

i **Front page**

Department of Mathematical Sciences

Examination paper for **TMA4268 Statistical learning**

Academic contact during examination: Mette Langaas

Phone: 988 47 649

Examination date: 23 May 2019

Examination time (from-to): 09:00 - 13:00

Permitted examination support material: C.

- One yellow A5 sheet with your own handwritten notes (stamped by the Department of Mathematical Sciences),
- specified calculator.

Other information:

- All answers must be justified, and relevant calculations provided.
- Observe that the maximal score is given for each problem, so you should not spend too much time and too many words on problems that have a low maximal score.
- The exam questions are only available in English since this is a course given in English.
- This exam is given in Inspira.
- For each problem you may write your answer into Inspira, or use the provided paper sheets. Remember to correctly specify the code for each problem (available in lower right corner of the problem) on your paper sheet, so that the scanned sheets will be matched with the correct problem. It is advisable that you write down the code when starting to write on the sheet.

Students will find the examination results in Studentweb. Please contact the department if you have questions about your results. The Examinations Office will not be able to answer this.

1 **Warming up with an overview**

[Maximal score: 3 points]

Choose the correct alternative for each drop-down menu.

i) In this course we have considered two types of learning, A: (neural networks, classification, supervised, regression) and B: (support vector machines, unsupervised, logistic, crossvalidation). In A both covariates and a known response is present, but in B only (responses, covariates) are present. In A we have spent time studying two problem types: C (dimensionality reduction, clustering, regression, regularization) and (matrix algebra, covariance, classification, crossvalidation).

Linear methods have been an important starting point for both problem types.

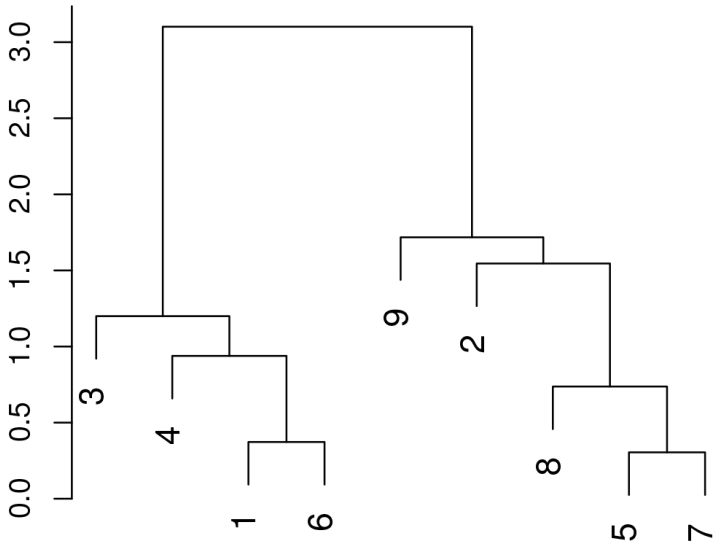
ii) In C we started with studying the relationship between one univariate response and one covariate in (K-nearest neighbour regression, simple linear regression, local regression, multiple linear regression), and then we added more covariates, and got (regression splines, principal component regression, multiple linear regression, simple linear regression). The least squares estimator was $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, which was (unbiased, biased, skewed) and had a covariance matrix that was (dependent on the responses, dependent on the covariates, a constant) in the multiple linear regression model. We looked at including non-linear effects in the covariates by adding (linear transformations, principal components, polynomials). Non-linear effects were

also the topic of (ridge regression, K-nearest neighbour regression, support vector classifier, lasso regression). Interactions (the fact that the effect of one covariate on the response is dependent on the value of another covariate) were easily included in (regression trees, ridge regression, lasso regression, step functions). Model selection was performed using (bagging, boosting, the AIC penalty, linear discriminant analysis) and (ridge regression, principal component regression, splines, lasso regression). More about regression in the last part of this exam.

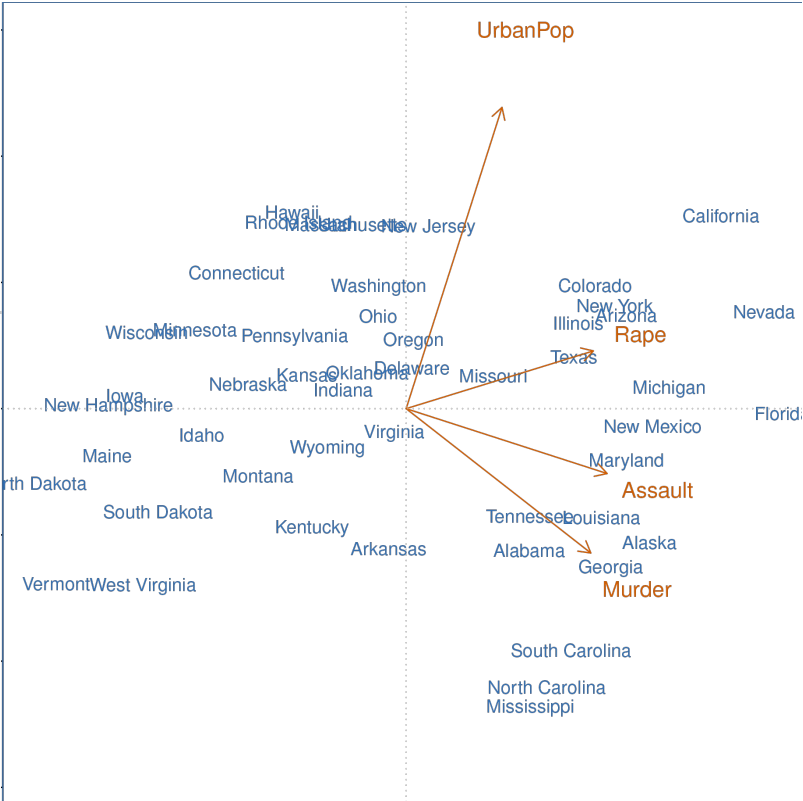
iii) The starting point for classification was the sampling and the (diagnostic, deviance, Gini, Bayes) paradigm. From the sampling paradigm we studied the (MMC, SVC and SVM, logistic regression, KNN, LDA and QDA), and from the other paradigm many more methods, which we will work with soon.

iv) In (crossvalidation, the bias-variance trade-off, supervised learning, unsupervised learning) the goal is to discover interesting aspects about the measurements x_1, x_2, \dots, x_p when y is not present. We have considered two types of methods: (principal component analysis, ridge regression, regularization, neural networks) and (lasso regression, clustering, R-programming, Python programming). Both methods have strong focus on finding (outliers, groups, residuals, errors) in the data, and on (the bias-variance trade-off, visualisation, crossvalidation, debugging).

The top figure on the right is an example of (hierarchical clustering, dimension reduction, response reduction, k-means clustering),



The bottom figure to the right gives a biplot from (ridge regression, k-means clustering, partial least squares, principal component analysis).



Maximum marks: 2.99

2 Match optimization criterion and method

[Maximal score: 2 points]

For each optimization criterion choose the correct method from the drop-down menu:

a) Maximize M by choosing $\beta_0, \beta_1, \dots, \beta_p$ subject to $\sum_{j=1}^p \beta_j^2 = 1$ and $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M$ for $i = 1, \dots, n$

(least squares regression, lasso regression, logistic regression, support vector classifier, regression tree, maximal margin classifier, ridge regression)

b) $\operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2$ (least squares regression, maximal margin classifier, support vector classifier, ridge regression, lasso regression, regression tree, logistic regression)

c) $\operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2$

(least squares regression, ridge regression, lasso regression, maximal margin classifier, support vector classifier, regression tree, logistic regression)

d) $\operatorname{argmin}_{R_1(j,s), R_2(j,s)} \left[\sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2 \right]$

(least squares regression, ridge regression, lasso regression, regression tree, maximal margin classifier, support vector classifier, logistic regression)

e) $\operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|$

(least squares regression, lasso regression, ridge regression, maximal margin classifier, regression tree, support vector classifier, logistic regression)

f) Maximize M by choosing $\beta_0, \beta_1, \dots, \beta_p, \varepsilon_1, \dots, \varepsilon_n$

subject to $\sum_{j=1}^p \beta_j^2 = 1$ and $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \varepsilon_i)$ where $\varepsilon_i \geq 0$ for all $i = 1, \dots, n$ and $\sum_{i=1}^n \varepsilon_i \leq C$

(least squares regression, regression tree, ridge regression, lasso regression, maximal margin classifier, support vector classifier, logistic regression)

Maximum marks: 2.04

3 Neural network – mathematical formula

[Maximal score: 4 points]

