

worksheet_19

April 9, 2024

1 Worksheet 19

Name: Daniyal Ahmed UID: U11469883

1.0.1 Topics

- Linear Model Evaluation

1.1 Linear Model Evaluation

Notice that R^2 only increases with the number of explanatory variables used. Hence the need for an adjusted R^2 that penalizes for insignificant explanatory variables.

```
[75]: import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures

SAMPLE_SIZE = 100
beta = [1, 5]
X = -10.0 + 10.0 * np.random.random(SAMPLE_SIZE)
Y = beta[0] + beta[1] * X + np.random.randn(SAMPLE_SIZE)

for i in range(1, 15):
    X_transform = PolynomialFeatures(degree=i, include_bias=False).
    ↪fit_transform(X.reshape(-1, 1))
    model = LinearRegression()
    model.fit(X_transform, Y)
    print(model.score(X_transform, Y))
```

```
0.9961938549419592
0.996238779946289
0.9962505068863905
0.9962834227607091
0.9962925331439183
0.9962935740928418
0.9963142839344125
0.9963544002855594
0.9963544068632645
0.9963851887600592
```

0.9964022179422242
0.9964125231063834
0.99641272334507
0.9964035706398343

a) Hypothesis Testing Sandbox (follow along in class) [Notes](#)

```
[76]: import numpy as np
      from scipy.stats import binom
      import matplotlib.pyplot as plt

      flips = [1, 0, 0, 1, 0]

      def num_successes(flips):
          return sum(flips)

      print(binom.pmf(num_successes(flips), len(flips), 1/2))

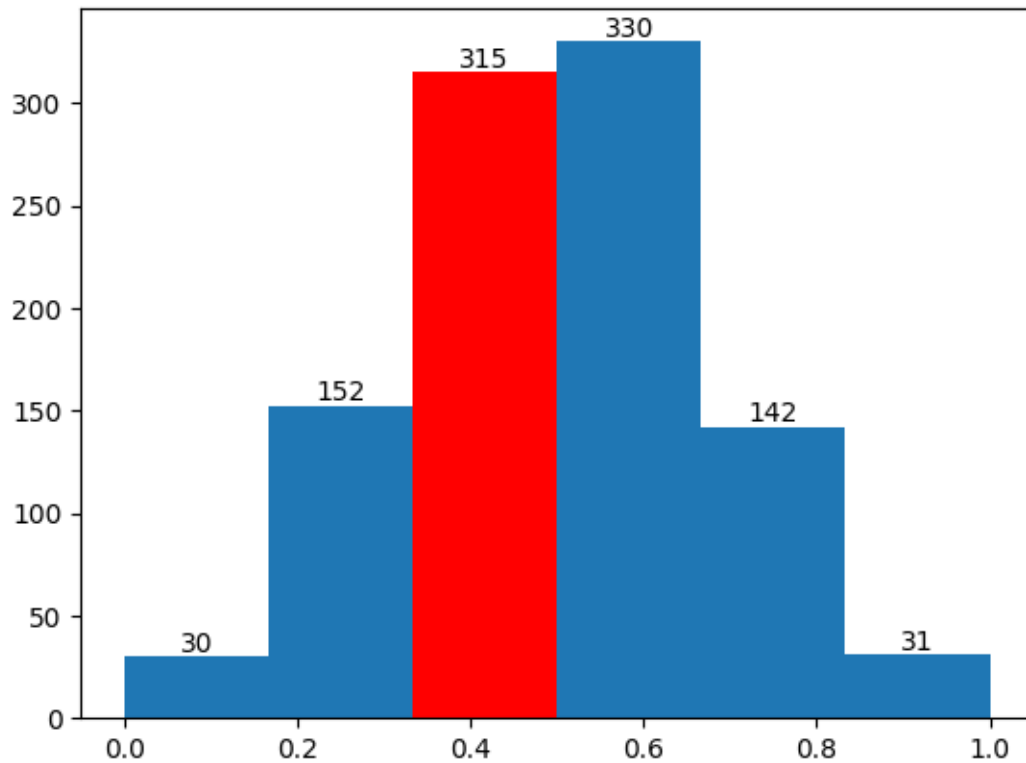
      SAMPLE_SIZE = 5
      flips = [np.random.choice([0, 1]) for _ in range(SAMPLE_SIZE)]
      print(flips)
      print(binom.pmf(num_successes(flips), SAMPLE_SIZE, 1/2))

      p_est = []

      for _ in range(1000):
          flips = [np.random.choice([0, 1]) for _ in range(SAMPLE_SIZE)]
          p_est.append(sum(flips) / SAMPLE_SIZE)

      fig, ax = plt.subplots()
      _, bins, patches = ax.hist(p_est, bins=SAMPLE_SIZE + 1)
      p = np.digitize([2/5], bins)
      patches[p[0]-1].set_facecolor('r')
      ax.bar_label(patches)
      plt.show()
```

