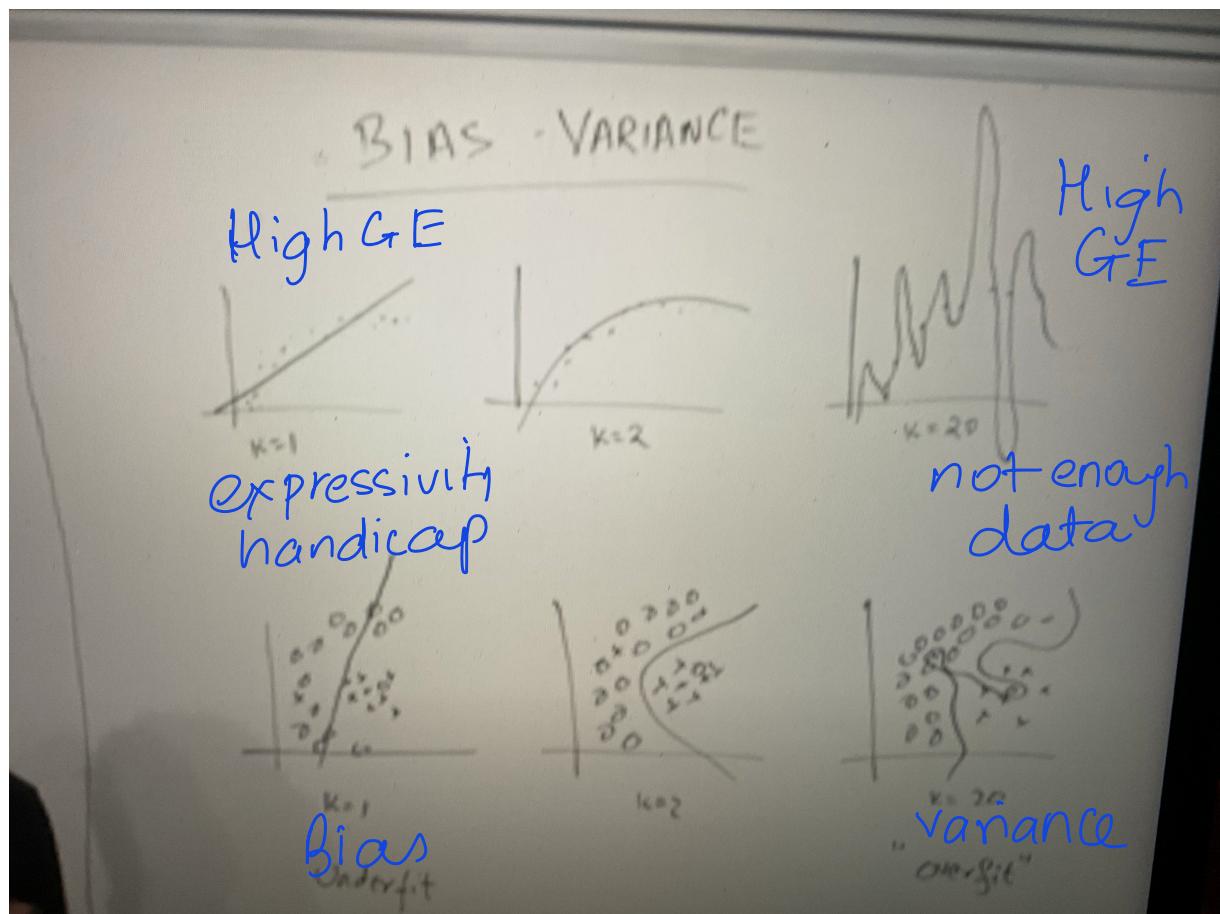


BIAS - VARIANCE



Generalisation Error : loss on unseen data, possibly ∞ , from whose distribution we got our training set.

Two components

- Bias - component of GE due to "expressivity handicap"
- Variance - component of GE due to finite sample of training set.

Squared Error

$$y^{(i)} = f(x^{(i)}) + \varepsilon^{(i)} \rightarrow \text{Data Generating Process}$$

$$E[\varepsilon] = 0 \quad V[\varepsilon] = \sigma^2$$

$$\text{Test Error}[G.E.] = E[(y - \hat{f}_n(x))^2]$$

[expectation over all ε in training set & test set]

↑ function obtained by training
 n eggs from Data.Clean.

$$= \underbrace{\sigma^2}_{\text{Irreducible error}} + \underbrace{\mathbb{E}[\hat{f}_n(x) - f(x)]^2}_{(\text{Bias})^2} + \underbrace{\text{Var}[\hat{f}_n(x)]}_{\text{Variance}}$$

process

Irreducible

error

due to noisiness
of data

[There is noise in
data ϵ , so for
same x too, you can
have diff. y]

$(\text{Bias})^2$

tells us how
systematically
wrong are
we.

