

April 2, 2025 - lecture #1 -

Notation - end of Chapter 1.

$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix}_{n \times p}$$

capital letter for matrices.
(bold in textbook)

of features
of obs

$$i = 1, \dots, n$$

$$j = 1, \dots, p$$

$$x_i = i^{\text{th}} \text{ row of } \mathbf{X} = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\mathbf{x}_j = \vec{x}_j = j^{\text{th}} \text{ column of } \mathbf{X} = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{nj} \end{pmatrix}$$

\vec{a} or a : vector of length n .

a : could be a vector not of length n , or it could be a scalar.

$\vec{a} \in \mathbb{R}^n$ } vector of length n with n in \mathbb{R}

all vectors are column vectors*

* vectors of length n are bold arrows have arrow on top *

$A \in \mathbb{R}^{w \times v}$ } matrix of dim. $w \times v$

$c \in \mathbb{R}$ } scalar

$a \in \mathbb{R}^q$ } vector of length q .

$a \in \{0, 1\}$] a can equal 0 or 1

$A \in \{1, 2, 3, 4, 5\}^{173 \times 12}$] 173×12 matrix, w/ entries that are 1, 2, 3, 4, or 5.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}_{2 \times 2}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 \\ 7 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}_{2 \times 1}$$

Chapter 2 - What is statistical learning?

Response Y^{ER} : what we want to predict.
↳ outcome.

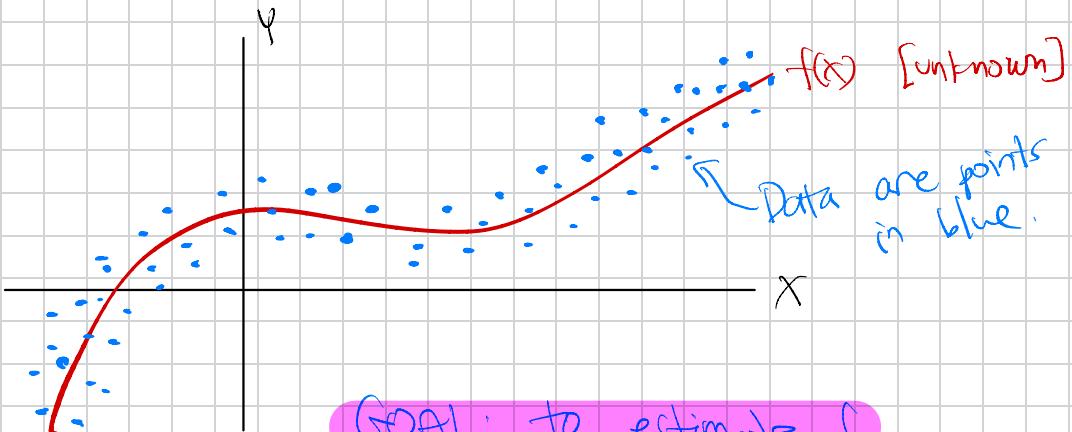
Predictors $X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^p$
↳ features, covariates

mean-zero
noise term
 $E(\epsilon) = 0$

$$Y = f(X) + \epsilon = f(x_1, x_p) + \epsilon$$

f ↴
 f is UNKNOWN.
Our goal is to est f .

$$f: \mathbb{R}^p \rightarrow \mathbb{R}$$



(\hat{f} is our estimate of f .)

Why do I want to est. f?

Prediction

Inference.

PREDICTION: we want a "good" est. of f , so that we can compute $\hat{y} = \hat{f}(x)$ for a new x .

e.g. I want to predict value of house based on # BR's, proximity to light rail, etc

INFERENCE: we want to understand the relationship btw $X_1, \dots, X_p \rightarrow Y$, if we are interested in it

e.g.: how much would this house's value increase if we added another BR?

How to estimate $f(\cdot)$?

Training data: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix}, \quad y \in \mathbb{R} \text{ (usually)}$$

Goal: $f(\cdot)$ such that $f(x) \approx y$
for all (x, y) .

→ Parametrically (Option 1)

→ non-parametrically (Option 2)

Option 1: Estimate f parametrically.

Step 1: "assume" some functional form for relationship b/w X & y . For example

$$y = f_0 + f_1 x_1 + \dots + f_p x_p + \epsilon$$

$f(x)$

Step 2: Use $\{(x_1, y_1), \dots, (x_n, y_n)\}$ to "fit the model", i.e. estimate f_0, f_1, \dots, f_p (call the estimates $\hat{f}_0, \hat{f}_1, \dots, \hat{f}_p$).

Examples: linear regression (L2)
logistic regression (L3)
:
:

Option 2: Estimate f non-parametrically.

⇒ Do NOT make explicit assumptions about the form of f .

⇒ Instead, find \hat{f} that is close to the training data & "not too wiggly".



Why would we prefer parametric versus non-parametric?

