Name: Jeong Yong Yang, Junho Son UID: U95912941, U64222022

## Topics

- Linear Model Evaluation

## 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.

```
1   import numpy as np
2   from sklearn.linear_model import LinearRegression
3   from sklearn.preprocessing import PolynomialFeatures
4
5   SAMPLE_SIZE = 100
6   beta = [1, 5]
7   X = -10.0 + 10.0 * np.random.random(SAMPLE_SIZE)
8   Y = beta[0] + beta[1] * X + np.random.randn(SAMPLE_SIZE)
9
10  for i in range(1, 15):
11      X_transform = PolynomialFeatures(degree=i, include_bias=False).fit_transform(X.reshape(-1, 1))
12      model = LinearRegression()
13      model.fit(X_transform, Y)
14      print(model.score(X_transform, Y))
```
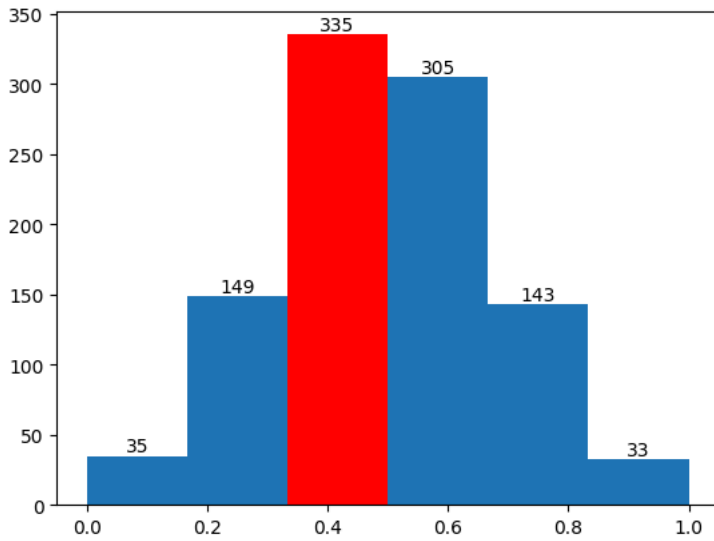
```
0.9954062656470016
0.9954159177864946
0.9955103152871462
0.9955994798701692
0.9957012434751283
0.9957335582161672
0.9958795013561041
0.9961007167628202
0.9961144218567842
0.9961354512007631
0.9961497434160319
0.9961650184057208
0.9961670316961133
0.9961347497238493
```

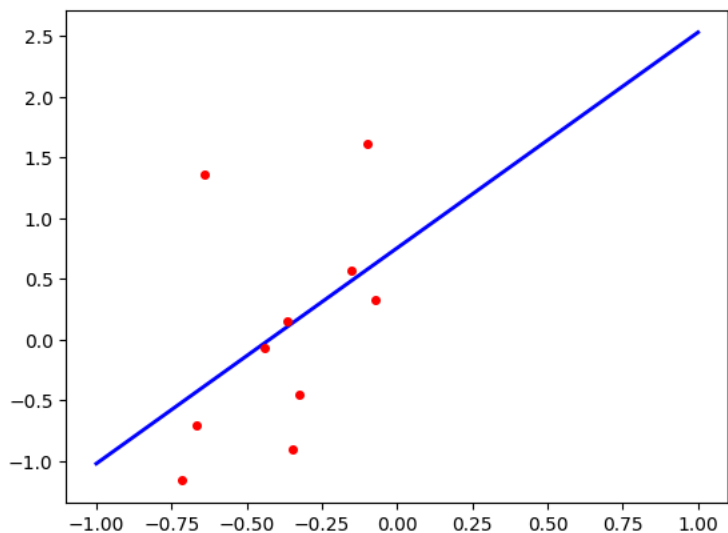a) Hypothesis Testing Sandbox (follow along in class) [Notes](#)

```
1 import numpy as np
2 from scipy.stats import binom
3 import matplotlib.pyplot as plt
4
5 flips = [1, 0, 0, 1, 0]
6
7 def num_successes(flips):
8     return sum(flips)
9
10 print(binom.pmf(num_successes(flips), len(flips), 1/2))
11
12 SAMPLE_SIZE = 5
13 flips = [np.random.choice([0, 1]) for _ in range(SAMPLE_SIZE)]
14 print(flips)
15 print(binom.pmf(num_successes(flips), SAMPLE_SIZE, 1/2))
16
17 p_est = []
18
19 for _ in range(1000):
20     flips = [np.random.choice([0, 1]) for _ in range(SAMPLE_SIZE)]
21     p_est.append(sum(flips) / SAMPLE_SIZE)
22
23 fig, ax = plt.subplots()
24 _, bins, patches = ax.hist(p_est, bins=SAMPLE_SIZE + 1)
25 p = np.digitize([2/5], bins)
26 patches[p[0]-1].set_facecolor('r')
27 ax.bar_label(patches)
28 plt.show()
```

```
0.31249999999999983
[0, 0, 0, 1, 1]
0.31249999999999983
```



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

```
 1 import numpy as np
 2 import matplotlib.pyplot as plt
 3
 4 SAMPLE_SIZE = 10
 5
 6 xlin = -1.0 + 1.0 * np.random.random(SAMPLE_SIZE)
 7 y = 0.0 + np.random.randn(SAMPLE_SIZE)
 8
 9 intercept = np.ones(np.shape(xlin)[0])
10 X = np.array([intercept, xlin]).T
11 beta = np.linalg.inv(X.T @ X) @ X.T @ y
12
13 xplot = np.linspace(-1,1,20)
14 yestplot = beta[0] + beta[1] * xplot
15 plt.plot(xplot, yestplot,'b-',lw=2)
16 plt.plot(xlin, y,'ro',markersize=4)
17 plt.show()
```
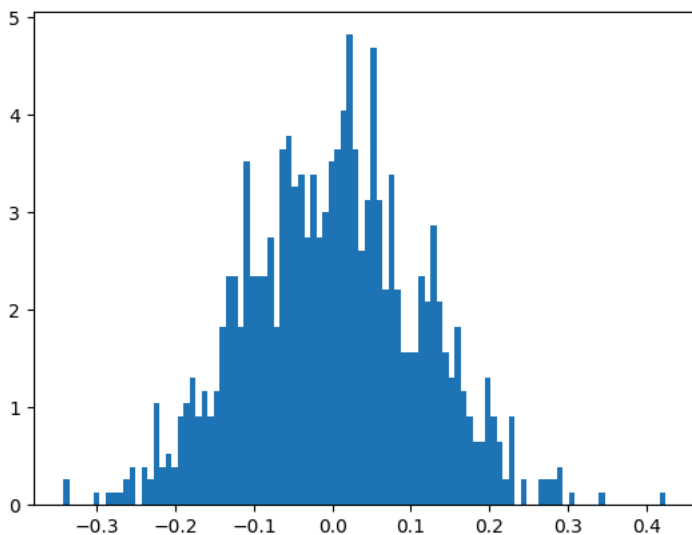


c) Using the above code, plot a histogram of the parameter estimates for the slope after generating $1000$ independent datasets. Comment on what the plot means. Increase the sample size to see what happens to the plot. Explain.

```
 1 SAMPLE_SIZE = 1000
 2 NUM_TRIALS = 1000
 3 beta_hist = []
 4 for _ in range(NUM_TRIALS):
 5     xlin = -1.0 + 1.0 * np.random.random(SAMPLE_SIZE)
 6     y = 0.0 + np.random.randn(SAMPLE_SIZE)
 7
 8     intercept = np.ones(np.shape(xlin)[0])
 9     X = np.array([intercept, xlin]).T
10     beta = np.linalg.inv(X.T @ X) @ X.T @ y
11     beta_hist.append(beta[1])
12
13 fig, ax = plt.subplots()
14 ax.hist(beta_hist, bins=100, density=True)
15 plt.show()
```



For a sample size of 100, most of the beta is 0, which confirms our estimation but there are some betas that are higher or lower than 0, indicating that there is a relationship between x and y. As the sample size decreases, the values of slopes get much steeper. As the sample size increases, the values of slopes get much flatter.

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 $k^{\text{th}}$ diagonal element of $(X^T X)^{-1}$. Thus, we know that
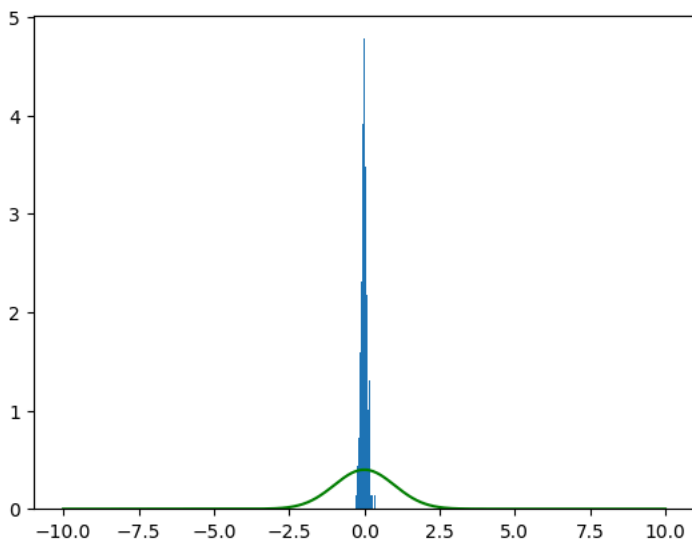
$$z_k = \frac{\hat{\beta}_k - \beta_k}{\sqrt{\sigma^2 S_{kk}}} \sim \mathcal{N}(0, 1)$$

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.

```
 1 from scipy.stats import norm
 2
 3 beta_hist = []
 4 for _ in range(1000):
 5     xlin = -1.0 + 1.0 * np.random.random(SAMPLE_SIZE)
 6     y = 0.0 + np.random.randn(SAMPLE_SIZE)
 7
 8     intercept = np.ones(np.shape(xlin)[0])
 9     X = np.array([intercept, xlin]).T
10     beta = np.linalg.inv(X.T @ X) @ X.T @ y
11     beta_hist.append(beta[1])
12
13 xs = np.linspace(-10,10,1000)
14 fig, ax = plt.subplots()
15 ax.hist(beta_hist, bins=100, density=True)
16 ax.plot(xs, norm.pdf(xs), color='green')
17 plt.show()
```
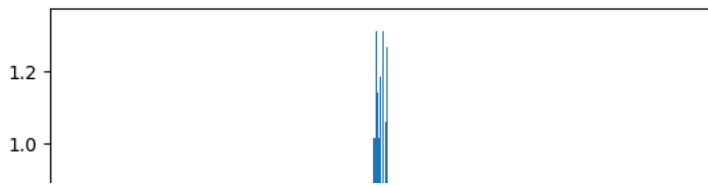
e) Above we normalized $\hat{\beta}$ by subtracting the mean and dividing by the standard deviation. While we know that the estimate of beta is an unbiased estimator, we don't know the standard deviation. So in practice when doing a hypothesis test where we want to assume that $\beta = 0$, we can simply 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.
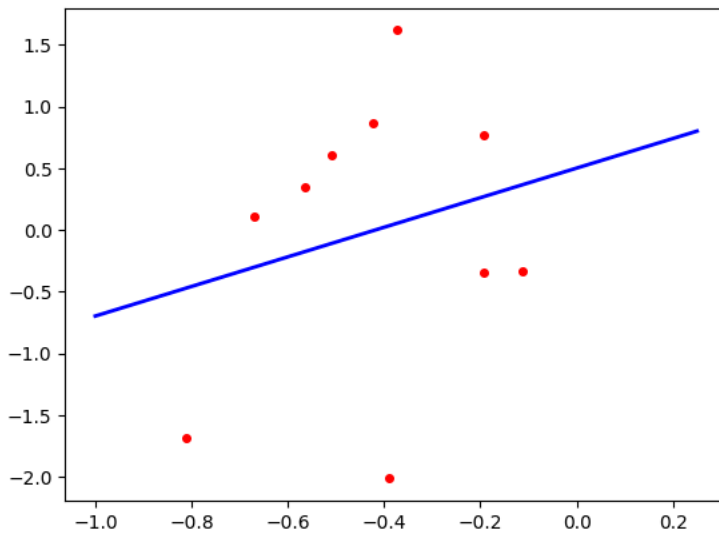
```
1 from scipy.stats import t
2
3 SAMPLE_SIZE = 100
4 NUM_TRIALS = 1000
5
6 def standard_error(ytrue, ypred):
7     length = len(ytrue)
8     residuals = (ytrue - ypred) ** 2
9     rss = 0
10    for x in residuals:
11        rss += x
12    return np.sqrt(rss / (length - 2))
13
14 beta_hist = []
15
16 for _ in range(NUM_TRIALS):
17     xlin = -1.0 + 1.0 * np.random.random(SAMPLE_SIZE)
18     y = 0.0 + np.random.randn(SAMPLE_SIZE)
19
20     intercept = np.ones(np.shape(xlin)[0])
21     X = np.array((intercept, xlin)).T
22     beta = np.linalg.inv(X.T @ X) @ X.T @ y
23     standardError = standard_error(y, X @ beta)
24     beta_hist.append(beta[1])
25
26 xs = np.linspace(-10,10,1000)
27 fig, ax = plt.subplots()
28 ax.hist(beta_hist, bins=100, density=True)
29 ax.plot(xs, t.pdf(xs, SAMPLE_SIZE - 2), color='red')
30 plt.show()
```

f) You are given the following dataset:



```
 1 import numpy as np
 2 import matplotlib.pyplot as plt
 3
 4 x = np.array([-0.1920605, -0.11290798, -0.56434374, -0.67052057, -0.19233284, -0.42403586, -0.8114285, -0.38986946, -0.37384161, -0.509302
 5 y = np.array([-0.34063108, -0.33409286, 0.34245857, 0.11062295, 0.76682389, 0.86592388, -1.68912015, -2.01463592, 1.61798563, 0.60557414])
 6
 7 intercept = np.ones(np.shape(x)[0])
 8 X = np.array([intercept, x]).T
 9 beta_hat = np.linalg.inv(X.T @ X) @ X.T @ y
10
11 print(beta_hat)
12
13 xplot = np.linspace(-1,.25,20)
14 yestplot = beta_hat[0] + beta_hat[1] * xplot
15 plt.plot(xplot, yestplot,'b-',lw=2)
16 plt.plot(x, y,'ro',markersize=4)
17 plt.show()
18
```

```
[0.50155603 1.19902827]
```



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

```
 1 length = len(x)
 2
 3 ypred = X @ beta_hat
 4 standardError = standard_error(y, ypred)
 5
 6 tScore = beta_hat[1] / standardError
 7
 8 p_value = 2 * (1 - t.cdf(np.abs(tScore), df = length - 2))
 9
10 print(f"The probability of observing a dataset at least as extreme as the above assuming beta = 0 is {p_value}")
```

```
    The probability of observing a dataset at least as extreme as the above assuming beta = 0 is 0.335417218459404
```

```
 1    코딩을 시작하거나 AI로 코드를 생성하세요.
```