# Introduction to Machine Learning
# Individual Laboration Report –1–

Erik Sven Vasconcelos Jansson
erija578@student.liu.se
Linköping University,  Sweden

November 27, 2016

## Assignment 1

Nobody likes *e-mail spam*, therefore methods for autonomously *predicting* if a given e-mail is probably *spam* or *not spam* is an important task. This is a classic example where *machine learning* is useful; given a set of *training data* and *testing data*, can we predict what is *spam* and *not spam* in the *testing set* (without knowing the answer) by deriving a *hypothesis function* built from the *training data*?

By using *k-nearest neighbor classification*, one can derive if an e-mail is spam or not by simply looking at *similar e-mails/messages*, and picking the most likely solution by doing a *"majority vote"*. First, a *distance function* needs to be implemented, which is the *cosine distance function* in Equation 1, whose implementation can be found in Listing 3, but with a optimized solution using only matrices.

$$d(X,Y) = 1 - \frac{X^T Y}{\sqrt{\sum_i X_i^2}\sqrt{\sum_i Y_i^2}} \qquad (1)$$

After defining the distance function $d(X,Y)$, one can find the *e-mail/message distance* for each $Y_j$ in respect to each $X_i$. Where $X$ is the *testing set* and $Y$ the *training set*. Each row of the resulting matrix contains the relative distance between $X_i$ and $\forall Y_j$. Therefore, sorting each row $X_i$ and picking the first $K$ elements gives the $K$ *closest messages* from the *training set* in respect to each *testing element*. By using this, the *k-nearest neighbors* can be found, and the prediction of $\hat{Y}$ (spam, not spam) is done by using Equation 2, where $K_i$ classify as being $C_i$.

$$\hat{Y} = \max_{\forall C_i} p(C_i|\boldsymbol{x}),\ p(C_i|\boldsymbol{x}) \propto K_i \div K \qquad (2)$$

The *k-nearest neighbor algorithm* is implemented in Listing 2, in the function knearest(t,k,t'). It works as previously described, where line 20 is calculating the *distance matrix* and line 21 sorting each row, so that all $Y$ distances are relative to $X_i$. Thereafter, in line 26-27 the classification is found for the $K$-nearest neighbors of $X_i$. The *mean* value is then taken, which is equivalent to $K_i \div K$ since only two classifications exist (spam and not spam), following a *Cover et al.* [CH67] K-NN descriptions.

Below follow *confusion matrices* & *ROC curves*, and it seems *knearest* gives better results than *kknn*. The *first confusion matrix* below belongs to the *testing data set* while the *second* one belongs to the *training data set*, and finally *testing* for *kknn*. The misslassification rates for each respective *confusion matrix* is: $(0.202, 0.347)$ for *training*, $(0.317, 0.347)$ for *testing* and finally $(0.345, 0.345)$ for *kknn*. Notice how predictions are more "accurate" in *training*.

| k=5 | false | true |
|---|---|---|
| false | 695 | 193 |
| true | 242 | 240 |

| k=1 | false | true |
|---|---|---|
| false | 639 | 178 |
| true | 298 | 255 |

| k=5 | false | true |
|---|---|---|
| false | 787 | 119 |
| true | 158 | 306 |

| k=1 | false | true |
|---|---|---|
| false | 939 | 2 |
| true | 6 | 423 |

| k=5 | false | true |
|---|---|---|
| false | 640 | 177 |
| true | 297 | 256 |

| k=1 | false | true |
|---|---|---|
| false | 640 | 177 |
| true | 297 | 256 |

The reason why this happens is because *training* is being used by the *predictor*, and therefore (especially when $k = 1$) will give very "accurate" predictions. This becomes a bit less apparent when we enforce majority voting (taking into account near neighbors), since the predictor doesn't become as biased towards the *training* data set as before (however it still seems $k = 5$ predicts better than *testing* $k = 5$). Note in the ROC curve below that *knearest* performs better than the *kknn*, at least in this case.



ROC Curve for the K–NN Spam Predictors

## Assignment 2

Knowing how to *infer* what parameter $\theta$ most likely produced a already given data vector $\boldsymbol{x}$ is useful to gain more information about underlying processes. In this case, that is *expected lifetime of machines*, which has been assumed to follow $p(\boldsymbol{x}|\theta) = \theta e^{-\theta x}$, where $\boldsymbol{x}$ are the *expected lifetimes* of $n$ machines, which are independent and identically distributed. The P.D.F. $\theta e^{-\theta x}$ is from a *exponential distribution*, as seen in *Jonsson et al.* [JN99]: $X \sim \mathrm{Exp}(\mu, \sigma^2)$...

Estimation of the parameter $\theta$ can be done with *MLE (Maximum Likelihood Estimation)*, which is usually done by using the *log-likelihood* of $\theta$ for a given data vector $\boldsymbol{x}$. The formula for *log-likelihood* is shown in Equation 3, where the parameter $\theta$ can then be estimated with $\hat{\theta}_{mle}$ by selecting the most probable $\theta$ from a given set of thetas $\Theta, \theta_i \in \Theta$, by *maximizing* the *average log-likelihood*, as in Eq. 4. This is seen in e.g. *Myung* [Myu03] and *Wikipedia*.

$$\mathcal{L}(\theta; \boldsymbol{x}) = \ln p(\boldsymbol{x}|\theta) = \sum_{i=1}^{n} \ln p(x_i|\theta) \qquad (3)$$
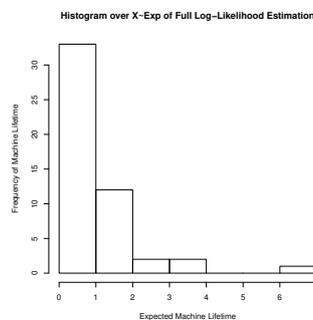
$$\hat{\theta}_{mle} = \max_{\forall \theta} \hat{\mathcal{L}}(\theta; \boldsymbol{x}), \ \hat{\mathcal{L}} = \frac{\ln p(\boldsymbol{x}|\theta)}{n} \qquad (4)$$

The actual implementation of these is found in Listing 5, in the $R$ functions `lnlikelihood` and `distribution` (gives $\theta e^{-\theta x_i}, \forall x_i \in \boldsymbol{x}$). For the MLE, lines `7-9` in Listing 4 calculate the *average log-likelihood* and then picks the *most probable* $\theta_i$. The MLE for the dataset shown in Listing 6 for when the *entire dataset* is used, and when only $|\boldsymbol{x}| = 6$ is taken, is shown in the table to the right. As can be seen in the graph to the right, the $\hat{\theta}_{mle}$ for when $|\boldsymbol{x}| = 6$ will overshoot the value for when the actual full dataset is used, therefore, it is less reliable then when $|\boldsymbol{x}| = 48$, not a good estimator.
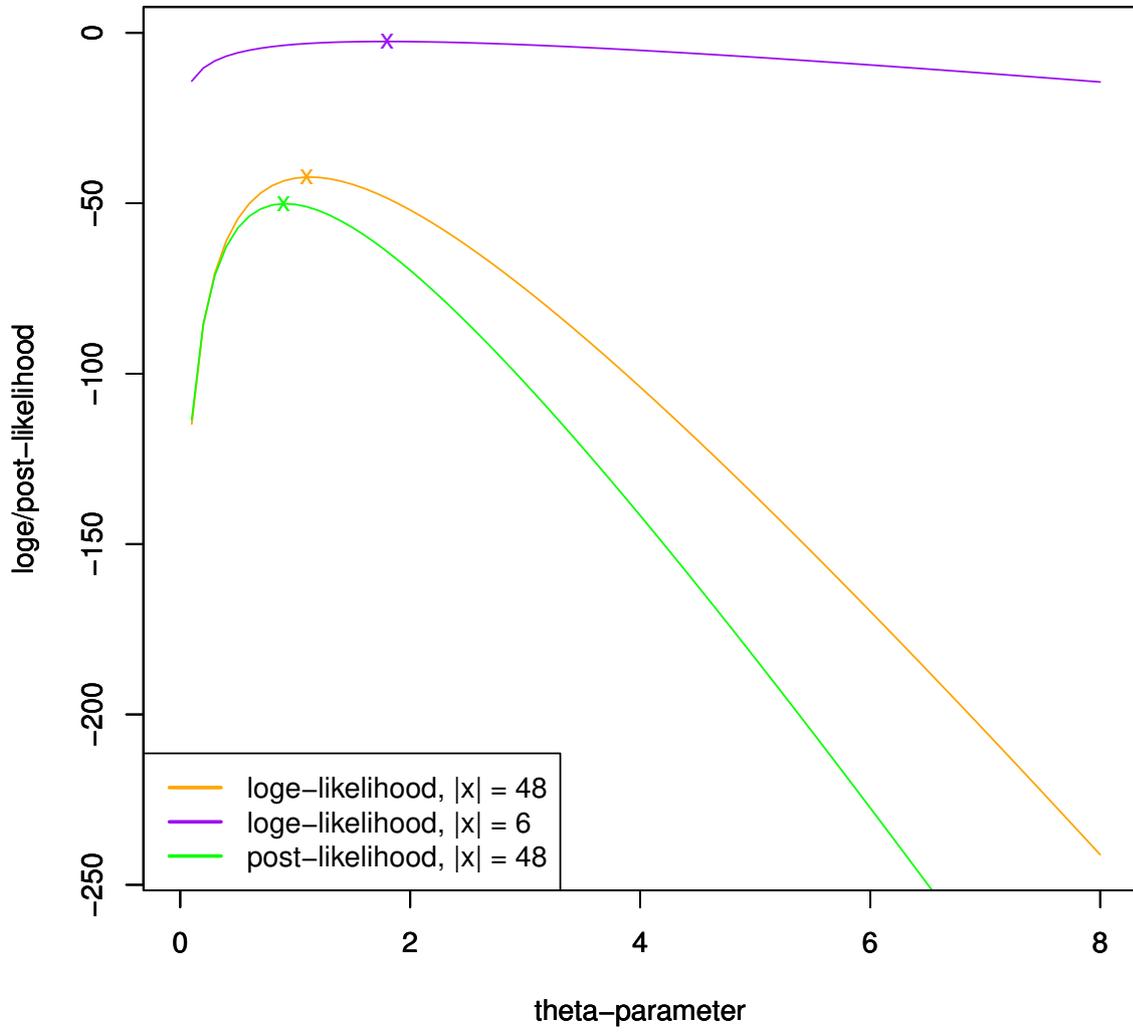
| Exp. Distribution | MLE/MAP |
|---|---|
| $p(\boldsymbol{x}|\theta), |\boldsymbol{x}| = 48$ | 1.1 |
| $p(\boldsymbol{x}|\theta), |\boldsymbol{x}| = 6$ | 1.8 |
| $p(\boldsymbol{x}|\theta)p(\theta), |\boldsymbol{x}| = 48$ | 0.9 |

If additional information is know about the distribution, in this case $p(\theta) = \lambda e^{-\lambda \theta}$ where $\lambda = 10$, the *MAP (Maximum a Posteriori)* estimation can be used instead. The implementation can be found in Listing 5, in the $R$ function `polikelihood`. The MAP for the $\boldsymbol{x}$ dataset is shown in the table, and the graph for it is the green line to the right. Finally, taking random samples from the distribution with $\hat{\theta}_{mle}$ produces the histogram to the right.

In conclusion, the parameter $\hat{\theta}_{mle}$ is a very good estimator for $\boldsymbol{x}$ given the Exp. distribution $p(\boldsymbol{x}|\theta)$.



Histogram over X–Exp of Full Log–Likelihood Estimation

**Max Log/Posteriori−Likelihood Estimation**

Legend:
- loge−likelihood, |x| = 48
- loge−likelihood, |x| = 6
- post−likelihood, |x| = 48

Axis labels: theta−parameter (x-axis), loge/post−likelihood (y-axis)

# References

[BGRS99] Kevin Beyer, Jonathan Goldstein, Raghu Ramakrishnan, and Uri Shaft. When is "nearest neighbor" meaningful? In *International conference on database theory*, pages 217–235. Springer, 1999.

[CH67] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27, 1967.

[JN99] Dag Jonsson and Lennart Norell. *Ett stycke statistik*. Studentlitteratur, 1999.

[Mur12] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT univeristy press, 1st edition, 2012.

[Myu03] In Jae Myung. Tutorial on maximum likelihood estimation. *Journal of mathematical Psychology*, 47(1):90–100, 2003.