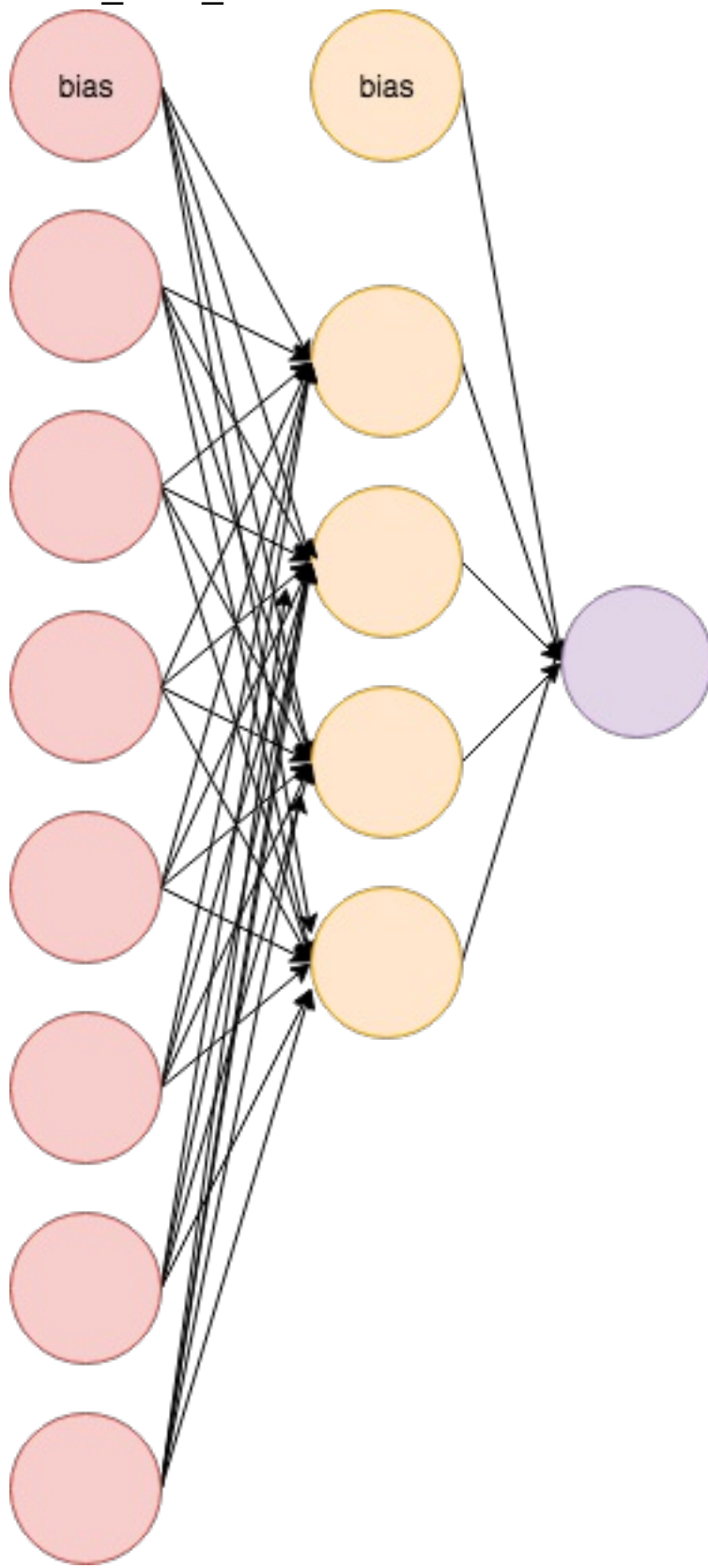
a) Write down the mathematical formula for a feedforward neural network with one input layer, one hidden layer and one output layer, with the following architecture

- input layer: seven input nodes
- one bias node for the hidden layer,
- hidden layer: four nodes,
- one bias node for the output layer,
- output layer: one node.

Let the nodes in the hidden layer have sigmoid activation function and the node in the output layer have linear activation function. A drawing of the feedforward neural network is given below.

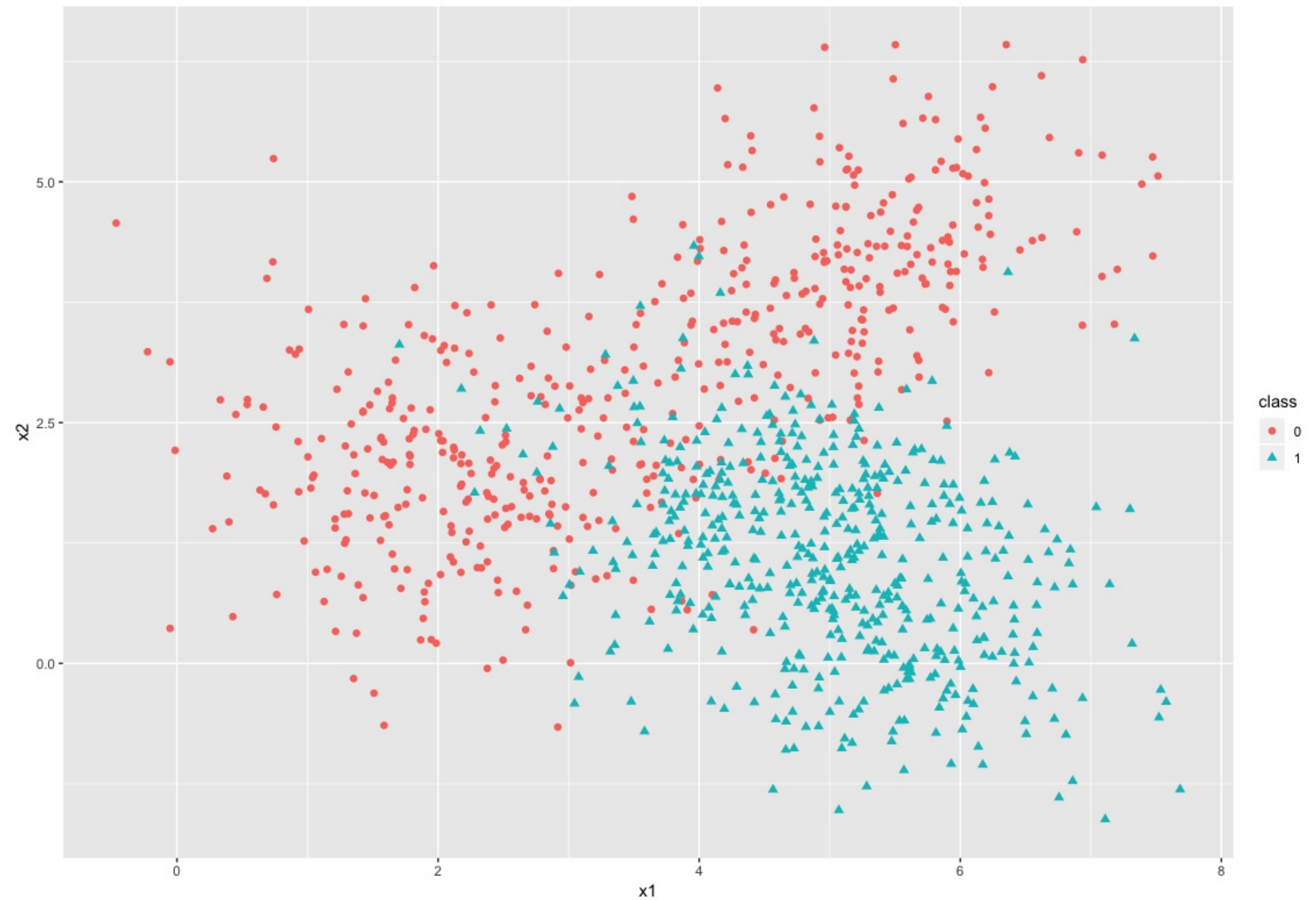
b) Which type of problem can this network solve?

c) What would be an appropriate loss function for this network (give mathematical formula)?



Fill in your answers here, and/or use the paper sheets provided.

Maximum marks: 4



[Maximal score: 3 points]

- a) Write down the mathematical formula for the probability of class 1 for the K -nearest neighbour classifier at a covariate value \mathbf{x}_0 , and explain your notation.
- b) What does the rule *majority vote* mean when applied to K -nearest neighbour classification? In our situation with two classes, what is the cut-off on the probability of class 1 equivalent of a majority vote?
- c) List one positive and one negative aspect of using K -nearest neighbour classification.

Fill in your answers here, and/or use the paper sheets provided.

Format | B | I | U | x₂ | x² | I_x | [copy] | [paste] | [undo] | [redo] | [list] | [table] | [math] | [sum] | ABC | [clear]

Words: 0

Maximum marks: 3

6 **Class boundaries for K-nearest neighbour**

[Maximal score: 2 points]

In the figure to the left we have used the K -nearest neighbour classification method, with K equal to **1, 3, 25** and **199**, to get a classification boundary between class 0 and 1 based on the training data. (To make the boundary stand out the training data are not shown in the plots.)

Match panels A–D to values of $K(1, 3, 25, 199)$ in the drop-down menus.

Panel A: (1, 3, 25, 199)

Panel B: (1, 3, 25, 199)

Panel C: (1, 3, 25, 199)

Panel D: (1, 3, 25, 199)

Maximum marks: 2

7 **Choosing K in K-nearest neighbour classification**

[Maximal score: 5 points]

Crossvalidation can be used to choose a good value for K, and is performed in the R-code and print-out shown in the pdf-file on the left.