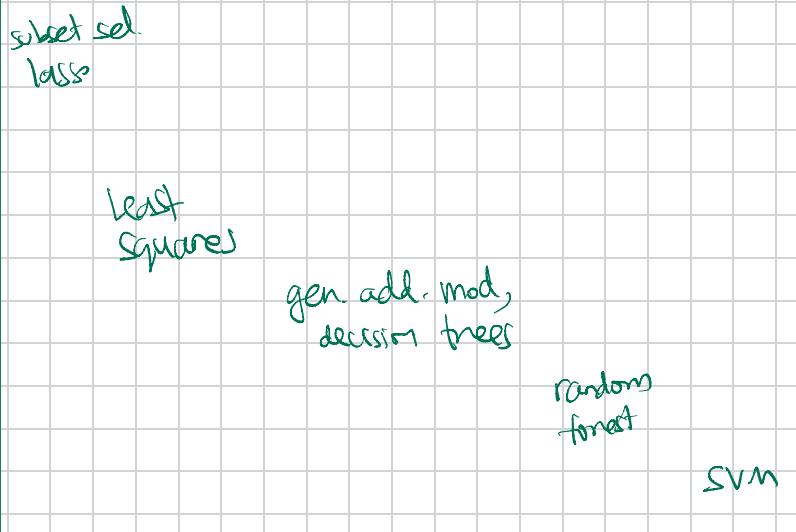
Parametric:

- ⊕ simple → easy to interpret -- useful for inference → easy computationally (?)
- ⊕ guard against overfitting: simplifies estimation task, which helps w.r.t. bias-variance trade-off. (B-v T.O.) especially when sample size (n) small.
- ⊕ we might "believe" the parametric model
- ⊖ strong assumption about f , & it might be wrong!

Non-Parametric

- (+) super flexible!! no assumption on f(.) is needed.-
- (-) Super flexible!! get hurt by B/r T/o.
(overfitting).
- (-) LARGE SAMPLE SIZE.
(+ relatively few features)
 $n \ggg p.$
- (-) not great for inference/interpretation,

↑
↑ Interpretability
↑



DL

2 Types of Statistical Learning

Supervised

(X, Y) pairs.

$$Y = f(X) + \epsilon$$

e.g., regression,

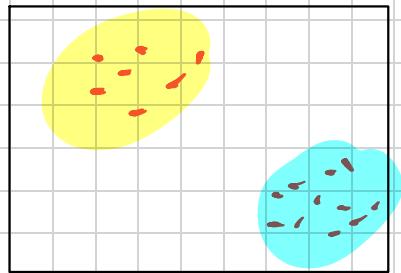
classification,

etc.,

DL, Xgboost, RFs, ...

Unsupervised

No Y ! Just X 's.



e.g. clustering, dimension reduction, etc

Supervised Learning

Regression: $y \in \mathbb{R}$

Classification

Quantitative
cont
real-valued

$y \in \{\text{blue, brown, green, black, hazel}\}$

$y \in \{0, 1\}$

discrete qualitative, categorical

$$y = f(x) + \epsilon \quad E(\epsilon) = 0$$

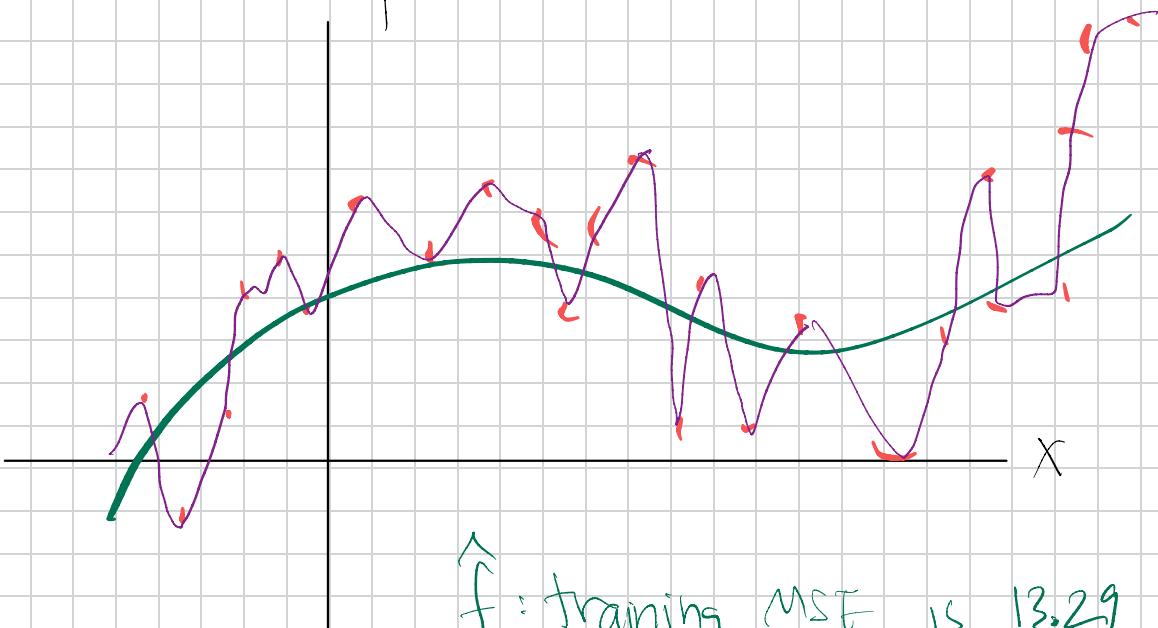
unknown.

