

Error bar:

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{\delta}{N} \log \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right)}$$

error bar

Example: $N = 100$, $\delta = 0.1$ (90% confidence)

$d_{\text{vc}} = 1$ (Too simple model)

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + 0.848$$

$$N = 1000, \quad E_{\text{out}}(g) \leq E_{\text{in}}(g) + 0.301$$

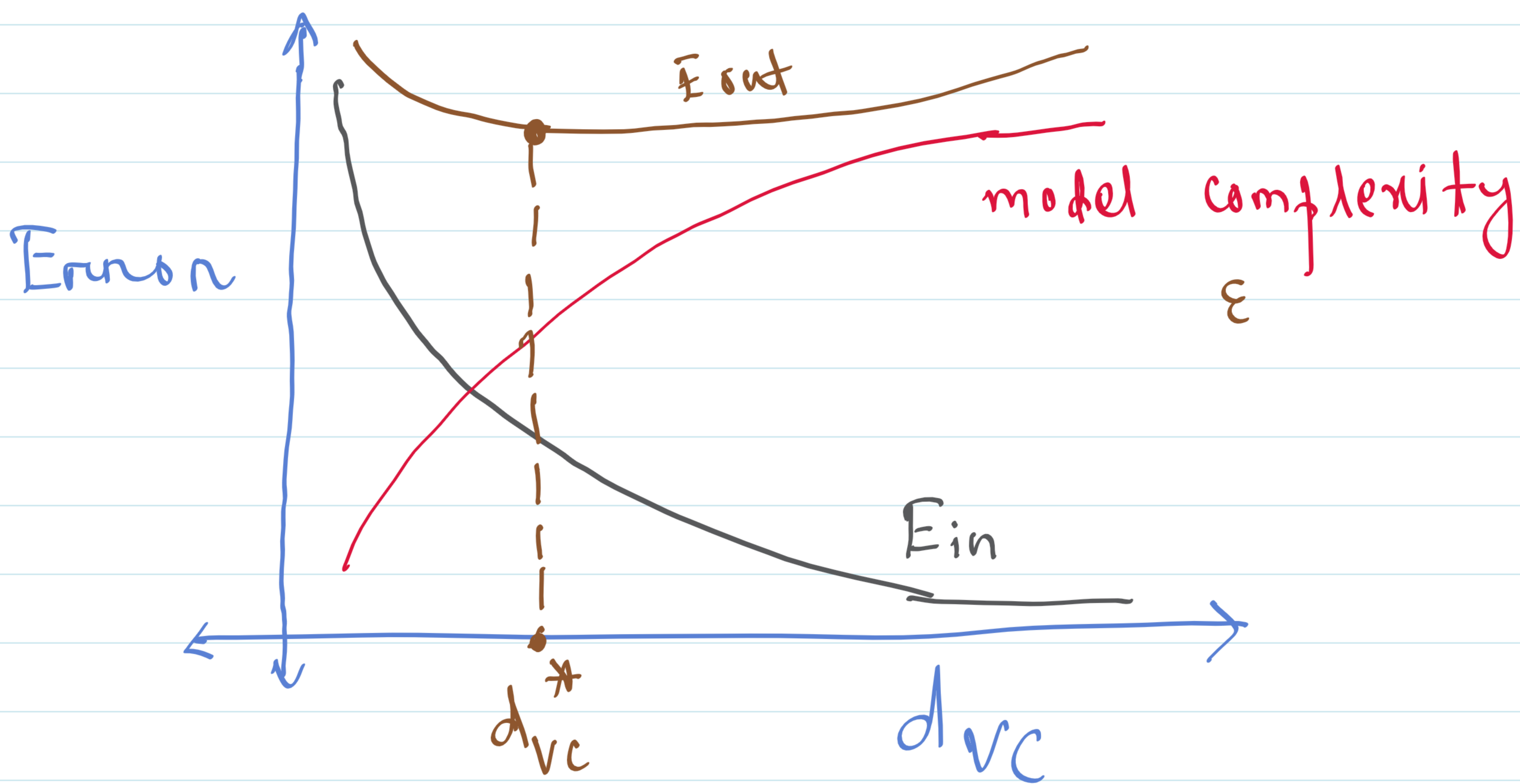
Model complexity:

$$\sqrt{\frac{\delta}{N} \log \left(\frac{4((2N)^{d_{\text{vc}}} + 1)}{\delta} \right)} := \mathcal{E}(N, \mathcal{H}, \delta)$$

penalty of the model complexity

→ If d_{vc} is large, then $\mathcal{E}(N, \mathcal{H}, \delta)$ is big.

