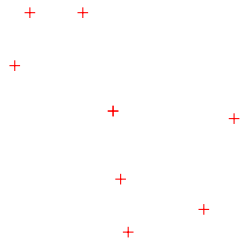


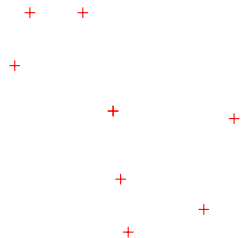
# Module 8.1 : Bias and Variance

We will begin with a quick overview of bias, variance and the trade-off between them.

- Let us consider the problem of fitting a curve through a given set of points

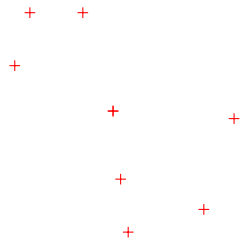


The points were drawn from a sinusoidal function (the true  $f(x)$ )



The points were drawn from a sinusoidal function (the true  $f(x)$ )

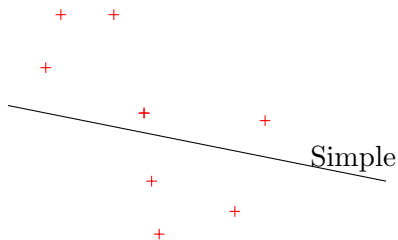
- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :



The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

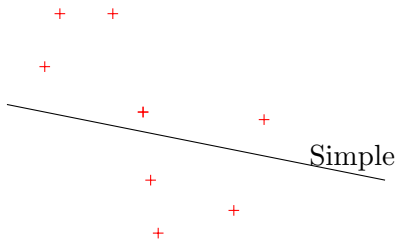
$$\begin{matrix} \textit{Simple} \\ \textit{(degree:1)} \end{matrix} \quad y = \hat{f}(x) = w_1x + w_0$$



The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

$$\begin{matrix} \textit{Simple} \\ \textit{(degree:1)} \end{matrix} \quad y = \hat{f}(x) = w_1x + w_0$$

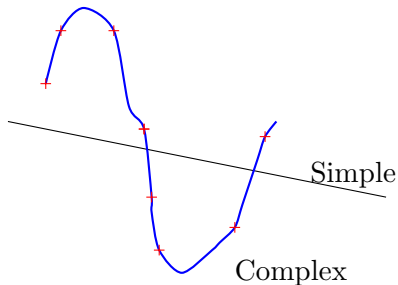


The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

$$\begin{array}{l} \text{Simple} \\ (\text{degree:1}) \end{array} \quad y = \hat{f}(x) = w_1x + w_0$$

$$\begin{array}{l} \text{Complex} \\ (\text{degree:25}) \end{array} \quad y = \hat{f}(x) = \sum_{i=1}^{25} w_i x^i + w_0$$



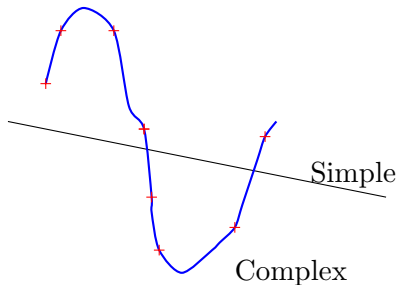
The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

$$\begin{array}{l} \text{Simple} \\ (\text{degree:1}) \end{array} \quad y = \hat{f}(x) = w_1x + w_0$$

$$\begin{array}{l} \text{Complex} \\ (\text{degree:25}) \end{array} \quad y = \hat{f}(x) = \sum_{i=1}^{25} w_i x^i + w_0$$





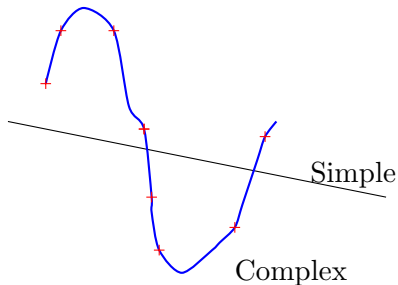
The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

$$\begin{array}{l} \text{Simple} \\ (\text{degree:1}) \end{array} \quad y = \hat{f}(x) = w_1x + w_0$$

$$\begin{array}{l} \text{Complex} \\ (\text{degree:25}) \end{array} \quad y = \hat{f}(x) = \sum_{i=1}^{25} w_i x^i + w_0$$

