

1. the law of total expectation  $E(x) = E(E(x|Y))$

$$\text{From } \text{Var}[x] = E[x^2] - [E[x]]^2$$

$$E(\text{Var}[x|Y]) = E(E[x^2|Y] - E[x|Y]^2) = E[x^2] - E(E[x|Y]^2)$$

$$\text{Var}(E[x|Y]) = E(E[x|Y]^2) - E(E[x|Y])^2 = E(E[x|Y]^2) - E[x]^2$$

$$\text{Then } \text{Var}[x] = E[x^2] - E[x]^2 = E(\text{Var}[x|Y]) + \text{Var}(E[x|Y])$$

2. for negative binomial distribution  $f_x(x) = \binom{x+r-1}{x} p^x (1-p)^r$

$$\text{likelihood function } L(p) = \prod_{i=1}^n f_x(x_i) = \prod_{i=1}^n \binom{x_i+r-1}{x_i} p^{x_i} (1-p)^{r_i}$$

$$\text{As only } x_1=2 \text{ is observed } L(p) = \binom{4}{2} p^2 (1-p)^3 = 6 p^2 (1-p)^3$$

$$\text{to have the largest } L(p) \quad \ln L(p) = \ln 6 + 2\ln p + 3\ln(1-p)$$

$$\text{Let } d\ln(L(p)) = \frac{2}{p} - \frac{3}{1-p} = 0 \quad \boxed{p = 0.4}$$

3. Each time tossing a coin is an independent Bernoulli trial, assume with head on probability  $p_0$ . For  $n=1000$  trials, follows a binomial distribution. Let  $X$  denote the number of heads out of  $n$ ,  $X \sim \text{Binomial}(1000, p_0)$

Set  $H_0: p_0 = 0.5$      $H_1: p_0 \neq 0.5$ , critical value 0.05

If  $H_0$  is true, denote the number of heads on occurred as  $x$

$$2 \times P[X \geq 560 | p_0 = 0.5] = 0.00016 < 0.05$$

so there is evidence to reject  $H_0$ ,

it's not reasonable to assume the coin is fair

4. As  $Y_1, \dots, Y_n$  are i.i.d. normal ( $\mu, 1$ ),

The sample mean  $\bar{Y}$  of  $Y_1, \dots, Y_n$  has distribution  $N(\mu, \frac{1}{n})$

$$\begin{aligned} \frac{\bar{Y}-\mu}{\sqrt{n}} &\sim N(0, 1), \quad P[-1.96 \leq \frac{\bar{Y}-\mu}{\sqrt{n}} \leq 1.96] = P[\frac{\bar{Y}-\mu}{\sqrt{n}} \leq 1.96] - P[\frac{\bar{Y}-\mu}{\sqrt{n}} \leq -1.96] \\ &= \phi(1.96) - \phi(-1.96) \end{aligned}$$

Then the 95% confidence interval  $(\bar{y} - 1.96/\sqrt{n}, \bar{y} + 1.96/\sqrt{n})$

and the future observation  $Y_{n+1}$  still has distribution  $N(\mu, 1)$

Because the confidence interval is obtained through the sample, a fixed interval.

$$\begin{aligned} P[\bar{y} - 1.96/\sqrt{n} \leq Y_{n+1} \leq \bar{y} + 1.96/\sqrt{n}] &= P[\bar{y} - \mu - 1.96/\sqrt{n} \leq Y_{n+1} - \mu \leq \bar{y} - \mu + 1.96/\sqrt{n}] \\ &= P[Y_{n+1} \leq \bar{y} + 1.96/\sqrt{n}] - P[Y_{n+1} \leq \bar{y} - 1.96/\sqrt{n}] \\ &= P[Y_{n+1} - \mu \leq \bar{y} - \mu + 1.96/\sqrt{n}] - P[Y_{n+1} - \mu \leq \bar{y} - \mu - 1.96/\sqrt{n}] \\ &= \phi(1.96/\sqrt{n} + \bar{y} - \mu) - \phi(-1.96/\sqrt{n} + \bar{y} - \mu) \end{aligned}$$