Variance

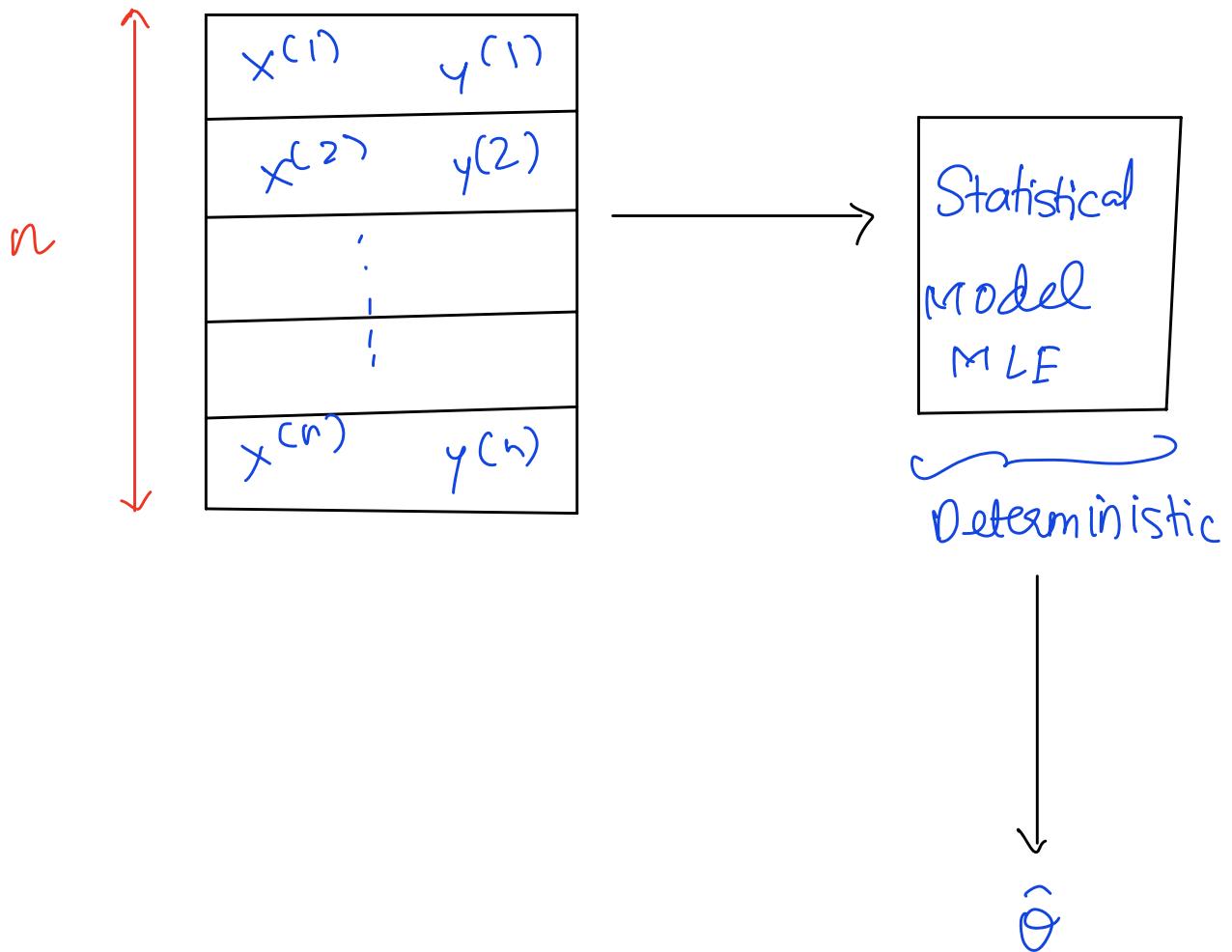
\downarrow

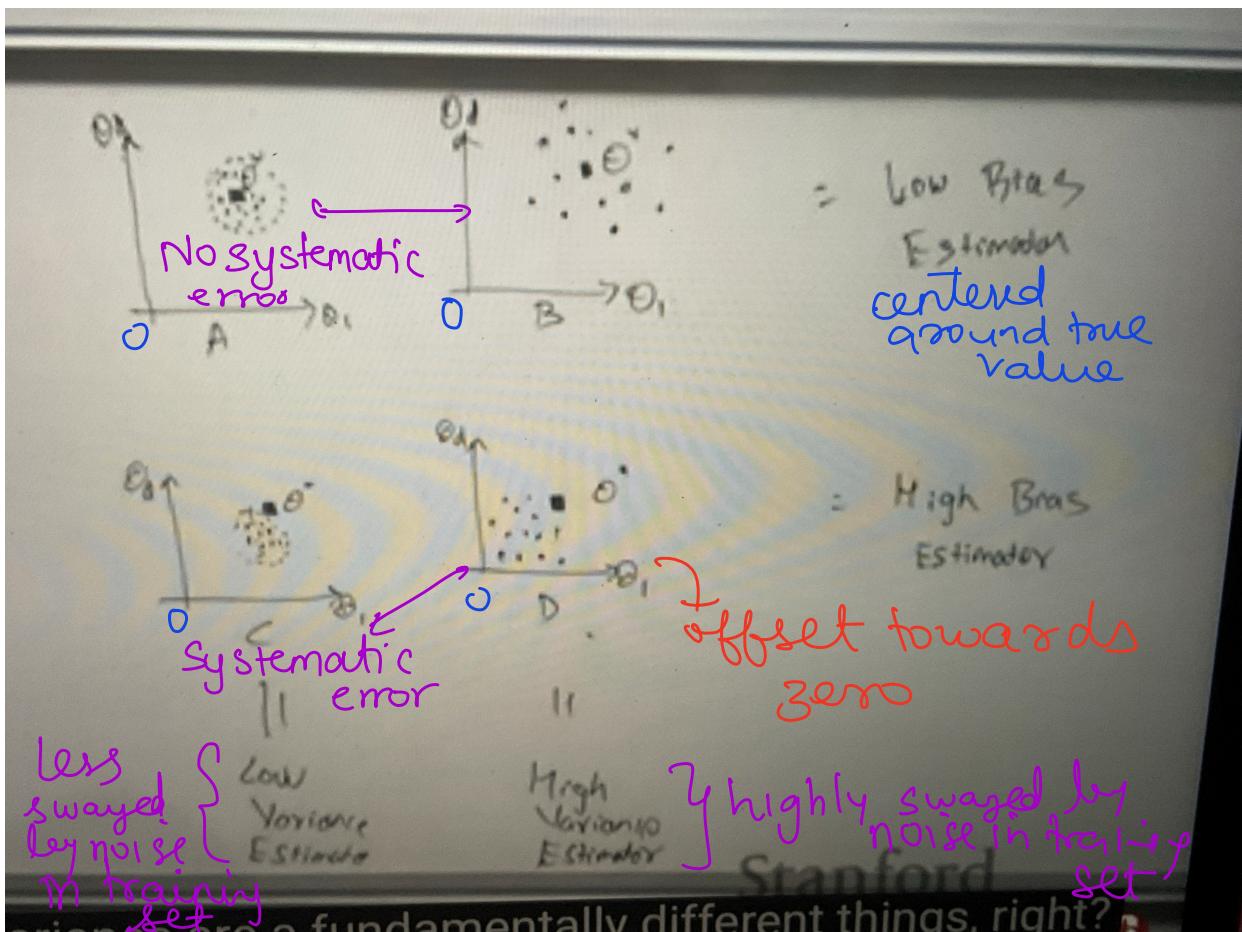
Repeat
by taking
new n

egs, fit
on it,
& then
prediction

You will get
a diff $f_n(x)$