Trade-off curve:



optimum value
of VC dimension

Generalisation bound for testing :

Testing set : $\mathcal{D}_{\text{test}} = \{x_1, \dots, x_L\}$

The final hypothesis is already determined.

So we do not need a union bound.

$$\text{Hoeffding : } \mathbb{P} \left[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq \underline{\underline{2e^{-2\epsilon^2 L}}}$$

Generalisation bound :

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{1}{2L} \log \frac{2}{\delta}}$$

$L \rightarrow \infty \rightarrow 0$

→ If you have big L , then $E_{\text{in}}(g)$ will be a good estimate for $E_{\text{out}}(g)$.

→ Independent of model complexity, dependent δ & L .

→ VC analysis → based on binary target fun's. (can be ~~be~~ extended to real values)

→ Alternative approach → real valued fun's based bias - Variance analysis.

→ VC analysis - E_{in} or E_{out} , $h(x) = f(x)$ or $h(x) \neq f(x)$

From VC analysis to Bias-Variance analysis

VC analysis is a decomposition. As we decompose E_{out} into E_{in} and ϵ i.e.,

$$E_{out} \leq E_{in} + \underbrace{\frac{8}{N} \log \frac{4((2N)^{dvc} + 1)}{5}}_{\text{penalty of the model complexity}}$$

\nwarrow
training error

→ Bias - Variance is another decomposition.

- 1) How well can H approximate f ?
- 2) How well can we zoom in a good h in H ?

→ VC analysis, $E_{out}(g) = P[g(x) \neq f(x)]$

Let $B = \{g(x) \neq f(x)\}$ (bad event) $B \in \{0, 1\}$

$$E_{out}(g) = P[B=1] = 1 \cdot P[B=1] + 0 \cdot P[B=0] = E[B]$$

$$E_{\text{out}}(g) = E_x \left[\mathbb{1}_{\{g(x) \neq f(x)\}} \right], \quad x \in \mathcal{P}(x)$$

(0-1 loss)

change the error measure: look at squared loss

$$E_{\text{out}}(g) = E_x \left[(g(x) - f(x))^2 \right] \quad (\text{Bias-Variance})$$

Dependancy on training set \mathcal{D} :

VC analysis $E_{\text{out}}(g^{(\mathcal{D})}) = E_x \left[\mathbb{1}_{\{g^{(\mathcal{D})}(x) \neq f(x)\}} \right]$

Bias-Variance analysis:

$$E_{\text{out}}(g^{(\mathcal{D})}) = E_x \left[(g^{(\mathcal{D})}(x) - f(x))^2 \right].$$

Q Why Not \mathcal{D} in VC analysis? As Hoeffding bound is uniform for all \mathcal{D} . Not true for \mathcal{D} (Bias-Variance)

Average over all \mathcal{D}

$$E_{\mathcal{D}} \left[E_{\text{out}}(g^{(\mathcal{D})}) \right] = E_{\mathcal{D}} \left[E_x \left[(g^{(\mathcal{D})}(x) - f(x))^2 \right] \right]$$

→ VC trade-off is a "worst-case" analysis, i.e., uniform on every \mathcal{D} .

→ Bias-Variance trade-off is an "average" analysis, i.e., average over different \mathcal{D} 's

Decomposing $E_{out}(g^{(D)})$

$$\begin{aligned} E_D [E_{out}(g^{(D)})] &= E_D [E_x [(g^{(D)}(x) - f(x))^2]] \\ &= E_x [E_D [(g^{(D)}(x) - f(x))^2]] \\ &= E_x [E_D [(g^{(D)}(x))^2 - 2g^{(D)}(x)f(x) + (f(x))^2]] \\ &= E_x [E_D [(g^{(D)}(x))^2] - 2E_D [g^{(D)}(x)]f(x) + (f(x))^2] \end{aligned}$$

$\bar{g}(x)$ - averaged final hypothesis is

The asymptotic limit of the estimate

$$\bar{g}(x) \approx \frac{1}{K} \sum_{k=1}^K g^{(D_k)}(x)$$

→ $g^{(D_k)}$ are inside the hypothesis set, but \bar{g} is not necessarily inside.