$$\textcircled{1} \text{ If } 1.96/\sqrt{n} + \bar{y} - \mu = 1.96 \text{ s.t. } \bar{y} = \mu + 1.96(1 - \frac{1}{\sqrt{n}}) \quad p = 0.95$$

$$\textcircled{2} \text{ If } 1.96/\sqrt{n} + \bar{y} - \mu > 1.96, \text{ s.t. } \bar{y} > \mu + 1.96(1 - \frac{1}{\sqrt{n}})$$

$$\text{or } \textcircled{3} \text{ If } 1.96/\sqrt{n} + \bar{y} - \mu < 1.96 \text{ s.t. } \bar{y} < \mu + 1.96(1 - \frac{1}{\sqrt{n}})$$

the interval will shift,  $p < 0.95$

5. a) As  $\bar{x}=10, s^2=1, n=10$

Set  $H_0: \mu=5$  critical value when  $\alpha=0.05 \quad z_{\alpha/2}=1.96$

$$Z = \frac{\bar{x}-\mu_0}{s/\sqrt{n}} = \frac{10-5}{1/\sqrt{10}} = 15.8 > 1.96$$

so we could reject  $H_0$  at significance level 0.05  
it is not reasonable to suggest  $\mu=5$

b) Set  $H_0: \mu=5$  critical value when  $\alpha=0.05 \quad t_{\alpha/2, 9}=2.26$

$$T_{n+1} = \frac{\bar{x}-\mu_0}{s/\sqrt{n}} = \frac{10-5}{2/\sqrt{10}} = 7.9$$

so we reject  $H_0$  at significance level 0.05  
it is not reasonable to suggest  $\mu=5$

$$6. \text{ As } E[(x_i - \mu)^2] = \text{Var}[x_i] = \sigma^2$$

$$E[(\bar{x} - \mu)^2] = \text{Var}[\bar{x}] = \frac{\sigma^2}{n}$$

As for  $\sigma^2$

$$\begin{aligned} E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] &= E\left[\sum_{i=1}^n (x_i - \mu + \mu - \bar{x})^2\right] = E\left[\sum_{i=1}^n (x_i - \mu)^2 - 2(x_i - \mu)(\bar{x} - \mu) + n(\mu - \bar{x})^2\right] \\ &= E\left[\sum_{i=1}^n (x_i - \mu)^2 - 2(\bar{x} - \mu)(n\bar{x} - n\mu) + n(\mu - \bar{x})^2\right] \\ &= \sum_{i=1}^n E[(x_i - \mu)^2] - nE[(\bar{x} - \mu)^2] = n\sigma^2 - n\frac{\sigma^2}{n} = (n-1)\sigma^2 \end{aligned}$$

$$\text{so } E\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right] = \sigma^2$$

$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$  is unbiased estimator for  $\sigma^2$

$$\text{Then } E\left[\frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 - \frac{\sigma^2}{n}\right] = 0$$

$\frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$  is an unbiased estimator of  $\text{Var}[\bar{x}]$

7. as  $\hat{\beta}_1$  is unbiased  $E[\hat{\beta}_1] = \beta_1$

from  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ ,  $\sum_{i=1}^n Y_i = N\beta_0 + \sum_{i=1}^n \beta_1 x_i + \sum_{i=1}^n \varepsilon_i$  so that  $\bar{Y} = \beta_0 + \beta_1 \bar{x} + \bar{\varepsilon}$

from  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$ ,  $\hat{\beta}_0 = \beta_0 + (\beta_1 - \hat{\beta}_1) \bar{x} + \bar{\varepsilon}$

$$E[\hat{\beta}_0] = E[\beta_0] + E[\bar{x}(\beta_1 - \hat{\beta}_1)] + E[\bar{\varepsilon}] = \beta_0 + \bar{x} E[(\beta_1) - E[\hat{\beta}_1]] = \beta_0$$

so  $\hat{\beta}_0$  is unbiased

Because of the law of large numbers and the consistency of  $\hat{\beta}_1$   
when  $n \rightarrow \infty$   $\hat{\beta}_1 \rightarrow \beta_1$ ,  $\hat{\beta}_1 \bar{x}$  converge to  $\beta_1 E[x]$