$$(x, y) \sim \text{Dist}(\theta^*)$$

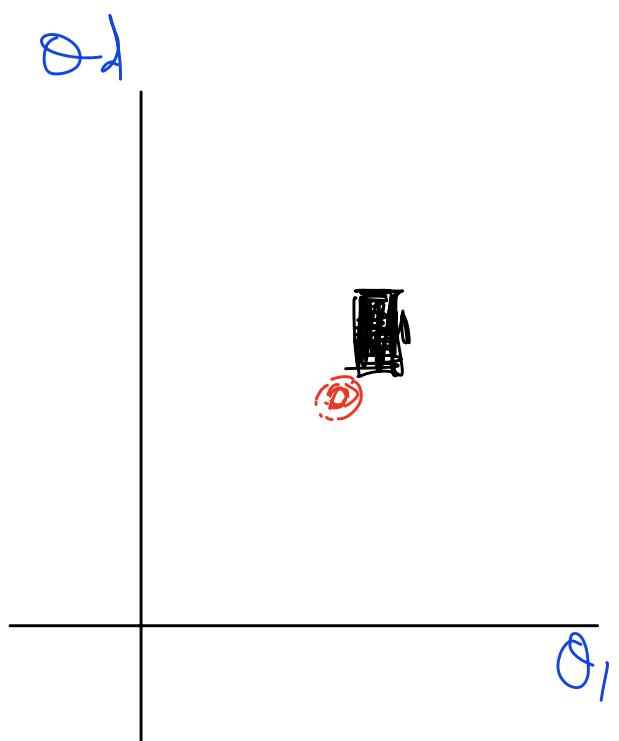
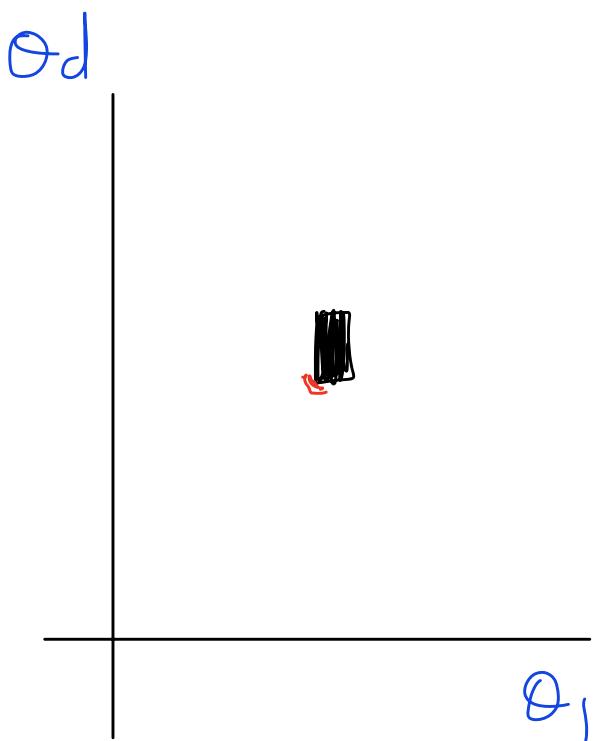
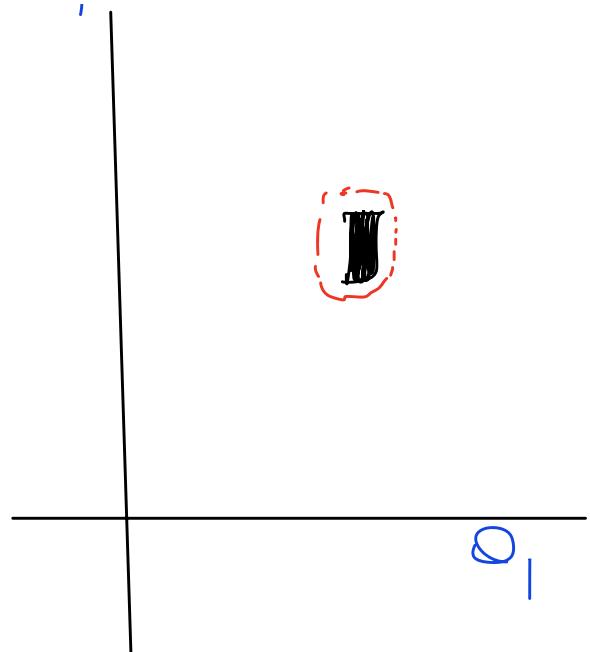
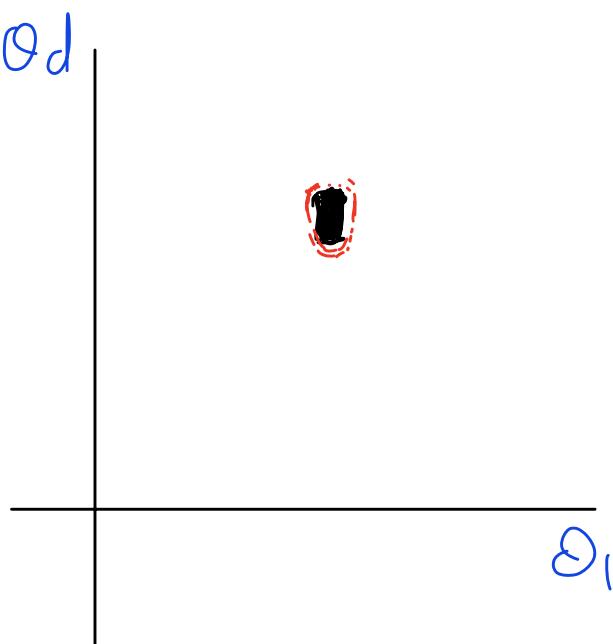