Bias-Variance decomposition:

$$\begin{aligned} E_D [E_{out}(g^{(D)})] &= E_x [E_D [g^{(D)}(x)^2] - 2\bar{g}(x)f(x) + f(x)^2] \\ &= E_x [E_D [g^{(D)}(x)^2] - \bar{g}(x)^2 + \bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2] \end{aligned}$$

$$\mathbb{E}_x \left[\underbrace{\mathbb{E}_D [g^{(D)}(x)^2 - \bar{g}(x)^2]}_{\text{Van}(x)} + \underbrace{\bar{g}(x)^2 - 2\bar{g}(x)f(x) + f(x)^2}_{(\bar{g}(x) - f(x))^2} \right]$$

$\mathbb{E}_D [g^{(D)}(x) - \bar{g}(x)]^2$

$\text{Van}(x)$

Random deterministic

bias(x)
 how far
 is \bar{g} (averaged hypothesis)
 from target fun^t

Overall,

$$\text{bias} = \mathbb{E}_x [\text{bias}(x)] = \mathbb{E}_x [(\bar{g}(x) - f(x))^2]$$

$$\text{Van} = \mathbb{E}_x [\text{Van}(x)] = \mathbb{E}_x [\mathbb{E}_D [(g^{(D)}(x) - f(x))^2]]$$