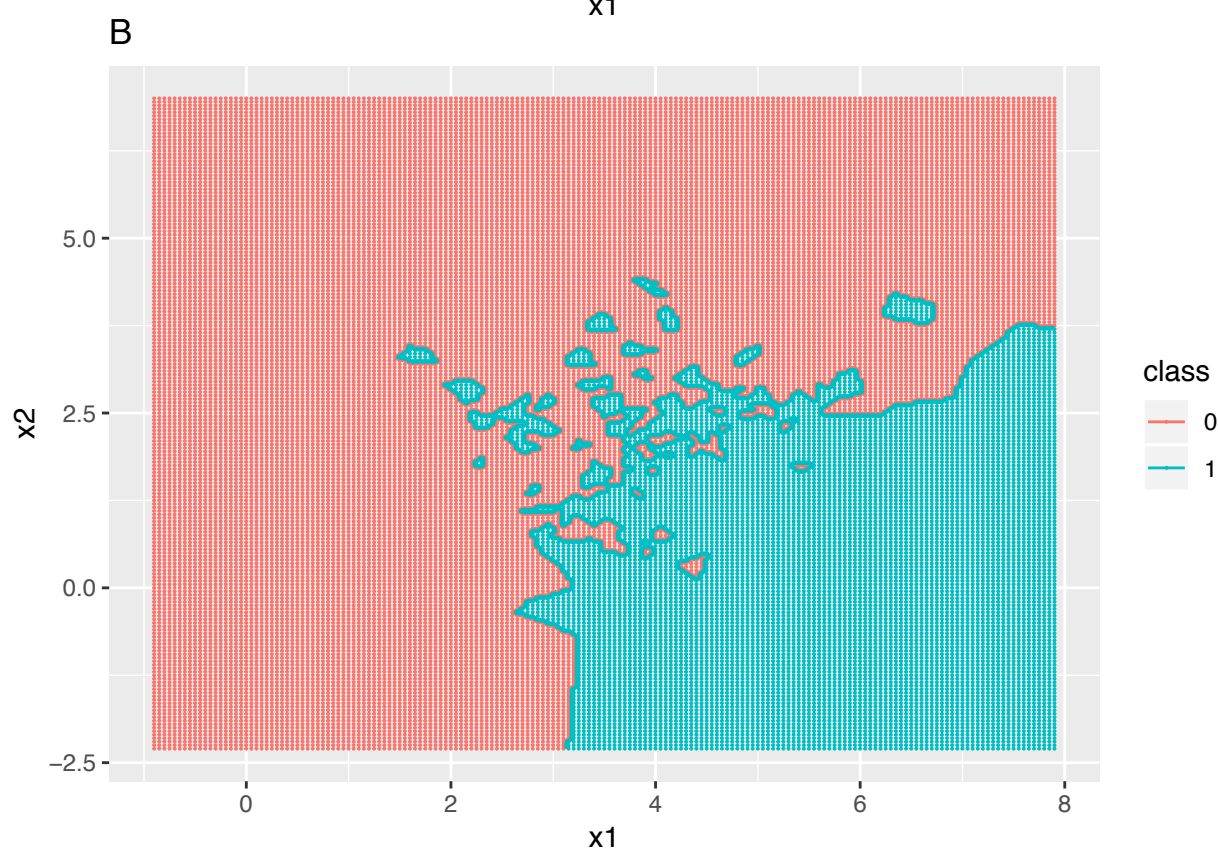
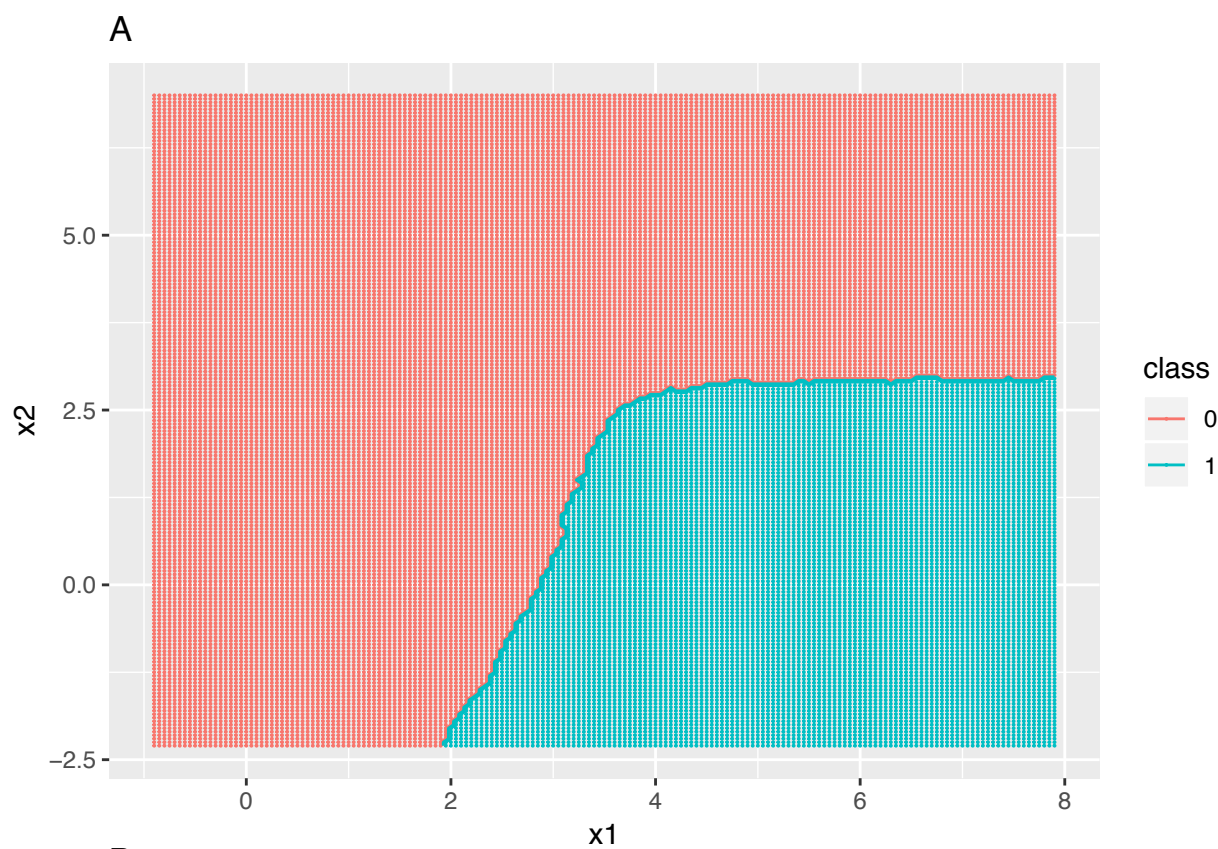
- a) Explain shortly what is done, and what the results are. Your explanation should include a drawing and the words: fold and error measure.
- b) Which value of K (number of neighbours) would you choose? Elaborate.
- c) What is the one-standard-error-rule? Can you from the print-out use this rule to choose K? If not, what can be added or changed in the R-code to be able to use this rule? You are not asked to rewrite the R-code.