# Appendix

Listing 1: Spam Prediction Script

```r
1   library("kknn")
2   source("distance.r")
3   source("knearest.r")
4
5   sensitivity <- function(x, y) {
6       tp <- sum(x == 1 & y == 1)
7       fn <- sum(x == 0 & y == 1)
8       return(tp / (tp + fn))
9   }
10
11  specificity <- function(x, y) {
12      tn <- sum(x == 0 & y == 0)
13      fp <- sum(x == 1 & y == 0)
14      return(tn / (tn + fp))
15  }
16
17  set.seed(12345) # For debugging.
18  data <- read.csv("spambase.csv")
19  # Pick randomly around half of the rows in dataset.
20  samples <- sample(1:nrow(data), floor(0.5*nrow(data)))
21  # Split given dataset evenly for training and tests.
22  learning <- data.matrix(data[samples,]) # Training.
23  testing <- data.matrix(data[-samples,]) # Testing.
24
25  # Predict spam for learning data K=5.
26  cat("\nknearest: predicting k = 5 (learning)\n")
27  k5 <- knearest(learning, 5, learning)
28  kr5 <- round(k5) # Classify >0.5 -> 1|0.
29  # Generate the confusion matrix for K=5.
30  cm5 <- table(kr5, learning[,ncol(learning)])
31  # Calculate given missclassification.
32  mc5 <- 1 - sum(diag(cm5)) / sum(cm5)
33  # Report confusion matrix and error.
34  print(cm5) ; print(mc5)
35
36  # Predict spam for learning data K=1.
37  cat("\nknearest: predicting k = 1 (learning)\n")
38  k1 <- knearest(learning, 1, learning)
39  kr1 <- round(k1) # Classify >0.5 -> 1|0.
40  # Generate the confusion matrix for K=1.
41  cm1 <- table(kr1, learning[,ncol(learning)])
42  # Calculate given missclassification.
43  mc1 <- 1 - sum(diag(cm1)) / sum(cm1)
44  # Report confusion matrix and error.
45  print(cm1) ; print(mc1)
46
47  # Predict spam for testing data K=5.
48  cat("\nknearest: predicting k = 5 (testing)\n")
49  k5 <- knearest(learning, 5, testing)
50  kr5 <- round(k5) # Classify >0.5 -> 1|0.
51  # Generate the confusion matrix for K=5.
52  cm5 <- table(kr5, testing[,ncol(testing)])
53  # Calculate given missclassification.
54  mc5 <- 1 - sum(diag(cm5)) / sum(cm5)
55  # Report confusion matrix and error.
56  print(cm5) ; print(mc5)
57
```

```
58  # Predict spam for testing data K=1.
59  cat("\nknearest: predicting k = 1 (testing)\n")
60  k1 <- knearest(learning, 1, testing)
61  kr1 <- round(k1) # Classify >0.5 -> 1|0.
62  # Generate the confusion matrix for K=1.
63  cm1 <- table(kr1, testing[,ncol(testing)])
64  # Calculate given missclassification.
65  mc1 <- 1 - sum(diag(cm1)) / sum(cm1)
66  # Report confusion matrix and error.
67  print(cm1) ; print(mc1)
68
69  cat("\nkknn: training and predicing with k = 1, 5\n")
70  m5 <- train.kknn(Spam ~ ., data = data.frame(learning), ks = c(5))
71  p5 <- predict(m5, data.frame(testing)) # Predict spam with k = 5.
72  pr5 <- round(p5) # Classify with the function >0.5 -> 1 else 0.
73  cm5 <- table(pr5, testing[,ncol(testing)])
74  # Calculate given missclassification.
75  mc5 <- 1 - sum(diag(cm5)) / sum(cm5)
76  # Report confusion matrix and error.
77  print(cm5) ; print(mc5)
78
79  cat("\nkknn: training and predicing with k = 1\n")
80  m1 <- train.kknn(Spam ~ ., data = data.frame(learning), ks = c(1))
81  p1 <- predict(m1, data.frame(testing)) # Predict spam with k = 1.
82  pr1 <- round(p1) # Classify with the function >0.5 -> 1 else 0.
83  cm1 <- table(pr1, testing[,ncol(testing)])
84  # Calculate given missclassification.
85  mc1 <- 1 - sum(diag(cm1)) / sum(cm1)
86  # Report confusion matrix and error.
87  print(cm1) ; print(mc1)
88
89  # Classify dataset by 0.05 steps...
90  response <- testing[,ncol(testing)]
91  classify <- seq(0.05, 0.95, by=0.05)
92  # Apply the classification rule for all.
93  kc5 <- sapply(k5, function(x) x > classify)
94  pc5 <- sapply(p5, function(x) x > classify)
95
96  # Find the sensitivity and specificity of knearest.
97  ksensitivity <- apply(kc5, 1, sensitivity, response)
98  kspecificity <- apply(kc5, 1, specificity, response)
99  psensitivity <- apply(pc5, 1, sensitivity, response)
100 pspecificity <- apply(pc5, 1, specificity, response)
101
102 plot(1 - kspecificity, ksensitivity, xlim=c(0.05,0.95), ylim=c(0.05,0.95), xlab="Specificity",
         ylab="Sensitivity", type='l')
103 lines(1 - kspecificity, ksensitivity, col="Orange") ; lines(1 - pspecificity, psensitivity,
         col="Purple")
104 legend(x = "bottomright", c("knearest", "kknn"), lty = c(1,1), lwd = c(2,2), col=c("Orange", "
         Purple"))
105 lines(0:1, 0:1, col="Red", xlim=c(0.05, 0.95), ylim=c(0.05,0.95))
106 title("ROC Curve for the K-NN Spam Predictors")
```

Listing 2: K-Nearest Neighbor Algorithm Implementation

```
1  source("distance.r") # cos-distance d.
2
3  # spam(i, t) - gets spam vector in i.
4  spam <- function(indices, training) {
5      spamid <- ncol(training) #  last.
6      return(mean(training[indices, spamid]))
```

```
 7  }
 8
 9  # knearest(t, k, t') - predicts values
10  # given in t', for training data in t.
11  # Done with using k-nearest neighbors.
12  knearest <- function(train, k, test) {
13      # Don't include the 'Spam' column, not feature.
14      test_features <- data.matrix(test[,-ncol(test)])
15      train_features <- data.matrix(train[,-ncol(train)])
16
17      # Compute the distance matrix between train and test
18      # using the cosine distance formula (see distance.r)
19      # which is then sorted, for picking the k-neighbors.
20      distances <- distance(train_features, test_features)
21      sorted_distance_ids <- as.matrix(t(apply(distances, 2, order))[,1:k])
22
23      # Finally, retrieve if the training data is spam or
24      # not, selecting the k-closest classifications, for
25      # later determining the most likely classification.
26      spamv <- apply(sorted_distance_ids, 1, spam, train)
27      kspam_vector <- spamv # Select only first K.
28      mean_spam <- data.matrix(kspam_vector)
29      # Still need to classify data by e.g. >0.5 -> 1.
30      return(mean_spam) # This step is done in spam.r.
31  }
```

Listing 3: Cosine Cost/Distance Formula

```
 1  # distance(X, Y) - distances between X, Y.
 2  # Uses the usual cosine distance function.
 3  # Batch operation into a matrix -> fast...
 4  distance <- function(matrix_x, matrix_y) {
 5      x_squared_sum <- rowSums(matrix_x^2)
 6      y_squared_sum <- rowSums(matrix_y^2)
 7      x_prime <- matrix_x / sqrt(x_squared_sum)
 8      y_prime <- matrix_y / sqrt(y_squared_sum)
 9      similarity_matrix <- x_prime %*% t(y_prime)
10      distance_matrix <- 1.0 - similarity_matrix
11      return(distance_matrix)
12  }
```

Listing 4: Inference Script for Machine Lifetime

```
 1  source("likelihood.r") # sum of ln(p(x|theta)).
 2  lifetimes <- read.csv("machines.csv") # Matrix?
 3  parameter <- seq(0.1, 8.0, by=0.1) # Testing...
 4  # Apply each parameter theta individually gives
 5  # the log-likelihood for each of the parameters
 6  p <- sapply(parameter,lnlikelihood,x=lifetimes)
 7  average_lnlikelihoods <- p / dim(lifetimes)[1];
 8  mle <- order(average_lnlikelihoods)[length(p)];
 9  mle <- mle * 0.1 ; cat("MLE(theta): ",mle,"\n")
10
11  # Plot the relation between theta = logl.
12  plot(parameter,p, xlab="theta-parameter",
13       ylab = "loge/post-likelihood", type = "l",
14       xlim=c(0,8), ylim=c(-242,-2), col="orange")
15  points(mle, lnlikelihood(lifetimes, mle),
```

```
16          col="orange", lwd=c(2, 2), pch="x");
17  title("Max Log/Posteriori-Likelihood Estimation")
18
19  lifetimes6 <- t(data.matrix(lifetimes)[1:6])
20  p6 <- sapply(parameter,lnlikelihood,x=lifetimes6)
21  average_lnlikelihoods6 <- p6 / dim(lifetimes6)[1]
22  mle6 <- order(average_lnlikelihoods6)[length(p6)]
23  mle6 <- mle6 * 0.1 ; cat("MLE(theta): ",mle6,"\n")
24
25  # Plot the relation between theta = ln-l.
26  par(new=TRUE) # Seems a little bit hacky.
27  plot(parameter,p6, xlab="theta-parameter",
28       ylab = "loge/post-likelihood", type = "l",
29       xlim=c(0,8), ylim=c(-242, -2), col="purple")
30  points(mle6, lnlikelihood(lifetimes6, mle6),
31          col="purple", lwd=c(2, 2), pch="x")
32
33  po <- sapply(parameter, polikelihood, x=lifetimes)
34  average_polikelihoods <- po / dim(lifetimes)[1]
35  mpe <- order(average_polikelihoods)[length(po)]
36  mpe <- mpe * 0.1 ; cat("MPE(theta): ",mpe,"\n")
37
38  # Plot the relation between theta = po-li.
39  par(new=TRUE) # Seems a little bit hacky.
40  plot(parameter,po, xlab="theta-parameter",
41       ylab = "loge/post-likelihood",type="l",
42       xlim=c(0,8), ylim=c(-242, -2),
43       col = "Green")
44  points(mpe, polikelihood(lifetimes, mpe),
45          col="Green", lwd=c(2, 2), pch="x")
46  legend(x = "bottomleft", c("loge-likelihood, |x| = 48",
47                             "loge-likelihood, |x| = 6",
48                             "post-likelihood, |x| = 48"),
49       lty = c(1,1), lwd = c(2,2),
50       col=c("Orange", "Purple", "Green"))
51
52  random_exponential <- rexp(50, mle)
53  hist(random_exponential, main="Histogram over X~Exp of Full Log-Likelihood Estimation",
54       xlab = "Expected Machine Lifetime", ylab="Frequency of Machine Lifetime")
```

Listing 5: Max Log- and Posteriori-Likelihood Estimation Formula

```
1  distribution <- function(x, theta) {
2      # An exponential distribution.
3      exponential <- exp((-theta)*x)
4      return(theta*exponential)
5  }
6
7  lnlikelihood <- function(x, theta) {
8      p <- log(distribution(x, theta))
9      return(sum(p)) # log-likelihood.
10 }
11
12 polikelihood <- function(x, theta) {
13     jp<-prod(distribution(x, theta))
14     posterior <- 10*exp(-10*theta)
15     return(log(jp*posterior))
16 }
```

Listing 6: The Given Machine Lifetime CSV Dataset (Excerpt)

```
 1  Length
 2  0.394761404022574
 3  1.17680263974688
 4  0.768466353536656
 5  0.126129498627658
 6  0.053941871673606
 7  0.839961192069958
 8  2.83505971839929
 9  1.2602572965046
10  4.41829496504448
11  0.737917928761447
12  0.282883279724047
13  0.405573239549994
14  1.38489578934096
15  1.41224693846584
16  1.81130825686632
17  1.58639749349854
18  1.0564000190196
```

# Introduction to Machine Learning
# Individual Laboration Report –2–

Erik Sven Vasconcelos Jansson
erija578@student.liu.se
Linköping University, Sweden

November 15, 2016

## Assignment 1

Reducing the amount of relevant features by using *feature selection* is an important task in *supervised machine learning*. Since many features $\mathcal{F}_i$ are more relevant than others, producing the optimal feature set $\hat{\mathcal{F}} \subseteq \mathcal{F}$, reducing $|\mathcal{F}|$, while also $\hat{\varepsilon}(\hat{\mathbf{y}}, \mathbf{y})$, the *Mean Squared Error (M.S.E)*. Here, we build a *brute-force feature selection* for *linear models*, using *k-folds cross validation* for ranking feature subsets.

First, every possible feature combination $\mathcal{F}_i \subseteq \mathcal{F}$ is generated. Thereafter, each $\mathcal{F}_i$ is tested through *k-fold cross-validation*, giving the *mean* $\hat{\varepsilon}_i(\hat{\mathbf{y}}, \mathbf{y})$ of the feature subset $\mathcal{F}_i$. By picking the feature subset $\mathcal{F}_i$ which produces $\min_i \hat{\varepsilon}_i(\hat{\mathbf{y}}, \mathbf{y})$, the best features are picked. In Listing 3 line 26 we generate all $\mathcal{F}_i$, which are then cross-validated in lines 30−31, then in lines 34−43 the best feature subset $\hat{\mathcal{F}}$ is given by evaluating the errors $\min_i \hat{\varepsilon}_i(\hat{\mathbf{y}}, \mathbf{y})$, where $\hat{\mathcal{F}} = \mathcal{F}_i$.