How do we assess model accuracy?

$$(\text{Training}) \text{ MSE} := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

{ evaluated on same n obs. used
to est f (training set).

We want (Training) MSE to
be small,
 y i.e., $y_i \approx \hat{f}(x_i)$, $i=1, \dots, n$



\hat{f} : training MSE is 13.29

\hat{f} : training MSE is 0.000000

What I actually care about is [naught not]

test MSE: for a "test" obs (x_0, y_0)

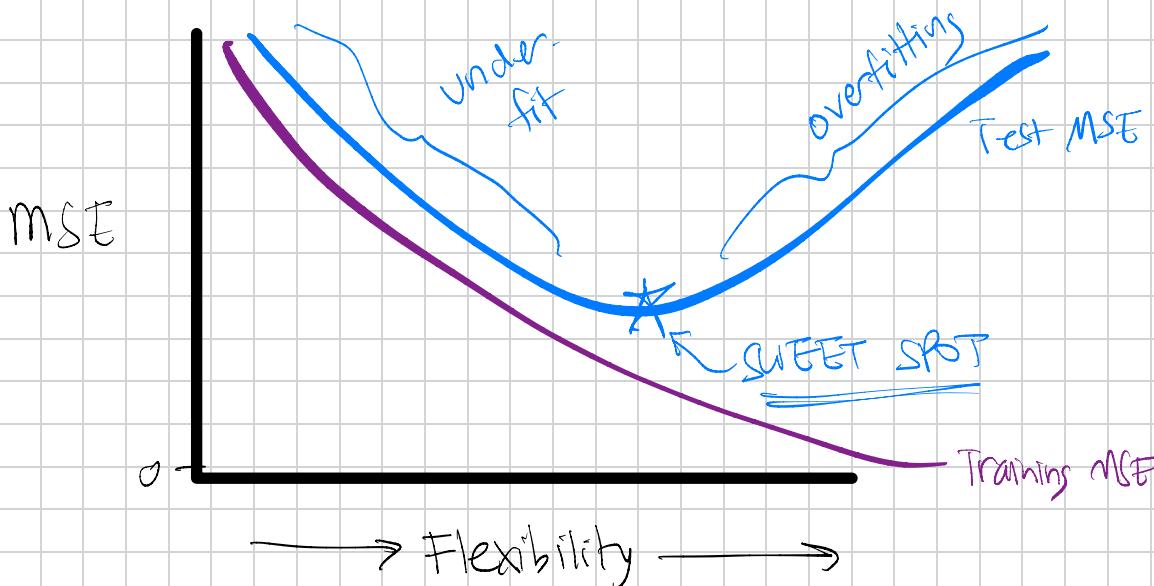
that wasn't in training set, +
that's drawn from same distribution
as training set,

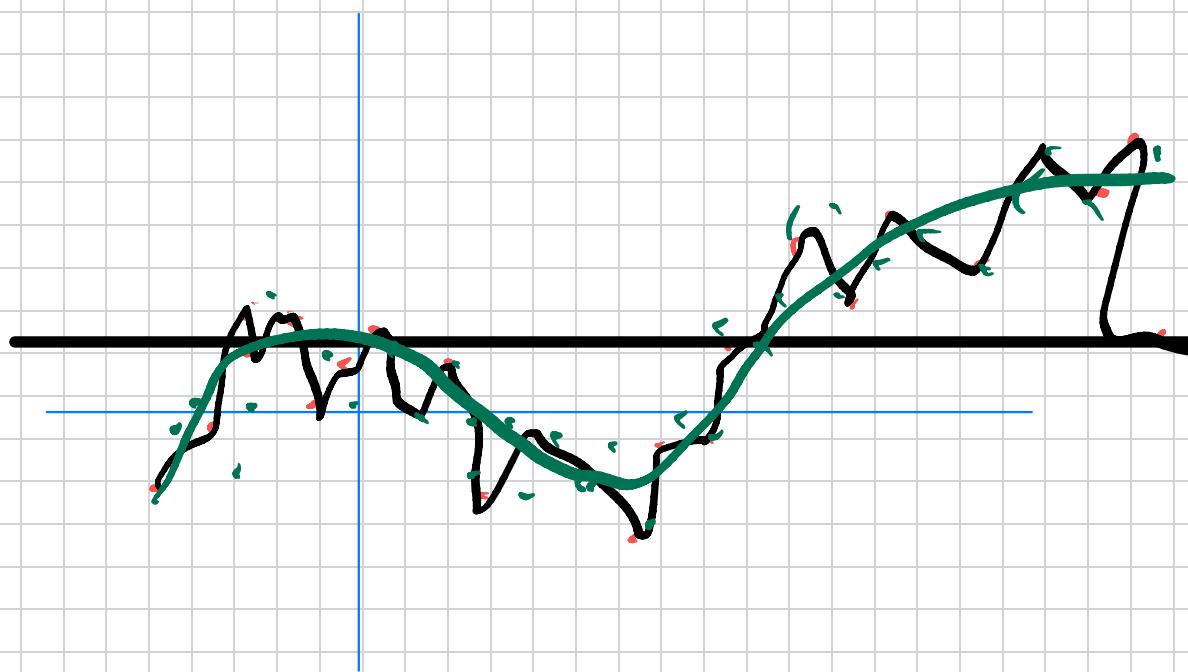
$$\text{Test MSE} = \text{Average} \left[(y_0 - \hat{f}(x_0))^2 \right]$$



average over a whole
batch of (x_0, y_0) 's.