- Note that in both cases we are making an assumption about how  $y$  is related to  $x$ . We have no idea about the true relation  $f(x)$



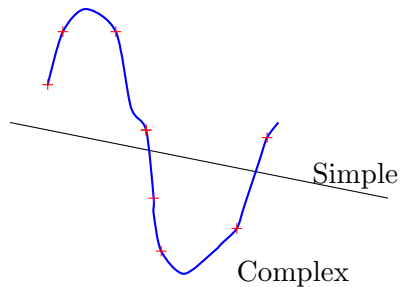
The points were drawn from a sinusoidal function (the true  $f(x)$ )

- Let us consider the problem of fitting a curve through a given set of points
- We consider two models :

$$\begin{array}{l} \text{Simple} \\ (\text{degree:1}) \end{array} \quad y = \hat{f}(x) = w_1x + w_0$$

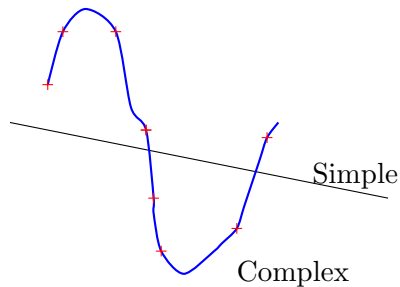
$$\begin{array}{l} \text{Complex} \\ (\text{degree:25}) \end{array} \quad y = \hat{f}(x) = \sum_{i=1}^{25} w_i x^i + w_0$$

- Note that in both cases we are making an assumption about how  $y$  is related to  $x$ . We have no idea about the true relation  $f(x)$
- The training data consists of 100 points



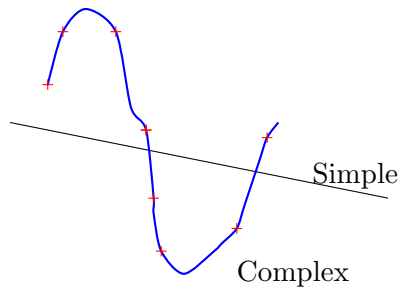
The points were drawn from a sinusoidal function (the true  $f(x)$ )

- We sample 25 points from the training data and train a simple and a complex model



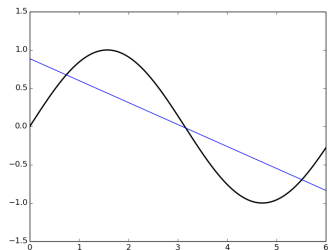
The points were drawn from a sinusoidal function (the true  $f(x)$ )

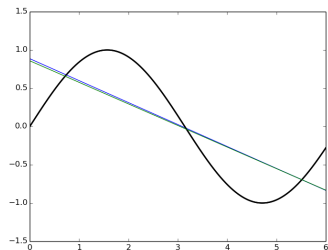
- We sample 25 points from the training data and train a simple and a complex model
- We repeat the process ' $k$ ' times to train multiple models (each model sees a different sample of the training data)

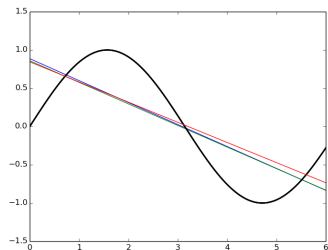


The points were drawn from a sinusoidal function (the true  $f(x)$ )

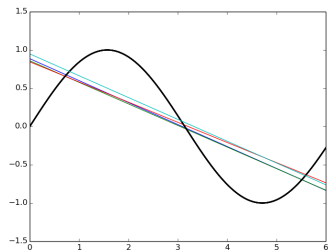
- We sample 25 points from the training data and train a simple and a complex model
- We repeat the process ' $k$ ' times to train multiple models (each model sees a different sample of the training data)
- We make a few observations from these plots