Finally,

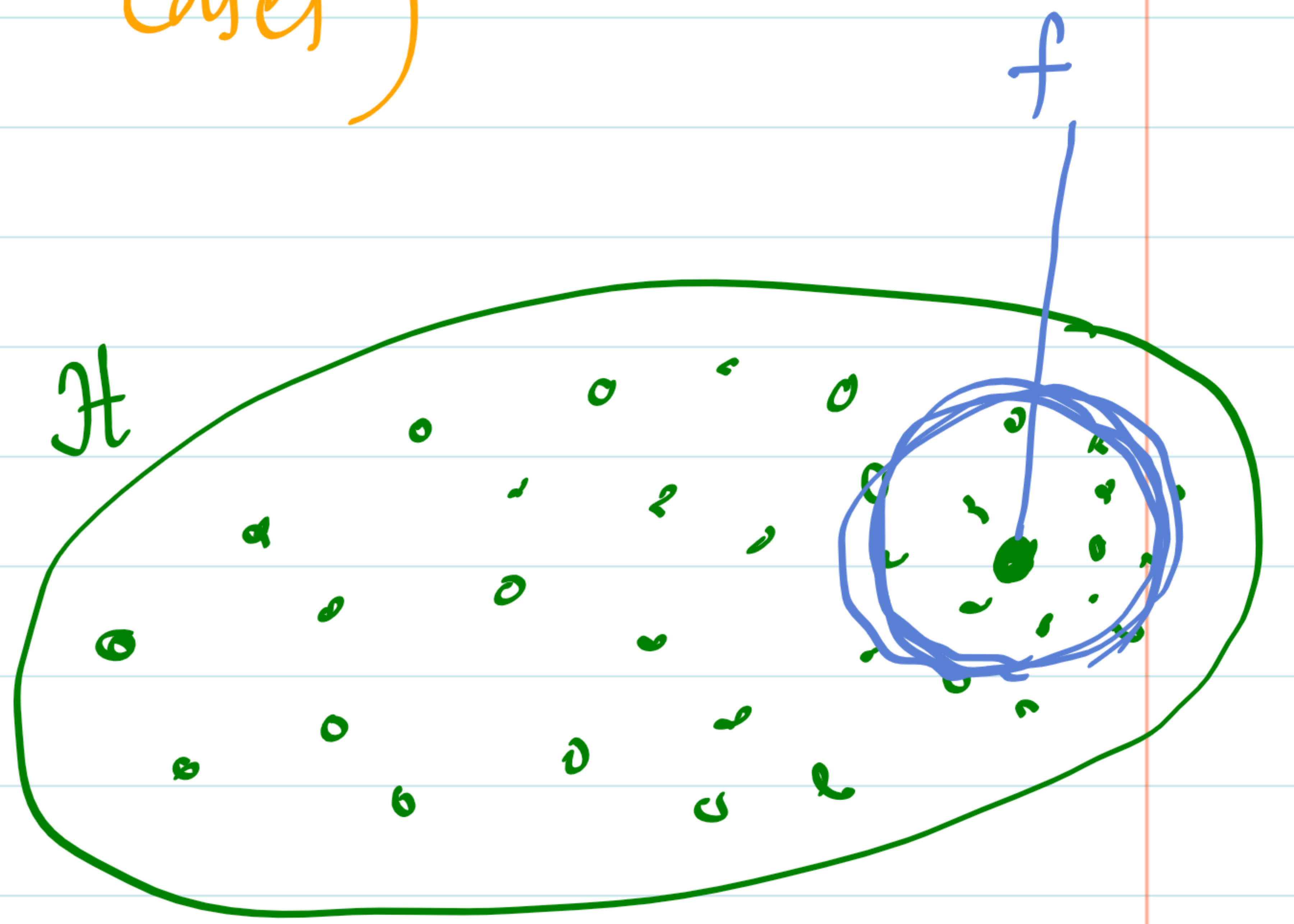
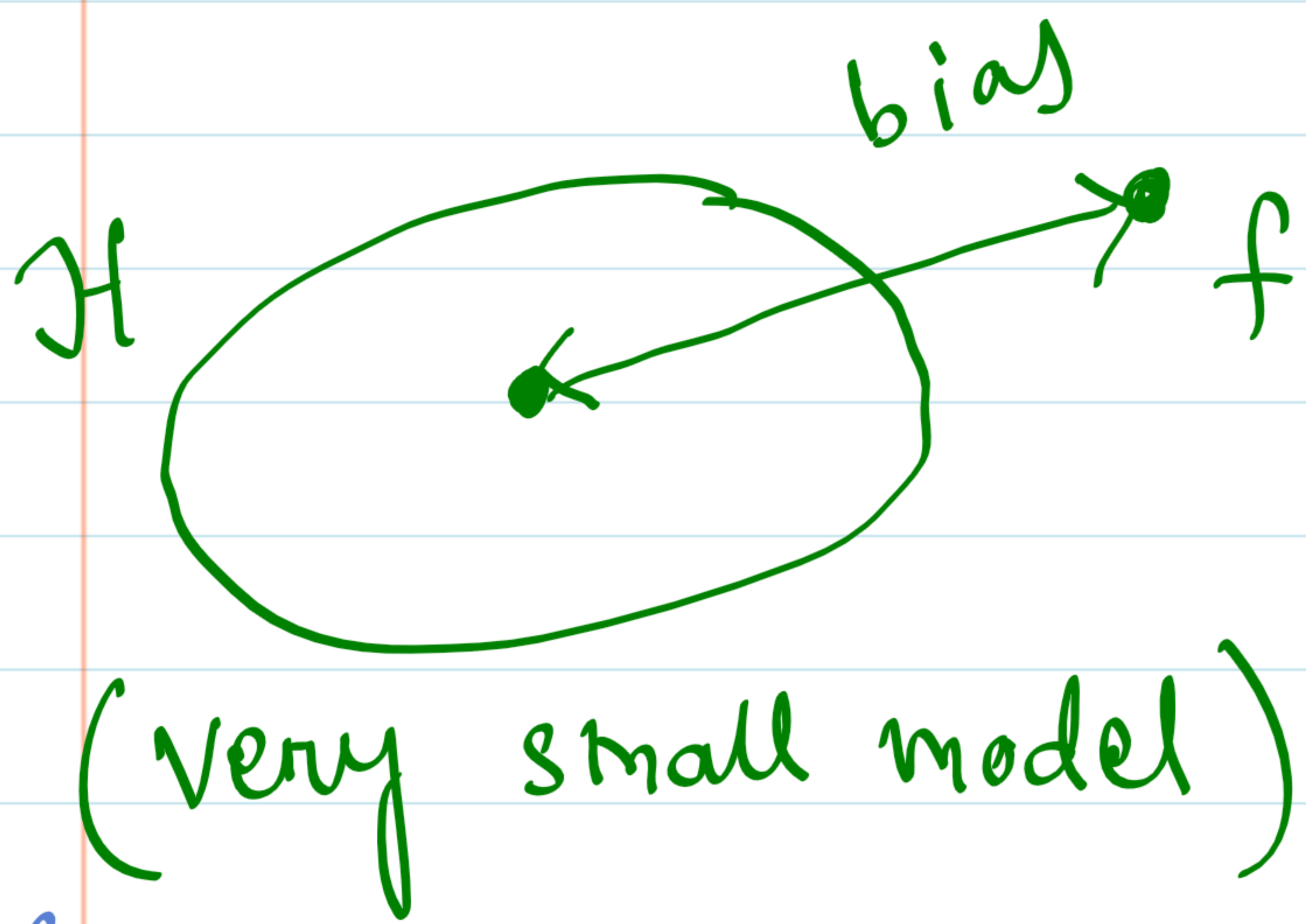
$$\mathbb{E}_D [\text{E}_{\text{out}}(g^{(D)})] = \mathbb{E}_x [\text{bias}(x) + \text{Van}(x)]$$

$$= \text{bias} + \text{Van}$$

bias(x): How close is the average fun^c \bar{g} compared to the target

Van(x): How much uncertainty you have around \bar{g} .

Picture (Two extreme cases)

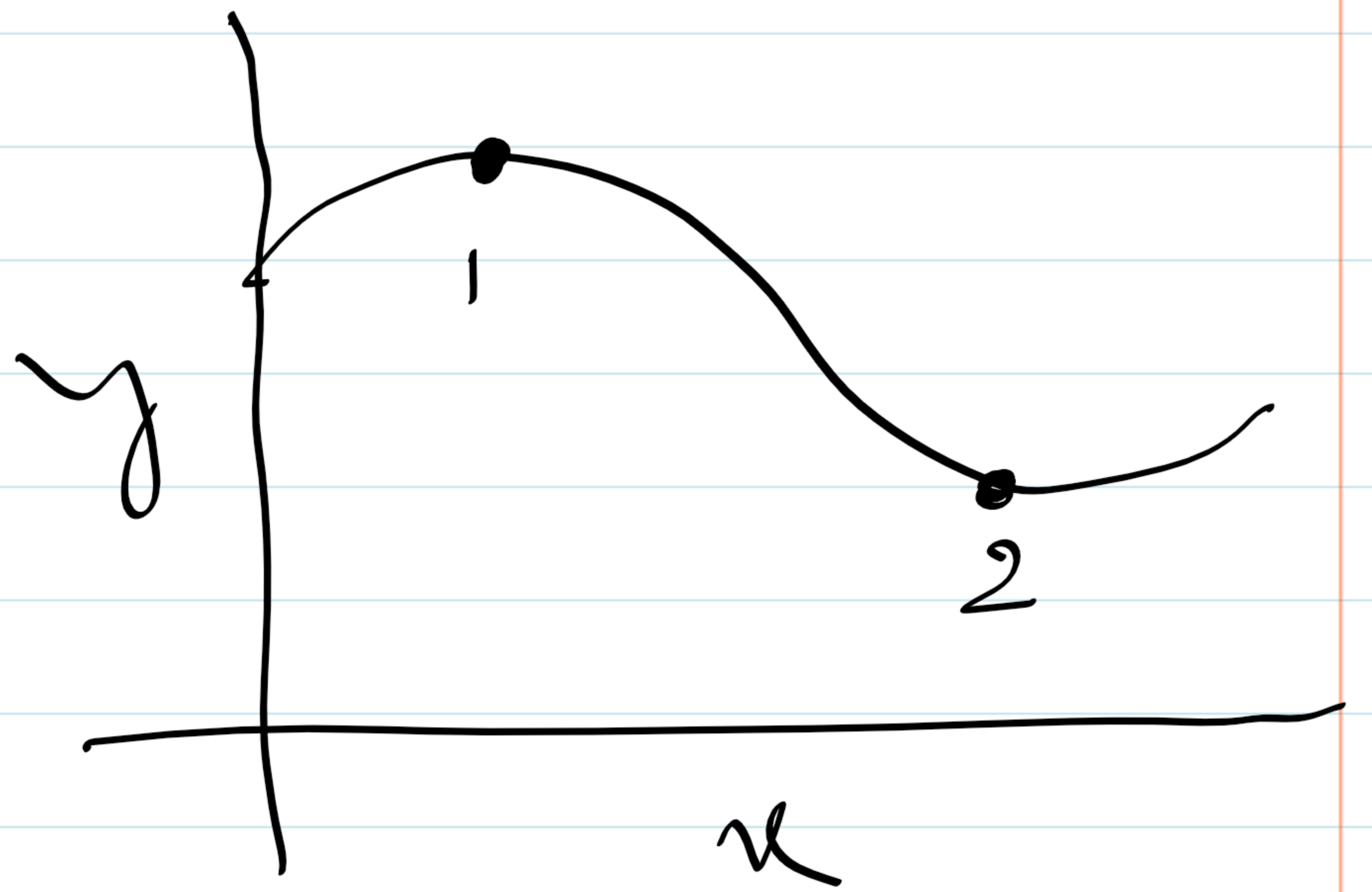


If you have a simple H then large bias, but small variance

If you have a complex H , then small bias but large variance

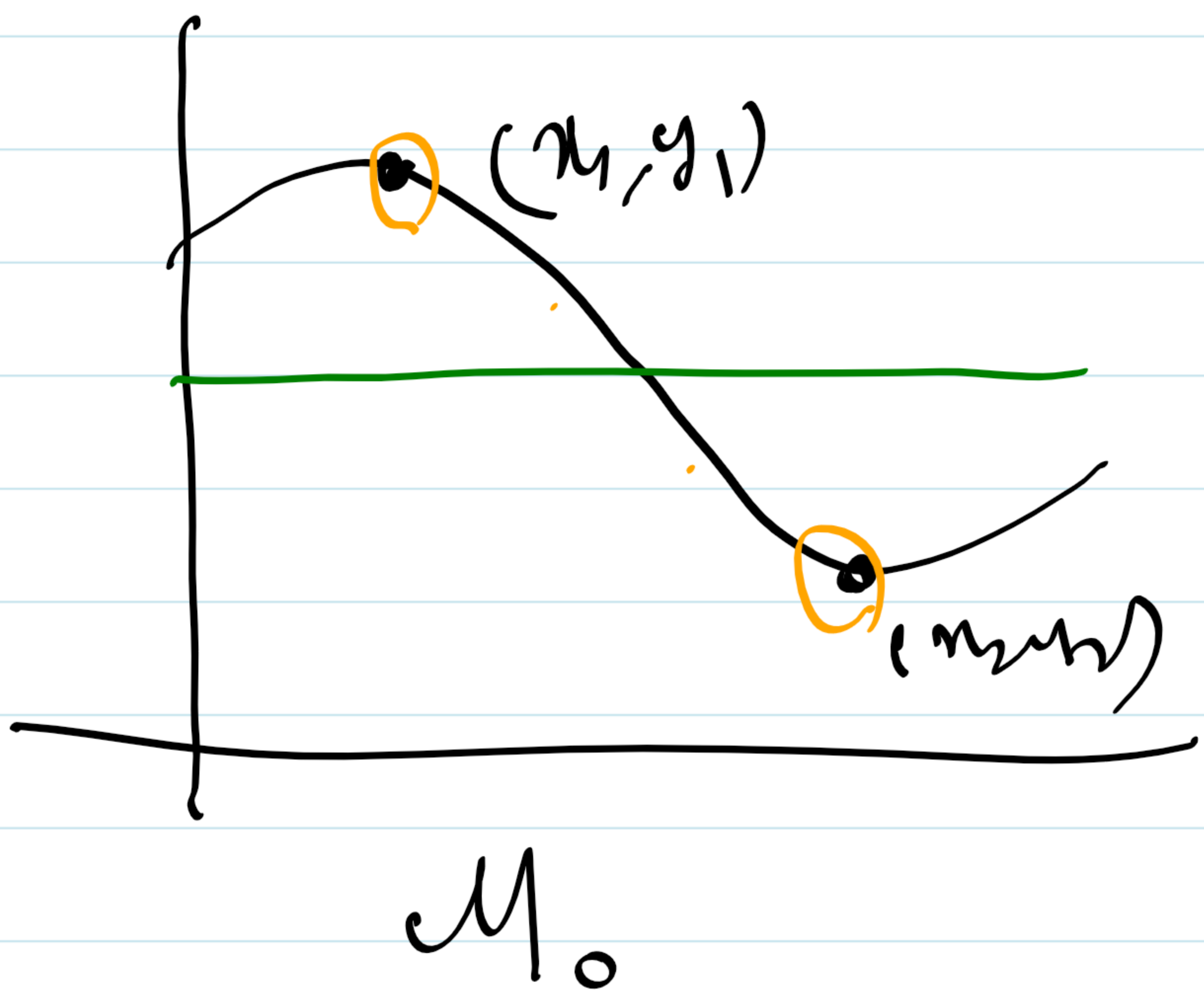
Example: Consider a $\sin(\cdot)$ func.
 $f(x) = \sin(\pi x)$.

$N=2$ training samples
 (x_1, y_1) & (x_2, y_2)
by sampling uniformly
on $[-1, 1]$



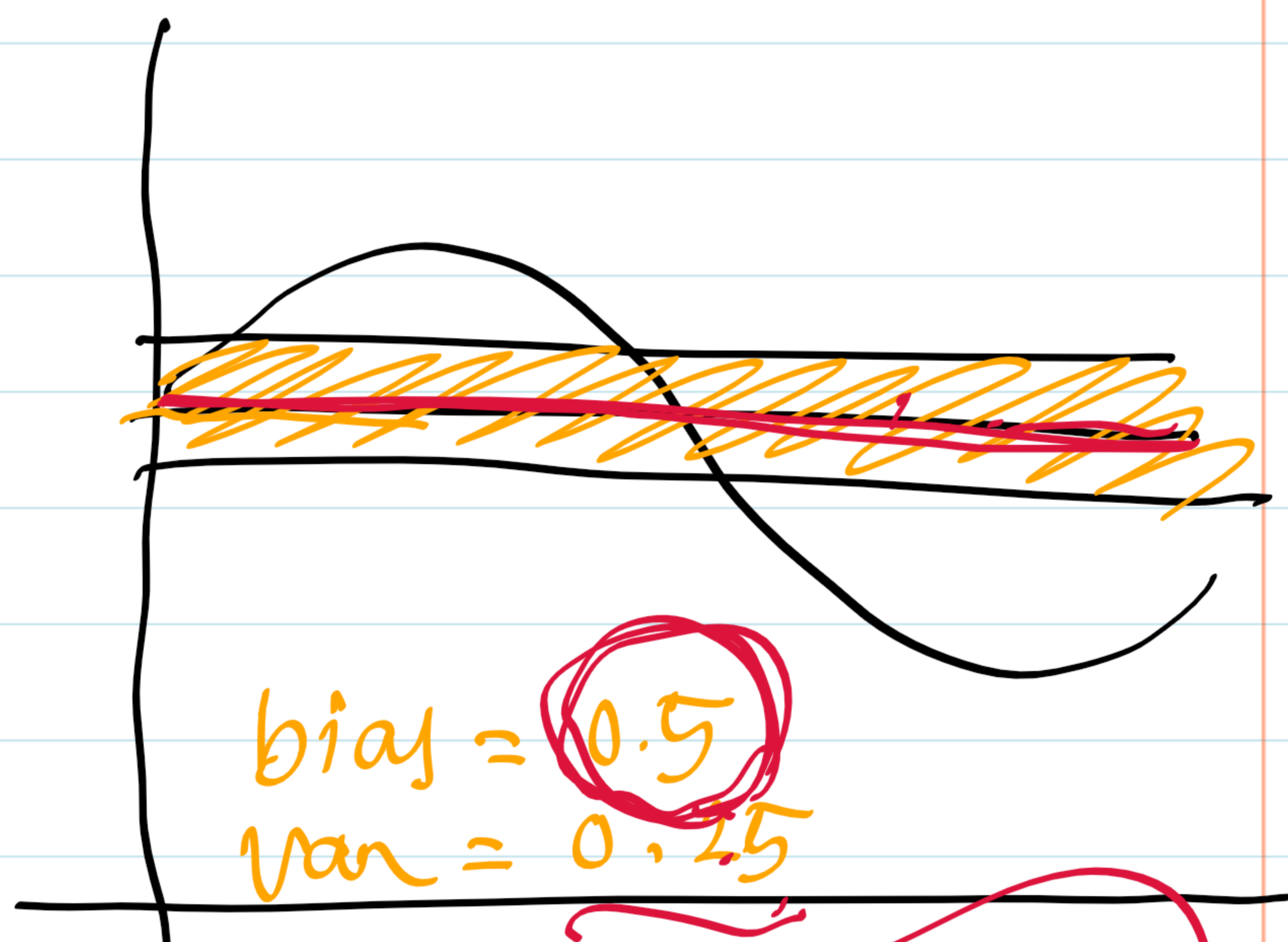
Hypothesis set 0: $\mathcal{M}_0 =$ set of all lines
of the form $h(x) = b$

Hypothesis set 1: $\mathcal{M}_1 =$ set of all lines
of the form $h(x) = ax + b$.



For μ_0 : $h(x) = \frac{y_1 + y_2}{2}$

For μ_1 : $h(x) = \left(\frac{y_2 - y_1}{x_2 - x_1} \right) x + (y_1 x_2 - y_2 x_1)$
(better)



μ_0 $E_{out} = 0.75$