0.31249999999999983
[0, 1, 1, 1, 1]
0.15625



b) Plot a data set and fitted line through the point when there is no relationship between X and y.

```
[77]: import numpy as np
import matplotlib.pyplot as plt

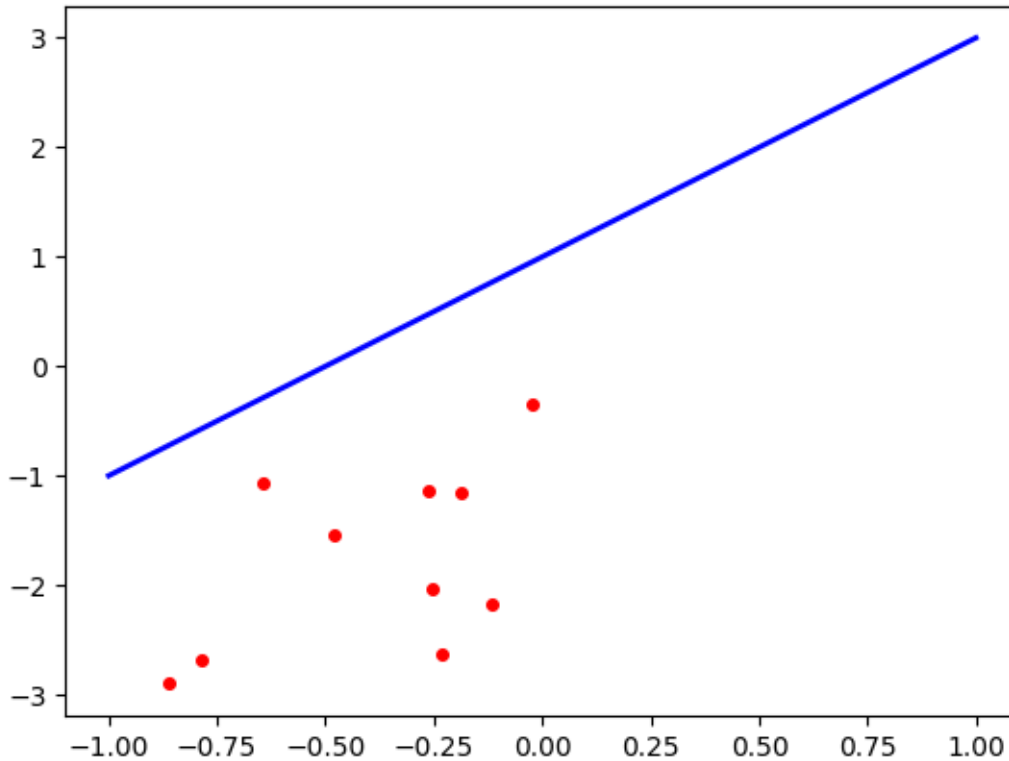
SAMPLE_SIZE = 10

xlin = -1.0 + 1.0 * np.random.random(SAMPLE_SIZE)
y = -1.0 + 1.0 * xlin + np.random.randn(SAMPLE_SIZE)

intercept = np.ones(np.shape(xlin)[0])
X = np.array([intercept, xlin]).T
beta = [1,2]

print(intercept)
xplot = np.linspace(-1,1,20)
yestplot = beta[0] + beta[1] * xplot
plt.plot(xplot, yestplot, 'b-', lw=2)
plt.plot(xlin, y, 'ro', markersize=4)
plt.show()
```

```
[1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
```



- c) Using the above code, plot a histogram of the parameter estimates for the slope after generating 1000 independent datasets. Comment on what the plot means. Increase the sample size to see what happens to the plot. Explain.