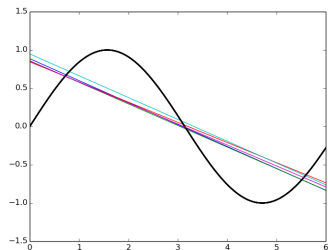


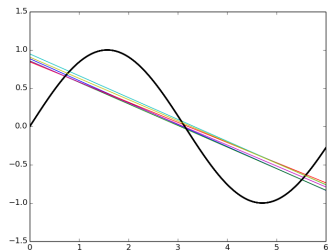


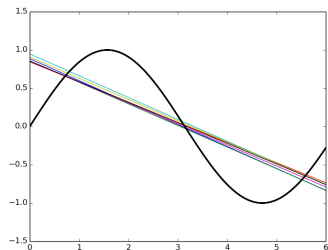


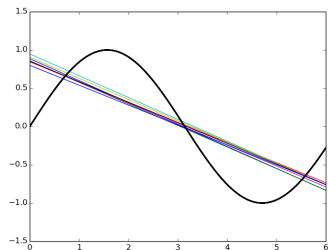


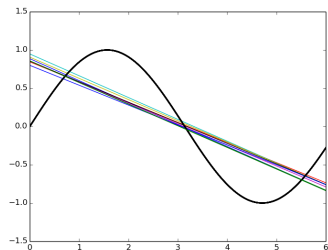


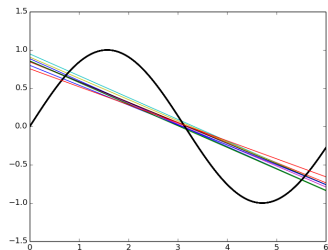


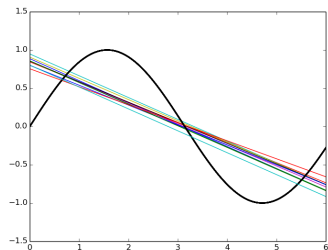




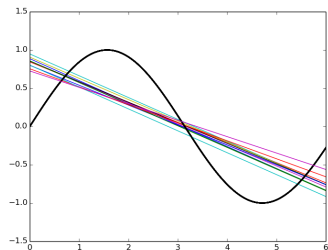


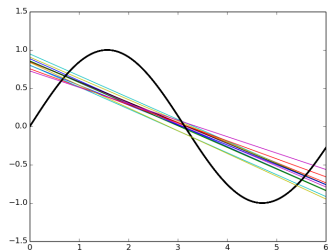


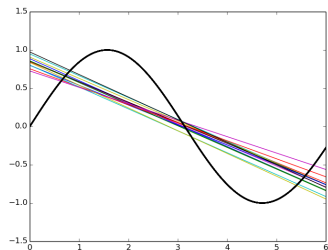


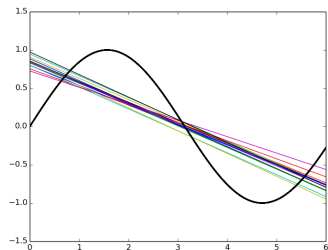


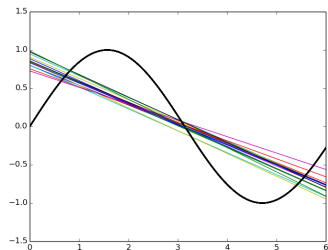


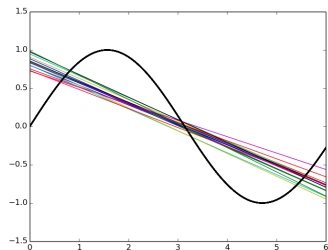


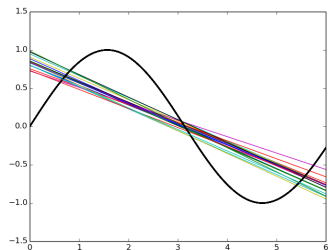


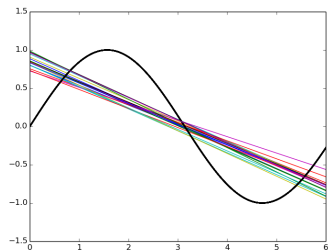




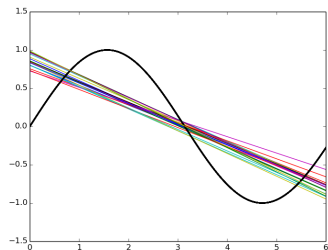


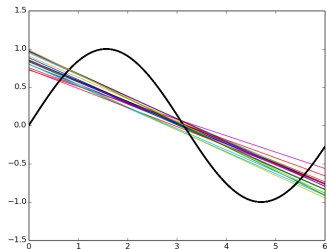




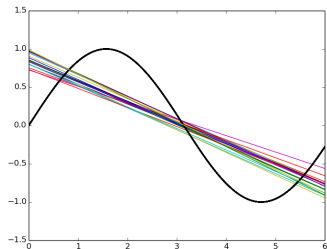




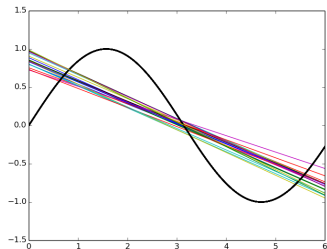




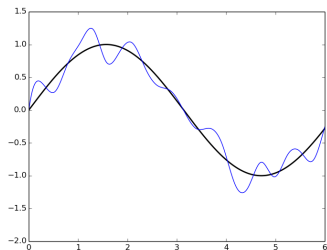
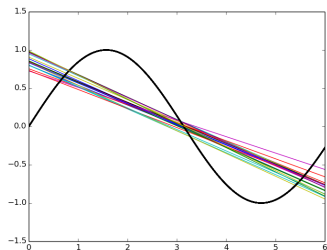