Bias doesn't tell us about variance & vice-versa

Suppose we take a larger data set with N examples (1000x bigger)

$$\hat{\theta}_d$$



In all 4 \rightarrow variance has come down although bias still remains

$$E[\hat{\theta} - \theta^*] = \text{Bias}$$

\uparrow True signal [unknown]
 Predicted

$$\text{Var}[\hat{\theta}] = \text{Variance}$$

$$n \rightarrow \infty$$

If Bias $\rightarrow 0$ \Rightarrow Consistent
 $n \rightarrow \infty$ Estimation

Gaussian

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{x}^{(i)} - \hat{\mu})^2$$

Biased Estimator

MLE

Underestimated bias

$$n \rightarrow \infty, \text{ bias} \rightarrow 0.$$

$\left[\begin{array}{l} \text{As you are noisy} \\ \hat{\mu} \rightarrow \text{estimated mean} \end{array} \right]$

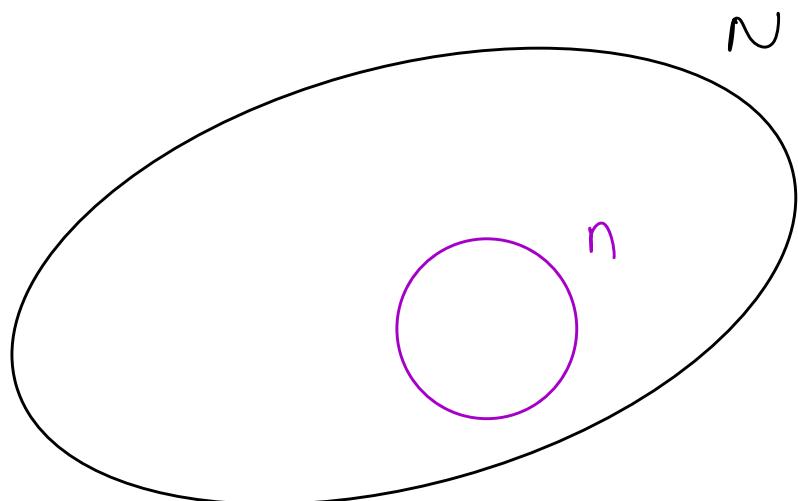
$$\mu = \text{mean}$$

Unbiased

$$= \frac{1}{n-1} \sum_{i=1}^n (x^{(i)} - \bar{x})^2$$

$$\text{OR } \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)^2$$

\downarrow Actual
 μ



Mean

Population
(parameter)

Sample
(statistic)