Now, how *k-fold cross validation* works is shown. Roughly, Algorithm 1 demonstrates these steps, by giving each individual $\mathcal{F}_i$ and respective $X_{\mathcal{F}_i}$, $\mathbf{y}_{\mathcal{F}_i}$ as arguments, the *feature matrix* and *target vector*. For each $\mathcal{F}_i$ and *k-fold iteration*, a *linear hypothesis function* is trained, predicting $\hat{\mathbf{y}}$ using Equations 1. Thereafter, the *M.S.E* of the prediction $\hat{\varepsilon}_i(\hat{\mathbf{y}}, \mathbf{y})$, is calculated by using Equation 3. Finally, the *mean* of these $\hat{\varepsilon}_i(\hat{\mathbf{y}}, \mathbf{y})$ is the result of $\mathcal{F}_i$ cross-validation.

$$\hat{\mathbf{w}}_t = (X_t^\mathsf{T} X_t)^{-1} X_t^\mathsf{T} \mathbf{y}_t \quad (1)$$

$$\hat{\mathbf{y}}_v = X_v \hat{\mathbf{w}}_t \quad (2)$$

$$\hat{\varepsilon}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2 \quad (3)$$

---

**Algorithm 1** K-Fold Cross-Validation (Linear $\mathcal{M}$)

---

**Require:** feature matrix $X_\mathcal{F}$ and target vector $\mathbf{y}_\mathcal{F}$, given a feature selection $\mathcal{F}$ with cardinality $|\mathcal{F}|$.

1: $(X_i, \mathbf{y}_i) \leftarrow \text{split}(X_\mathcal{F}, \mathbf{y}_\mathcal{F}, k)$ {Equally $|X_\mathcal{F}| \div k$}
2: **for** $i \leftarrow 1$ **to** $k$ **do** {Attempts every of $k$-folds}
3: $\quad X_t \leftarrow X_1 \cup \cdots \cup X_k - X_i$ {Except fold $i$}
4: $\quad \mathbf{y}_t \leftarrow \mathbf{y}_1 \cup \cdots \cup \mathbf{y}_k - \mathbf{y}_i$ {Except fold $i$}
5: $\quad \hat{\mathbf{w}}_t \leftarrow (X_t^\mathsf{T} X_t)^{-1} X_t^\mathsf{T} \mathbf{y}_t$ {Train model}
6: $\quad \hat{\mathbf{y}}_i \leftarrow X_i \hat{\mathbf{w}}_t$ {Predict target vector}
7: $\quad \hat{\varepsilon}_i(\hat{\mathbf{y}}_i, \mathbf{y}_i) \leftarrow \frac{1}{n} \sum_{j=1}^{n} (\hat{y}_j - y_j)^2$
8: **end for**
9: **return** $(\sum_{i=1}^{k} \hat{\varepsilon}_i(\hat{\mathbf{y}}_i, \mathbf{y}_i)) \div k$

---

The implementation of *feature selection* is found under Listing 3, while *k-fold cross-validation* should be found in Listing 2, *linear regression* in Listing 1. Most information was derived from *Andrew Ng's Regularization & Model Selection Handouts* [Ng16].

Finally, testing the *swiss dataset* on our *feature selection* implementation where $\mathcal{F} = \mathcal{U} -$ Fertility, gives $\hat{\mathcal{F}} = \{3, 4, 5\}$, also shown in the Table 1 below. Our $\hat{\mathcal{F}}$ seems reasonable, being *Catholic* is usually attributed with *low abortion rate*, *Infant Mortality* is directly linearly related as can be seen in the plot.

| Feature Selection | M.S.E. |
|---|---|
| $\{3, 4, 5\}$ | 90.6202 |
| $\{1, 3, 4, 5\}$ | 111.411 |
| $\{1, 2, 3, 5\}$ | 117.092 |

Table 1: where $1 \rightarrow$ *Agriculture*, $2 \rightarrow$ *Examination*, $3 \rightarrow$ *Education*, $4 \rightarrow$ *Catholic*, $5 \rightarrow$ *Infant Mortality*.
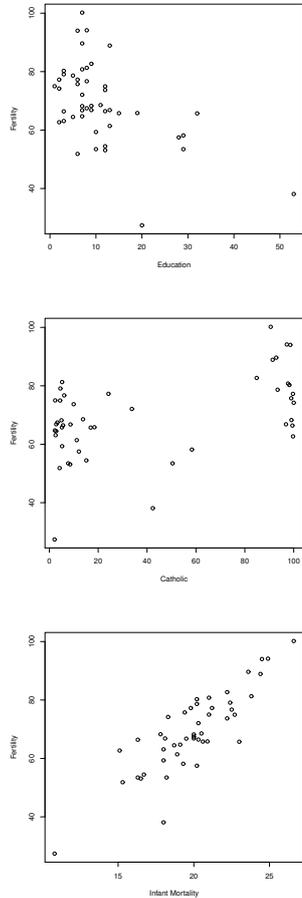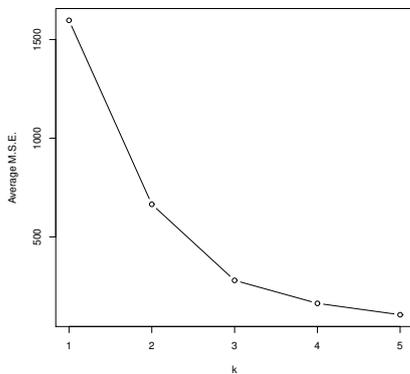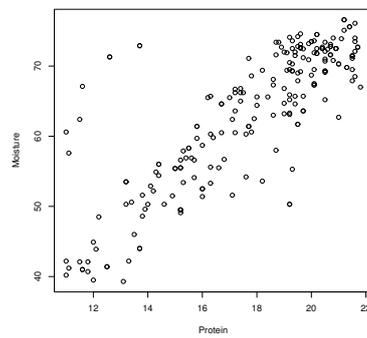
Figure 1: displays the relationships between feature subsets of *Size K* and their *Average M.S.E.* Notice that the *expected error decreases* by *increases in k*.



# Assignment 2

By using the *tecator dataset*, produced from observations on *predicting fat content* simply by using an *infrared absorbance spectrum* (w. several channels), we try to find relationships between their features.
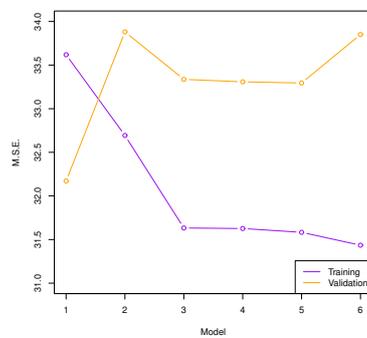
Figure 2: relationship seems very close to linear. *Protein* and *Moisture* seem to be related somehow, at least in the provided meat observation datasets.



Now we consider models $M_i$ where *Moisture* is *normally distributed* and related to *Protein* with a *polynomial function* of steadily increasing degree. Equation 4 displays this relation, MSE in Figure 3.

$$\text{Moisture} \sim \mathcal{N}(w_0 + \sum_{i=1}^{6} (w_i^i \cdot \text{Protein}_i), \sigma^2) \quad (4)$$

Figure 3: increasing complexity for the $i^{th}$ *Model* seems to generate *bias* towards the *training dataset*, motivating the higher *validation dataset* error rate.
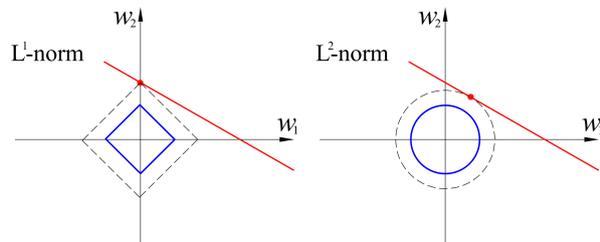
It seems the model *overfits* with more *complexity*. The models $M_i$ and their *M.S.E.* are calculated in lines 24-39 in Listing 4, the assignment 2 scripts. Also, according to the plot, the best model $M_i$ is located somewhere between $i = 1$ and 2, where the lines intersect, producing the lowest M.S.E. for *both datasets*, meaning there isn't *bias* towards either *D*.

Now we perform *feature/variable selection* on the *features Channel1-100* and *target/response* of *Fat*. By using *stepAIC* on a linear model generated with *R's lm*, a total of *64 variables* were selected here. This is done in lines 66-67 in Listing 4 (*w. MASS*).

By using the *Ridge regression model* on the same *features* and *response variables* with *glmnet library*. Figure 5 shows how *coefficients* relate to the *log $\lambda$*, *the log penalty*. Notice how all coefficients *converge uniformly* as *higher penalty* is added to the model. This is being done in line 72 with $\lambda = 0.0$ (Ridge).
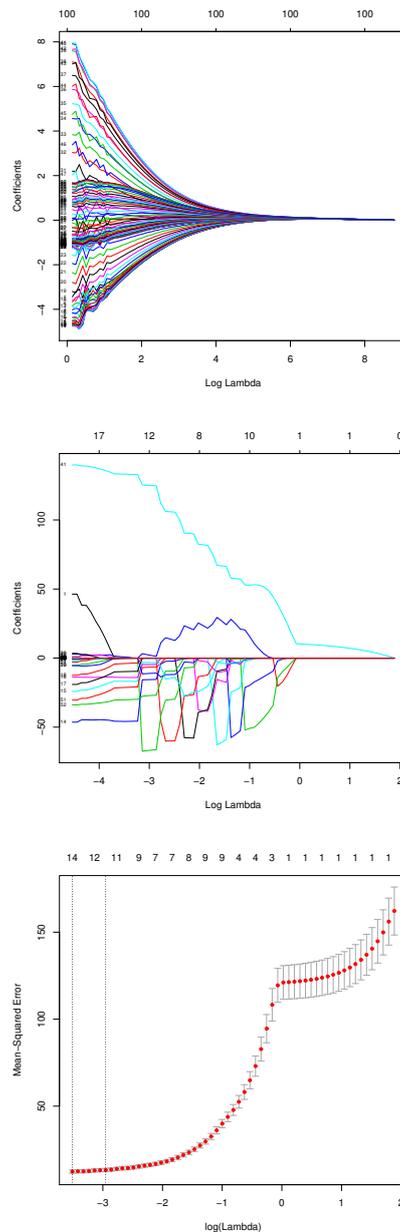
Thereafter, we use the *LASSO regression model* with the same *features and response variables*. See in Figure 5 how the *coefficients* converge *iteratively* instead of *simultaneously*, comparing toward *Ridge*. This can be explained graphically within Figure 4.

Figure 4: on the left we have *Lasso* and right *Ridge*. Notice, that *Lasso* can set $w_1 \rightarrow$ zero much faster. Distributed under CC by Rezamohammadighazi.

Finally, we *cross-validate* the *Lasso model* with *k-fold cross-validation* using the *glmnet library*. By using $k = 20$ and using *M.S.E.* as the metric, we find the last plot in Figure 5 which shows that the chosen features (the dotted interval) are those that produce the lowest M.S.E. Which are the 14 selected features. This is shown in the lines 96-97 of Listing 4. Lastly, we compare these results with those found with *stepAIC*. Notice how the *Lasso model* with *k-fold C.V.* produced a much lower number of selected features, 14, while *stepAIC()* produced a total of 64. *Lasso* gives us less features.

Figure 5: below are the relationships between the *coefficients* and the $log(\lambda)$ *penalty* for *Ridge* and *Lasso regression* (shown below in that exact order). Notice how *Ridge* converges all of the *coefficients simultaneously* while *Lasso* does this step *iteratively*, therefore, *Lasso regression* should converge "faster". Finally, the bottom plot displays how a *Lasso C.V.* relates to the increasing of *penalty factor* of $log(\lambda)$.

# References

[HQ14]   Trevor Hastie and Junyang Qian.   *Glm-net Vignette.* `stanford.edu/~hastie/glmnet/`, 2014. [Online; on 15/11/2016].

[Ng16]   Andrew Ng. *Handouts: Regularization and Model Selection.* Stanford University, 2016.