I want THIS to be small!





Bias-Variance Trade-off

$$Y = f(X) + \epsilon, \quad E(\epsilon) = 0$$

↑
unknown; est w/ training set

Test obs. (x_0, y_0)

In this class: X fixed (not random).

ϵ random,
so Y is random b/c of ϵ

Expected prediction error (exp test MSE).

$$\text{EPE} := E \left[[y_0 - \hat{f}(x_0)]^2 \right]$$

$$EA^2 = E(A^2) \neq (EA)^2 = [E(A)]^2$$

expectation over (x_0, y_0) 's
AND over training set

$$E[(Y_0 - \hat{f}(X_0))^2]$$

$Y = f(X) + \varepsilon$
 $E(\varepsilon) = 0$

$$= E\left[\underbrace{(Y_0 - f(X_0) + f(X_0) - E\hat{f}(X_0) + E\hat{f}(X_0) - f(X_0))^2}_{A+B+C}\right]$$

Note: $(A+B+C)^2 = A^2 + B^2 + C^2 + 2AB + 2AC + 2BC$

$$= E[A^2 + B^2 + C^2 + 2AB + 2AC + 2BC]$$

$$= E(A^2) + E(B^2) + E(C^2) + 2E[AB] + 2E[AC] + 2E[BC]$$

$$E(A^2) = E[(Y_0 - f(X_0))^2] = E(\varepsilon^2) = E(\varepsilon^2) - (E\varepsilon)^2$$

$= \text{Var}(\varepsilon)$

$$E(B^2) = E\left[\underbrace{(f(X_0) - E\hat{f}(X_0))^2}_{\text{not random}}\right] = \underbrace{[f(X_0) - E\hat{f}(X_0)]^2}_{\text{not random}}$$

$= [\text{Bias}(\hat{f}(X_0))]^2$

recall: parameter θ , estimator $\hat{\theta}$

$$(E\hat{\theta} - \theta)^2 = (\text{Bias}(\hat{\theta}))^2$$

$$E(C^2) = E[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2] = \text{Var}[\hat{f}(x_0)]$$

recall: param θ , estimator $\hat{\theta}$,

$$E[(\hat{\theta} - E(\hat{\theta}))^2] = \text{Var}(\hat{\theta})$$

$$E\left[\underbrace{(y_0 - f(x_0)}_A + \underbrace{(f(x_0) - E[\hat{f}(x_0)] + E[\hat{f}(x_0)] - \hat{f}(x_0))^2}_B\right] \quad (E(3x) = 3E(x))$$

$$E(AB) = E\left\{(y_0 - f(x_0))(f(x_0) - E[\hat{f}(x_0)])\right\}$$

not random

$$= (f(x_0) - E[\hat{f}(x_0)]) E\left[y_0 - f(x_0)\right]_{E_0}$$

$$= (f(x_0) - E[\hat{f}(x_0)]) E(\varepsilon) = 0$$

$$E(AC) = E\left[\underbrace{(y_0 - f(x_0))}_{E_0} (\underbrace{E[\hat{f}(x_0)] - \hat{f}(x_0)}_{\text{random b/c of } \varepsilon_0, \dots, \varepsilon_n})\right]$$

$$\text{if } y_0 - f(x_0) \underset{E[\hat{f}(x_0)] - f(x_0)}{\overset{\varepsilon_0}{\longrightarrow}} 0$$

$$= E(E_0) \cdot E[E[\hat{f}(x_0)] - \hat{f}(x_0)]$$

$$= 0$$

$$E[(y_0 - f(x_0) + \hat{f}(x_0) - f(x_0) - E[\hat{f}(x_0)] + E[\hat{f}(x_0) - f(x_0)])^2]$$