- Simple models trained on different samples of the data do not differ much from each other



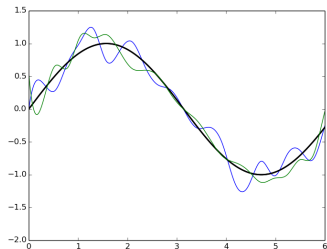
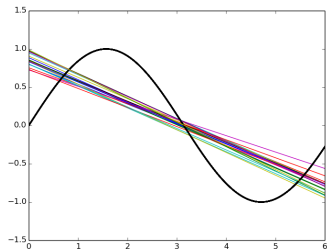
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



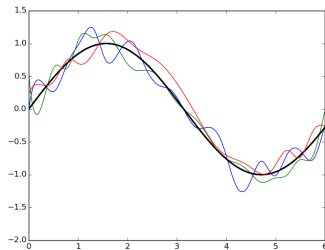
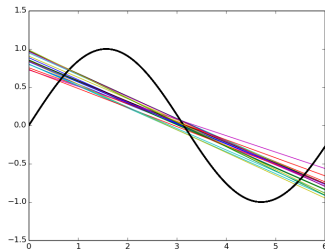
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



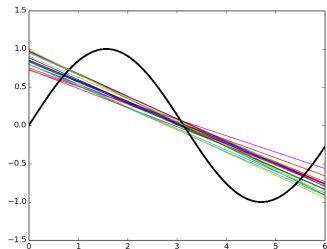
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



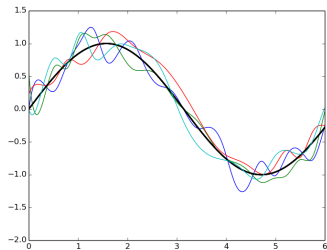
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



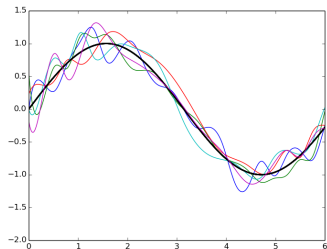
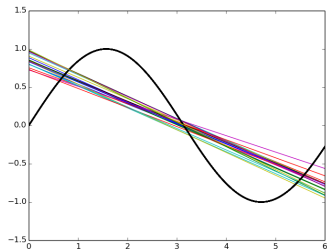
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



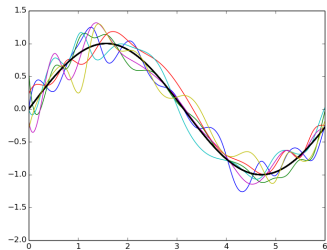
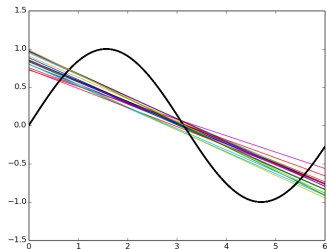
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



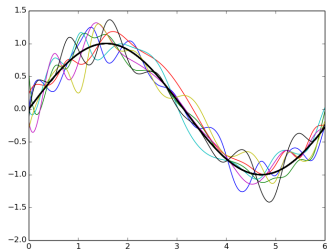
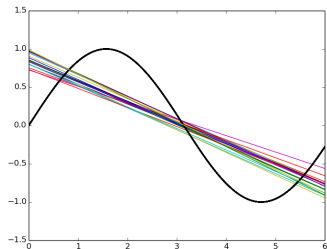