# Appendix

Listing 1: Estimation of the Linear Regression Model with a Hat Matrix

```
1   # ------linrhat(X, y)------
2   # Predicts the parameters w
3   # for the given features X,
4   # and the targets y through
5   # the use of an hat matrix.
6   linrhat <- function(X, y) {
7       # Nice hat good sir...
8       return(solve(t(X)%*%X)
9              %*% t(X)%*%y)
10  }
```

Listing 2: Implementation of a K-Fold Cross-Validation Algorithm for $\mathcal{M}$

```
1   source("linrhat.r")
2
3   # ------disjoin(X, k)------
4   # Produce a k disjoint sets
5   # of the training matrix X.
6   # Useful for kfoldcv below.
7   # Note: only returns index.
8   disjoin <- function(X, k) {
9       row <- 1:nrow(X)
10      folds <- nrow(X) / k
11      # Resulting disjoints.
12      S <- matrix(, k, folds)
13      U <- c() # Picked sets.
14      for (fold_row in 1:k) {
15          D <- setdiff(row, U)
16          Si <- sample(D, folds)
17          S[fold_row,] <- Si
18          U <- union(U, Si)
19      }
20
21      return(S)
22  }
23
24  # ------egerror(x, y)------
25  # Locate the generalization
26  # error within in our model
27  # by comparing the results:
28  # x and y targets. Returns:
29  # the difference for x & y.
30  egerror <- function(x, y) {
31      targets <- length(x)
32      # Using good old MSE...
33      sdiff <- sum((x - y)^2)
34      return(sdiff / targets)
35  }
36
37  # ------kfoldcv(X, y, k)-------
38  # Returns the estimated genera-
39  # lization error of the feature
40  # set according to a linear mo-
41  # del. It does this by applying
42  # an k-folding cross validation
```

```
43  # method, splitting X randomly,
44  # gives k disjoint subets of X.
45  kfoldcv <- function(X, y, k)  {
46      kfolding <- 1:k
47      sets <- disjoin(X, k)
48      ege <- c() # Empty set of errors.
49      for (i in kfolding) { # Every set.
50          kset <- sets[-i,] # Remove 'i'
51          iset <- sets[i,]  # Only 'i'.
52          # Pick dataset for all but 'i'
53          Xi <- X[kset,] ; yi <- y[kset]
54          # Estimated parameters w. Xi.
55          hypothesis <- linrhat(Xi, yi)
56          # Predict for the unused 'i'.
57          p <- X[iset,]%*%hypothesis
58          # Esimate error of this.
59          e <- egerror(p, y[iset])
60          # Add to list of these.
61          ege <- c(ege, abs(e))
62      }
63
64      # Oh yea baby,
65      # average errors.
66      return(mean(ege))
67  }
```

Listing 3: Brute-Force Feature Selection to Find Lowest M.S.E. Subset

```
1   library("ggplot2")
2   source("kfoldcv.r")
3
4   # ------kfold(i, X, y, k)------
5   # Wrapper for applying kfoldcv,
6   # since it needs to apply rows.
7   kfold <- function(i, X, y, k) {
8       Xi <- data.matrix(X[,i]);
9       return(kfoldcv(Xi, y, k))
10  }
11
12  # ------featsel(X, y, k)------
13  # Applies feature selection to
14  # X, giving the best subset of
15  # X which minimizes the error.
16  # This is done by applying the
17  # k-fold cross validation, for
18  # each possible subset of X's.
19  featsel <- function(X, y, k) {
20      features <- c() # Nothing.
21      lowerror <- Inf # Not good
22      relation <- data.frame(K   = 1:k,
23                             MSE = c(0))
24      for (i in 1:k) { # Wrapper
25          # Produce all combins.
26          fi <- t(combn(1:k, i))
27          # Apply combinations for
28          # each of the featureset
29          # giving the error value
30          ei <- apply(fi, 1, kfold,
31                      X, y, k)
32
33          # Order by error...
```

```
34          err <- ei[order(ei)]
35          fea <- fi[order(ei),]
36          fea<-data.matrix(fea)
37
38          relation$MSE[i]<-mean(err)
39          # Update best estimates.
40          if (lowerror > err[1]) {
41              lowerror <- err[1];
42              features <- fea[1,]
43          }
44      }
45
46      setEPS()
47      postscript("kvsmse.eps")
48      # Plot relation k v.s. M.S.E.
49      plot(relation$K, relation$MSE,
50          type="both", xlab="k",
51          ylab = "Average M.S.E.")
52      dev.off() # Write...
53      # Best features.
54      return(features)
55  }
56
57  set.seed(12345)
58  X <- data.matrix(swiss[,-1])
59  y <- data.matrix(swiss[,1]);
60  features <- featsel(X, y, 5)
61  X <- data.matrix(X[,features])
62  what <- linrhat(X, y)
63  yhat <- X %*% what
64
65  graph <- data.frame(X)
66  graph$Fertility <- yhat
67
68  setEPS()
69  postscript("education.eps")
70  plot(graph$Education, graph$Fertility,
71      xlab="Education", ylab="Fertility")
72  dev.off()
73
74  postscript("catholic.eps")
75  plot(graph$Catholic, graph$Fertility,
76      xlab="Catholic", ylab="Fertility")
77  dev.off()
78
79  postscript("mortality.eps")
80  plot(graph$Infant.Mortality, graph$Fertility,
81      xlab="Infant Mortality", ylab="Fertility")
82  dev.off()
```

Listing 4: Script for Assignment 2 on Linear Models, Ridge, Lasso and Cross-Validation

```
1  library("MASS")
2  library("glmnet")
3  data <- read.csv("tecator.csv");
4  data_matrix <- data.matrix(data)
5
6  setEPS()
7  postscript("linear.eps")
8  plot(data$Protein, data$Moisture,
9      xlab="Protein", ylab="Moisture")
```

```
10  dev.off()
11
12  # Do you think these data are described well by a linear model?
13  # Answer: yes, definitely. Both data are rising similar ratios.
14  # Moisture ~ N(mu = Xw, sigma^2), a linear model should work...
15  # Moisture_hat ~ Mi = polynomial function dependent on Protein:
16  #                    w0 + x1*w1 + x2*w2^2 + ... xi*wi^i = f(x)
17
18  set.seed(12345) # Set seed for getting same results...
19  indices <- sample(1:nrow(data), floor(nrow(data)*0.5))
20  training <- data[indices,] # Subset for training data.
21  validation <- data[-indices,] # Subset for validation.
22  validation <- validation[-nrow(validation),] # Shhh...
23
24  polynomials <- 1:6
25  imse <- data.frame(Model = polynomials,
26                     Training.MSE = c(0.0),
27                     Validation.MSE = c(0))
28  for (degree in polynomials) {
29      model <- lm(Moisture ~ poly(Protein, degree),
30                  training)
31      tprediction <- predict(model, training)
32      vprediction <- predict(model, validation)
33      vmse <- (vprediction - validation$Moisture)^2
34      tmse <- (tprediction - training$Moisture)^2
35      vmse <- mean(vmse) # MSE for validation set.
36      tmse <- mean(tmse) # MSE for training set...
37      imse$Validation.MSE[degree] = vmse # For i.
38      imse$Training.MSE[degree] = tmse # For i.
39  }
40
41  setEPS()
42  postscript("depends.eps")
43  plot(imse$Model, imse$Training.MSE,
44       xlab="Model", ylab="M.S.E.", "b", col="purple",
45       ylim=c(31.0, 34.0))
46  points(imse$Model, imse$Validation.MSE, col="orange",
47       type = "b", ylim=c(31.0, 34.0))
48  legend("bottomright", legend=c("Training", "Validation"),
49         col=c("purple", "orange"), lty=1)
50  dev.off()
51
52  # Which model is best according to the plot?
53  # Answer: model where i = 1, since the validation predictions
54  # are becoming more erronous as i gets larger, while training
55  # becomes more accurate (this is caused by overfitting data).
56  # How do the M.S.E. values change and why?
57  # Answer: the error for the training set seems to become less
58  # as i gets larger while validation becomes more error prone.
59  # This is caused because the model overfits,  with a bias to-
60  # wards the training set (since the model is based on these).
61  # Ppecifically, it overfits:  the models become more complex.
62  # Interpret this picture in terms of bias-variance tradeoffs.
63  # Answer: it seems to be biased towards the training dataset.
64  # Since the E[yhat(x0) - f(x)] isn't close to zero, biased...
65
66  model <- lm(Fat ~ . - Moisture - Protein, data)
67  feature_selection <- stepAIC(model, direction="both")
68  # How many were selected? Answer: 64,  coeffiecients.
69
70  y <- data_matrix[,102] # Select Fat as the only response...
71  X <- data_matrix[,2:101] # Select Channel1-Channel-100 features.
```

```
72  ridge <- glmnet(X, y, alpha = 0) # Fit with the Ridge regression.
73  lasso <- glmnet(X, y, alpha = 1) # Fit with the Lasso regression.
74
75  setEPS()
76  postscript("ridge.eps")
77  # Explicit about x-axis variable.
78  plot(ridge, xvar="lambda", label=TRUE)
79  dev.off()
80
81  # Report on how the coeffiecients change with lambda.
82  # Answer: ridge penalizes all features equally with lambda.
83  # Therefore, it will take longer to coverge all features...
84
85  setEPS()
86  postscript("lasso.eps")
87  # Explicit about x-axis variable.
88  plot(lasso, xvar="lambda", label=TRUE)
89  dev.off()
90
91  # Conclusions on the Ridge vs Lasso resulting plots?
92  # Answer: the Lasso regression seems to converge individual
93  # features, therefore converging faster, towards the values
94  # of some features. While Ridge converges simultaniously...
95
96  kfoldcv <- cv.glmnet(X, y,  type.measure="mse", nfolds=20)
97  feature_selection <- coef(kfoldcv, s = "lambda.min")
98
99  # Report the optimal lambda and how many variables were selected.
100 # Answer: the optimal lambda was 0.02985605 and 14 were selected.
101 # Conclusions: the interval area shown in the graph shows optimal
102 # number of feature selections,  basically with the lowest M.S.E.
103
104 setEPS()
105 postscript("kfold.eps")
106 plot(kfoldcv)
107 dev.off()
108
109 # Compare the results from steps (4) and (7). Basically, compare
110 # stepAIC() and the glmnet CV (using k-fold).    Answer: stepAIC
111 # gives a lot of features, 64 for the selection, while glmnet cv
112 # for Lasso gives 13 features. Therefore, stepAIC() seems to not
113 # penalize the features very well, and therefore chooses more...
```

# Introduction to Machine Learning
# Individual Laboration Report –3–

Erik Sven Vasconcelos Jansson
erija578@student.liu.se
Linköping University, Sweden

November 27, 2016

## Assignment 1

Now that *linear regression* and *cross-validation* have been studied for *regression problems* (involving continuous target variables), the question is how to apply the same concept to *classification problems* (which involve a discrete amount of target variables). Both *Linear Discriminant Analysis (LDA)* and *Logistic Regression* can be used to achieve this, and also to draw *decision boundaries*.

In this task the goal is to *classify* the *Sex* of *Australian Crabs* by using their *Carapace Length* and *Rear Width* using *LDA* and thereafter plotting the *decision boundary*. Plotting the data gives Figure 1. As can be seen, a decision boundary seems to exist.

Firstly, the *linear parameters* $w_0$ and $\boldsymbol{w}$ need to be found for each class $k$, essentially creating our *target hypothesis function* for each one of these $k$. Finding these parameters is done with Equations 1, essentially computing the mean $\hat{\boldsymbol{\mu}}_k$ for each class $k$ which is used to derive the *covariance matrices* $\Sigma_k$.

$$
\begin{aligned}
\hat{\pi}_k &= \frac{N_k}{N} \\
\hat{\mu}_k &= \frac{1}{N_k} \sum_{i \in k} \boldsymbol{x}_i \\
\Sigma_k &= \frac{1}{N_k} \sum_{i \in k} (\hat{\boldsymbol{\mu}}_k - \boldsymbol{x}_i)(\hat{\boldsymbol{\mu}}_k - \boldsymbol{x}_i)^{\mathsf{T}} \\
\hat{\Sigma} &= \frac{1}{N} \sum_{k \in K} N_k \Sigma_k = \sum_{k \in K} \frac{\operatorname{cov} X_k}{|X_k|}
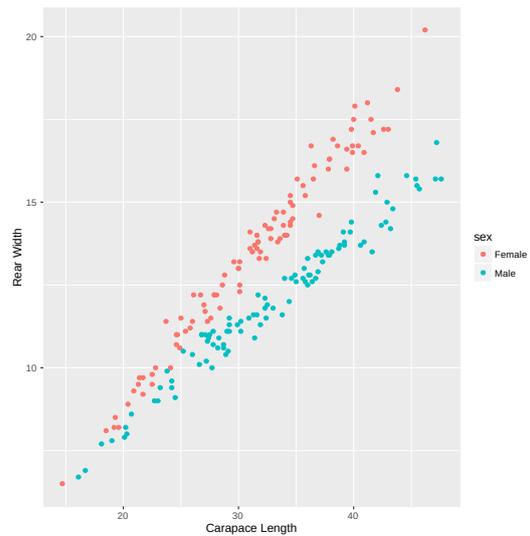\end{aligned}
\tag{1}
$$

$$
\begin{aligned}
w_k &= -\frac{\hat{\boldsymbol{\mu}}_k^{\mathsf{T}}}{2} \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_k + \log \hat{\pi}_k \\
\boldsymbol{w}_k &= \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_k
\end{aligned}
\tag{2}
$$

Finally, computing $\hat{\Sigma}$ is done by using these $\Sigma_k$, which is thereafter used to calculate $w_0$ and $\boldsymbol{w}$, $\forall k$. Implementation of these formulae in $R$ can be found in Listing 1 under lines 15-39 in the lda function. All equations are in *Friedman et al's* [FHT09] book.

$$
\delta_k(\boldsymbol{x}) = w_0 + \boldsymbol{w}
\tag{3}
$$

Using the *Discriminant Function* in Equation 3 above, one can plot the line through each class $k$. The *decision boundary* is in-between these 2 lines. In lines 58-61 we calculate the *intercept* and the *slope* of the *decision boundary*, which gives us all the information needed to produce a new Figure 2. The fit seems perfect, exactly as previous Figure 1.