A B C

$$E(BC) = E[(f(x_0) - E[\hat{f}(x_0)]) \cdot (E[\hat{f}(x_0)] - \hat{f}(x_0))]$$

not random

$$= (f(x_0) - E[\hat{f}(x_0)]). E[E[\hat{f}(x_0)] - \hat{f}(x_0)]$$

$\quad\quad\quad = E[f(x_0)] - E[\hat{f}(x_0)] = 0$

$$= 0$$

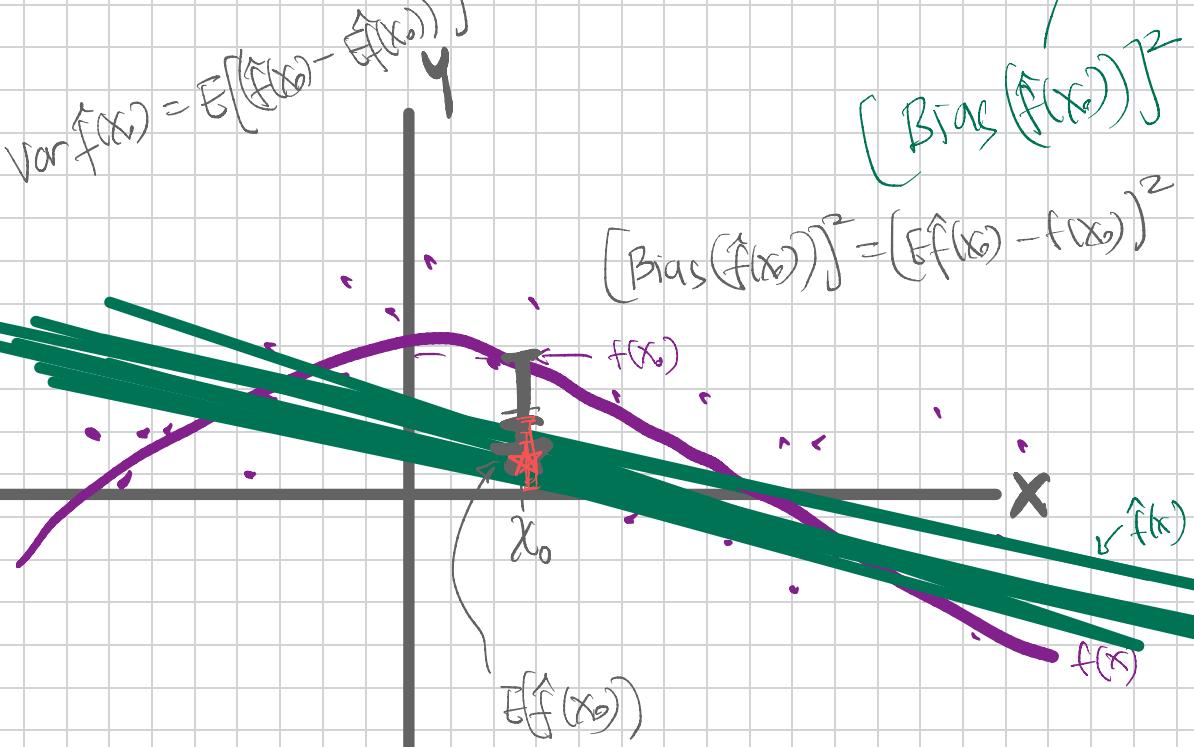
$$EPE(x_0) := E[(y_0 - \hat{f}(x_0))^2]$$

$$= \boxed{\text{Var}(\epsilon)} + \boxed{[\text{Bias}(\hat{f}(x_0))]^2} + \boxed{\text{Var}[\hat{f}(x_0)]}$$

BIAS-VARIANCE TRADE-OFF

irreducible error.

reducible error.