$\bar{Y}$  would converge to  $E(y) = \beta_0 + \beta_1 E(x)$

$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$  so  $\hat{\beta}_0$  would converge to  $\beta_0$ ,  $\hat{\beta}_0$  is consistent

b) Assume  $Y_i = Y|_{X_i}$ , random sample  $(x_1, Y_1), \dots, (x_n, Y_n)$

$Y|_X$  follows a normal distribution with variance  $\sigma^2$  and mean  $M_{Y|X} = \beta_0 + \beta_1 x$   
 random variables  $Y_i = Y|_{X_i}$  are i.i.d. normal

$$\text{For } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{j=1}^n (x_j - \bar{x})^2} \quad \text{as } \sum_{i=1}^n (x_i - \bar{x}) = 0$$

$$= \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i, \text{ which is a linear combination of the i.i.d. normally distributed } Y_i$$

so  $\hat{\beta}_1$  follows a normal distribution

to find its mean and variance,

$$E[\hat{\beta}_1] = E\left[\frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i\right] = \frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \cdot \sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i)$$

$$= \frac{\beta_0}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) + \beta_1 \frac{\sum_{i=1}^n x_i(x_i - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2} = 0 + \beta_1 \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{\sum_{j=1}^n x_j^2 + n\bar{x}^2 - 2n\bar{x}^2} = \beta_1$$

$$\text{Var}[\hat{\beta}_1] = \text{Var}\left[\frac{1}{\sum_{j=1}^n (x_j - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) Y_i\right] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}[Y_i]}{\left(\sum_{j=1}^n (x_j - \bar{x})^2\right)^2} = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

so  $\hat{\beta}_1$  follows a normal distribution with mean  $\beta_1$ , variance  $\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

7. c) For  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \frac{1}{n} \sum_{i=1}^n Y_i - \hat{\beta}_1 \cdot \frac{1}{n} \sum_{i=1}^n X_i$   
As  $\hat{\beta}_1$  and  $Y_i$  are both normally distributed  
 $\hat{\beta}_1$  is a linear combination  
so  $\hat{\beta}_1$  is normally distributed  
 $E[\hat{\beta}_0] = \beta_0$  has been proved in (a)

$$\text{Var}[\hat{\beta}_0] = \text{Var}[\bar{Y} - \hat{\beta}_1 \bar{X}] = \text{Var}[\bar{Y}] + (\bar{X})^2 \text{Var}[\hat{\beta}_1] - 2\bar{X} \text{Cov}[\bar{Y}, \hat{\beta}_1]$$

$$\text{As } \text{Var}[\bar{Y}] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n Y_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[Y_i] = \frac{6^2}{n}$$

$$\begin{aligned} \text{Cov}[\bar{Y}, \hat{\beta}_1] &= \text{Cov}\left\{\frac{1}{n} \sum_{i=1}^n Y_i, \frac{\sum_{j=1}^n (x_j - \bar{x}) Y_j}{\sum_{k=1}^n (x_k - \bar{x})^2}\right\} = \frac{1}{n} \frac{1}{\sum_{k=1}^n (x_k - \bar{x})^2} \text{Cov}\left\{\sum_{i=1}^n Y_i, \sum_{j=1}^n (x_j - \bar{x}) Y_j\right\} \\ &= \frac{1}{n \sum_{k=1}^n (x_k - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \sum_{j=1}^n \text{Cov}(Y_i, Y_j) \\ &= \frac{1}{n \sum_{k=1}^n (x_k - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) 6^2 = 0 \end{aligned}$$