Figure 1: Sex Classifications of Australian Crabs



1

Thereafter, by using *Logistic Regression* instead, Figure 3 can be obtained, which misses predictions.

| Class | $w_0$ | $w$ |
|--------|--------|-----|
| Female | -22.42 | (-2.161, 8.248) |
| Male | -12.42 | (-0.213, 2.565) |

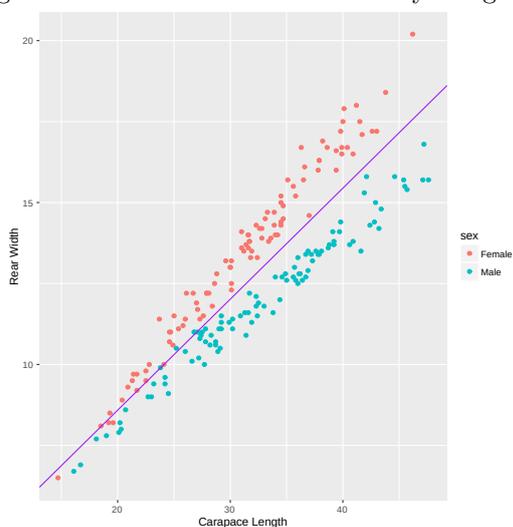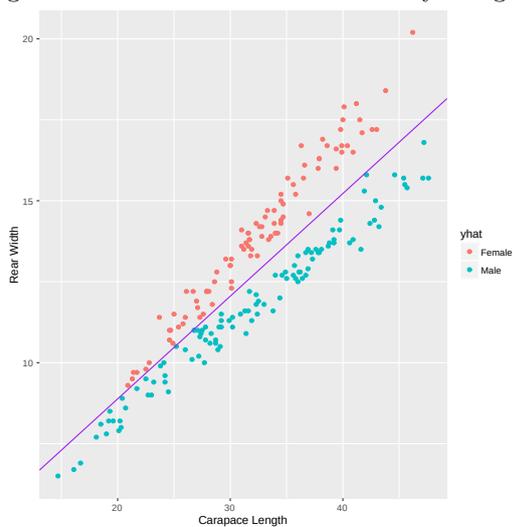Table 1: Predicted Parameters for Australian Crab

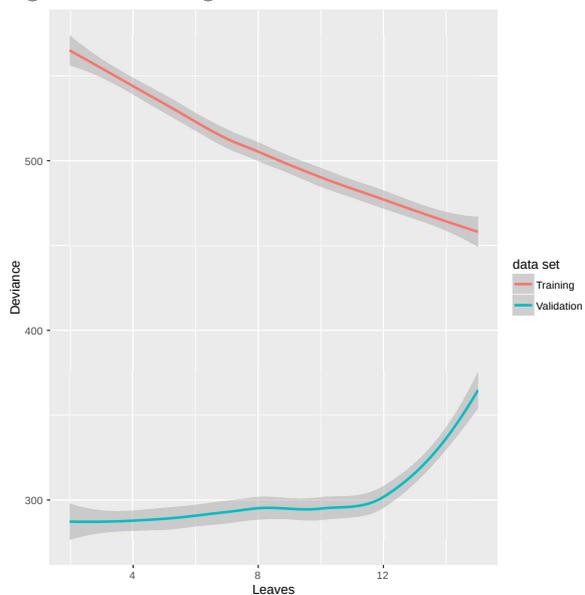Figure 2: Crab Sex Decision Boundary using LDA



# Assignment 2

Given several *observations* from different *customers* with certain *features* and *classifications* of whether they have *managed their loans good* or *bad*, we are tasked to *predict* if a new customer with certain set of features, will pay back their loans on time or not.

By using *decision tree learning*, which maps observations (the branches) to their conclusions (the leaves), we can predict this aforementioned model. In Listing 2 under the lines 26-45, we use `tree` to fit our model for usage on *training & testing* sets. These models can be fit using different measures of *impurity*, where we only consider *deviance* and the *gini index* here. In Table 2 are the results for these. It seems that *gini and/or deviance* classifies better.

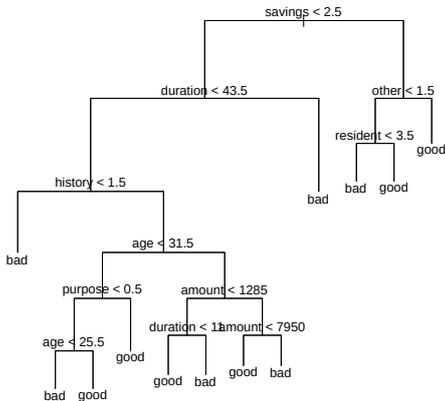| Data Set | Impurity Measure | Miss Rate |
|----------|------------------|-----------|
| Training | Deviance | 0.212 |
| Testing | Deviance | 0.236 |
| Training | Gini | 0.230 |
| Testing | Gini | 0.282 |
| Training | Deviance & Gini | 0.212 |
| Testing | Deviance & Gini | 0.236 |

Table 2: Decision Tree Missclassification Rates

Figure 4: Training and Validation Deviance Values



Figure 3: Crab Sex Decision Boundary using LR

Afterwards, we use the *training* and *validation data sets* to choose the *optimal tree depth*. Following lines `47-66`, we iteratively *prune* the decision tree, essentially adding more *leaves*, and thereafter plot the dependence on the *number of leaves* and the *estimated deviance of the model*. This plot can be seen in Figure 4. As can be seen, when more than 12 leaves are selected, the *validation set* has a pretty large increase in *deviance*, which isn't the desired behaviour. On the other hand, the *training data set* has a constant decrease in deviance as the number of leaves increase. Therefore, the best number of leaves seems to be 12, giving best of both worlds. Additionally, the depth of this optimal tree is 6 and has a misclassification rate of 0.24 with the following features selected: *savings, duration, history, age, purpose, amount, other and resident*, all of which can be found easy with `summary(tree)`.

Finally, the decision tree is plotted in Figure 5. Most of the branches and their outcomes seem quite reasonable (with common sense), a person that has very few savings and has a history of usually delays loans, will probably be a bad customer who doesn't pay back them. While a person which has savings and is a resident will probably be more responsible and pay back his loans. It seems that the tree has been able to select the most important features as the branches of the decision tree, which are highly correlated with either classification result (the leaves). The depth can be seen here too.

Figure 5: Visualization of the Decision Tree



Now we proceed to fit another model with the same situation, but using *Naïve Bayes* instead, which basically assumes that all *features* are independent, like so: $p(C_k|x_1, x_2, ..., x_n) \propto p(x_i|C_k)$.

Under lines `82-90` we fit the model with `e1071`'s `naiveBayes` *classification*, and produce the *confusion matrices* and *missclassifications* below for both the *training* and the *testing data sets*.

*Missclassification: train* = 0.300 & *test* = 0.306.

| train | bad | good | | test | bad | good |
|---|---|---|---|---|---|---|
| **bad** | 95 | 98 | | **bad** | 142 | 147 |
| **good** | 52 | 255 | | **good** | 83 | 378 |

Table 3: Normal Naïve Bayes Confusion Matrix

It seems that in this case, *Naïve Bayes* performs worse than *decision trees* since the missclassification rate is quite a bit higher. Now, let's assume that we are given a *loss matrix*, a way of *penalizing* the *predictor* to apply weight to certain negatives, which in this case is $L = \left( \begin{smallmatrix} 0 & 1 \\ 10 & 0 \end{smallmatrix} \right)$. Doing this in lines `92-109` results in a prediction giving the *confusion matrix* below and missclassifications of train = 0.274 and test = 0.278. These missclassifications are a lot lower than those in the normal "lossless" model. The reason why this happens is because the cost of doing a faulty classification is highly penalized in one case (while the other isn't penalized as much), this is a underlying reason for these results.

| train | bad | good | | test | bad | good |
|---|---|---|---|---|---|---|
| **bad** | 27 | 17 | | **bad** | 40 | 24 |
| **good** | 120 | 336 | | **good** | 185 | 501 |

Table 4: "Lossy" Naïve Bayes Confusion Matrix

# References

[FHT09] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning.* Springer series in statistics, Berlin, second (11th) edition, 2009.

# Appendix

Listing 1: Linear Discriminant Analysis Assignment

```
1   library("glmnet")
2   library("ggplot2")
3   library("grDevices")
4
5   mu <- function(X)  { return(colMeans(X)) }
6   softmax <- function(X, wi, wj) {
7       X <- data.matrix(X)
8       ihypothesis <- exp(X %*% wi[-1] + wi[1])
9       jhypothesis <- exp(X %*% wj[-1] + wj[1])
10      jhypothesis <- jhypothesis +
11                     exp(X %*% wi[-1] + wi[1])
12      return(ihypothesis / jhypothesis)
13  }
14
15  lda <- function(X, y) {
16      classes <- levels(y) # Only c = 2
17      class1  <- which(y == classes[1])
18      class2  <- which(y == classes[2])
19      X1      <- data.matrix(X[class1,])
20      y1      <- data.matrix(y[class1]);
21      X2      <- data.matrix(X[class2,])
22      y2      <- data.matrix(y[class2]);
23
24      mu1 <- mu(X1) ; mu2 <- mu(X2)
25      pi1 <- length(y1) / length(y)
26      pi2 <- length(y2) / length(y)
27      sigma <- cov(X1)*nrow(X1)+
28              cov(X2)*nrow(X2)
29      sigma <- sigma/nrow(X)
30
31      w01 <- -0.5 * mu1 %*% solve(sigma) %*% mu1 + log(pi1)
32      wx1 <- solve(sigma) %*% mu1 # Some sort of weird magic.
33      w1 <- matrix(c(w01, wx1), 1, 3)
34
35      w02 <- -0.5 * mu2 %*% solve(sigma) %*% mu2 + log(pi2)
36      wx2 <- solve(sigma) %*% mu2 # Some more magic here too.
37      w2 <- matrix(c(w02, wx2), 1, 3)
38      return(rbind(w1, w2)) # w1, w2.
39  }
40
41  classify <- function(X, d) {
42      return(d[1] + d[2]*X[,1] +
43                    d[3]*X[,2])
44  }
45
46  crabs <- read.csv("crabs.csv")
47  X <- crabs[,c("CL", "RW")]
48  y <- crabs[,c("sex")]
49
50  setEPS()
51  cairo_ps("crabs.eps")
52  print(qplot(CL, RW, data = crabs, color = sex,
53      geom = c("point"),
54      xlab = "Carapace Length",
55      ylab = "Rear Width"))
56  dev.off()
57
```

```
58  parameters <- lda(X, y)
59  difference <- parameters[1,]-parameters[2,]
60  intercept <- difference[1] / difference[3]
61  slope <- difference[2] / difference[3]
62
63  sex <- classify(X, difference) > 0.0
64  sex[sex == FALSE] = "Female"
65  sex[sex == TRUE] = "Male"
66
67  setEPS()
68  cairo_ps("boundarylda.eps")
69  print(qplot(CL, RW, data = crabs, color = sex,
70        geom = c("point"),
71        xlab = "Carapace Length",
72        ylab = "Rear Width") +
73        geom_abline(intercept = -intercept,
74                    slope = -slope, colour='purple'))
75  dev.off()
76
77  fit <- cv.glmnet(data.matrix(X), data.matrix(y),
78      family = "binomial", type.measure = "class")
79  yhat <- predict(fit, data.matrix(X), type="class")
80
81  setEPS()
82  cairo_ps("boundarylr.eps")
83  print(qplot(X$CL, X$RW, color = yhat,
84        geom = c("point"),
85        xlab = "Carapace Length",
86        ylab = "Rear Width") +
87        geom_abline(intercept = -coef(fit)[1] / coef(fit)[3],
88                    slope = -coef(fit)[2] / coef(fit)[3],
89                    colour='purple'))
90  dev.off()
91
92  cat("Decision boundary with linear discriminant analysis:",
93      -intercept, "+", -slope, "* k\n")
94  cat("Decision boundary with linear regression:",
95      -coef(fit)[1] / coef(fit)[3], "+",
96      -coef(fit)[2] / coef(fit)[3], "* k\n")
```

Listing 2: Decision Trees and Naïve Bayes Assignment

```
1   library("tree")
2   library("ggplot2")
3   library("reshape2")
4   library("grDevices")
5   library("e1071")
6
7   set.seed(12345) # As always.....
8   scores <- read.csv("scores.csv")
9   n <- nrow(scores) # Observation.
10  samples <- sample(1:n,  n / 2.0)
11  others <- setdiff(1:n, samples)
12  halves <- sample(others, n/4.0)
13
14  training <- scores[samples,]
15  trainingX <- training[,-ncol(training)]
16  trainingy <- training[,ncol(training)]
17
18  validation <- scores[halves,]
19  validationX <- validation[,-ncol(validation)]
```

```
20  validationy <- validation[,ncol(validation)]
21
22  testing  <- scores[-halves,]
23  testingX <- testing[,-ncol(testing)]
24  testingy <- testing[,ncol(testing)]
25
26  fit <- tree(good_bad ~ ., data = training, split = c("deviance"))
27  training_prediction <- predict(fit, trainingX, type= "class")
28  testing_prediction <- predict(fit, testingX, type="class")
29  cat("Missclassifications only with deviance impurity: (",
30      mean(training_prediction != trainingy), "," ,
31      mean(testing_prediction != testingy), ")\n")
32
33  fit <- tree(good_bad ~ ., data = training, split = c("gini"))
34  training_prediction <- predict(fit, trainingX, type= "class")
35  testing_prediction <- predict(fit, testingX, type="class")
36  cat("Missclassifications only with the gini impurity: (",
37      mean(training_prediction != trainingy), "," ,
38      mean(testing_prediction != testingy), ")\n")
39
40  fit <- tree(good_bad ~ ., data = training, split = c("deviance", "gini"))
41  training_prediction <- predict(fit, trainingX, type= "class")
42  testing_prediction <- predict(fit, testingX, type="class")
43  cat("Missclassifications only with deviance and gini: (",
44      mean(training_prediction != trainingy), "," ,
45      mean(testing_prediction != testingy), ")\n")
46
47  max_depth <- 15
48  training_deviance <- rep(0, max_depth)
49  validation_deviance <- rep(0, max_depth)
50  for (depth_level in 2:max_depth) {
51      pruned <- prune.tree(fit, best = depth_level)
52      pred <- predict(pruned, validation, type="tree")
53      training_deviance[depth_level] <- deviance(pruned)
54      validation_deviance[depth_level] <- deviance(pred)
55  }
56
57  deviances <- data.frame(2:max_depth, training_deviance[-1], validation_deviance[-1])
58  colnames(deviances) <- c("Leaves", "Training", "Validation")
59  collapsed_deviances <- melt(deviances, id="Leaves")
60
61  setEPS()
62  cairo_ps("deviance.eps")
63  # It seems depth 12 is good, since validation goes hayware after that...
64  print(ggplot(data=collapsed_deviances, aes(x=Leaves, y=value, color=variable)) +
65      geom_smooth() + labs(x="Leaves", y="Deviance", color="data set"))
66  dev.off()
67
68  # The final tree has depth 6, see output of `final_tree`.
69  # The parameters chosen are: savings, duration, history, age,
70  # purpose, amount, other, resident, in `summary(final_tree)`.
71  final_tree <- prune.tree(fit, best = 12) # The best choice...
72  prediction <- predict(final_tree, testing, type = "class")
73  cat("Missclassification for the optimal tree depth: (",
74      mean(prediction != testingy), ")\n")
75
76  fit <- naiveBayes(good_bad ~ ., data = training)
77  training_prediction <- predict(fit, training, type = "class")
78  testing_prediction <- predict(fit, testing, type = "class")
79  cat("Missclassifications using Naive Bayes method:  (",
80      mean(training_prediction != trainingy), "," ,
81      mean(testing_prediction != testingy), ")\n")
```