$$\text{Bias}(\hat{f}(x_0))^2 = 0$$

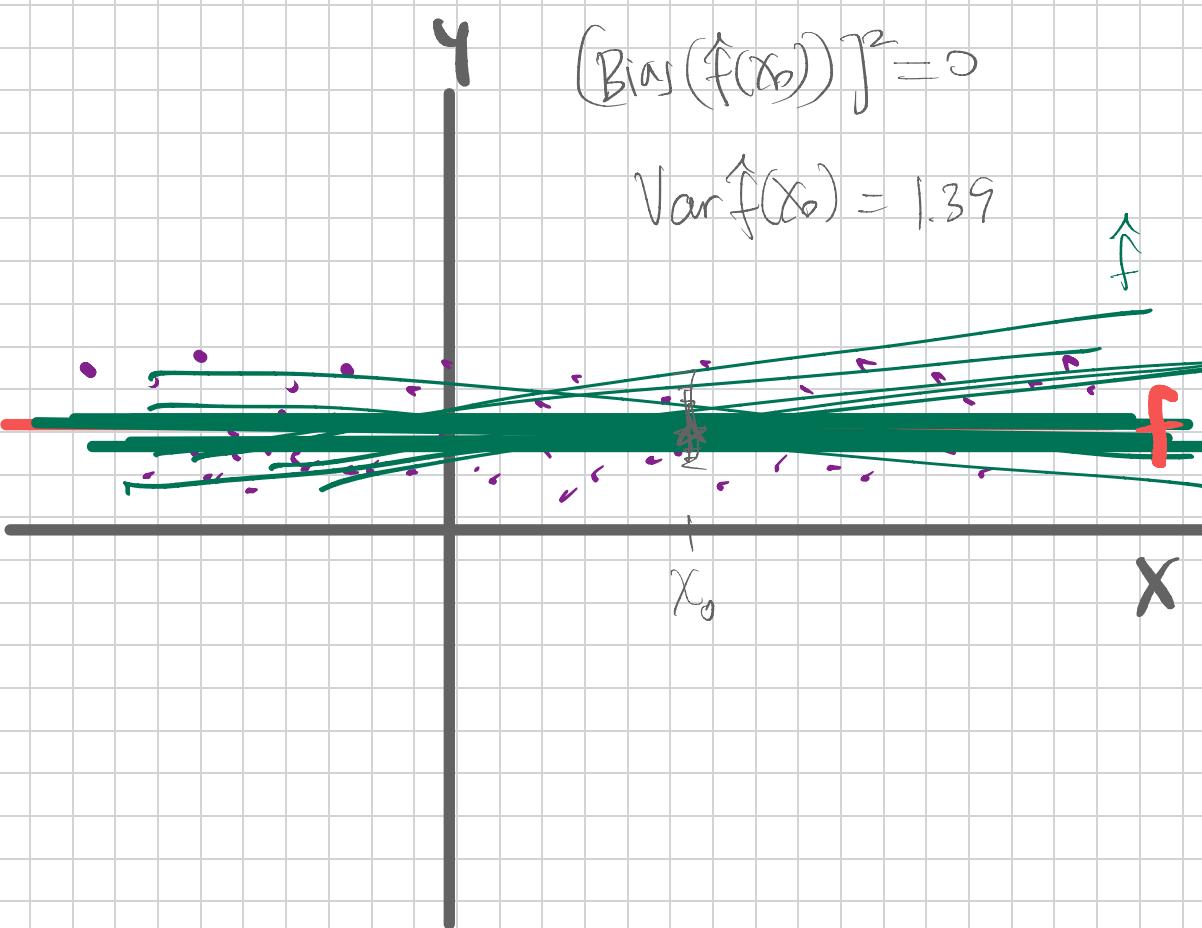
$$\text{Var} \hat{f}(x_0) = 1.39$$

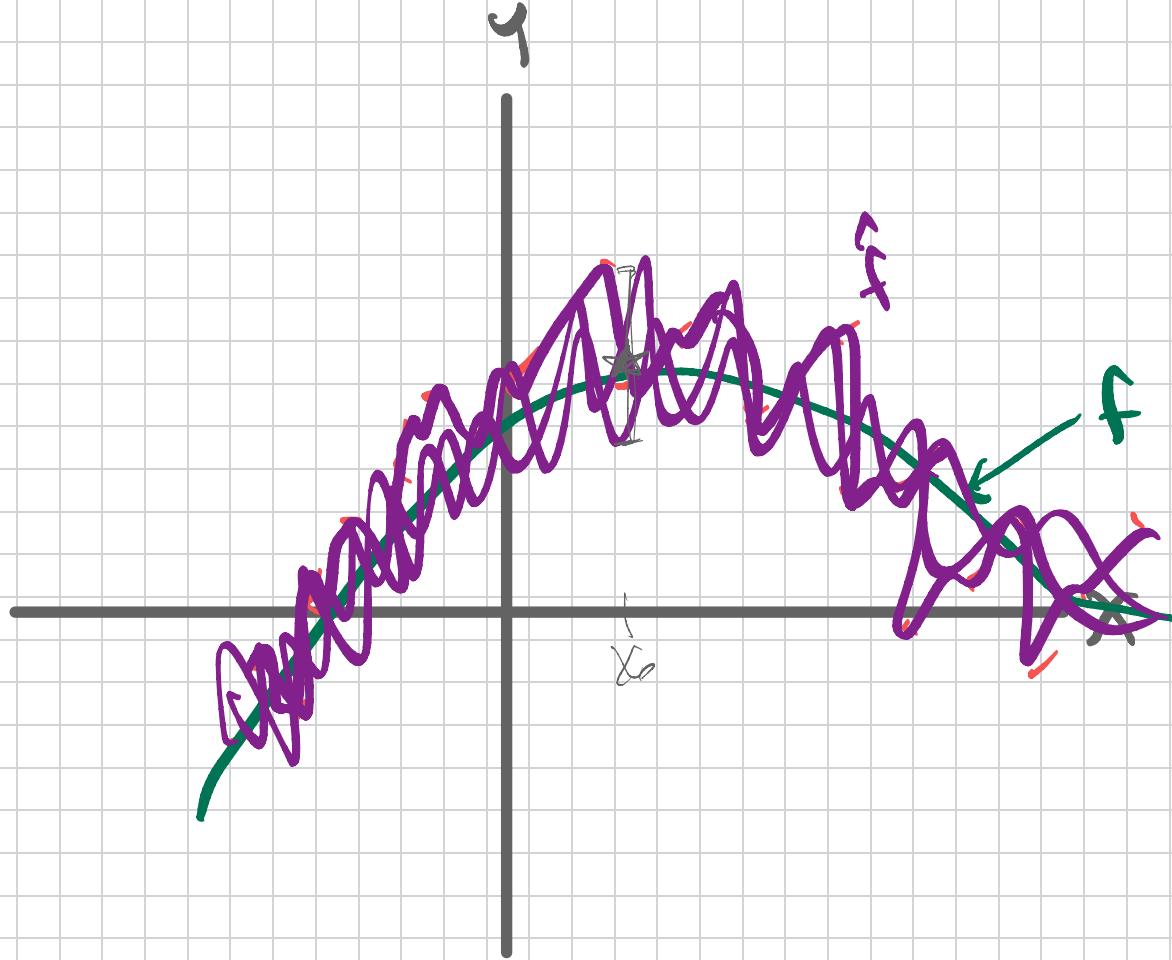
↑
f

f

X

x₀





$$(\text{Bias}(\hat{f}(x_0)))^2 \approx 0$$

$\text{Var}(\hat{f}(x_0))$ is big

Classification:

$y \in \{1, 2, \dots, K\}$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

(training) (mis) classification error:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(y_i \neq \hat{f}(x_i))} \rightarrow \hat{f}(x_i)$$

$$\mathbb{1}_A = \begin{cases} 1 & \text{if } A \text{ holds} \\ 0 & \text{else} \end{cases}$$

$$\mathbb{1}_{(y_i \neq \hat{f}(x_i))} = \begin{cases} 1 & \text{if } y_i \neq \hat{f}(x_i) \\ 0 & \text{otherwise} \end{cases}$$

test (mis)classification error

$$\text{Average } (\mathbb{1}_{(y_0 \neq \hat{f}(x_0))}) \rightarrow \hat{f}(x_0)$$

for test obs y_0 .