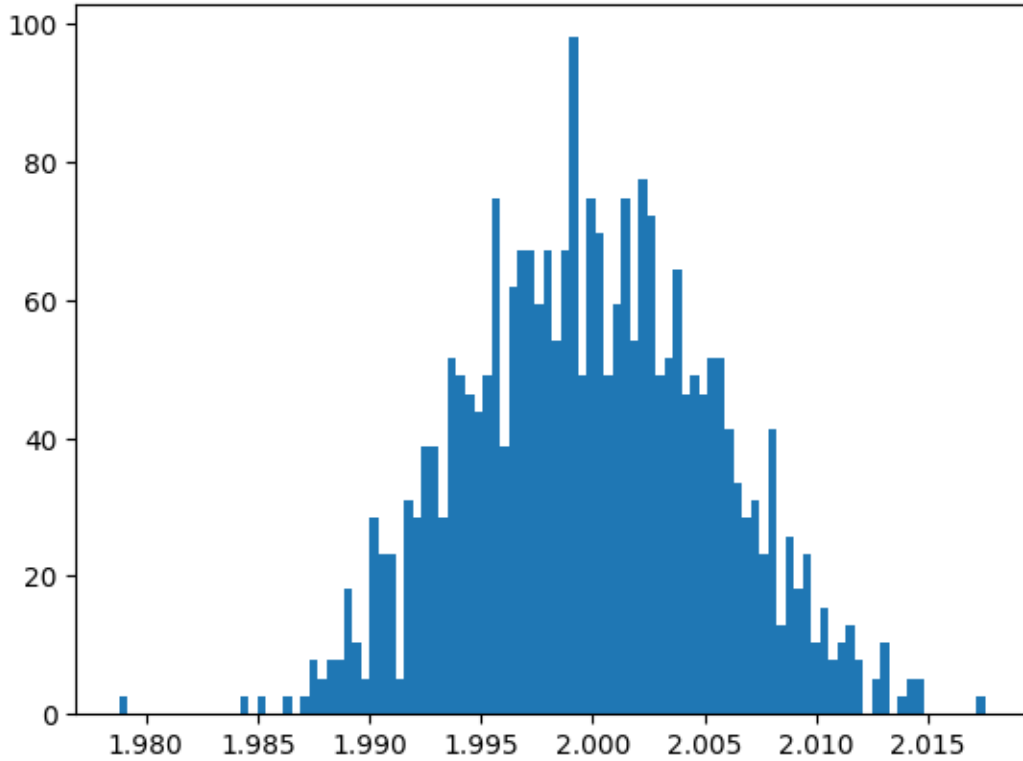
```
[143]: beta_hist = []
SAMPLE_SIZE = 100000
for _ in range(1000):

    xlin = 1.0 + 2.0 * np.random.random(SAMPLE_SIZE)
    y = 1 + 2*xlin + np.random.randn(SAMPLE_SIZE)

    intercept = np.ones(np.shape(xlin)[0])
    X = np.array([intercept, xlin]).T

    beta_hist.append((np.linalg.inv(X.T @ X) @ X.T @ y)[-1])

fig, ax = plt.subplots()
ax.hist(beta_hist, bins=100, density=True)
plt.show()
```



The plot seems to be a normal distribution of points, what happens when we increase the sample size is that the distribution becomes more and more spread apart, which makes sense this we are adding more and more samples.

d) We know that:

$$\hat{\beta} - \beta \sim \mathcal{N}(0, \sigma^2 (X^T X)^{-1})$$

thus for each component k of $\hat{\beta}$ (here there are only two - one slope and one intercept)

$$\hat{\beta}_k - \beta_k \sim \mathcal{N}(0, \sigma^2 S_{kk})$$

where S_{kk} is the k^{th} diagonal element of $(X^T X)^{-1}$. Thus, we know that

$$z_k = \frac{\hat{\beta}_k - \beta_k}{\sqrt{\sigma^2 S_{kk}}} \sim \mathcal{N}(0, 1)$$