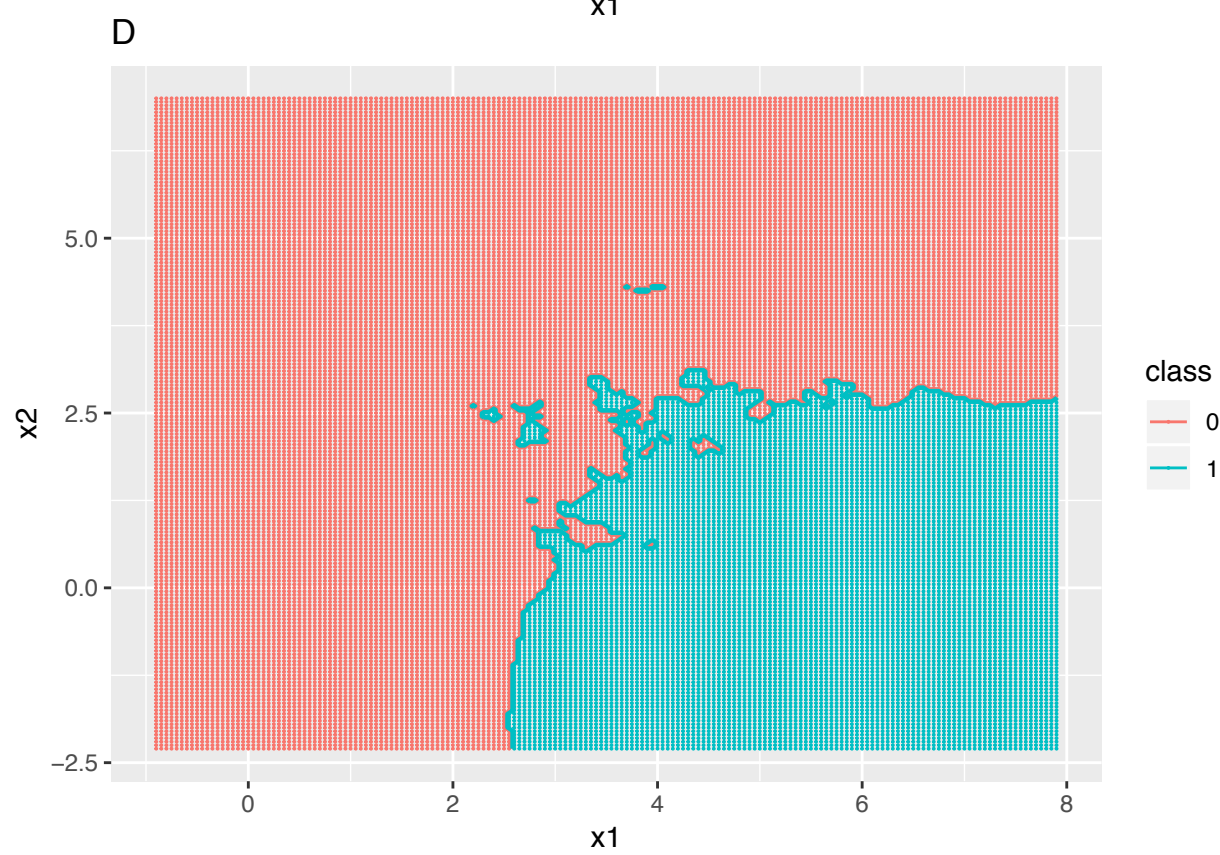
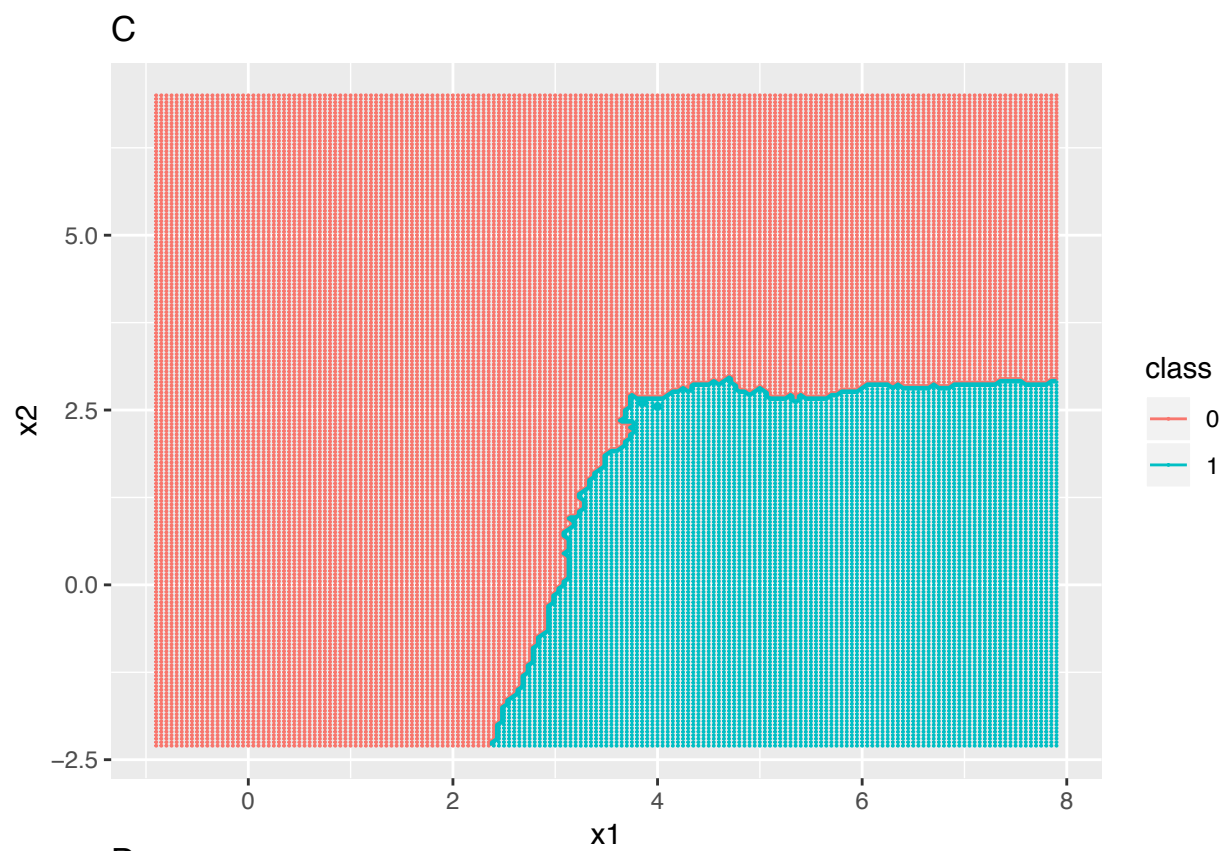
Fill in your answer here and/or use the paper sheets provided.

Format - **B** *I* U \times_2 \times^2 I_x Σ ABC

Words: 0

Maximum marks: 5





Choosing K in K -nearest neighbour classification

```
# the training data contains 500 observations from each class
# train.x: a 1000 times 2 matrix with training covariates x1 and x2
# train.y: a 1000 vector with the class (0 or 1) as a factor for the training data
```

```
Kgrid=seq(1,199,by=2)
Kgrid
```

```
##   [1]  1  3  5  7  9 11 13 15 17 19 21 23 25 27 29 31 33
##  [18] 35 37 39 41 43 45 47 49 51 53 55 57 59 61 63 65 67
##  [35] 69 71 73 75 77 79 81 83 85 87 89 91 93 95 97 99 101
##  [52] 103 105 107 109 111 113 115 117 119 121 123 125 127 129 131 133 135
##  [69] 137 139 141 143 145 147 149 151 153 155 157 159 161 163 165 167 169
##  [86] 171 173 175 177 179 181 183 185 187 189 191 193 195 197 199
```

```
nmodels=length(Kgrid)
nfolds=5
```

```
set.seed(4268)
folds <- createFolds(train.y,k=nfolds)
# the random sampling is done within the levels of y when y is a factor
# in an attempt to balance the class distributions within the splits.
# folds[[1]] gives indecies for training observations in fold 1,
# folds[[2]] the same for fold 2 etc.
```

```
cv.errors=rep(0,nmodels)
```

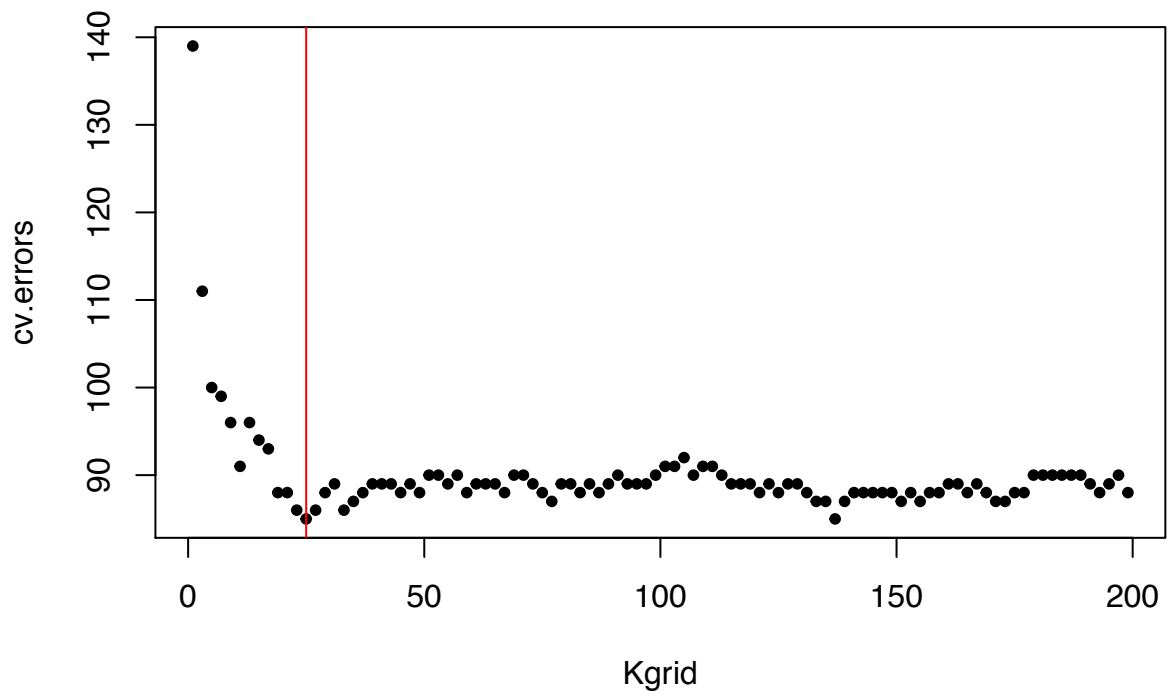
```
for(i in 1:nfolds)
{
  for(j in 1:nmodels)
  {
    pred=class::knn(train=train.x[-folds[[i]],],
                    test=train.x[folds[[i]],],
                    cl=train.y[-folds[[i]]], k=Kgrid[j])
    # gives the predicted class 0/1 for the validation fold
    cv.errors[j]=cv.errors[j]+sum(pred!=train.y[folds[[i]]])
  }
}
cv.errors
```

```
##   [1] 139 111 100 99 96 91 96 94 93 88 88 86 85 86 88 89 86
##  [18] 87 88 89 89 89 88 89 88 90 90 89 90 88 89 89 89 88
##  [35] 90 90 89 88 87 89 89 88 89 88 89 90 89 89 89 90 91
##  [52] 91 92 90 91 91 90 89 89 89 88 89 88 89 89 88 87 87
##  [69] 85 87 88 88 88 88 87 88 87 88 88 89 89 88 89 88
##  [86] 87 87 88 88 90 90 90 90 90 90 89 88 89 90 88
```

```
Kgrid[which.min(cv.errors)]
```

```
## [1] 25
```

```
plot(x=Kgrid,y=cv.errors,pch=20)
abline(v=25,col=2)
```



9 Bayes decision rule

[Maximal score: 10 points]

a) What is a Bayes classifier, Bayes decision boundary and Bayes error rate?

