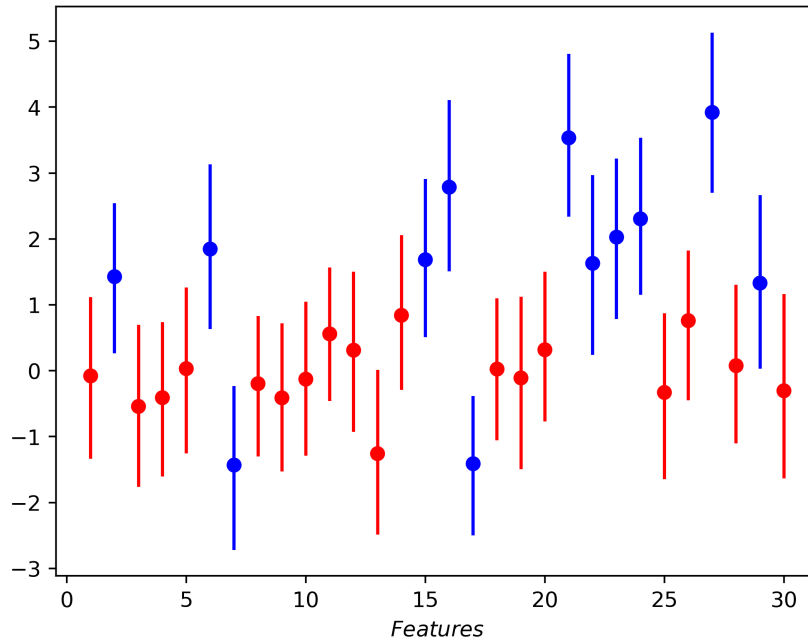


Question 1

a



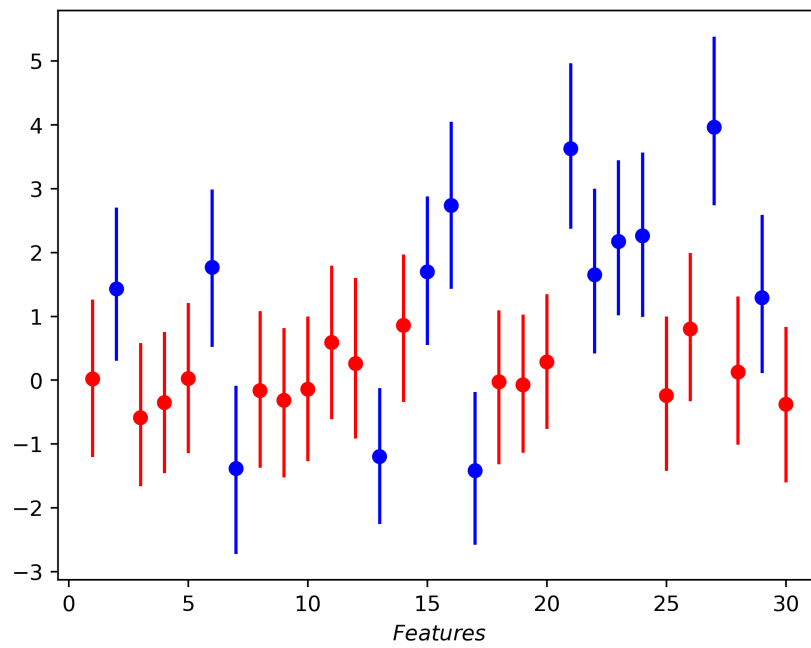
Screenshot of code here: 0.1.

i. C is a hyperparameter used for regularization to prevent overfitting of the Logistics Regression model. This is inversely proportional to the penalty constant λ and penalises models that have a lot of features and therefore the effect is that it reduces each feature importance.

ii. The effect that this has on the Bootstrap graph is that it increases the variance of each feature, therefore making the 90% confidence intervals larger and more reliable since their largeness means they are more likely to include the true value.

At $C = 0.1$, most of the features have their average at 0 with a very small confidence interval. Because I know how the data was generated and I know that the mean for each feature *should* be 0, I know that this is a reliable estimate. However in the real world when I don't know how the data was generated, a model with $C = 0.1$ would probably be overfit on the data and not contain a true mean.

b



Screenshot of code here: 0.2.

Appendix

0.1 q1a

```
23 def q1a():      You, 2 hours ago • initial
24     # Import data
25     data = import_data(Q1_DATA_DIR)
26
27     Xtrain = data.iloc[:, :30].to_numpy()
28     Ytrain = data.Y.to_numpy()
29
30     # Set random seed
31     np.random.seed(12)
32
33     # Set B and C and p
34     B = 500
35     C = 1000
36     p = Xtrain.shape[1]
37
38     # Find confidence
39     coef_mat = np.zeros(shape=(B,p))
40     for b in range(B):
41         b_sample = np.random.choice(np.arange(Xtrain.shape[0]), size=Xtrain.shape[0])
42         Xtrain_b = Xtrain[b_sample]
43         Ytrain_b = Ytrain[b_sample]
44         mod = LogisticRegression(penalty="l1", solver="liblinear", C=C).fit(Xtrain_b, Ytrain_b)
45         coef_mat[b,:] = mod.coef_
46
47     means = np.mean(coef_mat, axis=0)
48     lower = np.quantile(coef_mat, 0.10, axis=0)
49     upper = np.quantile(coef_mat, 0.90, axis=0)
50
51     colors = ["red" if lower[i] <= 0 and upper[i] >= 0 else "blue" for i in range(p)]
52
53     plt.vlines(x=np.arange(1,p+1), ymin=lower, ymax=upper, colors=colors)
54     plt.scatter(x=np.arange(1,p+1), y=means, color=colors)
55     plt.xlabel("$Features$")
56     plt.savefig("../outputs/NPBootstrap.png", dpi=400)
```

0.2 q1b

```
58 def q1b():
59     # Import data
60     data = import_data(Q1_DATA_DIR)
61
62     Xtrain = data.iloc[:, :30].to_numpy()
63     Ytrain = data.Y.to_numpy()
64
65     # Set random seed
66     np.random.seed(20)
67
68     # Set B and C and p
69     B = 500
70     C = 1000
71     p = Xtrain.shape[1]
72
73     # Train initial model
74     mod = LogisticRegression(penalty="l1", solver="liblinear", C=C).fit(Xtrain, Ytrain)
75     B0 = mod.intercept_[0]
76     coefs = mod.coef_[0]
77
78     # Find confidence
79     coef_mat = np.zeros(shape=(B,p))
80     for b in range(B):
81         # Get X values randomly as before. Here, n = Xtrain.shape[0].
82         b_sample = np.random.choice(np.arange(Xtrain.shape[0]), size=Xtrain.shape[0])
83         Xtrain_b = Xtrain[b_sample]
84
85         Ytrain_b = np.zeros(shape=(1,Xtrain.shape[0]))
86         idx = 0
87         for i in Xtrain_b:
88             # Calculate "p" for bernoulli for this instance
89             prob = math.exp(B0 + coefs.T @ i) / (1 + math.exp(B0 + coefs.T @ i))
90
91             prob = math.exp(B0 + coefs.T @ i) / (1 + math.exp(B0 + coefs.T @ i))
92
93             # The Ytrain_b must come from a bernoulli distribution
94             response = np.random.binomial(n=1, p=prob, size=1)
95
96             # Add to ytrain data
97             Ytrain_b[0, idx] = response[0]
98             idx += 1
99
100         # Train new model fit to the b_set
101         mod_b = LogisticRegression(penalty="l1", solver="liblinear", C=C).fit(Xtrain_b, Ytrain_b[0])
102
103         # store coefficients
104         coef_mat[b,:] = mod_b.coef_
105
106     means = np.mean(coef_mat, axis=0)
107     lower = np.quantile(coef_mat, 0.10, axis=0)
108     upper = np.quantile(coef_mat, 0.90, axis=0)
109
110     colors = ["red" if lower[i] <= 0 and upper[i] >= 0 else "blue" for i in range(p)]
111
112     plt.vlines(x=np.arange(1,p+1), ymin=lower, ymax=upper, colors=colors)
113     plt.scatter(x=np.arange(1,p+1), y=means, color=colors)
114     plt.xlabel("$Features$")
115     plt.savefig("../outputs/NPParameterisedBootstrap.png", dpi=400)
```