$$\begin{aligned} \text{So } \text{Var}[\hat{\beta}_0] &= \frac{6^2}{n} + (\bar{X})^2 \frac{6^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{6^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 + n \bar{x}^2 \right] \\ &= \frac{6^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \left( \sum_{i=1}^n x_i^2 + n \bar{x}^2 - \sum_{i=1}^n 2x_i \bar{x} + n \bar{x}^2 \right) = \frac{6^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

so  $\hat{\beta}_0$  follows a normal distribution with mean  $\beta_0$  and variance  $\frac{6^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}$

# Bonus Question

How bias and variance would change when model complexity changes?

## 1. Analyze from definition

- Bias: how much the prediction value obtained by the model differs from the real value.

*sampling* and *estimation* could both introduce bias. As the discussion is mainly about the impact of model complexity, more focus will be on *estimation error* (for sampling error, resampling and repeat the model building process then derive the average of prediction values could help, however, no change in model complexity).

So the causes for a larger bias could be a simpler model to do a simpler approximation, as a more complicated model usually could do a better fit. The risk of underfitting exists. Because with simpler model, less assumptions are made. For example, with the simplest model I could think (simple linear regression)

- Variance: with the same model, if a different training set is used, how much the estimate of the function will be changed. Intuitively, when a more complicated model is selected, more parameters need to be calculated. It would be much easier to obtain a result with a large variance compared to the former results. High-variance might fit the training data well (as with different training data, large change in estimate), but the risk of overfitting exists (fit too good to predict the future value). Large noise might be fitted by the complicated model.

So, there exists bias-variance tradeoff. From the lecture slides,

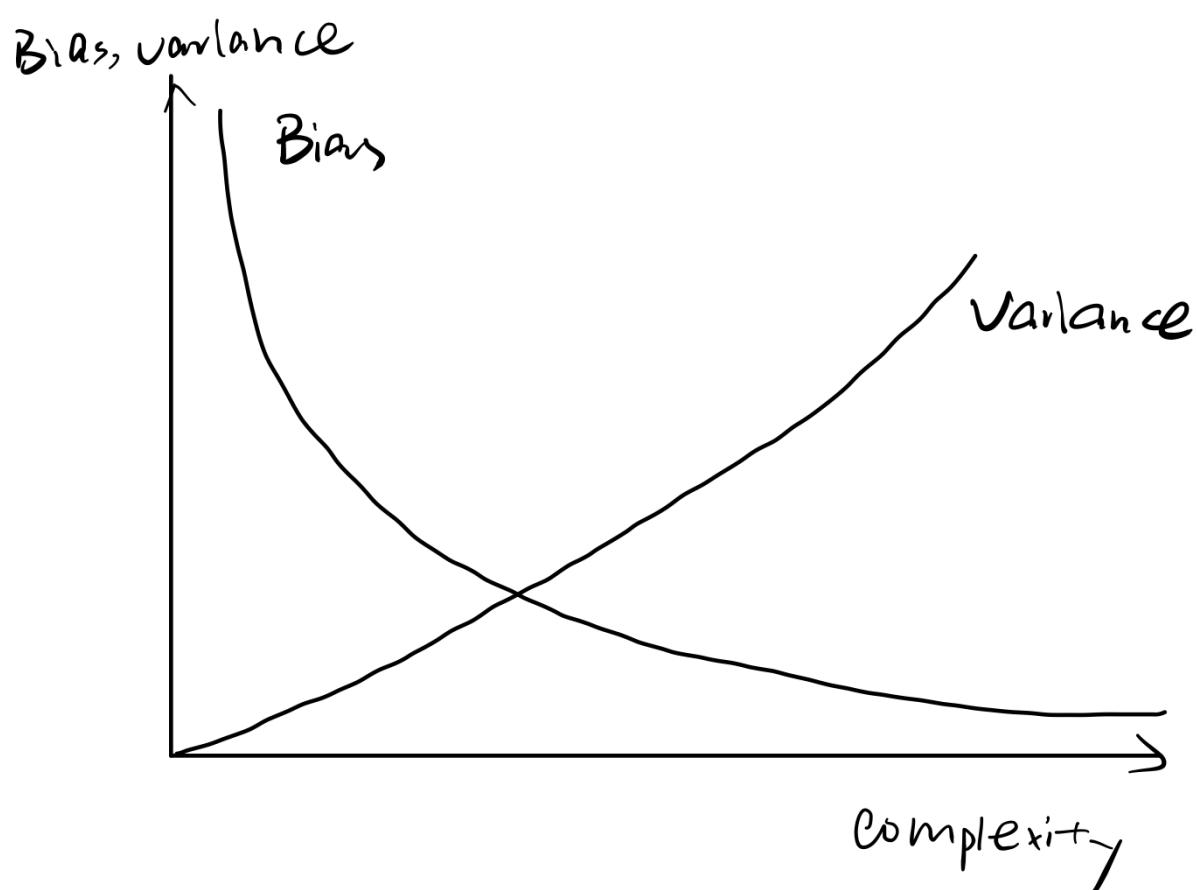
$$\text{MSE}(m) = \mathbb{E}[(Y - m)^2] = \text{Var}[Y] + (\mathbb{E}[Y - m])^2,$$