Verify that this is the case through a simulation and compare it to the standard normal pdf by plotting it on top of the histogram.

```
[103]: from scipy.stats import norm
```

```

beta_hist = []
for k in range(1000):
    xlin = 1.0 + 2.0 * np.random.random(SAMPLE_SIZE)
    y = 1 + 2*xlin + np.random.randn(SAMPLE_SIZE)
    betas = [1,2]

    intercept = np.ones(np.shape(xlin)[0])
    X = np.array([intercept, xlin]).T
    beta_hats = np.linalg.inv(X.T @ X) @ X.T @ y

    y_pred = beta_hats[0]+beta_hats[1]*xlin

    #RSS
    RSS = np.sum((y-y_pred)**2)

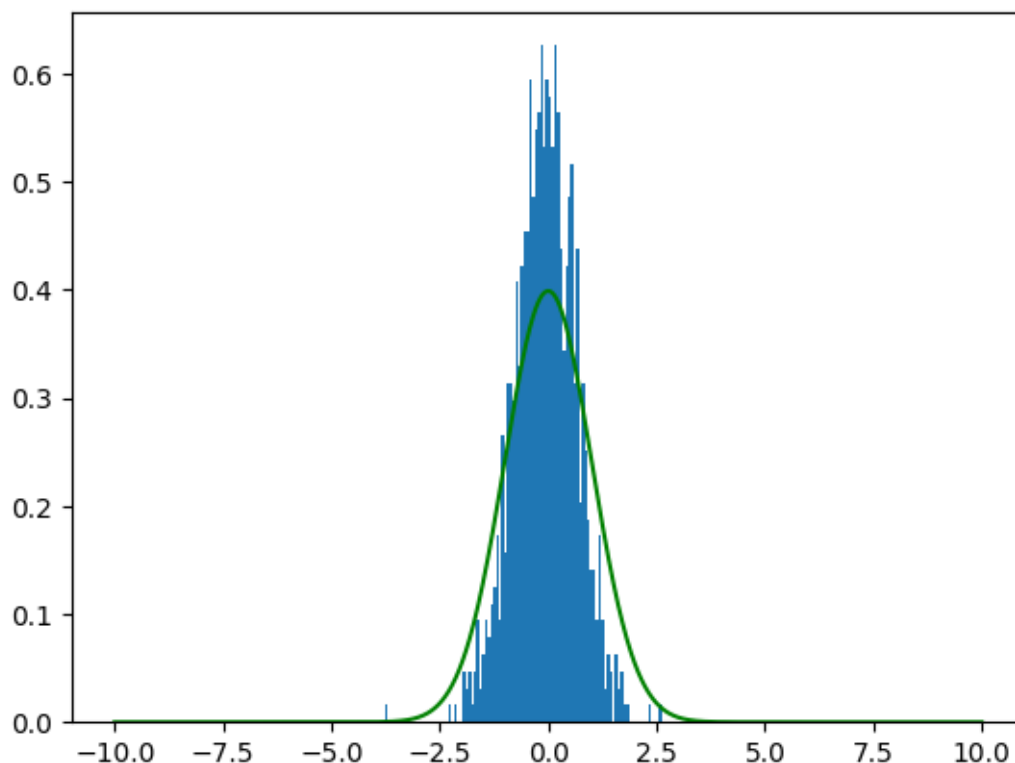
    standard_error = np.sqrt(RSS/(SAMPLE_SIZE-2))

    skk = np.diag((X.T @ X)**-1,k)
    result = beta_hats[1] -betas[1]/ np.sqrt(standard_error)

    beta_hist.append(result)

result = np.array(result)
xs = np.linspace(-10,10,1000)
fig, ax = plt.subplots()
ax.hist(beta_hist, bins=100, density=True)
ax.plot(xs, norm.pdf(xs, 0,1), color='green')
plt.show()