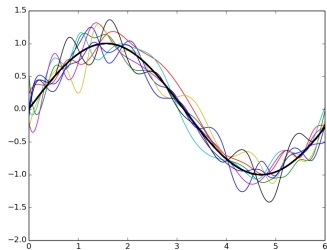
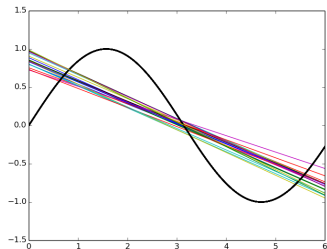
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



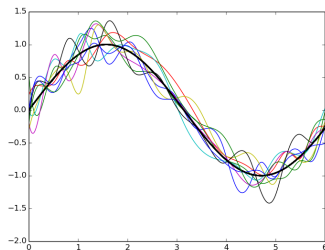
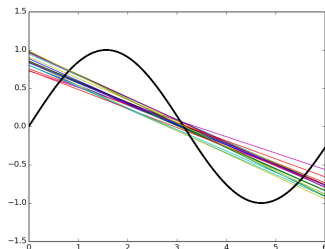
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



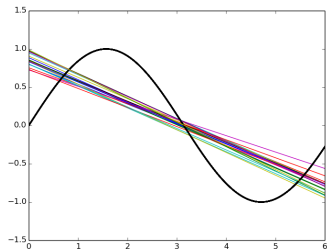
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



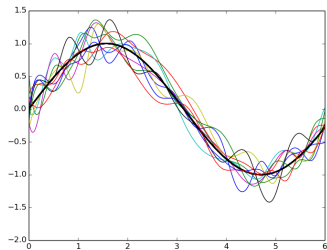
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

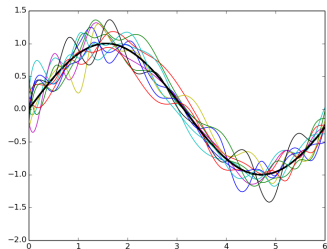
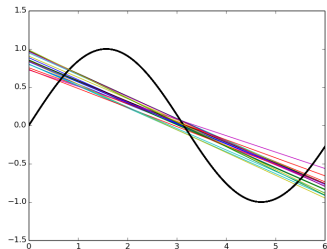


- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

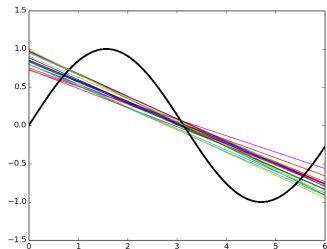


- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

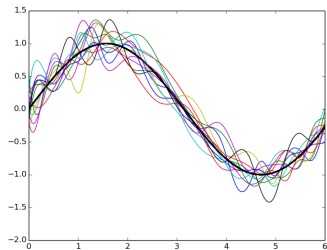




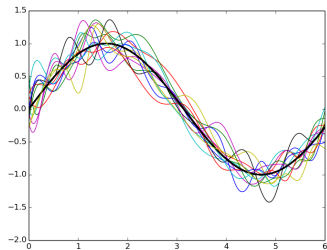
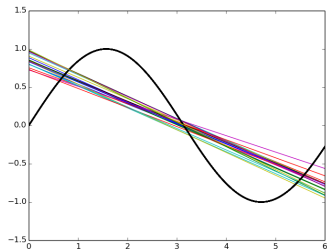
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



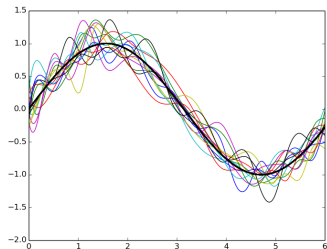
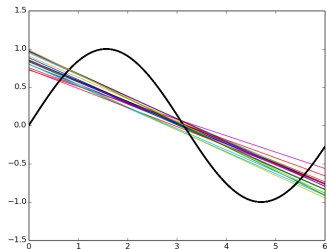
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



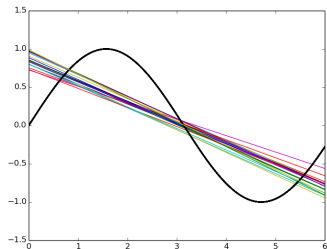




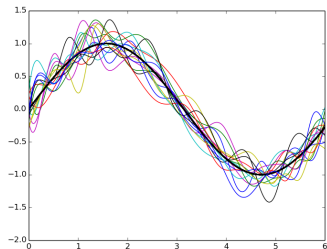
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