Model with small variance and high bias *underfit* the truth target, while with high variance and small bias *overfit* the truth target.

## 2. Answer for the trend

This answer is mainly based on the intuition in 1. and the fitted results shown in 3.

	<b>Bias</b>	<b>Variance</b>
monotone	yes	yes
change in rate	first change fast, then slower	roughly the same rate, do not differ much
Increase / decrease	Decrease	Increase

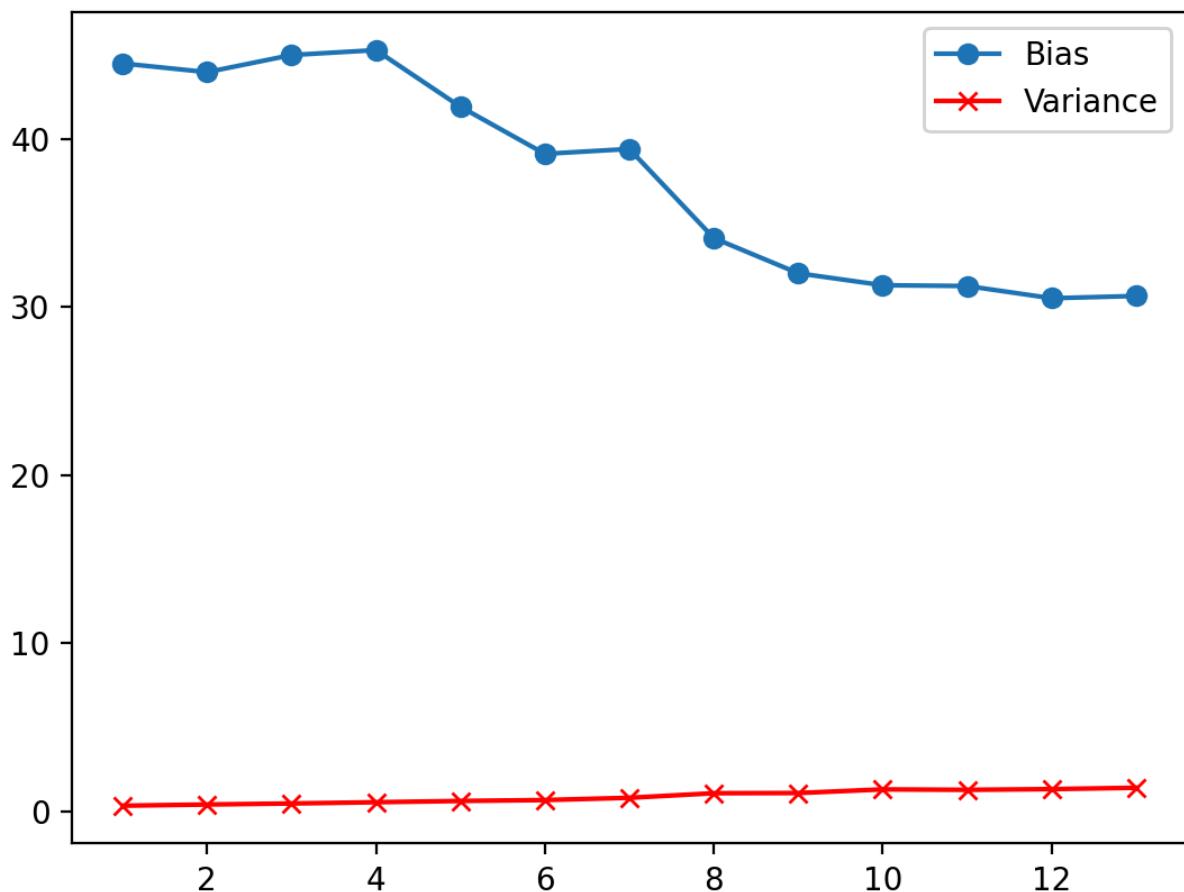


3. The fit is performed through Python and the raw data obtained by url = '<https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv>'. The head of the imported data set is

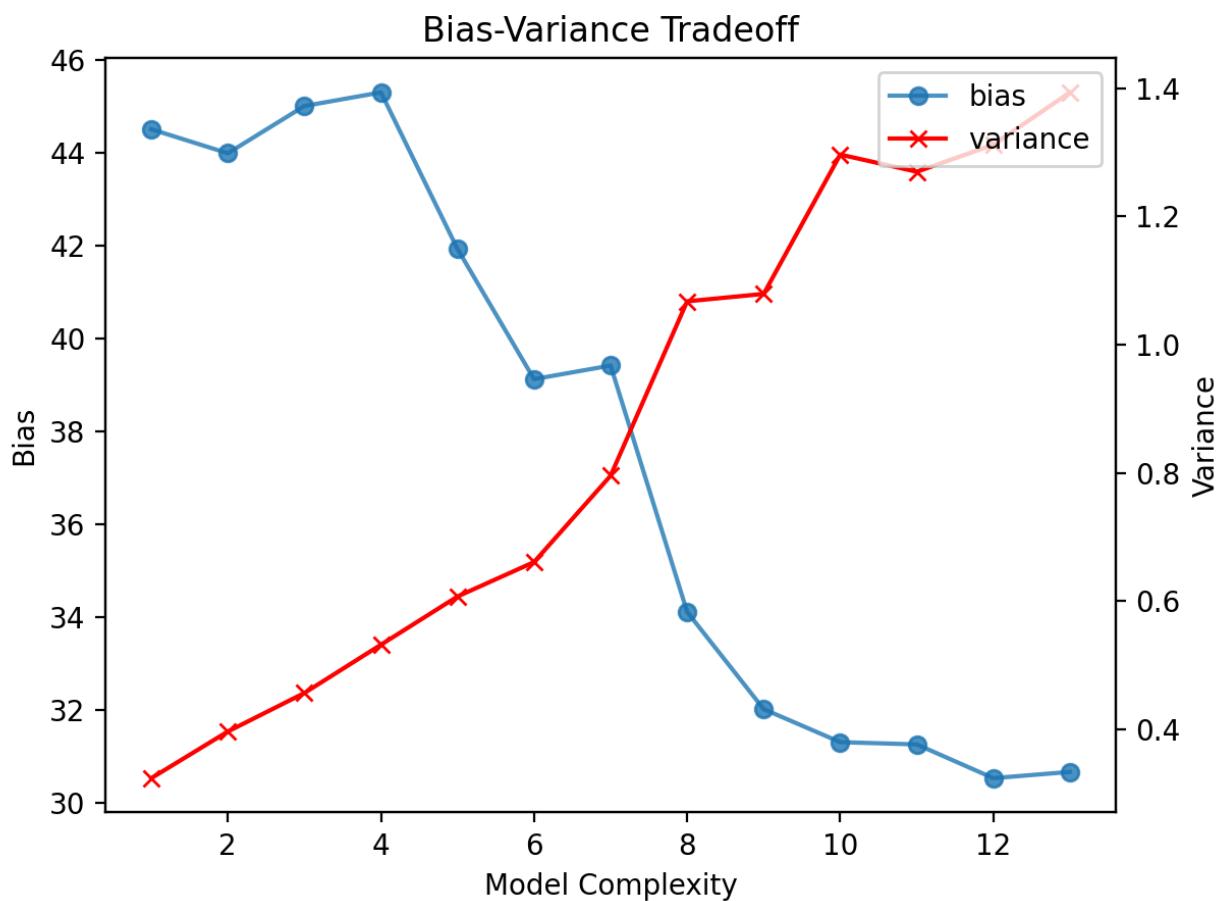
	0	1	2	3	4	5	6	7	8	9	10	11	12	13
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222.0	18.7	396.90	5.33	36.2