Mean

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$N \quad n$$

Variance

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{u})^2}{N}$$

$$S_n = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

\hat{T}_b biased
underestimate

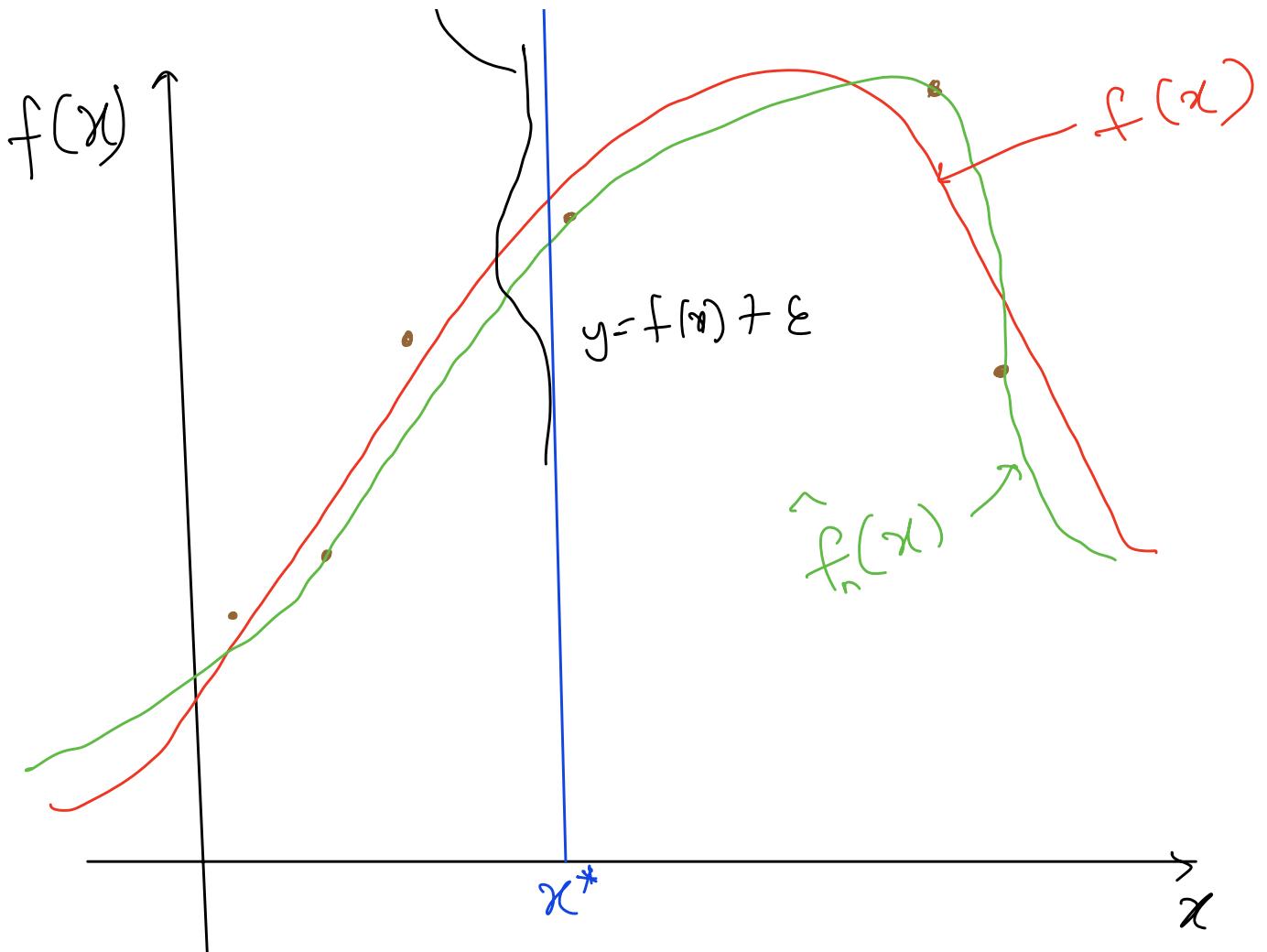
unbiased

$$S_{n-1} = \frac{\sum_{i=1}^{n-1} (x_i - \bar{x})^2}{n-1}$$

$$\lim_{n \rightarrow \infty} \text{var}[\hat{\theta}] \xrightarrow{\text{rate}} 0$$

[statistical efficiency
 $\hookrightarrow \frac{1}{\eta}, \frac{1}{\sqrt{n}}$

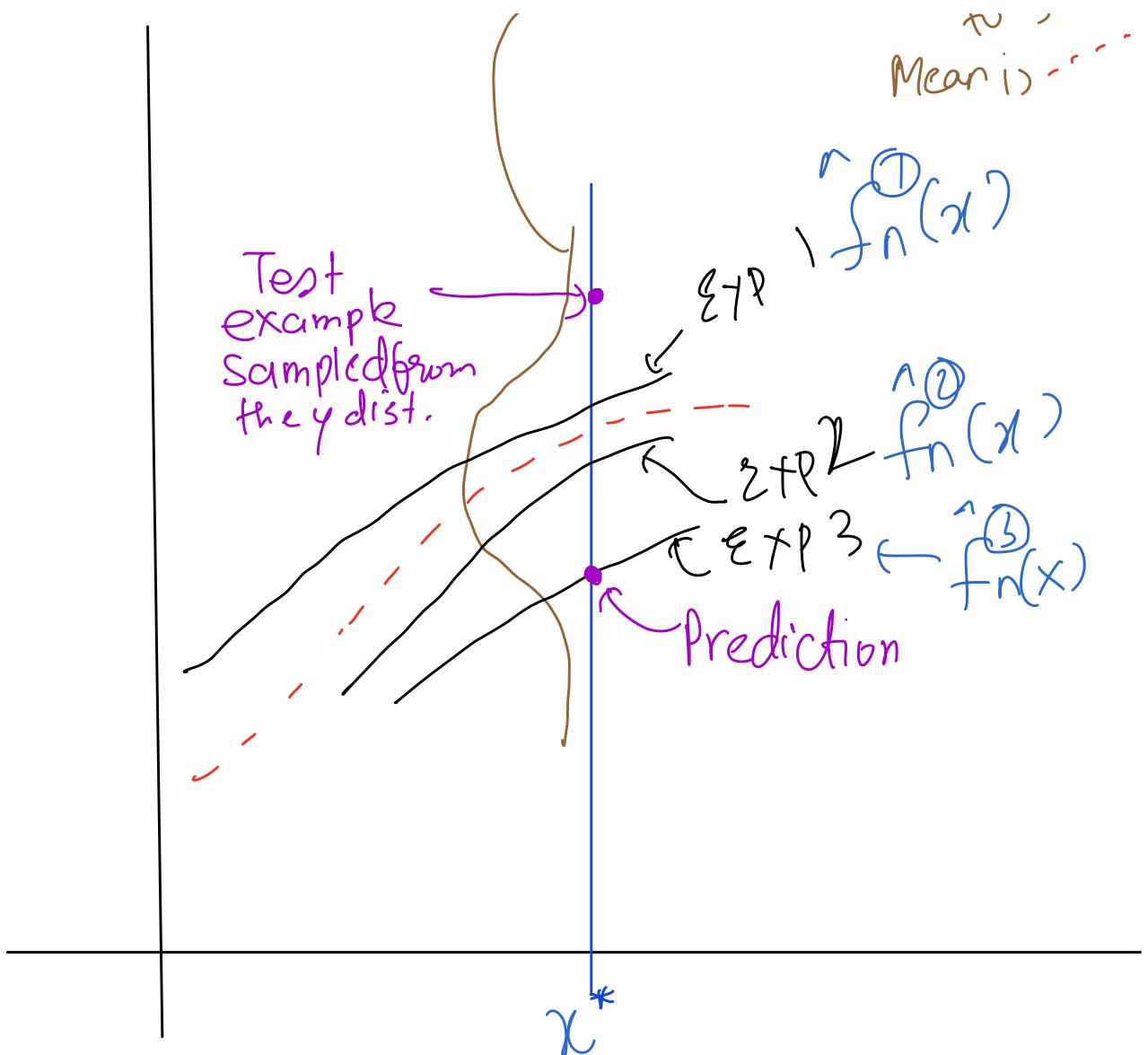
mean \rightarrow red
 Variance $\rightarrow \sigma^2$



