# COMP9417: Homework Set #2

### z5113817

University of New South Wales — July 18, 2021

All code for this homework set is available here.
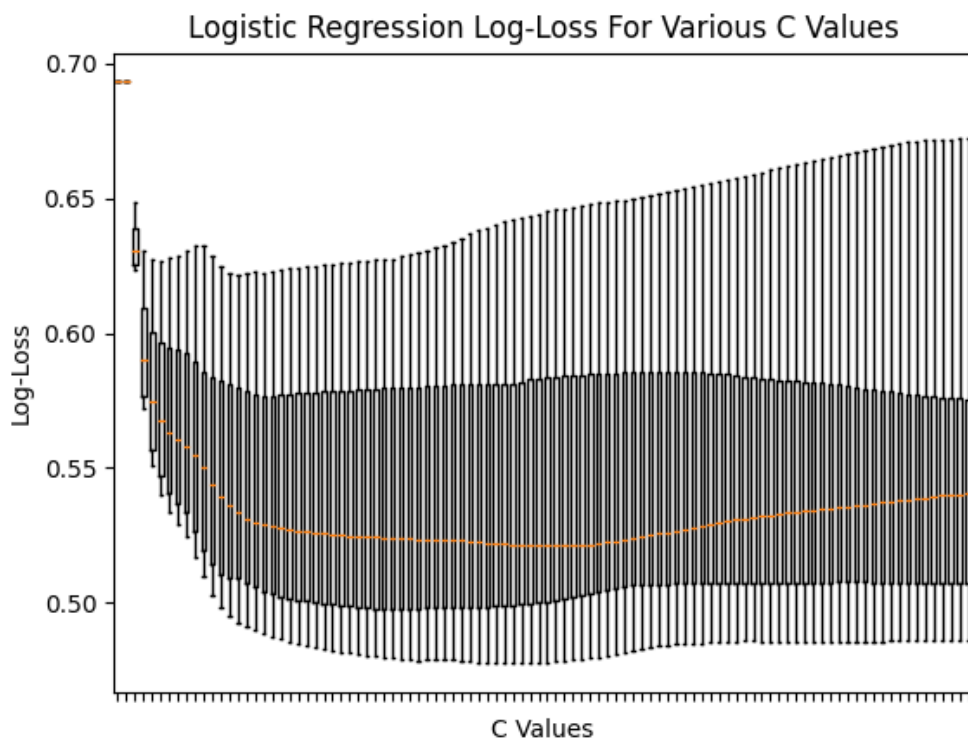
# Question 1

**a**

The possible values for both $y_i$ and $\tilde{y}_i$ are binary. Even though they have different values ($y_i \in \{0,1\}$ and $\tilde{y}_i \in \{-1,1\}$), the objective of each logistic regression implementation is to divide the dataset into 2 classifications. Because of this, the actual value that each classification has will not affect the parameters that the regression is attempting to optimise (($\hat{\beta}_0, \hat{\beta}$) and ($\hat{w}, \hat{c}$)). Therefore the solutions for the parameters being minimised by each regression will be the same.

$C$ is a hyper-parameter that adjusts the sensitivity that the model has to its coefficients. Compared with the standard LASSO parameter $\lambda$, $C$ is a multiple of the Loss function whereas $\lambda$ is a multiple of the Penalty.

**b**

Boxplot of testing accuracy for each value of C:



The value of C returning best results: **0.18794747474747472**
The testing accuracy of this model: **76%**

**c**

**From GridSearchCV:**
The value of C returning best results: **0.012219191919191918**
The testing accuracy of this model: **75.2%**
In our answer for b , we determined the "best" value of C as the value which corresponded to the average lowest *log-loss* value across all folds. The value from the *GridSearchCV* are different because, by default, it determines the "best" value of C as the one which corresponds to the average highest *score* across all folds [1].

We can modify the *GridSearchCV* class by providing our own metric for *scoring*. The following code is a scorer that uses the smallest *log-loss* value as its scoring metric.

```
scoring = make_scorer(
    log_loss,                # The sklearn implementation of log_loss
    greater_is_better=False, # A smaller log_loss is a better value
    needs_proba=True         # Calculating log_loss needs the probability predictions
)
```

This yields the following results:
*grid_lr.best_estimator_*: LogisticRegression(C=0.18188787878787877, penalty='l1', solver='liblinear')
*grid_lr.best_score_* (lowest log-loss): -0.5374067706086373
This value of *C* matches my value in b .

---

[1] See *scoring* parameter in documentation. If the estimator provided exposes a *score* method and a value for *scoring* is not provided, then *score* is used to determine the "best" value of C

# Code For Questions

## q1b

```python
def q1b():
    # Import data and take first 500 rows
    data = import_data(Q1_DATA_DIR)
    data = data.head(500)

    # Create folds
    folds = []
    fold_size = 50
    for i in range(0,10):
        folds.append(data[i*fold_size:i*fold_size+fold_size])

    # Create grid of 100 C values
    C_grid = linspace(0.0001, 0.6, 100)

    # Iterate over C values and train Logistic Model
    scores = []
    for C in C_grid.tolist():
        # Create model
        model = LogisticRegression(penalty="l1", solver="liblinear", C=C)

        # Iterate over each fold to train and test on each
        fold_scores = []
        for i in range(0, len(folds)):
            # Set up training and test data
            folds_train = folds.copy()  # Create a temp so that we don't mutate original 'folds'
            test_df = folds_train.pop(i)
            train_df = concat(folds_train)

            # Fit model to data
            clf = model.fit(train_df.drop("Y", axis=1), train_df["Y"])

            # Score data
            clf_probs = clf.predict_proba(test_df.drop("Y", axis=1))
            fold_score = log_loss(test_df["Y"], clf_probs)
            fold_scores.append(fold_score)

        scores.append(fold_scores)
```

```
58
59        # Save boxplot
60        fig, ax = plt.subplots()
61        ax.boxplot(scores)
62        ax.set_xticklabels("")        You, 21 hours ago • Finish q1c
63        ax.set_title("Logistic Regression Log-Loss For Various C Values")
64        ax.set_xlabel("C Values")
65        ax.set_ylabel("Log-Loss")
66        plt.savefig("./outputs/q1b_boxplot.png")
67
68        # Record the value of C that returns the "best" CV score
69        #
70        # "Best" is taken to be the lowest average log-loss
71        averages = list(map(lambda x : mean(x), scores))
72
73        # Get index of lowest average
74        i = averages.index(min(averages))
75
76        # Map this to a C value
77        C_best = C_grid[i]
78        print(C_best)
79
80        # Retrain this model and report its accuracy
81        model = LogisticRegression(penalty="l1", solver="liblinear", C=C_best)
82
83        # Let's test on the first fold and train on the remainder
84        test_df = folds.pop(0)
85        train_df = concat(folds)
86        clf = model.fit(train_df.drop("Y", axis=1), train_df["Y"])
87
88        print(clf.score(test_df.drop("Y", axis=1), test_df["Y"]))
```

## q1c

```
91    def q1c():
92        # Import data
93        data = import_data(Q1_DATA_DIR)
94        data = data.head(500)
95        Xtrain = data.drop("Y", axis=1)
96        Ytrain = data["Y"]
97
98        # Create grid of 100 C values
99        C_grid = linspace(0.0001, 0.6, 100)
100       param_grid = { "C": C_grid }
101
102       # Set our own scoring metric
103       scoring = make_scorer(
104           log_loss,
105           greater_is_better=False,
106           needs_proba=True
107       )
108
109       # Assignment code
110       grid_lr = GridSearchCV(
111           estimator=LogisticRegression(penalty='l1', solver='liblinear'),
112           cv=10,
113           param_grid=param_grid,
114           scoring=scoring
115       )
116       grid_lr.fit(Xtrain, Ytrain)
117       print(grid_lr.best_estimator_)
118       print(grid_lr.best_score_)
```