The last column is the y value we need to fit. So there are 13 parameters. The model complexity is increased by adding more features in the linear regression mode. After fitting, we have 13 models with the corresponding bias and variance.

The line plot of bias and variance that have the same y-axis shows below.



To have a better view, different axes are used. Then the different pattern of bias and variance are obvious. Because of the limitation of the real dataset and simple method to increase the model complexity, it has some fluctuations. But the overall pattern meets the assumptions in 1 and 2. Bias will decrease and variance will increase when the model complexity increases.



The code its attached at the end.

```
%matplotlib notebook
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Ridge
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import cross_val_score
from mlxtend.evaluate import bias_variance_decomp

# estimate the bias and variance for a regression model
from pandas import read_csv
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from mlxtend.evaluate import bias_variance_decomp
# load dataset
url =
'https://raw.githubusercontent.com/jbrownlee/Datasets/master/housing.csv'
```

```

dataframe = read_csv(url, header=None)
# separate into inputs and outputs
data = dataframe.values
y = data[:, -1]
biasList = []
varList = []
for i in range(12,0, -1):
    X= data[:, i:-1]
    X_train, X_test, y_train, y_test = train_test_split(X, y,
random_state=0)
    linreg = LinearRegression().fit(X_train, y_train)
    mse, bias, var = bias_variance_decomp(
        linreg, X_train, y_train, X_test, y_test, loss = 'mse')
    biasList.append(bias)
    varList.append(var)

x_num = [i for i in range(1,13)]
x_num

fig = plt.figure()
plt.plot(x_num, biasList, '-o', x_num, varList, '-xr')
plt.legend(['Bias', 'Variance'])

plt.show()

fig = plt.figure()
ax1 = fig.add_subplot(111)
ax1.plot(x_num, biasList, '-o', alpha=0.8, label = 'bias')
ax1.set_ylabel('Bias')
ax1.set_xlabel('Model Complexity')
ax1.set_title('Bias-Variance Tradeoff')

ax2 = ax1.twinx()
ax2.plot(x_num, varList, '-xr', label = 'variance')
ax2.set_ylabel('Variance')
fig.legend(loc=1, bbox_to_anchor=(1,1),
bbox_transform=ax1.transAxes)
plt.show()

```

Also, residuals vs. fitted values are plotted for further analysis. We could see that it has been more close to the mean 0, which meets our assumption.