6

```
82  cat("\nConfusion matrices for using Naive Bayes:\n")
83  print(table(training_prediction, trainingy))
84  print(table(testing_prediction, testingy))
```

# Introduction to Machine Learning
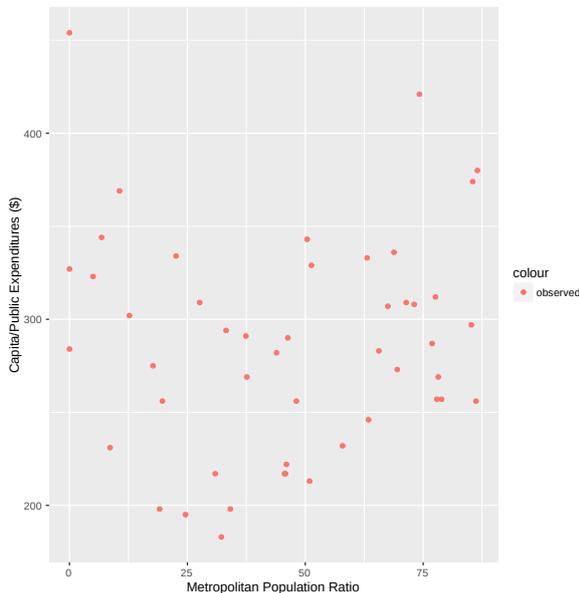# Individual Laboration Report –4–

Erik Sven Vasconcelos Jansson
erija578@student.liu.se
Linköping University,  Sweden
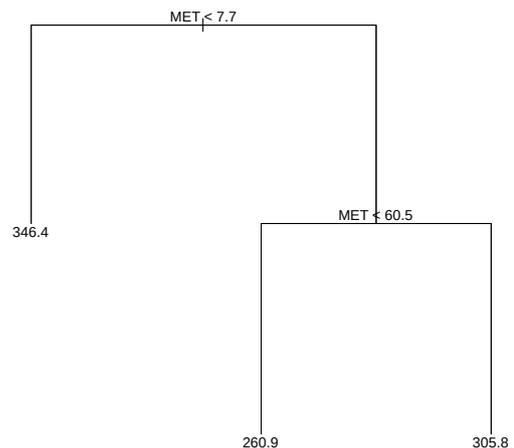
December 4, 2016

## Assignment 1

We are given the data set *State*, where observations regarding the *metropolitan ratio* and the *local public expenditure* for several *states* exist. Plotting these against each other gives Figure 1. Notice that the data is quite spread out, and that there doesn't seem to be any easily visible pattern unfortunately.

Figure 1: Metropolitan Ratio vs Expenditures ($)



For this task we are recommended to use *regression trees* as the *predictor* for the above situation. We fit the model using all observations, and thereafter *cross-validate* the model to find the optimal number of leaves for the *decision tree*. See Figure 2.
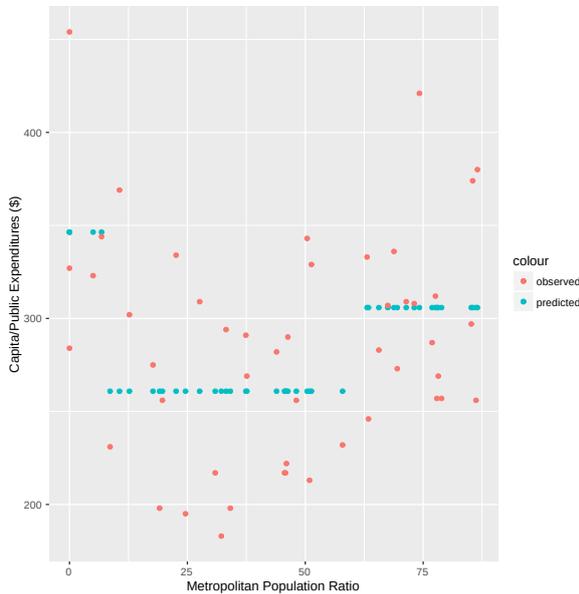
Figure 2: Optimal Regression Tree Predictor



Notice that *three leaves* have been selected, therefore, the *optimal regression tree* is found by *pruning* the original tree and obtaining the above. This is done in Listing 1 under source lines 20−30.
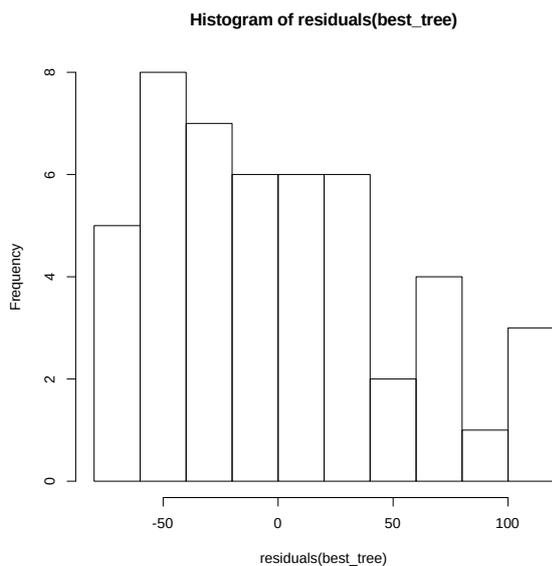
Thereafter, we predict the observed data using our optimal regression tree. Plotting these against the original observations gives Figure 3. Notice how the predictor has given estimates that are roughly the *mean* of each "bucket" selected by the regression tree. Therefore, the selected model is not expected to perform very good since it does not account for the noise present in the model very well, but will predict some correctly, since it's around the mean.

Figure 3: Regression Tree Predicted Expenditure



Finally, we plot the *residuals* from the *regression tree model*, which are basically the *distance* between the *observed responses* and the *predicted responses*. See Figure 4, where we have plotted a histogram of the residuals. Notice that it *doesn't* resemble a *bell curve*, meaning the *error isn't normally distributed*.
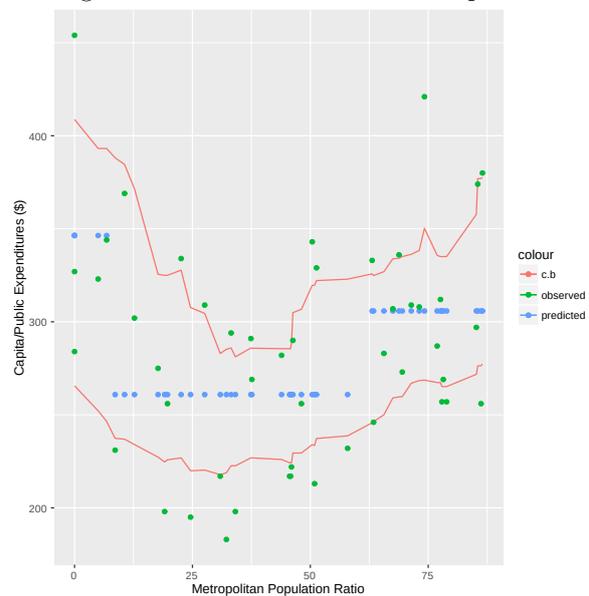
Figure 4: Histogram of Prediction Residuals



Since the chosen model doesn't seem particularly good, at least for these observations, we want to find *more information* about our estimator. Using *bootstrap* seems like a good idea. *Bootstrapping* is used to *estimate the properties* of an estimator by *re-sampling* from a given *approximated distribution*.

Because we don't know the underlying distribution, we want to use *non-parametric bootstrapping*. It works by *re-sampling observations* with *replacement*, the distribution is then calculated by using: $\hat{f}(D_1), \hat{f}(D_2), ..., \hat{f}(D_B)$, where $\hat{f}$ is our estimator. We use the `boot` function in $R$, re-sampling 1024 times. See Listing 1 lines `49-64` for the algorithm.
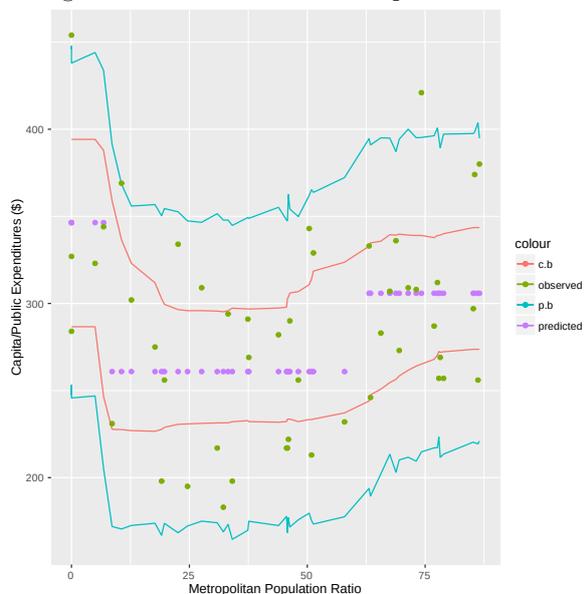
Figure 5: Non-Parametric Bootstrap C.B.



According to Figure 5 which displays the plot of the *confidence bands* of our estimator, several of our observations fall outside the 95% confidence level, which means that our estimator isn't good at all. The curve also seems to be quite bumpy, because the predictor only jumps between the three leaves.

Now, we are given a precondition: $Y \sim \mathcal{N}(\mu_i, \sigma^2)$ which means the *target* is now *normally distributed*. In light of this, we can use *parametric bootstrapping*, since the distribution has now been assumed to be known. It functions by re-sampling from the given distribution, which is different from before since we were only taking samples from the original data set, while we now generate entirely new, fresh, samples.

By doing this in Listing 1 under lines `82-114` we retrieve both the *confidence* and *prediction bands*. Both of these bands are plotted in Figure 6 below.

Figure 6: Parametric Bootstrap C.B. + P.B.



Notice, that the *prediction band* covers most of the observations, as opposed to the *confidence band*. *Prediction bands* account for the *noise* in the data. It seems reasonable that around 5% of the observations fall outside the *prediction band*, since we are targeting a 95% level of confidence (or by $\alpha = 5\%$). In this case, the model which we have predicted seems reasonable, since most predictions are right. *However*, this is probably not true, since Figure 4 shows a distribution that is most likely not normal, therefore *non-parametric bootstrapping* seems more accurate in this particular case (also by intuition).
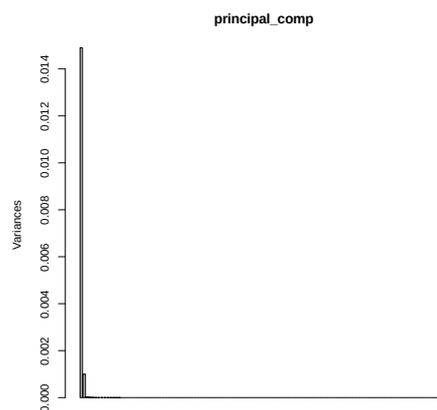
## Assignment 2

Here, we are given a data set containing measures of *near-infrared spectra* and their *viscosity level* for a collection of *diesel fuels*. There are a lot of *features*.

*Principal Component Analysis (PCA henceforth)* is a *dimensionality reduction* technique, where the goal is to find a set of *principal components* where the given *observations* are might to be correlated. Our task is to find the *principal components* of the given data set accounting for 99% of the variance.

By using the built-in `prcomp` function, we apply *PCA* in Listing 2 under lines `7-17`. Afterwards in Figure 7 we produce a *scree plot*, which tells us that *two features* seem to account for all of the variance. More accurately, a bit above 99% of the variance...

Figure 7: Scree Plot for PCA



Additionally, the chosen features *X750* and *X752* are plotted against each other, giving the Figure 8. Notice that there are a couple of outliers in $[1.0, 1.5]$ which can be classified as being *unusual diesel fuels*.

Figure 8: Score for PC1 and PC2



3

Plotting the so called *loadings* of these principal components gives Figure 9. A high *loading* implies that there is a *strong correlation* for a given feature.

Now we apply *Independent Component Analysis* instead, which assumes the given components are statistically independent. We do this in Listing 2 under lines 45-68, and produce the *trace plot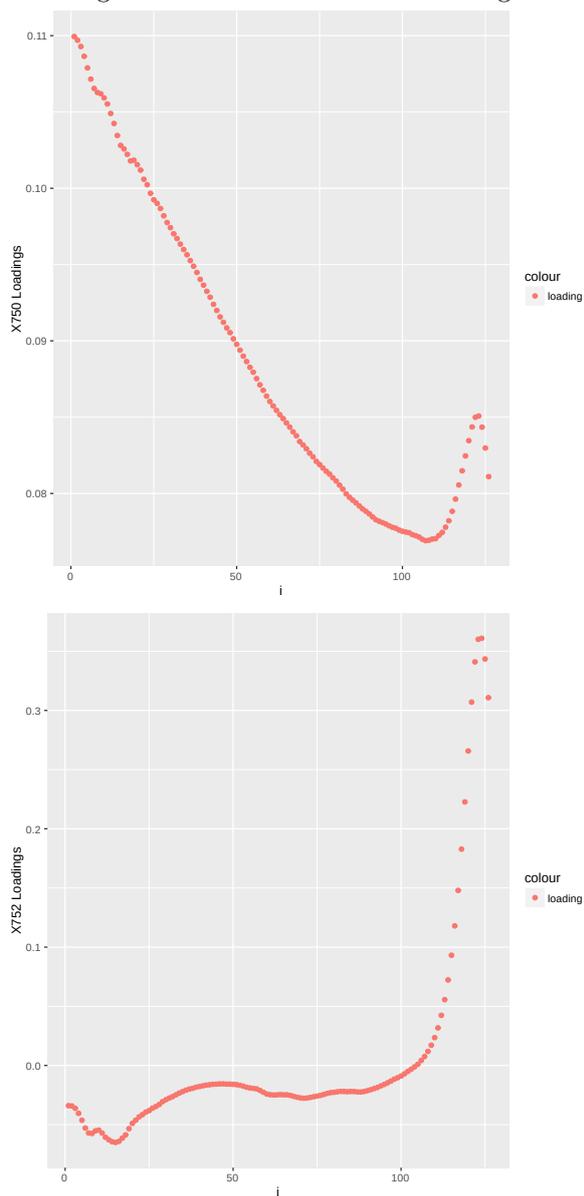* for these in Figure 10. The contents of these plots are similar to those in Figure 9, which are *loadings*, but are inverted, since we are measuring *independence* instead of *correlation*, the opposite PCA measure. Finally, we plot the score(s) of these in Figure 11.
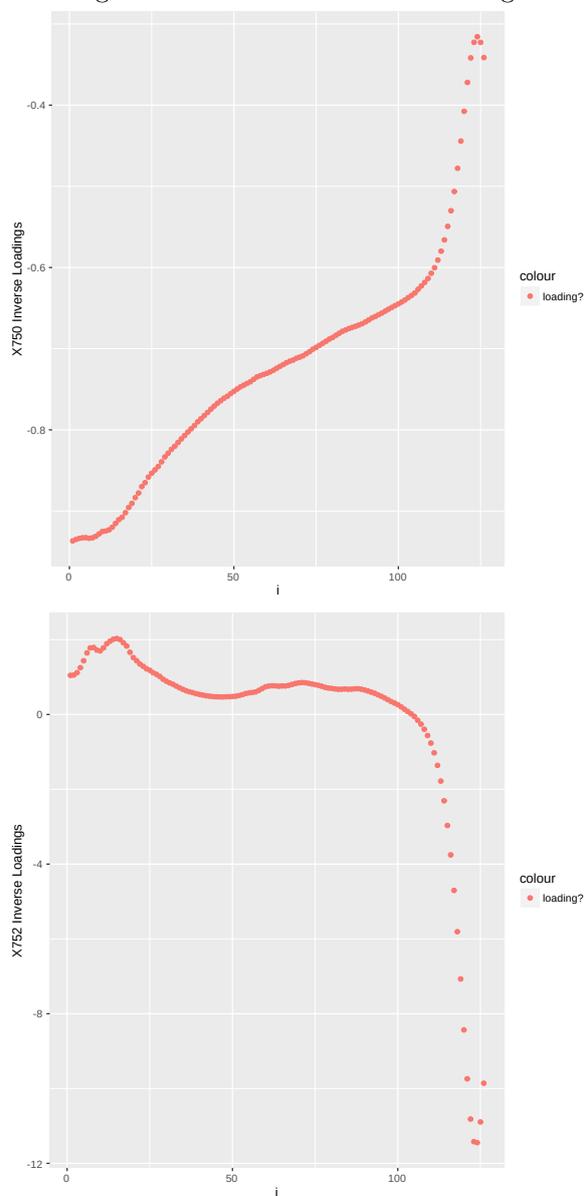
Figure 9: Trace Plot of PCA Loadings



Figure 10: Trace Plot of ICA Loadings



Notice that both PC1 and PC2 seem to have a correlated feature(s) that spike in loading amount, while the rest seem to less significant in comparison. This might be a good candidate for a third principal component if we required to have higher variance.

Notice that the direction of the score has been swapped, the underlying reason for this is similar to those previously mentioned regarding loadings. It's interesting to note that the units are different.

# References

[FHT09] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning.* Springer series in statistics, Berlin, second (11th) edition, 2009.

Figure 11: Score for PCA and ICA



Finally, we apply *cross-validation* in Figure 11, and note that the ranges of $[10, 18]$ p.c. components produces least *MSE*, in relation to the p.c. amount.

Figure 12: Cross-Validation for PCA



5

# Appendix

Listing 1: Script for Assignment 1 on Bootstrapping Regression Tree Model

```r
1   library("tree")
2   library("boot")
3   library("ggplot2")
4   library("grDevices")
5
6   setEPS() # Enable save to eps.
7   state <- read.csv2("state.csv")
8   # Reorder data according to $MET.
9   state <- state[order(state$MET),]
10
11  cairo_ps("state.eps")
12  # Plotting MET vs EX. See file.
13  print(qplot(MET, EX, data = state,
14      xlab = "Metropolitan Population Ratio",
15      ylab = "Capita/Public Expenditures ($)",
16      geom = c("point")))
17  dev.off()
18
19  set.seed(12345) # Required for cross-validation.
20  # Make sure that there are at least 8 leaves now.
21  control <- tree.control(nrow(state), minsize = 8)
22  # Fit our model by using regression trees (8 leaves).
23  fit <- tree(EX ~ MET, data = state, control = control)
24  optimal <- cv.tree(fit) # Do k-fold cross-validation.
25  least_deviance_index <- which.min(optimal$dev)
26  leaves <- optimal$size[least_deviance_index]
27  best_tree <- prune.tree(fit, best = leaves)
28  # Gives "best" tree according to k-fold cv.
29  yhat <- predict(best_tree, newdata = state)
30  # Predict EX by using best regression tree.
31
32  cairo_ps("tree.eps")
33  plot(best_tree)
34  text(best_tree)
35  dev.off()
36
37  cairo_ps("predicted_state.eps")
38  # Plotting MET vs EX. See the file.
39  print(qplot(MET, yhat, data = state,
40      xlab = "Metropolitan Population Ratio",
41      ylab = "Capita/Public Expenditures ($)",
42      geom = c("point")))
43  dev.off()
44
45  cairo_ps("histogram.eps")
46  hist(residuals(best_tree))
47  dev.off()
48
49  bootstrap_predictor <- function(data, indices) {
50      sample <- data[indices,] # Pick a subset of data.
51      control <- tree.control(nrow(sample), minsize = 8)
52      # Fit our model by using regression trees (8 leaves).
53      fit <- tree(EX ~ MET, data = sample, control = control)
54      leaves <- optimal$size[least_deviance_index]
55      best_tree <- prune.tree(fit, best = leaves)
56      # Gives "best" tree according to k-fold cv.
57      yhat <- predict(best_tree, newdata = data)
```
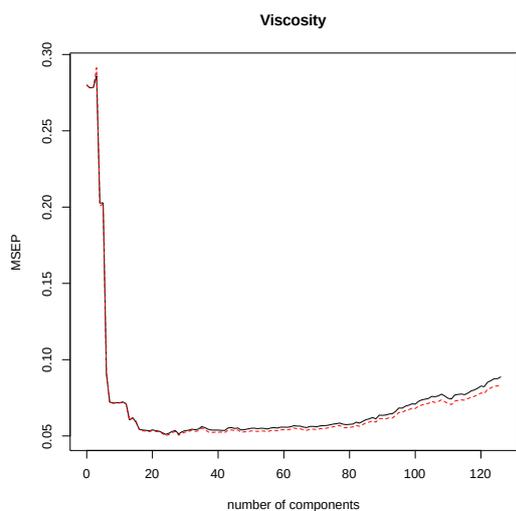
```
58        return(yhat) # Prediction from subset.
59    }
60
61    set.seed(12345) # Required for the bootstrapping.
62    # Apply non-parametric bootstrap to our regression tree
63    # model, picking out 1024 different indices from state.
64    bootstrap <- boot(state, bootstrap_predictor, R = 1024)
65
66    cairo_ps("bootstrap.eps")
67    plot(bootstrap)
68    dev.off()
69
70    # Find the confidence bands.
71    bands <- envelope(bootstrap)
72
73    cairo_ps("bands.eps")
74    # Plotting MET vs EX. See the file.
75    print(qplot(MET, yhat, data = state,
76        xlab = "Metropolitan Population Ratio",
77        ylab = "Capita/Public Expenditures ($)",
78        geom = c("point")) + geom_line(data = state, aes(x = MET, y = bands$point[1,], col = "c.b"
         )) +
79                            geom_line(data = state, aes(x = MET, y = bands$point[2,], col = "c.b"
         )))
80    dev.off()
81
82    bootstrap_prediction <- function(data) {
83        control <- tree.control(nrow(data), minsize = 8)
84        # Fit our model by using regression trees (8 leaves).
85        fit <- tree(EX ~ MET, data = data, control = control)
86        leaves <- optimal$size[least_deviance_index]
87        best_tree <- prune.tree(fit, best = leaves)
88        # Gives "best" tree according to k-fold cv.
89        yhat <- predict(best_tree, newdata = data)
90        sample <- rnorm(nrow(data), yhat, sd(resid(fit)))
91        return(yhat) # Prediction from random subset.
92    }
93
94    bootstrap_confidence <- function(data) {
95        control <- tree.control(nrow(data), minsize = 8)
96        # Fit our model by using regression trees (8 leaves).
97        fit <- tree(EX ~ MET, data = data, control = control)
98        leaves <- optimal$size[least_deviance_index]
99        best_tree <- prune.tree(fit, best = leaves)
100       # Gives "best" tree according to k-fold cv.
101       yhat <- predict(best_tree, newdata = data)
102       return(yhat) # Prediction from subset.
103   }
104
105   bootstrap_random <- function(data, mle) {
106       data$EX <- rnorm(nrow(data), predict(mle, newdata = data), sd(resid(mle)))
107       return(data) # MLE is basically the best tree model from the prediction.
108   }
109
110   set.seed(12345) # Required for the paramatric bootstrapping using confidence bands instead.
111   bootstrapc <- boot(state, bootstrap_confidence, R = 1024, mle = best_tree, ran.gen =
          bootstrap_random, sim = "parametric")
112   set.seed(12345) # Required for the paramatric bootstrapping using prediction bands instead.
113   bootstrapp <- boot(state, bootstrap_prediction, R = 1024, mle = best_tree, ran.gen =
          bootstrap_random, sim = "parametric")
114   confidence_bands <- envelope(bootstrapc) ; prediction_bands <- envelope(bootstrapp) # Bands
          for the Parametric Boostrap...
```

```
115
116   cairo_ps("npcbootstrap.eps")
117   plot(bootstrapc)
118   dev.off()
119   cairo_ps("nppbootstrap.eps")
120   plot(bootstrapp)
121   dev.off()
122
123   cairo_ps("npbands.eps")
124   # Plotting MET vs EX. See the file.
125   print(qplot(MET, yhat, data = state,
126       xlab = "Metropolitan Population Ratio",
127       ylab = "Capita/Public Expenditures ($)",
128       color = "predicted",
129       geom = c("point")) + geom_line(data = state, aes(x = MET, y = prediction_bands$point[1,],
      col = "p.b")) +
130                            geom_line(data = state, aes(x = MET, y = prediction_bands$point[2,],
      col = "p.b")) +
131                            geom_line(data = state, aes(x = MET, y = confidence_bands$point[1,],
      col = "c.b")) +
132                            geom_line(data = state, aes(x = MET, y = confidence_bands$point[2,],
      col = "c.b")) +
133                            geom_point(data = state, aes(x = MET, y = EX, color="real")))
134   dev.off()
```

