

Bias Variance decomposition

- Non - Linear Regression

$x \in \mathbb{R}$

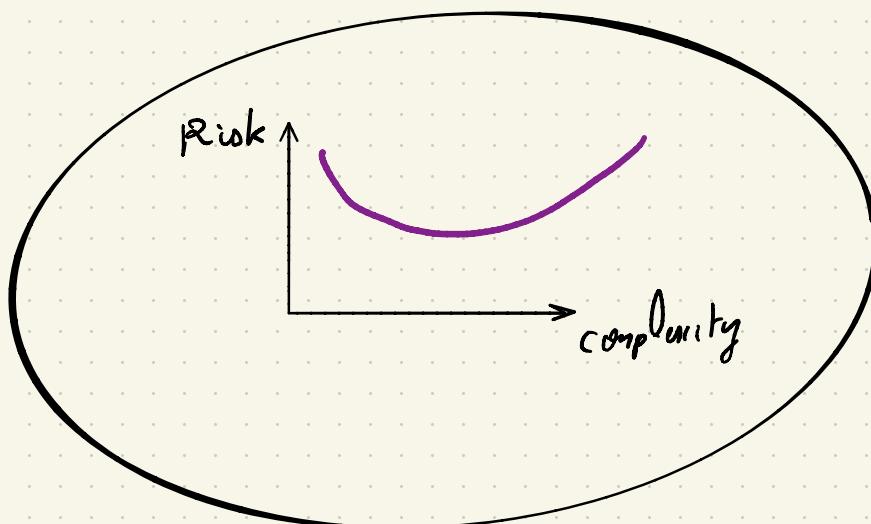
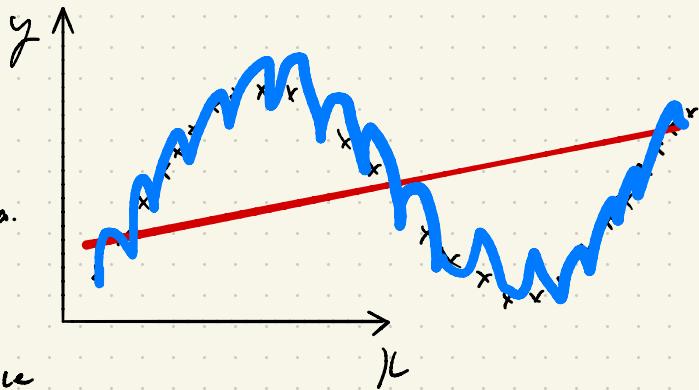
$$(x, x^2, x^3, \dots, x^d)$$

Complexity

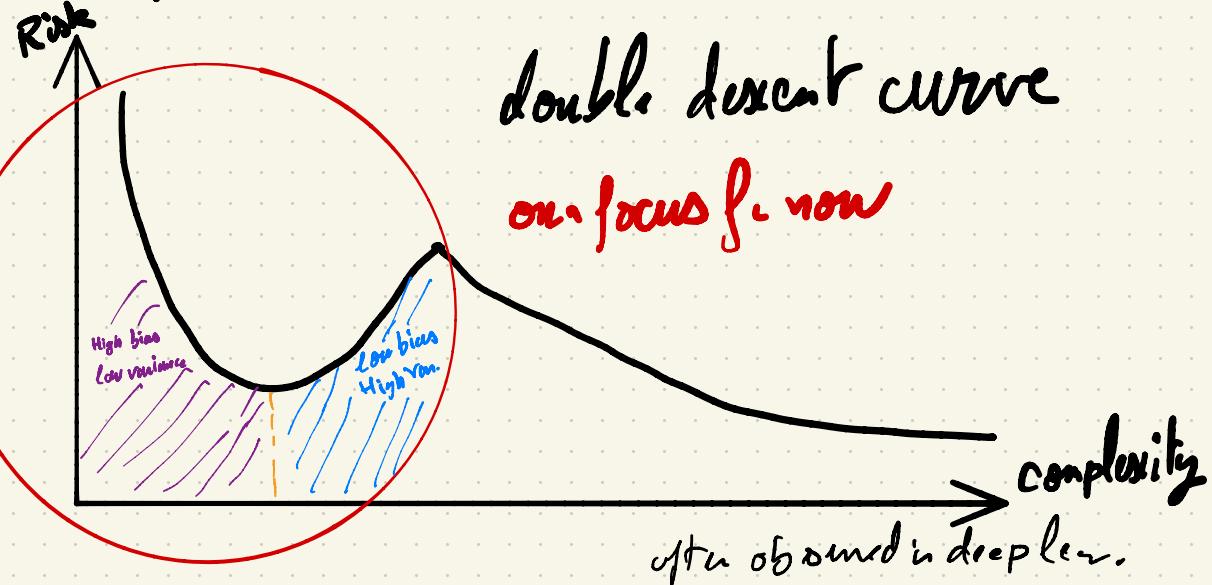
complexity of the class

- low complexity \rightarrow large bias
• small varia.

