}_i|\mathbf{y}_i)$ is also hard to calculate due to similar reason. <span style="color:red">(Recall curse of dimensionality)</span>
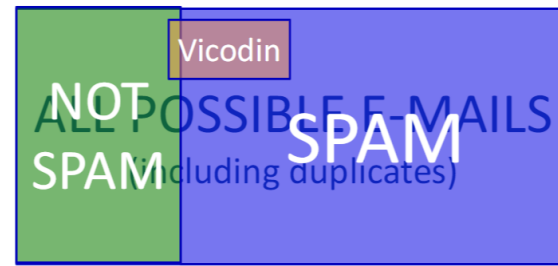  We then <span style="color:red">assume the independence</span> (might introduce bias, but generally OK)

$$p(\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iB}|\mathbf{y}_i) \approx \prod_{b=1}^{B} p(\mathbf{x}_{ib}|\mathbf{y}_i)$$

# Naive Bayes (Algorithm and implementation)

- Now, the task is to estimate $p(x|y)$ for each possible x and y; and the margin prob $p(y)$ for each y

$$p(hello=1, vicodin=0, 340=1 | spam) \approx p(hello=1|spam) \, p(vicodin=0|spam) \, p(340=1|spam)$$

HARD     easy     easy     easy

ALL POSSIBLE E-MAILS (including duplicates)

NOT SPAM

SPAM

Vicodin

- Easy to estimate:

$$p(vicodin=1 | spam=1) = \frac{\# \text{spam messages w/ vicodin}}{\# \text{spam messages}}$$

- Label smoothing: what happen if any term in $\prod_{b=1}^{B} p(\mathbf{x}_{ib}|\mathbf{y}_i)$ is zero?

$$\frac{(\# \text{spam messages with lactase}) + 1}{(\# \text{spam messages}) + 2}$$

- Avoid probability underflow: use **log-prob** instead

$$p(y_i=c | x_i) \propto \prod_{j=1}^{d} \left[ p(x_{ij} | y_i = c) \right] p(y_i = c)$$

→ All these are $< 1$ so the product gets very small.

# Thanks for your time! Questions?