Bayes classifier ← unattainable dream

↳ assigns a test obj. x_0 to the class for which $\Pr(Y=k | X=x_0)$ is largest.

don't know this:

need to est. this,

The Bayes classifier minimizes the expected test (mis)classification error

= Bayes error

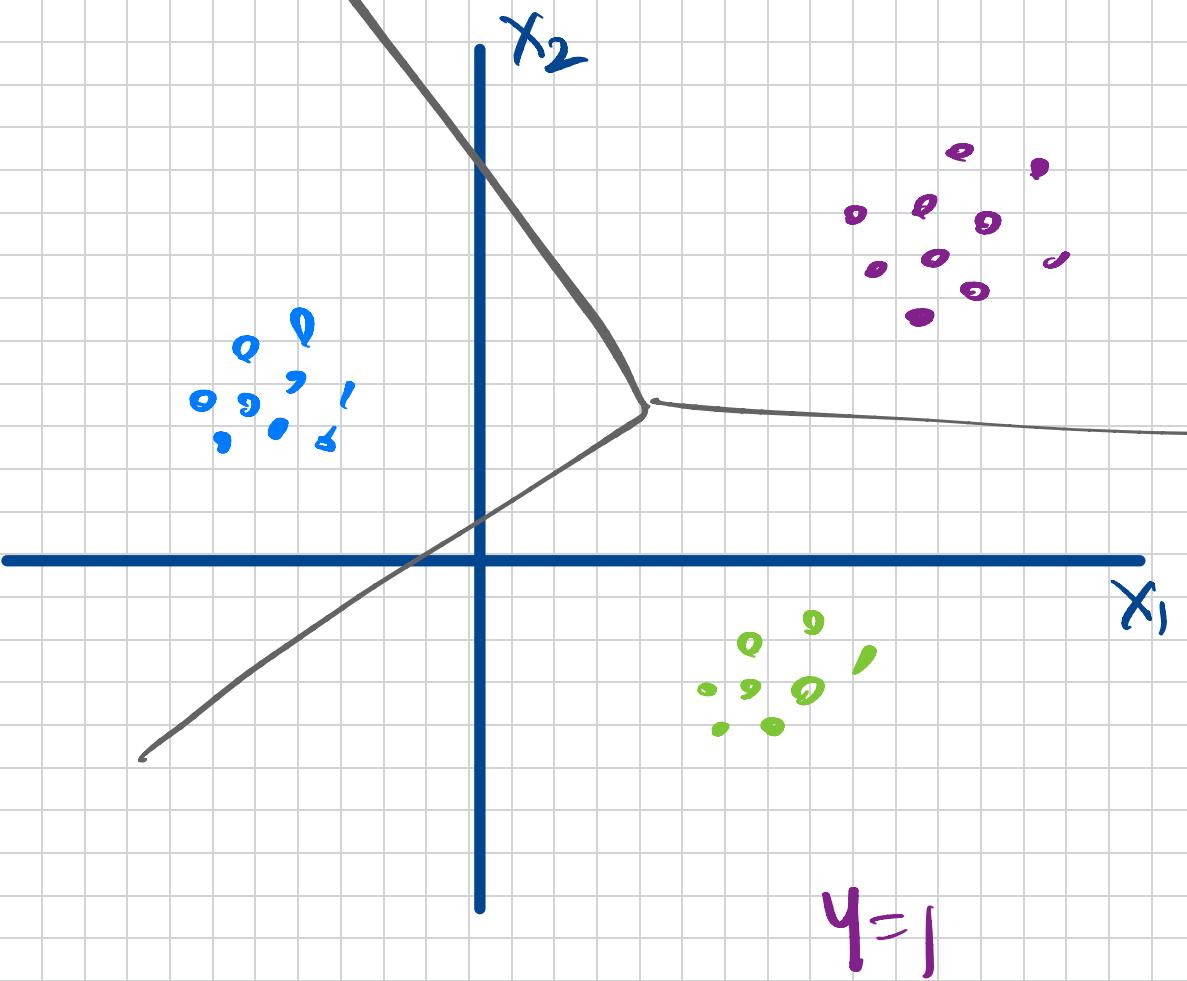
Error of Bayes classifier:

$$\Pr(Y=\text{cat} | X=x_0) = 0.7$$

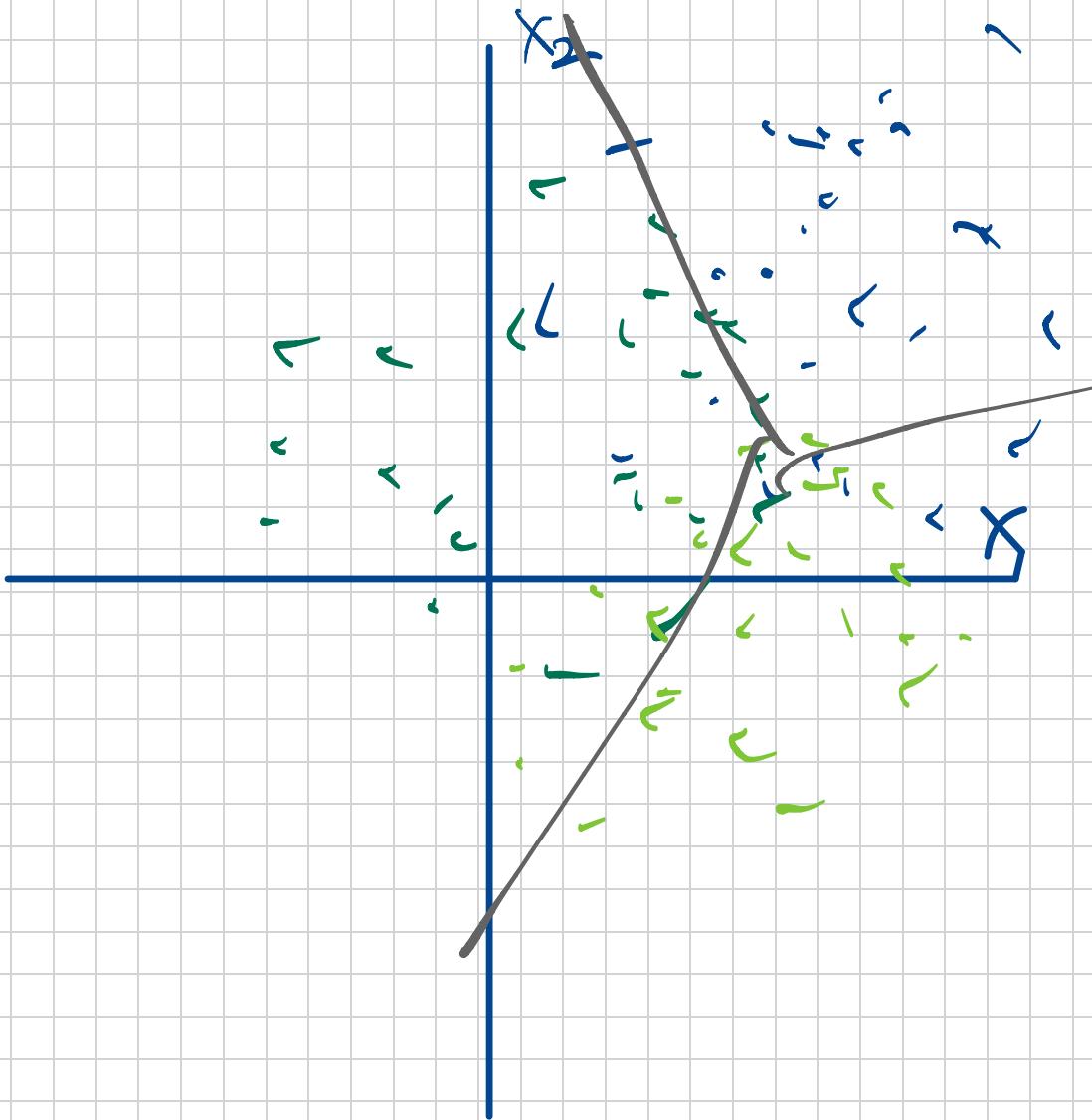
$$\Pr(Y=\text{dog} | X=x_0) = 0.2$$

$$\Pr(Y=\text{monkey} | X=x_0) = 0.1$$

$$1 - E \left[\max_k \Pr(Y=k | X=x_0) \right]$$



Bayes error is TINY



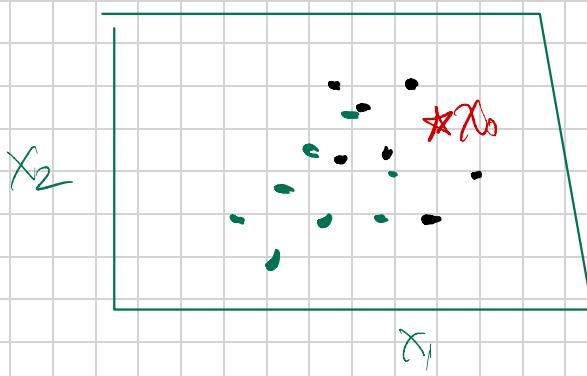
K-nearest neighbors.

Non-parametric approach to est.

$$\Pr(Y=j \mid X=x_0)$$

$$\approx \frac{1}{K} \sum_{i \in N(x_0)} I(y_i = j)$$

where $N(x_0)$ are the indices
of the K "nearest neighbors"
to x_0 in the training set.



$K=1$: zero misclass.
error on training set

$K=5$: nonzero
misclass.
error on test.