```



- e) Above we normalized $\hat{\beta}$ by subtracting the mean and dividing by the standard deviation. While we know that the estimate of beta is an unbiased estimator, we don't know the standard deviation. So in practice when doing a hypothesis test where we want to assume that $\beta = 0$, we can simply use $\hat{\beta}$ in the numerator. However we don't know the standard deviation and need to use an unbiased estimate of the standard deviation instead. This estimate is the standard error s

$$s = \sqrt{\frac{RSS}{n - p}}$$

where p is the number of parameters beta (here there are 2 - one slope and one intercept). This normalized $\hat{\beta}$ can be shown to follow a t-distribution with $n-p$ degrees of freedom. Verify this is the case with a simulation.

```
[126]: from scipy.stats import t

def standard_error(ytrue, ypred):

    #RSS
    RSS = np.sum((ytrue-ypred)**2)
```

```

standard_error = np.sqrt(RSS/(SAMPLE_SIZE-2))
return standard_error

beta_hist = []
for _ in range(1000):
    xlin = 0.0 + 2.0 * np.random.random(SAMPLE_SIZE)
    y = 0 + 2*xlin + np.random.randn(SAMPLE_SIZE)
    betas = [0,2]

    intercept = np.ones(np.shape(xlin)[0])
    X = np.array([intercept, xlin]).T
    beta_hats = np.linalg.inv(X.T @ X) @ X.T @ y
    y_pred = beta_hats[0]+beta_hats[1]*xlin

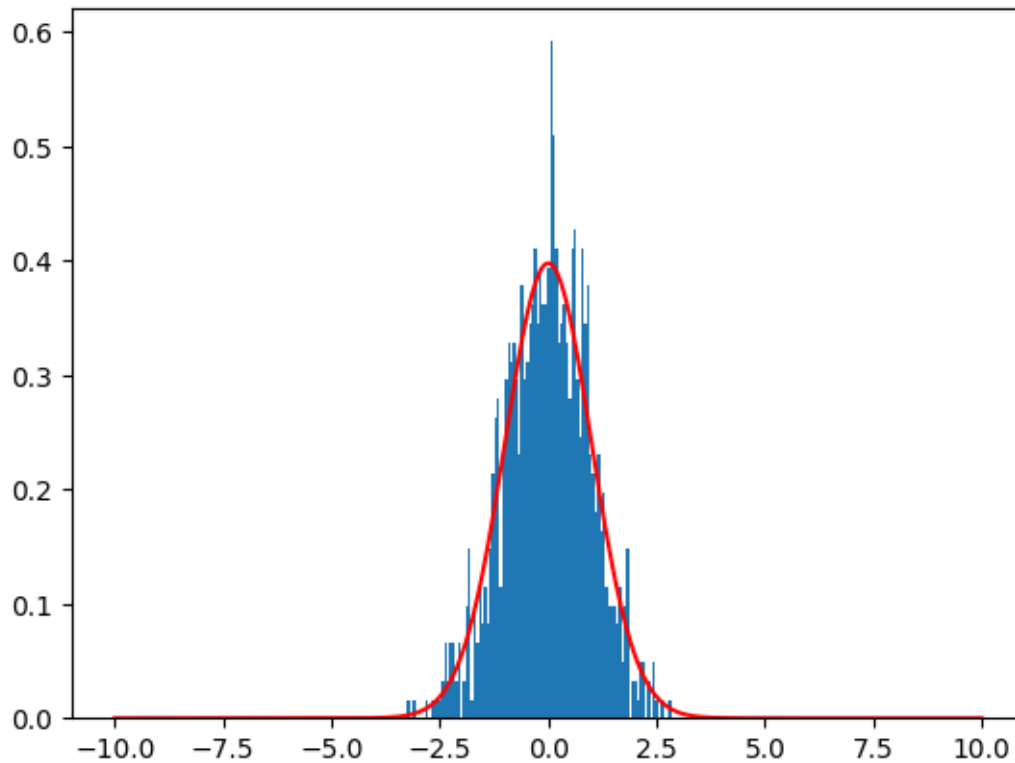
    se = standard_error(y,y_pred)

    #to show it can follow a t-distribution you have to use the T stat formula, ↵
    ↪I had to google this part, I had no idea what this was before
    #but I assumed that a T stat would get me the results I wanted, since we ↵
    ↪wanna follow a T distribution
    t_stat = (beta_hats[1] - 2) / (se / np.sqrt(np.var(xlin) * SAMPLE_SIZE))

    beta_hist.append(t_stat)

xs = np.linspace(-10,10,1000)
fig, ax = plt.subplots()
ax.hist(beta_hist, bins=100, density=True)
ax.plot(xs, t.pdf(xs, SAMPLE_SIZE - 2), color='red')
plt.show()