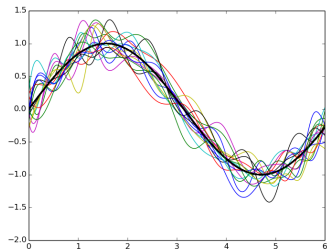
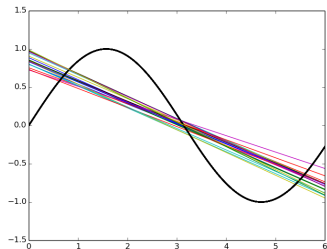


- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

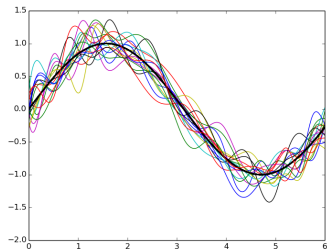
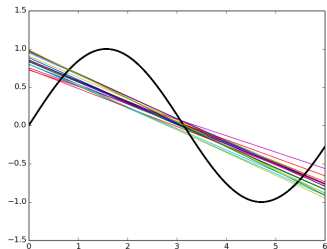


- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

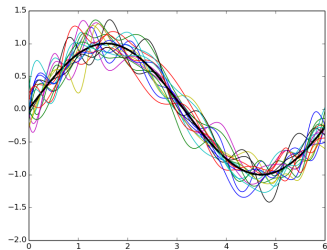
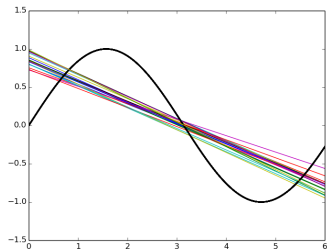




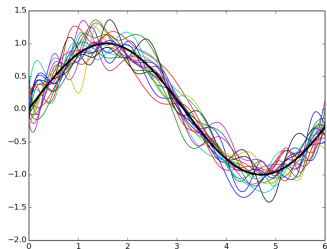
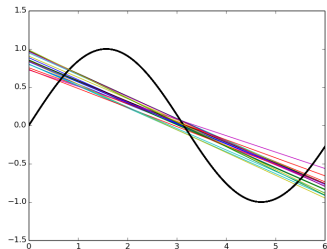
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



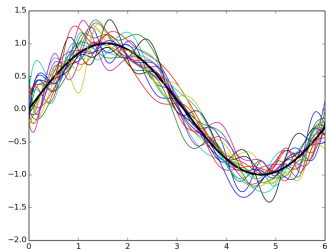
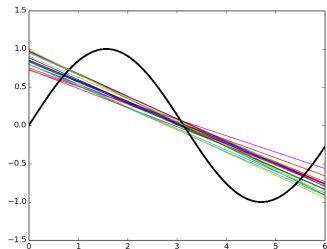
- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

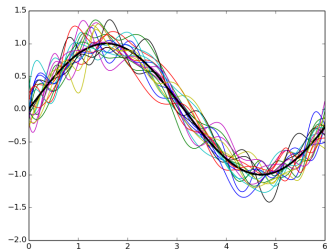
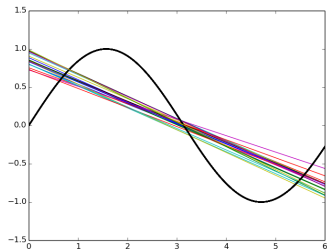


- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)



- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)

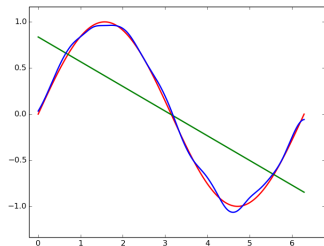




- Simple models trained on different samples of the data do not differ much from each other
- However they are very far from the true sinusoidal curve (under fitting)
- On the other hand, complex models trained on different samples of the data are very different from each other (high variance)

- Let  $f(x)$  be the true model (sinusoidal in this case) and  $\hat{f}(x)$  be our estimate of the model (simple or complex, in this case) then,

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$



Green Line: Average value of  $\hat{f}(x)$   
for the simple model

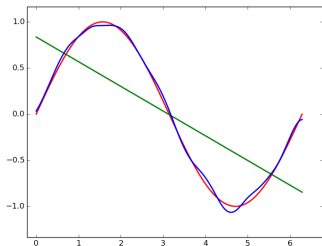
Blue Curve: Average value of  $\hat{f}(x)$   
for the complex model

Red Curve: True model ( $f(x)$ )

- Let  $f(x)$  be the true model (sinusoidal in this case) and  $\hat{f}(x)$  be our estimate of the model (simple or complex, in this case) then,

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

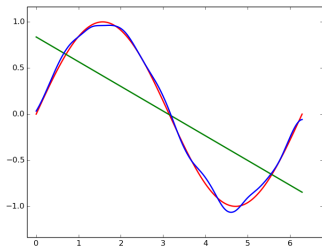
- $E[\hat{f}(x)]$  is the average (or expected) value of the model



Green Line: Average value of  $\hat{f}(x)$   
for the simple model

Blue Curve: Average value of  $\hat{f}(x)$   
for the complex model

Red Curve: True model ( $f(x)$ )



Green Line: Average value of  $\hat{f}(x)$   
for the simple model

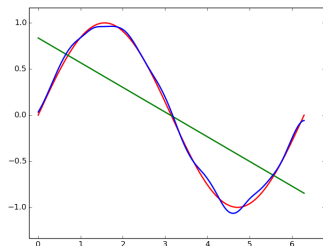
Blue Curve: Average value of  $\hat{f}(x)$   
for the complex model

Red Curve: True model ( $f(x)$ )

- Let  $f(x)$  be the true model (sinusoidal in this case) and  $\hat{f}(x)$  be our estimate of the model (simple or complex, in this case) then,

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

- $E[\hat{f}(x)]$  is the average (or expected) value of the model
- We can see that for the simple model the average value (green line) is very far from the true value  $f(x)$  (sinusoidal function)



Green Line: Average value of  $\hat{f}(x)$  for the simple model

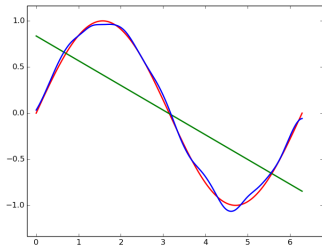
Blue Curve: Average value of  $\hat{f}(x)$  for the complex model

Red Curve: True model ( $f(x)$ )

- Let  $f(x)$  be the true model (sinusoidal in this case) and  $\hat{f}(x)$  be our estimate of the model (simple or complex, in this case) then,

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

- $E[\hat{f}(x)]$  is the average (or expected) value of the model
- We can see that for the simple model the average value (green line) is very far from the true value  $f(x)$  (sinusoidal function)
- Mathematically, this means that the simple model has a high bias



Green Line: Average value of  $\hat{f}(x)$  for the simple model

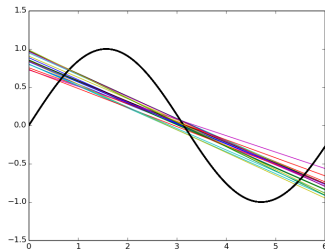
Blue Curve: Average value of  $\hat{f}(x)$  for the complex model

Red Curve: True model ( $f(x)$ )

- Let  $f(x)$  be the true model (sinusoidal in this case) and  $\hat{f}(x)$  be our estimate of the model (simple or complex, in this case) then,

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

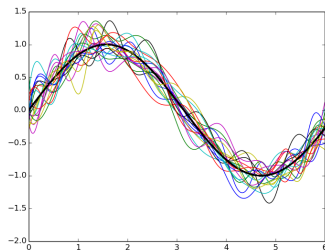
- $E[\hat{f}(x)]$  is the average (or expected) value of the model
- We can see that for the simple model the average value (green line) is very far from the true value  $f(x)$  (sinusoidal function)
- Mathematically, this means that the simple model has a high bias
- On the other hand, the complex model has a low bias

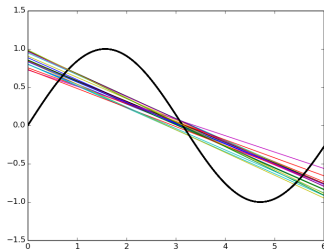


- We now define,

$$\text{Variance } (\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

(Standard definition from statistics)