- High complexity
 - \rightarrow Low bias
 - High variance



Important note



Data model:

$$y = \underbrace{f(x)}_{\text{True model}} + \underbrace{\varepsilon}_{\text{noise}}$$

noise: $\varepsilon \sim D_\varepsilon$ iid, $\perp D_x$
 $E[\varepsilon] = 0$

- In general, the model is not realizable.
i.e. f is not in our model class

$$\mathbb{E}_D [(y - f(x))^2]$$

f from

$$S_{\text{Train}} \rightarrow \text{Algo. to get } f_{S_{\text{Train}}} \in \mathcal{F}(S_{\text{Train}})$$

x_0 fixed

$$(y - \int_{S_{\text{Train}}} f(x_0))^2 = (f(x_0) + \epsilon - \int_{S_{\text{Train}}} f(x_0))^2$$

$$\begin{aligned} \mathbb{E}_{\substack{\mathcal{D} \sim D \\ S_{\text{Train}} \sim \mathcal{D}}} [(y - \int_{S_{\text{Train}}} f(x_0))^2] &= \mathbb{E}_{\substack{\mathcal{D} \sim D \\ S_{\text{Train}} \sim \mathcal{D}}} [(f(x_0) + \epsilon - \int_{S_{\text{Train}}} f(x_0))^2] \\ &= \mathbb{E}_{\substack{S_{\text{Train}} \sim \mathcal{D}}} [(f(x_0) - \int_{S_{\text{Train}}} f(x_0))^2] + \mathbb{E}_{\substack{\mathcal{D} \sim D \\ S_{\text{Train}} \sim \mathcal{D}}} [\epsilon^2] + \text{Cross.T.} \end{aligned}$$

$$\text{Since } \mathbb{E}[\epsilon] = 0, \quad \mathbb{E}[\epsilon^2] = \text{Var}(\epsilon)$$

$$\text{now take: } \mathbb{E} [2\mathbb{E}[(f(x_0) - \int_{S_{\text{Train}}} f(x_0))]] = 2\mathbb{E}_0 [\mathbb{E}[f(x_0) - \int_{S_{\text{Train}}} f(x_0)]]$$

$$\mathcal{D}_\epsilon \perp\!\!\!\perp \mathcal{D} = 0$$

$$\mathbb{E}_{\substack{\mathcal{D} \sim D \\ S_{\text{Train}} \sim \mathcal{D}}} [(y - \int_{S_{\text{Train}}} f(x_0))^2] = \text{Var}(\epsilon) + \mathbb{E}_{S_{\text{Train}}} [(f(x_0) - \int_{S_{\text{Train}}} f(x_0))^2]$$

$$\begin{aligned} \mathbb{E}_{S_{\text{Train}}} [(f(x_0) - \int_{S_{\text{Train}}} f(x_0))^2] &= \mathbb{E}_{S_{\text{Train}} \sim D} \left[(f(x_0) - \mathbb{E}_{S_{\text{Train}} \sim D} [\int_{S_{\text{Train}}} f(x_0)]) + \mathbb{E}_{S_{\text{Train}} \sim D} [\int_{S_{\text{Train}}} f(x_0)] - \int_{S_{\text{Train}}} f(x_0) \right]^2 \\ &= (f(x_0) - \mathbb{E}_{S_{\text{Train}} \sim D} [\int_{S_{\text{Train}}} f(x_0)])^2 + \mathbb{E}_{S_{\text{Train}} \sim D} [(f(x_0) - \mathbb{E}_{S_{\text{Train}} \sim D} [f(x_0)])^2] + \text{Cross.T.} \end{aligned}$$

$$\begin{aligned} \text{C.T.} &= \mathbb{E}_{S_{\text{Train}}} \left[(f(x_0) - \mathbb{E}_{S_{\text{Train}} \sim D} [\int_{S_{\text{Train}}} f(x_0)]) \right] \cdot \left(\mathbb{E}_{S_{\text{Train}} \sim D} [\int_{S_{\text{Train}}} f(x_0)] - \int_{S_{\text{Train}}} f(x_0) \right) \\ &= (f(x_0) - \mathbb{E}_{S_{\text{Train}}} [\int_{S_{\text{Train}}} f(x_0)]) \cdot (\mathbb{E}_{S_{\text{Train}}} [\int_{S_{\text{Train}}} f(x_0)] - \mathbb{E}_{S_{\text{Train}}} [\int_{S_{\text{Train}}} f(x_0)]) = 0 \end{aligned}$$

$$\underset{S_{\text{Train}}}{\mathbb{E}_E} \left[(f(x_0) + E - f_{S_{\text{Train}}}(x_0))^2 \right]$$

$$= \text{Var}(E) + \left(f(x_0) - \underbrace{\mathbb{E}_{S_{\text{Train}}} [f_{S_{\text{Train}}}(x_0)]}_{\text{bias}} \right)^2$$

$$+ \underbrace{\mathbb{E}_{S_{\text{Test}}} \left[\left(f_{S_{\text{Test}}}(x_0) - \mathbb{E}_{S_{\text{Train}}} [f_{S_{\text{Train}}}(x_0)] \right)^2 \right]}_{\text{Variance}}$$

\Rightarrow Decomposition in 3 positive terms

- Noise Term, noise data
- Bias Term
- Variance Term

Bias Term: $f(x_0)$ = actual function

$\mathbb{E}_{S_{\text{Train}}} [f_{S_{\text{Train}}}(x_0)]$: expectation of prediction
(average predict.)

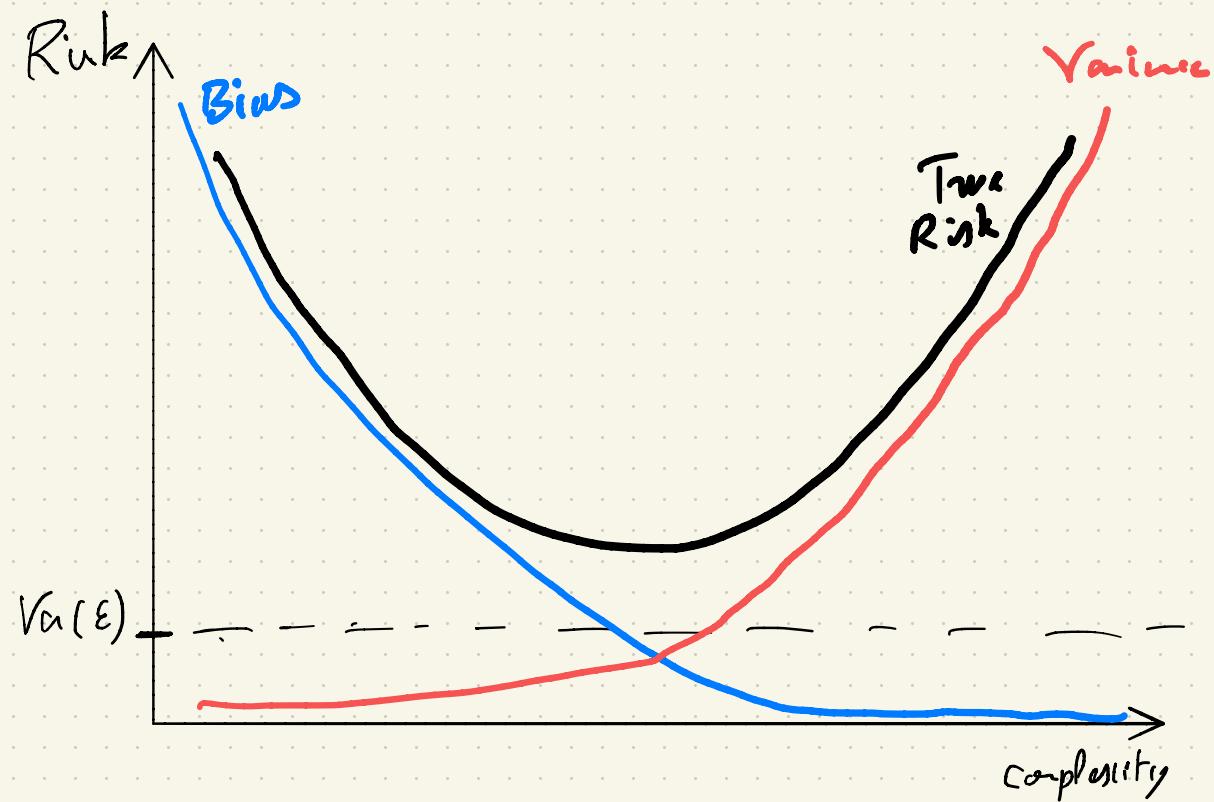
\rightarrow if complexity small \rightarrow high bias
 High \rightarrow small bias

Variance Term :

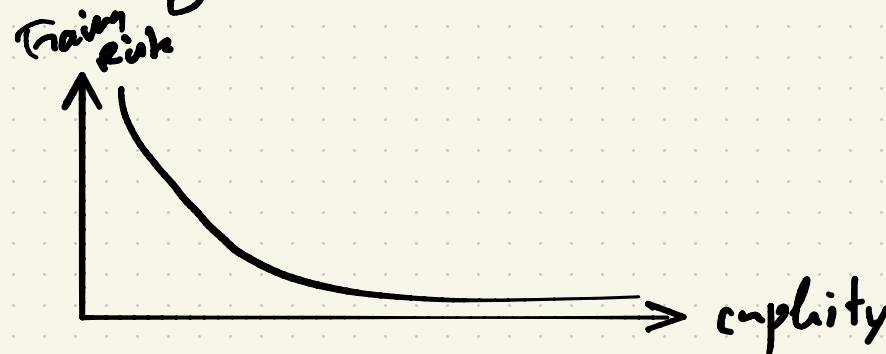
consistency of prediction

Complexity \nearrow Variance \uparrow

In Summary :



Training risk:



For neural networks = Train with a pert. alg.
(SGD)