- → training set points (noise added)

→ invisible to us

→ Distribution which tells us which value of $x^{(i)}, y^{(i)}$ sample



y is a sample for \mathcal{D}

$$\mathbb{E}[(\text{distance})^2] = \text{test error}$$

$$= \sigma^2 + \mathbb{E}[f(x) - \hat{f}_n(x)]^2 + \text{Var}[\hat{f}_n(x)]$$

how
far apart
are our
predicted
 y in
diff exp -

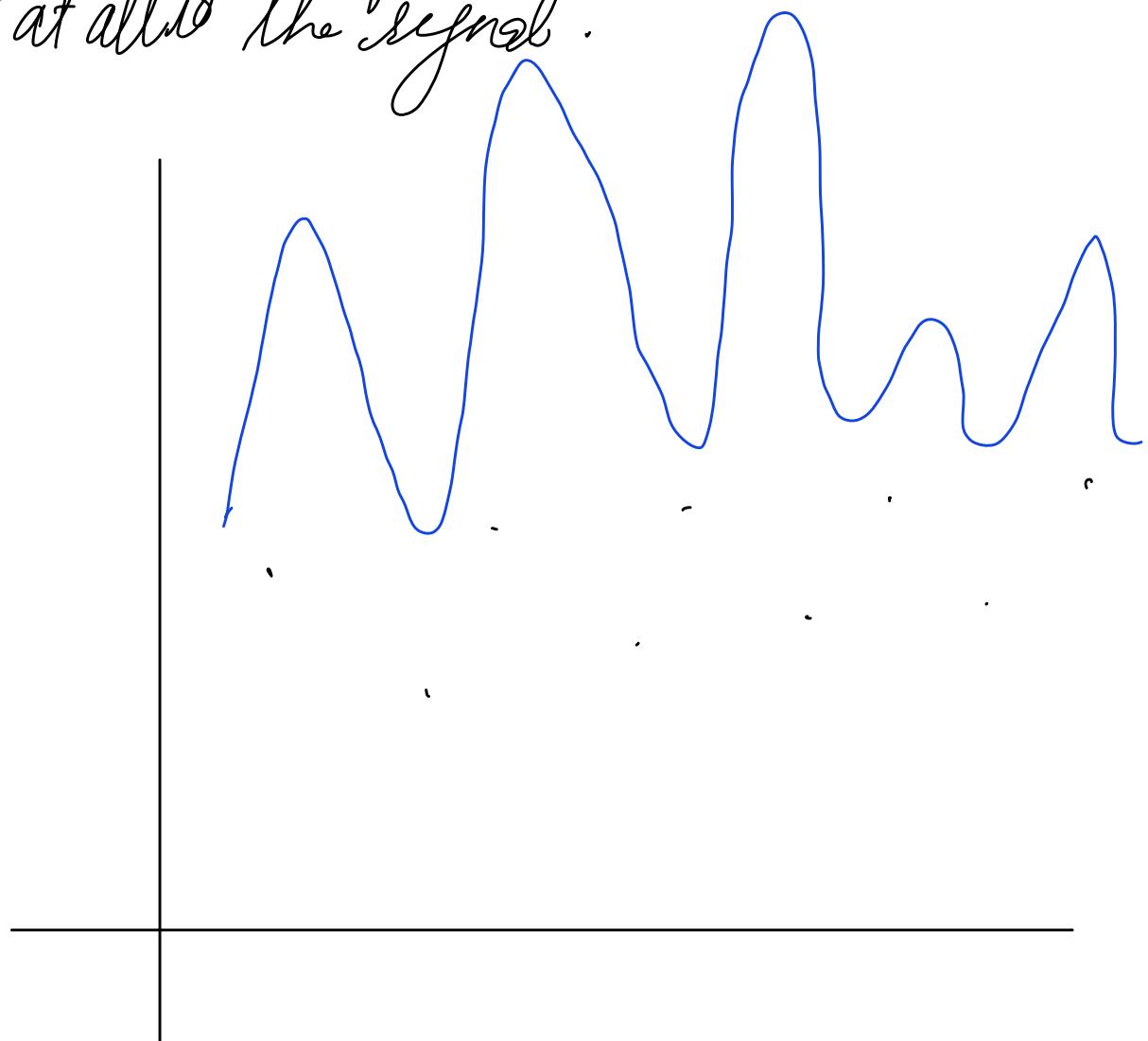
Irreducible
error
(due to variance
of ϵ)

(bias error)²
If for our
diff experiment,
the predicted
 y for x^* is
below only/
above only, ~~or~~
or is well
centered

Underfitting \approx High bias

Overfitting \approx High Variance

You can have both at same time like if your model is overfitting/underfitting at same time if model fits to the noise & doesn't fit at all to the signal.