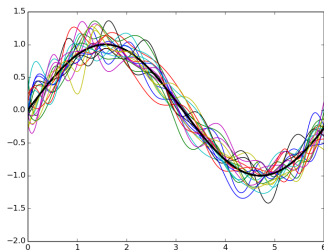




- We now define,

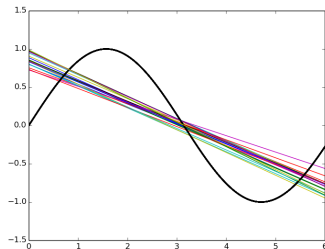
$$\text{Variance } (\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

(Standard definition from statistics)



- Roughly speaking it tells us how much the different  $\hat{f}(x)$ 's (trained on different samples of the data) differ from each other

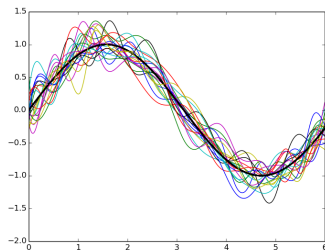




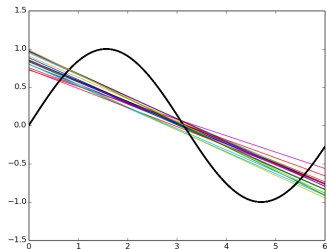
- We now define,

$$\text{Variance } (\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

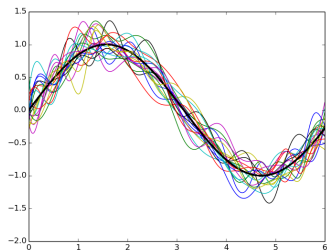
(Standard definition from statistics)

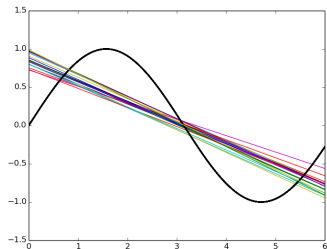


- Roughly speaking it tells us how much the different  $\hat{f}(x)$ 's (trained on different samples of the data) differ from each other
- It is clear that the simple model has a low variance whereas the complex model has a high variance

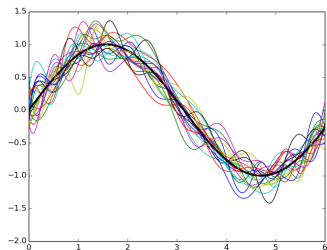


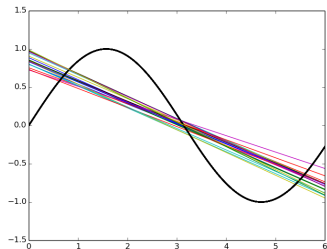
- In summary (informally)



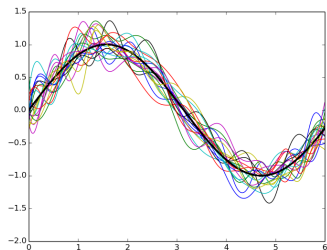


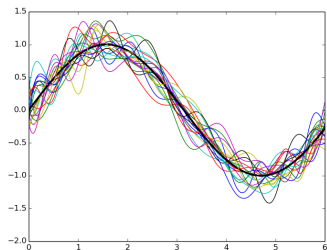
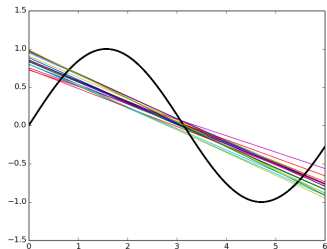
- In summary (informally)
- Simple model: high bias, low variance



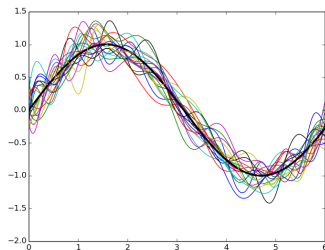
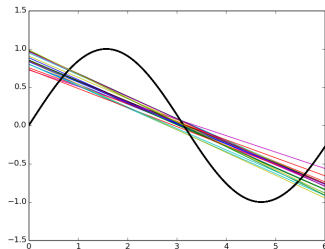


- In summary (informally)
- Simple model: high bias, low variance
- Complex model: low bias, high variance





- In summary (informally)
- Simple model: high bias, low variance
- Complex model: low bias, high variance
- There is always a trade-off between the bias and variance



- In summary (informally)
- Simple model: high bias, low variance
- Complex model: low bias, high variance
- There is always a trade-off between the bias and variance
- Both bias and variance contribute to the mean square error. Let us see how