```

f) You are given the following dataset:

```
[133]: import numpy as np
import matplotlib.pyplot as plt

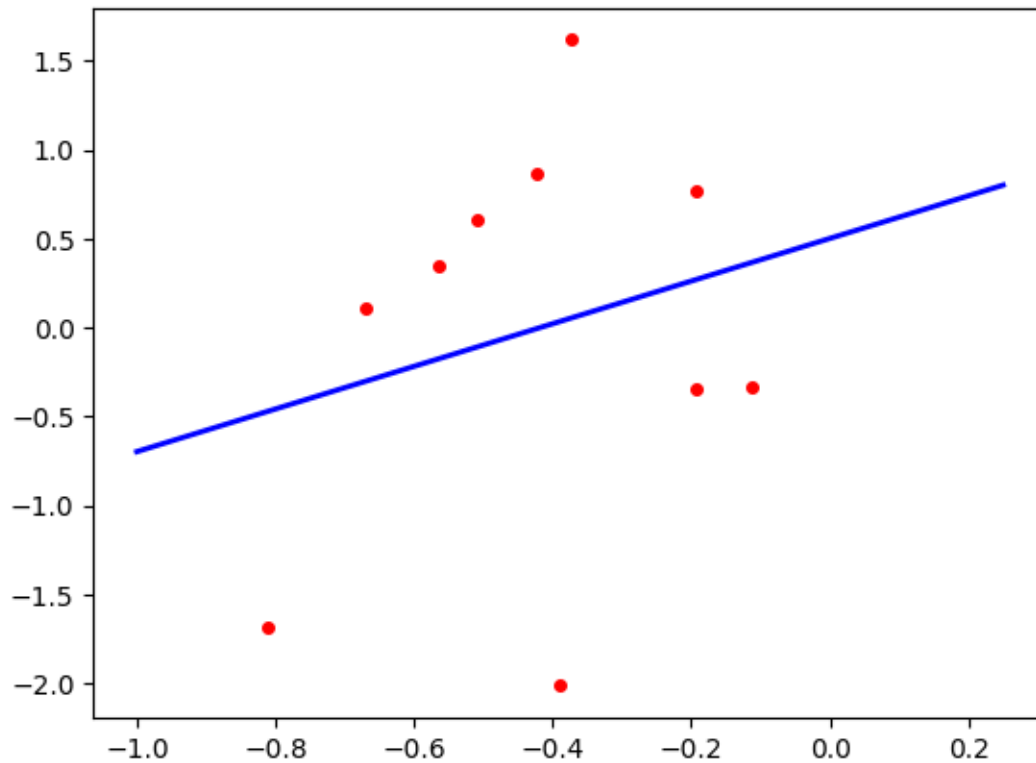
x = np.array([-0.1920605, -0.11290798, -0.56434374, -0.67052057, -0.19233284,
↪ -0.42403586, -0.8114285, -0.38986946, -0.37384161, -0.50930229])
y = np.array([-0.34063108, -0.33409286, 0.34245857, 0.11062295, 0.76682389, 0.
↪ 86592388, -1.68912015, -2.01463592, 1.61798563, 0.60557414])

intercept = np.ones(np.shape(x)[0])
X = np.array([intercept, x]).T
beta_hat = np.linalg.inv(X.T @ X) @ X.T @ y

print(beta_hat)

xplot = np.linspace(-1, .25, 20)
yestplot = beta_hat[0] + beta_hat[1] * xplot
plt.plot(xplot, yestplot, 'b-', lw=2)
plt.plot(x, y, 'ro', markersize=4)
plt.show()
```

```
[0.50155603 1.19902827]
```



what is the probability of observing a dataset at least as extreme as the above assuming $\beta = 0$?

[138]: *#Confidence Interval for the slope*

```
y_mean = np.mean(y)
```

```
se = standard_error(y, y_mean)
```

```
#FROM THE SLIDES
```

```
lowerbound = y_mean - 1.96 * se
```

```
upperbound = y_mean + 1.96 * se
```

```
print(lowerbound, upperbound)
```

```
#The probability is about 1/100 since the lower and upperbound are 100 units  
↪ apart, THIS IS CALCULATED WITH 95% CONFIDENCE INTERVAL
```

```
-0.6809770637074094 0.6671588737074093
```

[]:

[]: