## Question 1:

A situation that would benefit from a classification model is, "Should I mow my lawn or not?" Predictors for this model would include the height of the grass, the probability of rain, time of day, my energy level, and air temperature.

## Question 2.1:

| Value of C | Training Error | Applications | | Error Interpretation |
| --- | --- | --- | --- | --- |
| | | Positive (exp. 354) | Negative (exp. 300) | |
| 0.0005 | 35.02% | 73 | 581 | Much worse than default |
| 0.005 | 13.61% | 351 | 303 | Same as default |
| 0.05 | 13.61% | 351 | 303 | Same as default |
| 0.5 | 13.61% | 351 | 303 | Same as default |
| Default of 1 | 13.61% | 351 | 303 | Default |
| 10 | 13.61% | 351 | 303 | Same as default |
| 100 | 13.61% | 351 | 303 | Same as default |
| 1000 | 13.76% | 350 | 304 | Slightly worse than default |
| 10000 | 13.76% | 350 | 304 | Slightly worse than default |

Model ←
-0.0010065348*A1 -0.0011729048*A2 -0.0016261967*A3 + 0.0030064203*A8 +
1.0049405641*A9 -0.0028259432*A10 + 0.0002600295*A11 -0.0005349551*A12
-0.0012283758*A14 + 0.1063633995*A15 + 0.08158492

Given that we're analyzing whether to extend credit to a person; I assumed that the potential
loss of someone defaulting was greater than the marginal revenue of extending them credit.
With this assumption, I think that since the training error is the same for values of C between
0.005 and 100, we should place more importance on not extending credit to a potential
defaulter.  As a result, I think that C = 100 is the optimal choice.

*R-script is also attached and named question-2.1.R*

```
4   # Load credit card data and convert to matrix
5   data_df <- read.table("credit_card_data-headers.txt", header=TRUE)
6   data_mx <- data.matrix(data_df)
7
8   # Train the model
9   model <- ksvm(
10    data_mx[,1:10],          # predictors
11    data_mx[,11],            # response
12    kernel="vanilladot",     # kernel selector
13    type="C-svc",            # ksvm type
14    C=100,                   # cost of constraints violation
15    scaled = TRUE            # scale flag
16    )
17
18  # Print model error
19  model@error
20
21  # Score data_mx against the model
22  results <- predict(model, data_mx[,1:10], type='response')
23
24  # Calculate number of predicted Positive applications (actual is 354)
25  sum(results)
26
27  # calculate a1...am for model
28  a <- colSums(model@xmatrix[[1]] * model@coef[[1]])
29  a
30
31  # calculate a0
32  a0 <- -model@b
33  a0
```

## Question 2.2

| Value of k | Accuracy |
|---|---|
| 1 | 94.19% |
| 2 | 94.19% |
| 3 | 94.19% |
| 4 | 94.19% |
| 5 | 92.20% |
| 6 | 90.98% |
| 7 | 91.28% |
| 8 | 91.74% |
| 9 | 90.67% |
| 10 | 90.06% |
| 15 | 89.76% |
| 20 | 89.60% |

For k values between one and four, the model's accuracy held constant.  Increases in the value of k past four resulted in accuracy losses.  As a result, I believe that four would be an optimal choice for k.  This k value choice is due to the belief that when applying the model to a new data point, comparing it to four neighbors instead of one, two, or three, the likelihood of misclassifying the new data point will be lessened.

*R-script is also attached and named question-2.2.R*

```
1   #Load the kknn library
2   require('kknn')
3
4   # Load credit card data into a data frame
5   data_df <- read.table("credit_card_data-headers.txt", header=TRUE)
6
7   # set.seed(42) - comment in to repeat results across different values of k
8
9   # Identifies 70% of the records that will be used to determine the train_df
10  sample <- sample(nrow(data_df), round(nrow(data_df)*.7))
11
12  # Create training (80pct) and test (20pct) datasets
13  train_df <- data_df[sample, ]
14  test_df <- data_df[-sample, ]
15
16  # Train the model - test using full dataset per homework instructions
17  model <- kknn(R1 ~ ., train_df, data_df, k = 4, scale = TRUE)
18
19  # Compare data_df R1 to model fit
20  results_df <- data.frame(data_df$R1, model$fit)
21
22  # Create a confusion matrix - round model$fit since fit is continuous and R1 is 0 or 1
23  CM <- table(data_df[, 11], round(model$fit))
24
25  # Calculate accuracy using CM
26  accuracy <- (sum(diag(CM)))/sum(CM)
27
28  # Print accuracy
29  accuracy
```