In our classification problem the training and test data sets were simulated as follows:

- Prior class probabilities: $\pi_0 = P(Y = 0) = 0.5$ and $\pi_1 = P(Y = 1) = 0.5$.
- Class conditional probabilities:
 - $f_0(\mathbf{x}) = 0.5 \cdot \frac{1}{2\pi|\Sigma_{01}|} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{01})^T \Sigma_{01}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{01})\right] + 0.5 \cdot \frac{1}{2\pi|\Sigma_{02}|} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{02})^T \Sigma_{02}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{02})\right]$
 - $f_1(\mathbf{x}) = \frac{1}{2\pi|\Sigma_1|} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right]$
- with mean vectors $\boldsymbol{\mu}_{01} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$, $\boldsymbol{\mu}_{02} = \begin{bmatrix} 5 \\ 4 \end{bmatrix}$ and $\boldsymbol{\mu}_1 = \begin{bmatrix} 5 \\ 1 \end{bmatrix}$, and
- covariance matrices $\Sigma_{01} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $\Sigma_{02} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$, and $\Sigma_1 = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$

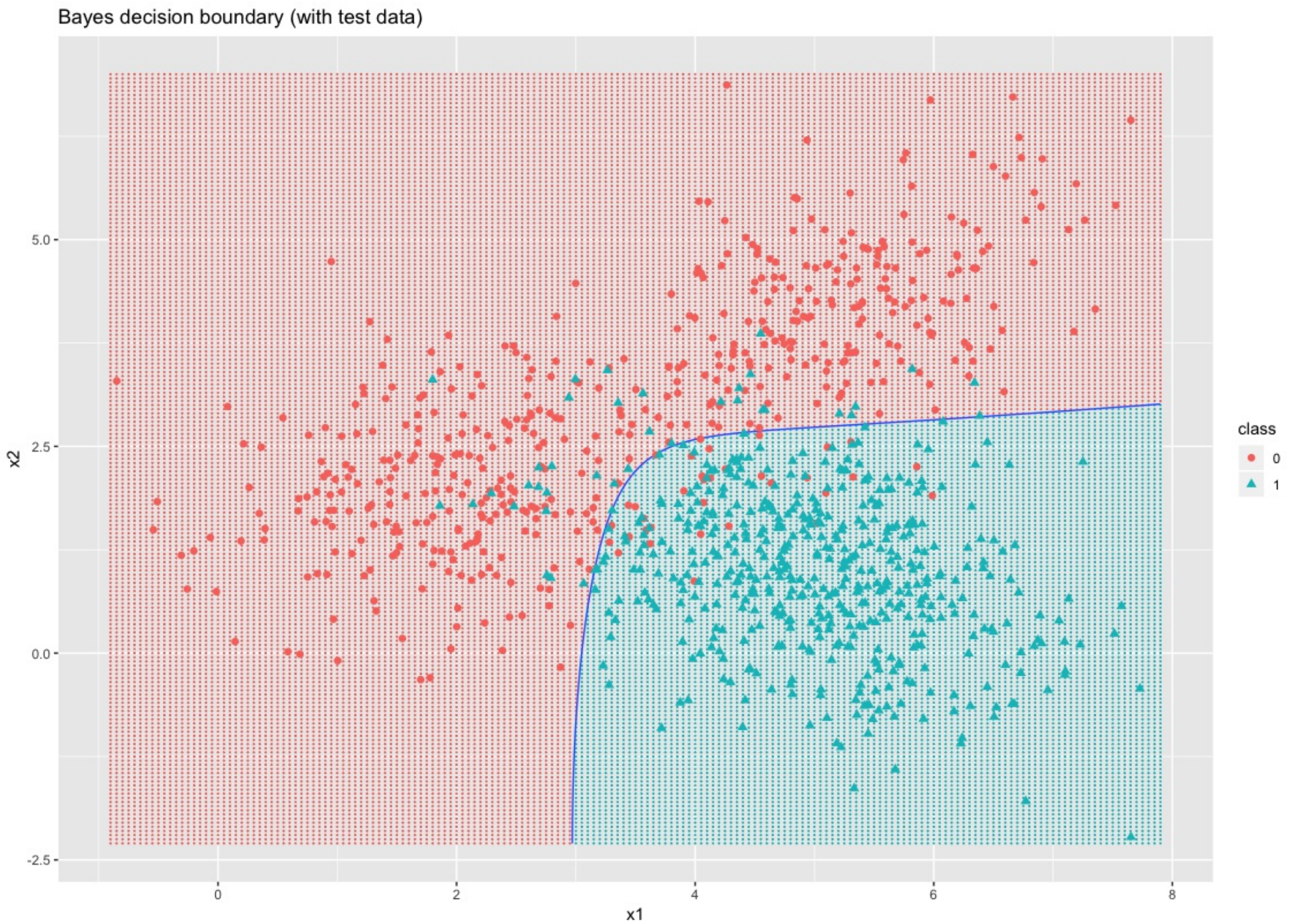
b) How would you proceed to find the Bayes decision boundary? Write down the equation to be solved, explain what are the unknowns. You are *not* asked to solve the equation.

c) The Bayes decision boundary is shown as a (blue) curve in the plot below. The Bayes error rate was found (by simulation) to be 8%. Explain what this means.

d) What would happen to the Bayes decision boundary and Bayes error rate if we change the prior probabilities of the two classes to $\pi_0 = 0.9$ and $\pi_1 = 0.1$? Elaborate.

e) The methods we have considered in this course are *K*—nearest neighbour classification, linear discriminant analysis, quadratic discriminant analysis, logistic regression, classification tree, bagged trees, random forest, maximal margin classifier, support vector classifier, support vector machine, and feedforward neural network.

Now that you know what the Bayes decision boundary looks like, which of the classification methods (select only one method) do you think would provide the best solution to our classification problem? Justify your answer.



Fill in your answer here and/or use the paper sheets provided.

Format | B | I | U | x₂ | x² | I_x | [copy] | [paste] | [undo] | [redo] | [list] | [list] | [link] | [table] | [edit] | [sum] | ABC | [close]

Words: 0

Maximum marks: 10

10 **Regression with bike rentals**

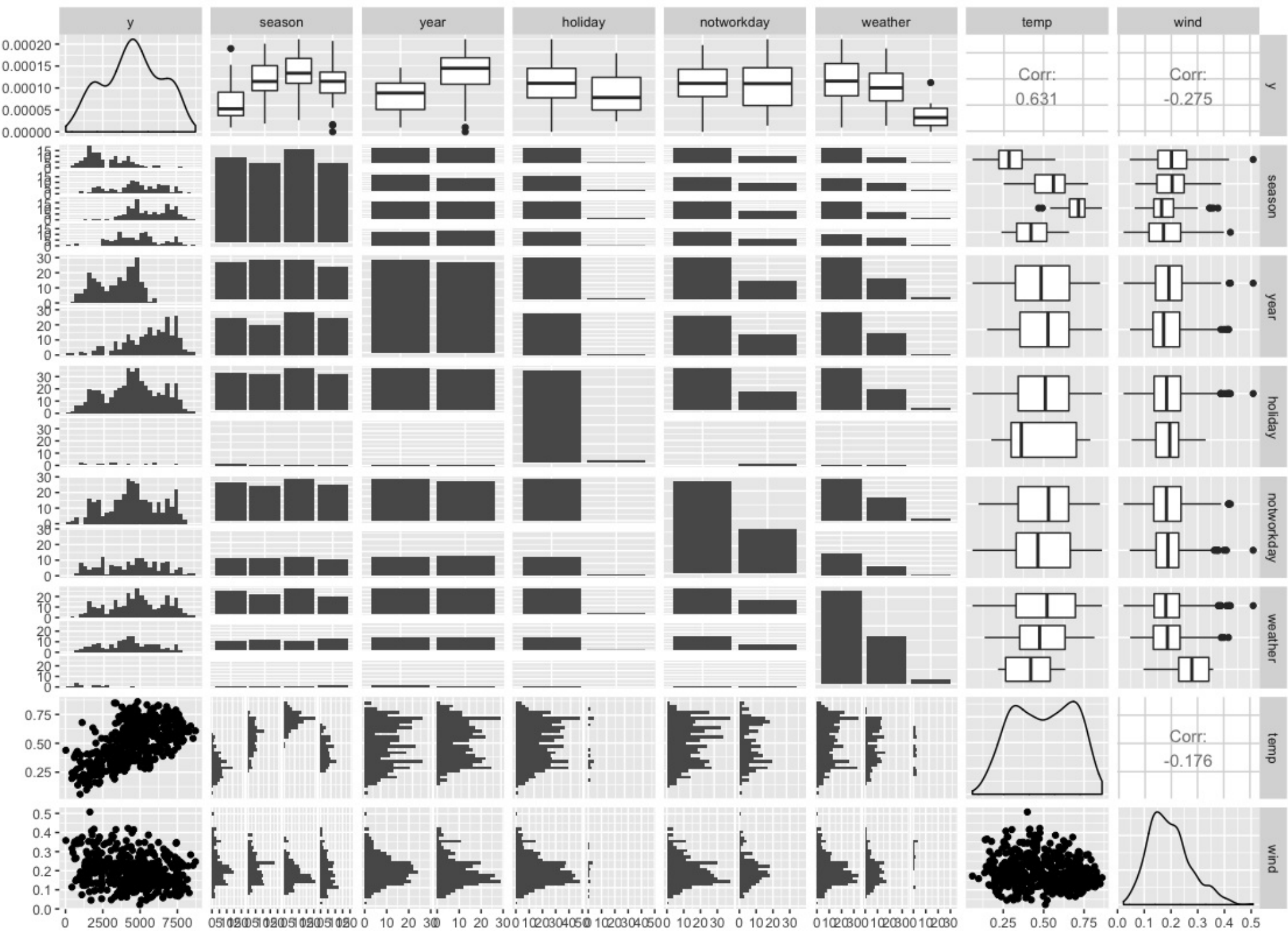
Comment: We will now have 4 problems on regression with the same data set-up.