Goal: Do well on Generalization Error

- Cross Validation

Do well on training set/min. loss

Split into	Train	70	60	80%
	Valid/Dev	20	20	10
	Test	10	20	10
Dowell on gen. error		↪ Estimate generalization error		

We check loss on valid set. Then do hyperparameter tuning

By repeatedly doing cross validation,
the model is looking at valid set,
it might note it.

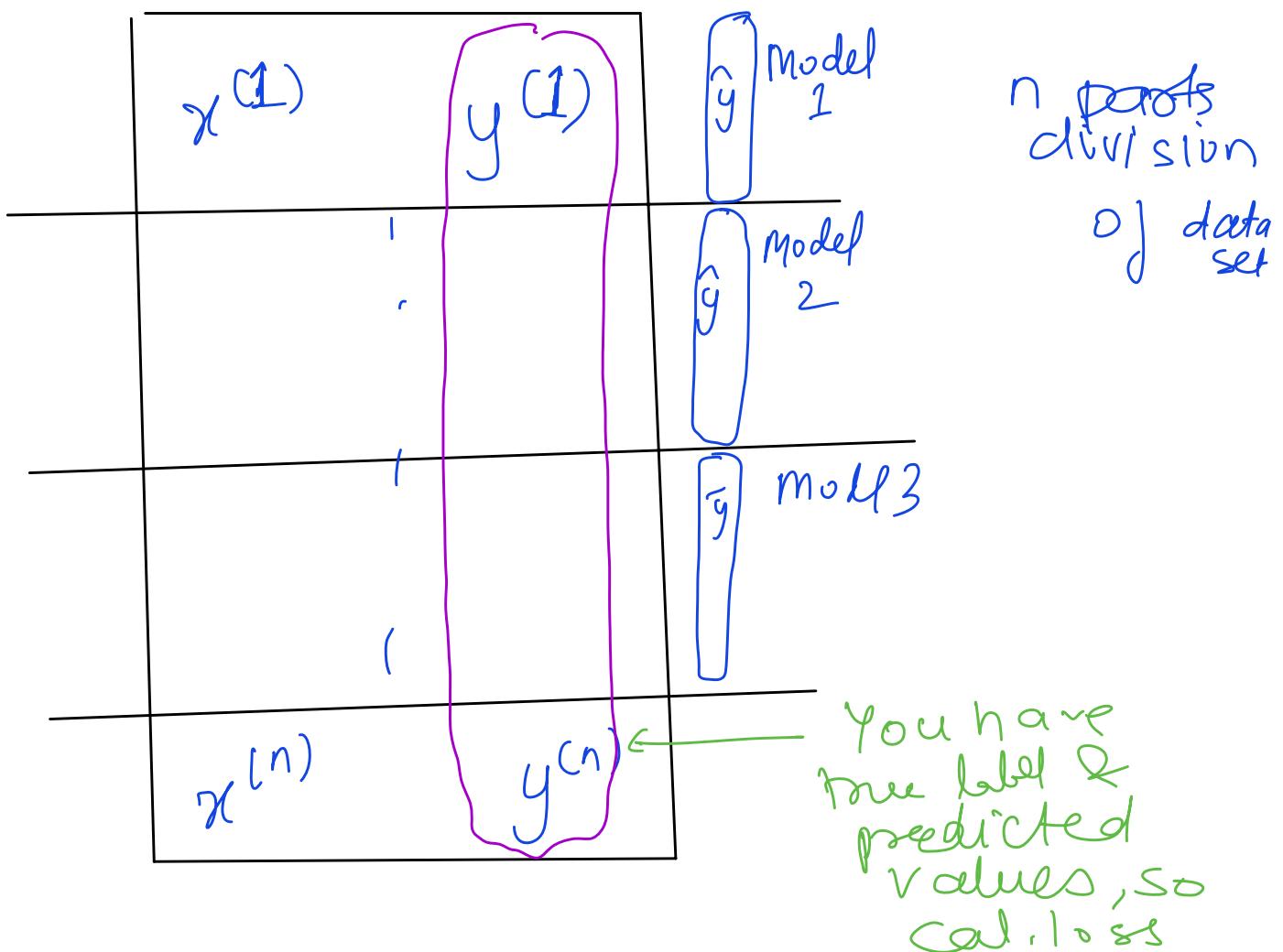
To counter that we have test set

The above method is HoldOut Cross Validation

K-Fold Cross Validation

Small Datasets

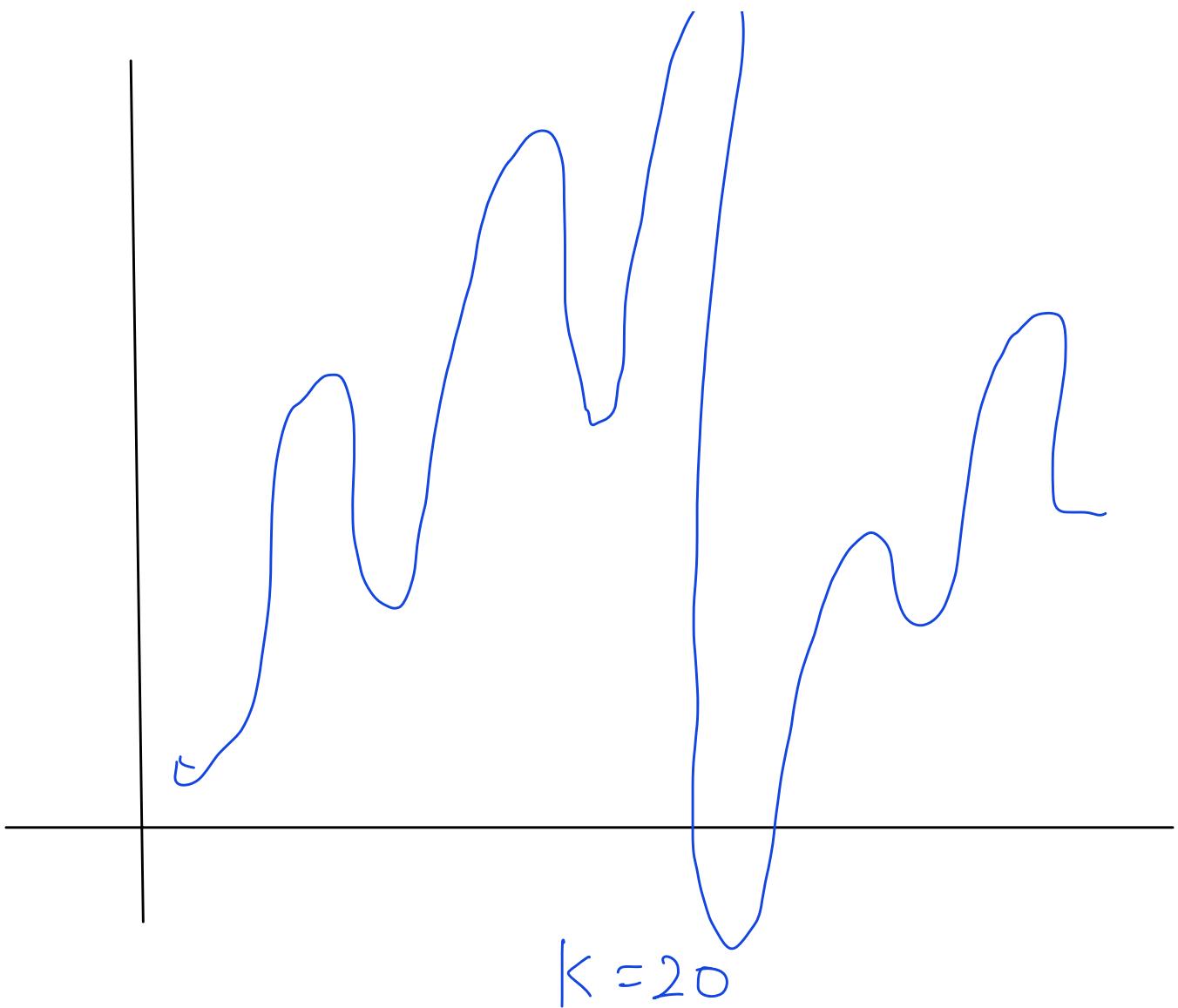
leave one out Cross Validation
 $K = n$



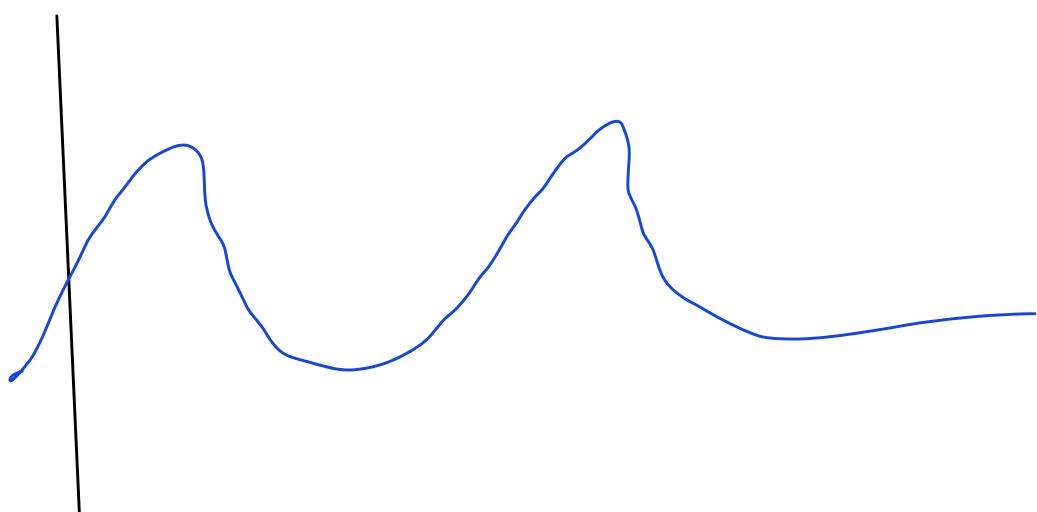
Model 1, fold 1 is val set
2-n → train set
Prediction on fold 1

Model 2, fold 2 is val set
1, 3-n → train set
pred. on fold 2
⋮
⋮ So on

Then, you can either ensemble all k-model & average their pred. at test time
Use this technique for hyperpar. tuning & then refit on entire set



$| \theta | >> 0 \rightarrow \text{bad}$



You want less wiggly but
still want it to be expressive

|| Regularisation

Encouraging small $\|\theta\|$
is Regularization

$$J(\theta) = \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)}))^2 + \lambda \|\theta\|_2^2$$

$\lambda \rightarrow$ relative
weighting factor

$$h_\theta(x) \rightarrow \theta^T x$$