We will look at a data set of daily counts of bike rentals in a bike sharing system in Washington DC in the first two years of operation, matched with data on climate, season and type of day.

(Source: <http://archive.ics.uci.edu/ml/datasets/bike+sharing+dataset>)

- y: daily count of number of bike rentals (our response)
- year: (0: 2011, 1: 2012)
- season: factor with four levels (1: spring, 2: summer, 3: autumn, 4: winter)
- holiday: (0: not holiday, 1: holiday)
- notworkday: (0: neither weekend nor holiday, 1: weekend or holiday)
- weather: factor with three levels (1: clear to partly cloudy, 2: mist and/or clouds, 3: snow, rain, thunderstorm)
- temp: normalized temperature in Celsius
- wind: normalized wind speed

The data was divided randomly into a training set of size 500 and a test set of size 231.
The training data is presented in the pairs-plot below.



Smoothing spline and polynomials
[Maximal score: 6 points]

We use the daily count of number of bike rentals (y) as response and temperature (temp) as the only covariate.

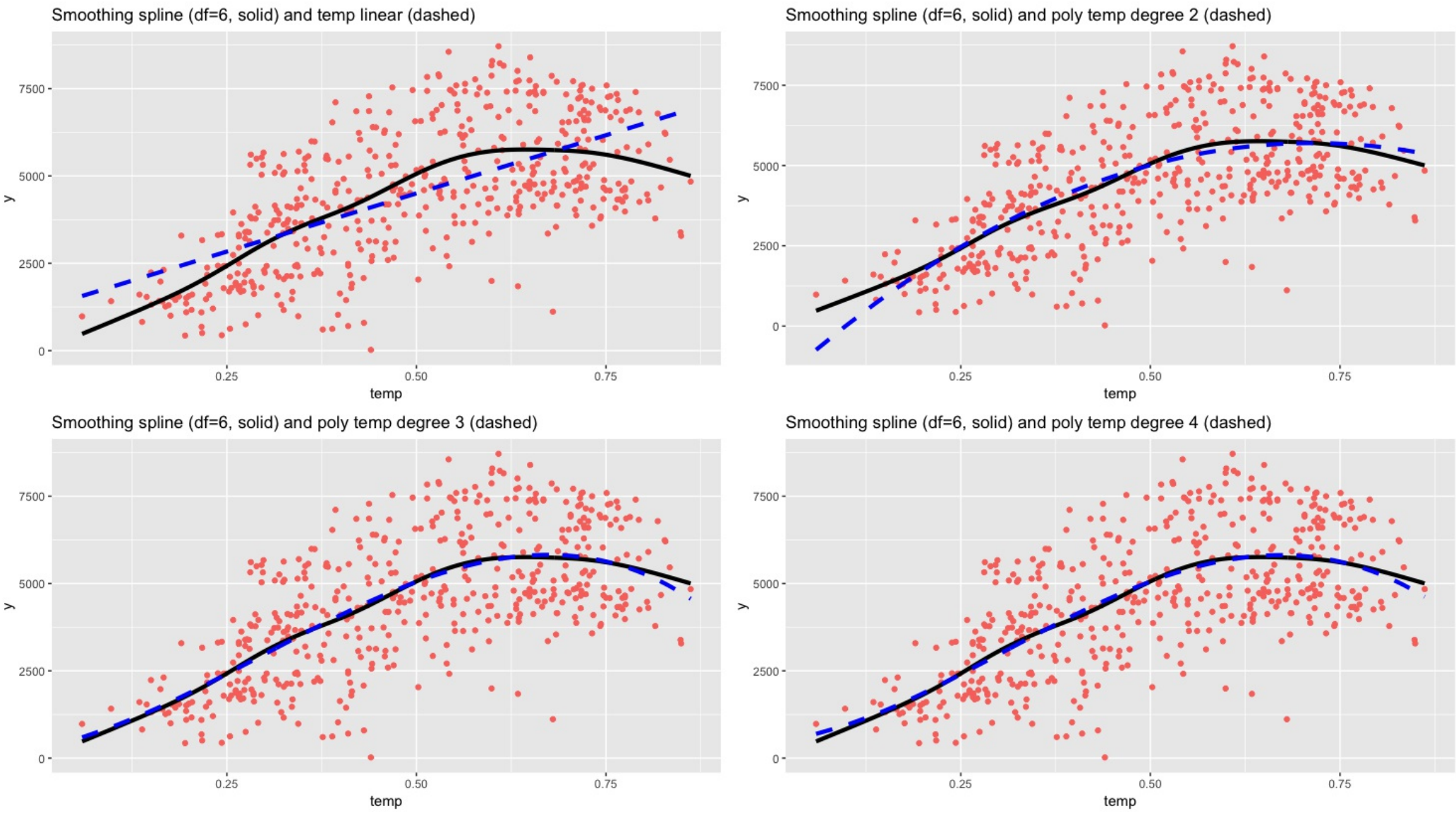
a) What is the idea behind a smoothing spline, and which criterion is the smoothing spline minimizing?
What is the connection between a smoothing spline and a cubic spline?

A total of five curves were fitted to the training data:

- a smoothing spline with 6 effective degrees of freedom (solid, black)
- polynomial of degrees 1, 2, 3 and 4 (dashed, blue)

See the results in the four panels below.

b) What does it mean that the smoothing spline has 6 effective degrees of freedom?
How would you evaluate the different model fits in the four panels?
We will next fit a multiple linear regression with all available covariates. What is your suggestion for how to handle the temp covariate?



Fill in your answer here and/or use the paper sheets provided.

Format | B | I | U | x₂ | x² | I_x | [copy] | [paste] | [undo] | [redo] | [list] | [bulleted] | Ω | [table] | [pencil] | Σ | ABC | [clear]

Words: 0

Maximum marks: 6

A multiple linear regression was fitted to the data, with linear effects of all covariates, and in addition also a quadratic and a cubic effect of temperature. Dummy variable coding (treatment contrast) was used for the factors, with the lowest level (season=1 and weather=1) as the reference category. A pdf with R-code and print-out is found in the left panel.

- year=1
- season=1
- holiday=0
- notworkday=0
- weather=3
- temp=0.3
- wind=0.2

Fill in your answer here and/or use the paper sheets provided.

13/15

Multiple linear regression

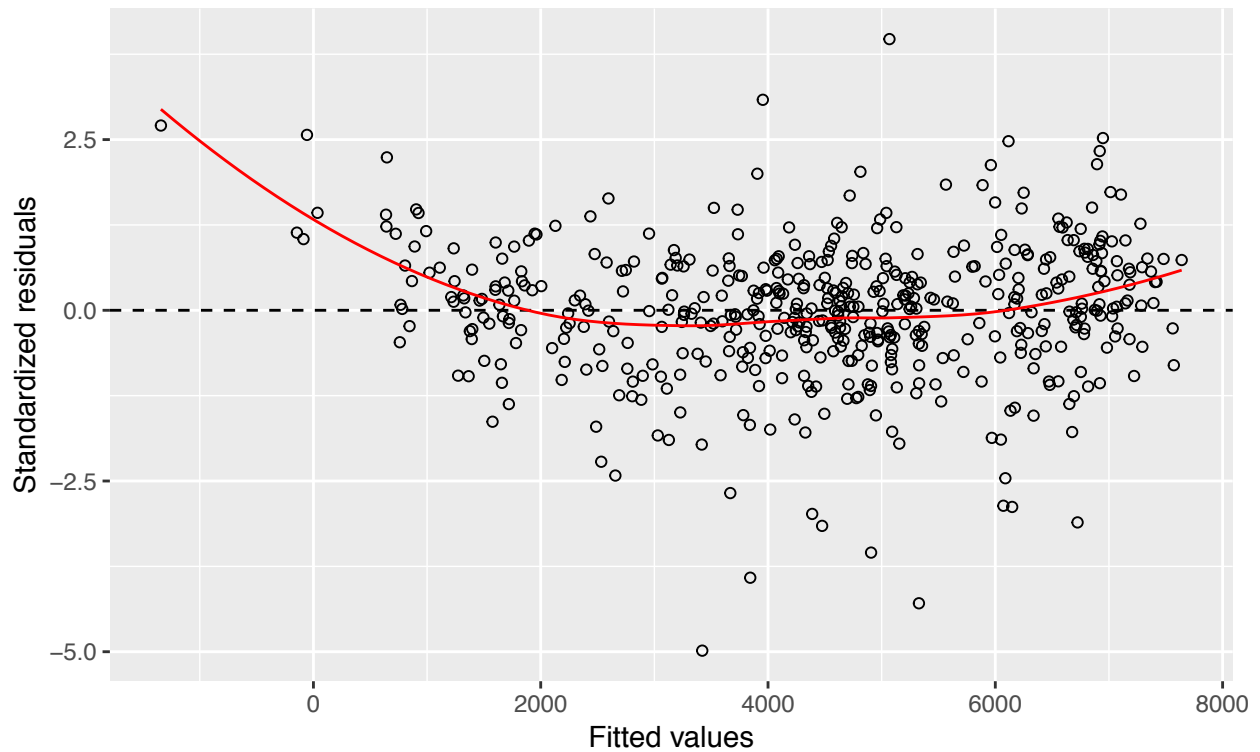
```
fitlm=lm(y~season+year+holiday+notworkday+weather+temp+I(temp^2)+I(temp^3)+wind,data=train)
# I(temp^2) just means temp^2, but in a model formula this extra I() is needed
summary(fitlm)
```

```
##
## Call:
## lm(formula = y ~ season + year + holiday + notworkday + weather +
##      temp + I(temp^2) + I(temp^3) + wind, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3399.3  -344.9    36.9   429.2  2767.6
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2902.20     481.36   6.029 3.26e-09 ***
## season2         828.76     121.20   6.838 2.41e-11 ***
## season3       1093.71     155.38   7.039 6.62e-12 ***
## season4       1286.18     104.99  12.251 < 2e-16 ***
## year          2052.39      64.82  31.663 < 2e-16 ***
## holiday        -479.99     199.13  -2.410  0.0163 *
## notworkday     -100.66      70.22  -1.434  0.1524
## weather2       -705.99      69.99 -10.087 < 2e-16 ***
## weather3     -2495.95     186.56 -13.379 < 2e-16 ***
## temp        -14643.01    3430.61  -4.268 2.37e-05 ***
## I(temp^2)     57717.61    7464.37   7.732 6.10e-14 ***
## I(temp^3)    -47873.21    5071.47  -9.440 < 2e-16 ***
## wind          -2726.69     425.89  -6.402 3.61e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 709.4 on 487 degrees of freedom
## Multiple R-squared:  0.8707, Adjusted R-squared:  0.8675
## F-statistic: 273.2 on 12 and 487 DF,  p-value: < 2.2e-16
```

```
ggplot(fitlm, aes(.fitted, .stdresid)) + geom_point(pch = 21) + geom_hline(yintercept = 0,
  linetype = "dashed") + geom_smooth(se = FALSE, col = "red", size = 0.5,
  method = "loess") + labs(x = "Fitted values", y = "Standardized residuals",
  title = "Fitted values vs standardized residuals", subtitle = deparse(fitlm$call))
```

Fitted values vs standardized residuals

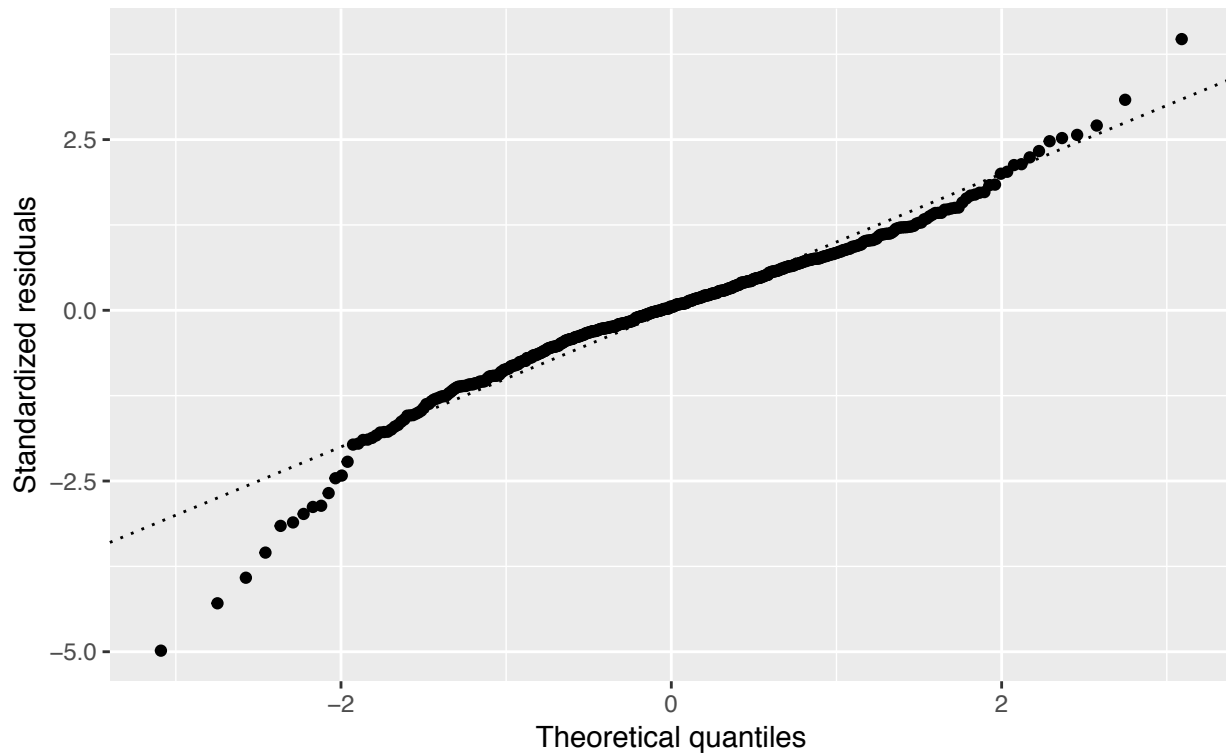
lm(formula = y ~ season + year + holiday + notworkday + weather +



```
ggplot(fitlm, aes(sample = .stdresid)) + stat_qq(pch = 19) + geom_abline(intercept = 0,  
  slope = 1, linetype = "dotted") + labs(x = "Theoretical quantiles",  
  y = "Standardized residuals", title = "Normal Q-Q", subtitle = deparse(fitlm$call))
```


Normal Q-Q

lm(formula = y ~ season + year + holiday + notworkday + weather +



```
library(nortest)
ad.test(rstudent(fitlm))

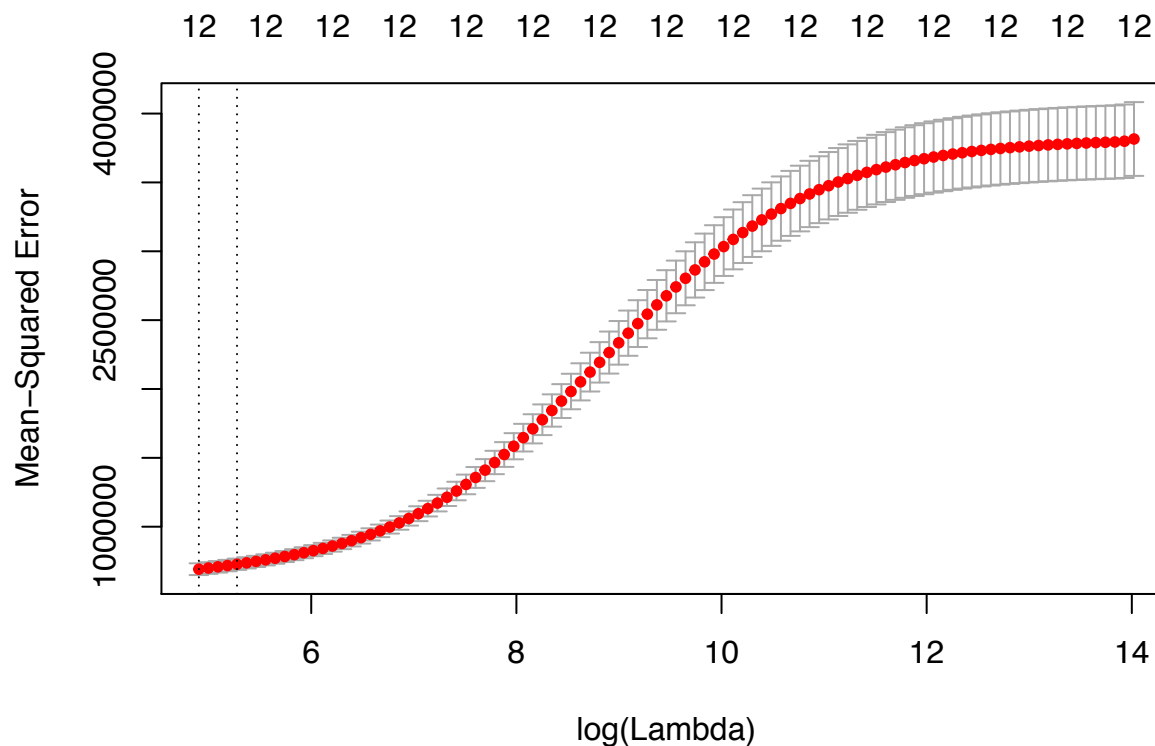
##
##  Anderson-Darling normality test
##
## data:  rstudent(fitlm)
## A = 3.9514, p-value = 7.525e-10

trainerrorlm=mean((predict(fitlm)-train$y)^2)
predlm = predict(fitlm, newdata=test)
testerrorlm=mean((predlm-test$y)^2)
c(trainerrorlm,testerrorlm)

## [1] 490198.7 568906.1
```


Ridge regression

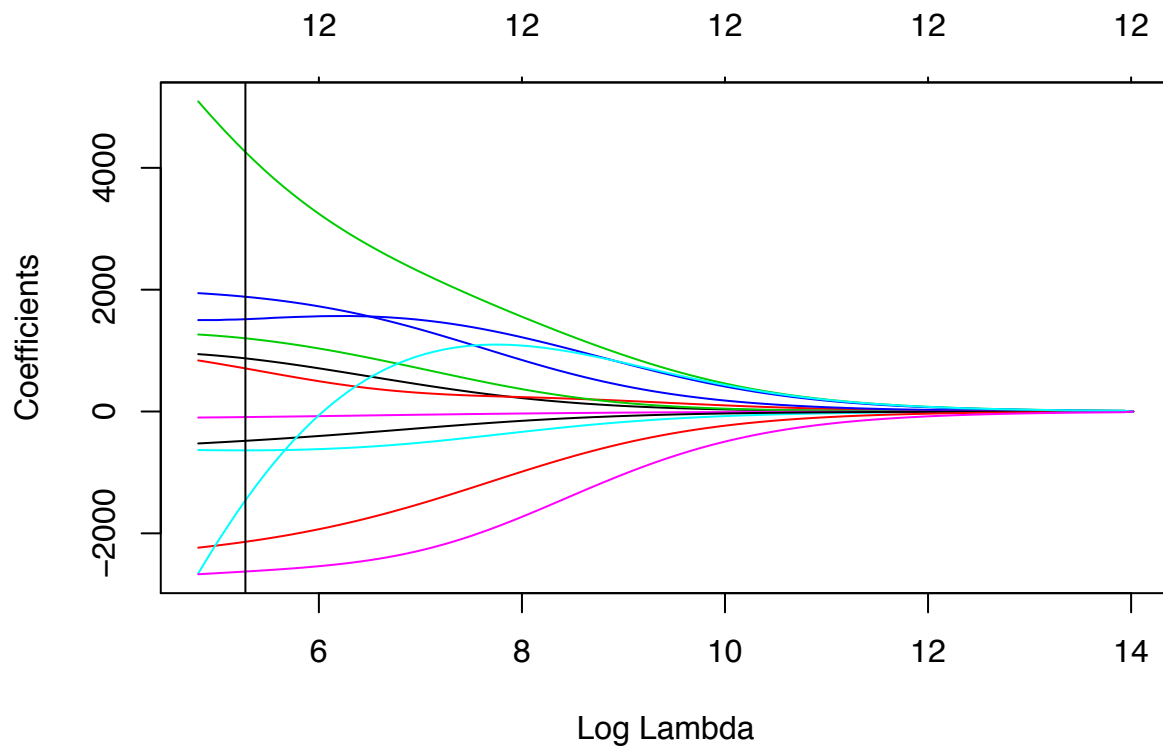
```
library(glmnet)
# set up model matrix, not include the intercept term because glmnet doesn't want that
model = formula(~season+year+holiday+notworkday+weather+temp+I(temp^2)+I(temp^3)+wind)
mm = model.matrix(model,data=train)[-1]
set.seed(4268) #for reproducibility
cvfitridge = cv.glmnet(mm,train$y,alpha=0,standardize=TRUE,nfolds=10)
plot(cvfitridge)
```



```
print(cvfitridge$lambda.1se)
```

```
## [1] 195.6469
```

```
# also plot how coefficients change with (log)lambda
plot(cvfitridge$glmnet.fit,xvar="lambda")
abline(v=log(cvfitridge$lambda.1se))
```



```
# look at final model
coef(cvfitridge)

## 13 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept) 1363.37674
## season2     872.55853
## season3     707.40757
## season4    1201.65306
## year        1883.97691
## holiday     -637.68536
## notworkday  -91.38161
## weather2    -481.23220
## weather3   -2139.03652
## temp        4259.71239
## I(temp^2)   1515.42979
## I(temp^3)  -1458.51139
## wind       -2626.81004

# mse on training and test set
trainererrorridge=mean((predict(cvfitridge, s="lambda.1se", newx=mm)-train$y)^2)
predridge =predict(cvfitridge, s="lambda.1se", newx=model.matrix(model,data=test)[,-1])
testerrorridge=mean((predridge-test$y)^2)
c(trainererrorridge,testerrorridge)

## [1] 687947.1 730840.7

# estimated regression coefficients
# for least squares (lm) and ridge regression (glmnet)
res=cbind(fitlm$coefficients,coef(cvfitridge))
colnames(res)=c("lm","ridge")
res
```

```
## 13 x 2 sparse Matrix of class "dgCMatrix"
##           lm           ridge
## (Intercept) 2902.1995 1363.37674
## season2      828.7597  872.55853
## season3     1093.7146  707.40757
## season4     1286.1808 1201.65306
## year        2052.3907 1883.97691
## holiday      -479.9857 -637.68536
## notworkday   -100.6592 -91.38161
## weather2     -705.9897 -481.23220
## weather3    -2495.9517 -2139.03652
## temp        -14643.0085 4259.71239
## I(temp^2)    57717.6057 1515.42979
## I(temp^3)   -47873.2054 -1458.51139
## wind         -2726.6905 -2626.81004

# and remembering the errors for lm and ridge
c(trainererrorlm,testerrorlm)

## [1] 490198.7 568906.1

c(trainererrorridge,testerrorridge)

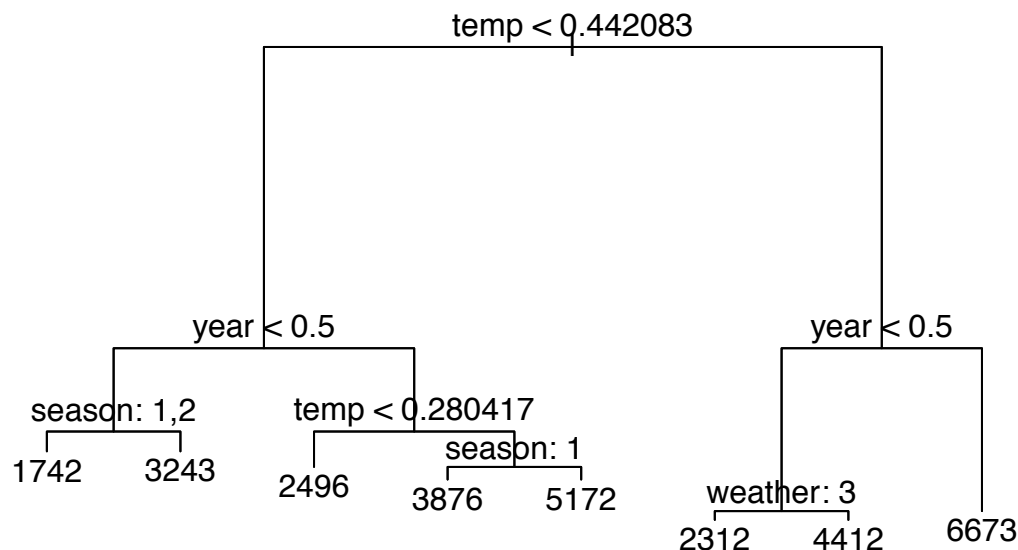
## [1] 687947.1 730840.7
```


Regression tree

```
library(tree)
fitT=tree(y~., data=train)
summary(fitT)
```

```
##
## Regression tree:
## tree(formula = y ~ ., data = train)
## Variables actually used in tree construction:
## [1] "temp"      "year"      "season"    "weather"
## Number of terminal nodes: 8
## Residual mean deviance: 774800 = 381200000 / 492
## Distribution of residuals:
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -5150.0  -430.7   126.0     0.0   504.4   2110.0
```

```
plot(fitT)
text(fitT,pretty=TRUE)
```



```
fitT
```

```
## node), split, n, deviance, yval
##      * denotes terminal node
##
## 1) root 500 1.895e+09 4493
##    2) temp < 0.442083 210 5.322e+08 3073
##      4) year < 0.5 110 1.009e+08 2138
##        8) season: 1,2 81 3.150e+07 1742 *
##        9) season: 4 29 2.128e+07 3243 *
##      5) year > 0.5 100 2.293e+08 4101
##        10) temp < 0.280417 25 1.599e+07 2496 *
##        11) temp > 0.280417 75 1.274e+08 4636
##          22) season: 1 31 1.819e+07 3876 *
##          23) season: 2,4 44 7.866e+07 5172 *
##    3) temp > 0.442083 290 6.319e+08 5522
```

```
##      6) year < 0.5 143 9.863e+07 4339
##      12) weather: 3 5 5.695e+05 2312 *
##      13) weather: 1,2 138 7.677e+07 4412 *
##      7) year > 0.5 147 1.382e+08 6673 *
```