Listing 2: Script for Assignment 2 on Principal/Individual Component Analysis

```
 1   library("pls")
 2   library("ggplot2")
 3   library("fastICA")
 4   library("reshape2")
 5
 6   setEPS() # Enables saving EPS format.
 7   spectra <- read.csv2("NIRSpectra.csv")
 8   xspectra <- spectra[,-ncol(spectra)]
 9   yspectra <- spectra[,ncol(spectra)]
10   principal_comp <- prcomp(xspectra)
11   lambda <- principal_comp$sdev^2
12
13   # Notice both X750, X752.
14   cairo_ps("screeplot.eps")
15   screeplot(principal_comp,
16             ncol(xspectra))
17   dev.off()
18   cairo_ps("biplot.eps")
19   biplot(principal_comp)
20   dev.off()
21
22   cairo_ps("score.eps")
23   print(qplot(principal_comp$x[,1],
24               principal_comp$x[,2],
25               xlab = "X750",
26               ylab = "X752"))
27   dev.off()
28
29   x750loadings <- principal_comp$rotation[,1]
30   x752loadings <- principal_comp$rotation[,2]
31
32   cairo_ps("x750loadings.eps")
33   print(qplot(1:length(x750loadings),
34               x750loadings, xlab="i",
```

```
35            ylab="X750 Loadings"))
36  dev.off()
37
38  cairo_ps("x752loadings.eps")
39  print(qplot(1:length(x752loadings),
40              x752loadings, xlab="i",
41              ylab="X752 Loadings"))
42  dev.off()
43
44  set.seed(12345) # But WHY?!?!?!??!?!?!?!
45  independent_comp <- fastICA(xspectra, 2)
46
47  W <- independent_comp$K %*% independent_comp$W
48  x750whitening <- W[,1] # Un-mixed and whitened
49  x752whitening <- W[,2] # Un-mixed and whitened
50
51  cairo_ps("x750traceplot.eps")
52  print(qplot(1:length(x750whitening),
53              x750whitening, xlab="i",
54              ylab="X750 Inverse Loadings"))
55  dev.off()
56
57  cairo_ps("x752traceplot.eps")
58  print(qplot(1:length(x752whitening),
59              x752whitening, xlab="i",
60              ylab="X752 Inverse Loadings"))
61  dev.off()
62
63  cairo_ps("icascore.eps")
64  print(qplot(independent_comp$S[,1],
65              independent_comp$S[,2],
66              xlab = "X750",
67              ylab = "X752"))
68  dev.off()
69
70  set.seed(12345)
71  principal_compcv <- pcr(Viscosity ~ ., data = spectra,
72                          validation = "CV")
73  cairo_ps("pcacv.eps")
74  validationplot(principal_compcv,
75                 val.type = "MS")
76  dev.off()
```

# Introduction to Machine Learning
# Individual Laboration Report –5–

Erik Sven Vasconcelos Jansson

`erija578@student.liu.se`

Linköping University, Sweden

December 12, 2016

Increasing the accuracy of *weather forecasts* is an important task. We propose an estimator which produces the *air temperature forecast* in *Sweden*, given a *latitude/longitude coordinate* and also *date*. Some observations by *SMHI*, taken from weather stations, have been given for training our estimator.

By using a *Nadaraya–Watson regression kernel*, we can estimate the temperatures $\boldsymbol{y}'$. This is done by taking the *kernels* $k_\sigma(\boldsymbol{x}^{(i)}, \boldsymbol{x}')$ for each $i^{th}$ data from the training set and using it as a *weight* when considering the response variable $\boldsymbol{y}^{(i)}$. Essentially, the kernel $k_\sigma(\boldsymbol{x}^{(i)}, \boldsymbol{x}')$ will reduce $\boldsymbol{y}^{(i)}$'s significance in the *total contribution* by giving less weight when the $\boldsymbol{x}^{(i)}$ and $\boldsymbol{x}'$ are further away (in some measure).

We have used a *Gaussian Radial Basis Function* as our *kernel*, which is defined in Equation 1 below. Note the parameter $\sigma$, which can be considered as the *spread* or *width* of the kernel, and also $\boldsymbol{x}^{(i)} - \boldsymbol{x}'$ which is the *distance function*; giving our kernel the property of a *similarity function* (because of $e^{(\cdots)}$).

By using $k_\sigma(\boldsymbol{x}^{(i)}, \boldsymbol{x}')$ in *Nadaraya–Watson's* $\boldsymbol{y}'$ estimator, shown in Equation 2, we are essentially *weighing* how important the contributions from $\boldsymbol{y}^{(i)}$ are to $\boldsymbol{y}'$, because *similar* $\boldsymbol{x}^{(i)}$ will give higher $k_\sigma$.

$$k_\sigma(\boldsymbol{x}, \boldsymbol{x}') = \exp\left( \frac{-\left\| (\boldsymbol{x} - \boldsymbol{x}') \right\|^2}{2\sigma^2 \left\{ \sigma \approx h \right\}} \right) \qquad (1)$$

$$\boldsymbol{y}'(\boldsymbol{x}, \boldsymbol{x}') = \frac{\sum_n \boldsymbol{y}^{(i)} k_\sigma(\boldsymbol{x}^{(i)}, \boldsymbol{x}')}{\sum_n k_\sigma(\boldsymbol{x}^{(i)}, \boldsymbol{x}')} \qquad (2)$$

Practically, the *kernel* is calculated in Listing 1 under `gaussian_kernel` and the *estimation* is being done in the function `forecast`. However, note that the final contributions use `forecast_kernel`, which will be described now.

Below follows the applied *distance functions*, which give the measured distance between a pair of *locations*, *times of the day*, and also *dates of year*. These are used in `forecast_kernel` for each respective `gaussian_kernel` invocation. Additionally, these distances are *normalized* to range in-between 0.0 - 1.0. See Listing 1 for these values.

$$d_l = r \, \mathrm{hav}^{-1}(h), \ \mathrm{hav}(\varphi) = \frac{1 - \cos\varphi}{2}$$

$$d_t = \begin{cases} |x - y| & |x - y| < (x + y) \bmod 24 \\ (x + y) \bmod 24 & |x - y| \geq (x + y) \bmod 24 \end{cases}$$

$$d_d = \begin{cases} |x - y| & |x - y| < (x + y) \bmod 365 \\ (x + y) \bmod 24 & |x - y| \geq (x + y) \bmod 365 \end{cases}$$

Therefore, the final `forecast_kernel` is being calculated as seen below, where $k_l$ uses the *location* distance, $k_d$ the *date distance* and $k_t$ *time distance*.

$$k_f(\boldsymbol{x}, \boldsymbol{x}') = k_l(\boldsymbol{x}, \boldsymbol{x}') + k_d(\boldsymbol{x}, \boldsymbol{x}') + k_t(\boldsymbol{x}, \boldsymbol{x}')$$

Within Table 1 are our chosen $\sigma/h$ *spread/width*.

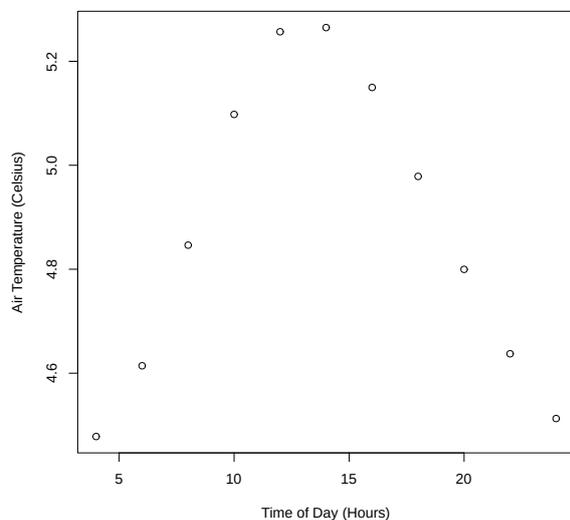| Feature | Spread |
|---------|--------|
| Location | 0.192 |
| Day | 0.256 |
| Time | 0.256 |

Table 1: Kernel Width

These have been chosen to decrease contribution, so for example *locations* are not easily influenced...

Finally, we estimate $\boldsymbol{y}$ for *2013-10-04* during the entire day in *58.4274 latitude* and *14.826 longitude*. See Table 2 and Figure 1 for the predicted results.

| Time | Temperature (°C) |
|------|------------------|
| 04:00 | 4.478177 |
| 06:00 | 4.614255 |
| 08:00 | 4.846348 |
| 10:00 | 5.097822 |
| 12:00 | 5.256865 |
| 14:00 | 5.264849 |
| 16:00 | 5.149687 |
| 18:00 | 4.978578 |
| 20:00 | 4.799737 |
| 22:00 | 4.637349 |
| 24:00 | 4.512750 |

Table 2: Forecast for 2013-11-04

Figure 1: Air Temperature Forecast Graphs



Notice how the plot above produces a *bell curve*, which is to be expected from a temperature forecast according to previous data given by *SMHI*. Also, these values don't seem to be that far off the truth, however they seem to be slightly colder than usual. One possible cause for this can be motivated by the *independence* of each *kernel*, outweighing the other.

For example, assume both *location* and *time of day* are highly correlated to our $\boldsymbol{x'}$, therefore, they will *contribute highly* with their $\boldsymbol{y^{(i)}}$. Now, for the sake of argument, assume *day of the year* is *not highly correlated* with $\boldsymbol{x'}$, therefore, the contribution $\boldsymbol{y^{(i)}}$ will not be significant, at least not compared to *location* and *time of day*. Therefore, even if we are taking an observation which is far away from the requested date, the contribution will still be high, which leads to most predictions being influenced with the "mean" temperatures in Sweden. Therefore, our hypothesis why most predictions are colder than expected is because the three kernels are being accounted independently of each other...

Additionally, some words need to be said regarding our choice for the *kernel spread/width* values in Table 1. These were chosen on the assumption that *locations* further away than *350 kilometers* are not very good contributors, as are *dates* with a distance further away than *45 days* and between $\approx 5$ *hours*. The $\sigma$ where chosen such that these values were reached, and only correlated less than 10%, thereafter proceeding with normal Gaussian falloff.

# References

[FHT09] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning*. Springer series in statistics, Berlin, second (11th) edition, 2009.

# Appendix

Listing 1: Nadaraya–Watson Gaussian Radial Basis Function Kernel Forecast Estimator

```r
1  library("ggplot2")
2  library("reshape2")
3  library("grDevices")
4  library("geosphere")
5
6  # Fulhax for finding day diff.
7  epoch <- as.Date("1970-01-01")
8  epocht <- as.POSIXlt(paste(epoch,
9                       "00:00:00"))
10 stations <- read.csv("stations.csv")
11 temps50k <- read.csv("temps50k.csv")
12 weather <- merge(stations, temps50k,
13             by = "station_number")
14
15 # Language R has a very shitty library for doing dates/times.
16 weather$time <- as.POSIXlt(paste("1970-01-01", weather$time))
17 weather$hour <- round(difftime(weather$time, epocht,
18                             units="hours"))
19 weather$day <- sub("^\\d{4}", "1970", weather$date)
20 weather$day <- as.Date(weather$day) - epoch
21 weather$day <- as.numeric(weather$day)
22
23 # Find faulty convert (because of leap year).
24 working <- which(complete.cases(weather$day))
25 weather <- weather[working,] # Remove errors.
26 indx <- sample(1:nrow(weather),nrow(weather))
27
28 gaussian_kernel <- function(distance, sdspread)  {
29     # The Gaussian Radial Basis Function.
30     gamma_spread   <-  1 / (2*sdspread^2)
31     return(exp(-gamma_spread*distance^2))
32 }
33
34 date_time_modulo_distance <- function(x, y, n) {
35     numeric_distance <- abs(x - y)
36     modulo_distance <- n - numeric_distance
37     if (numeric_distance < modulo_distance)
38         return(numeric_distance)
39     else return(modulo_distance)
40 }
41
42 forecast_kernel <- function(longitudes, latitudes,
43                             days,   time_points) {
44     time_points <- as.numeric(time_points) # I dunno why. FUCKING R.
45     location_distance <- distHaversine(c(longitudes[1], latitudes[1]),
46                                     c(longitudes[2], latitudes[2]))
47     location_distance <- (location_distance / 1000) # Kilometers.
48     day_distance <- date_time_modulo_distance(days[1], days[2], 365)
49     day_kernel <- gaussian_kernel(as.numeric(day_distance / 182.5), 0.256)
50     time_point_distance <- date_time_modulo_distance(time_points[1], time_points[2], 24)
51     time_point_kernel <- gaussian_kernel(as.numeric(time_point_distance / 12), 0.256)
52     location_kernel <- gaussian_kernel(as.numeric(location_distance / 1572), 0.192)
53     # cat("(", location_distance, day_distance, time_point_distance, ")",
54     #     "=> (", location_kernel, day_kernel, time_point_kernel, ") \n")
55     return((day_kernel + location_kernel + time_point_kernel))
56 }
57
```

```
58  forecast <- function(date, location, weather) {
59      clock <- paste0(seq(4,24,by=2))
60      date <- sub("^\\d{4}", "1970", date)
61      date <- as.Date(date) # Kill me.....
62      day <- abs(as.numeric(date - epoch))
63      # See lecture slides for predicting yhat-kernels.
64      total_kernel_sum <- vector(length = length(clock))
65      weighted_temperatures <- vector(length = length(clock))
66      similarities <- matrix(0, nrow(weather), length(clock))
67      stp <- 512 # Print steps for showing the progress...
68
69      # Loop through each observation and desired time.
70      for (i in 1:nrow(weather)) { # Shitty apply fails...
71          if (i%%stp==0) cat("Progress",(i/nrow(weather))*100,"%\n")
72          for (j in 1:length(clock)) { # Need to predict for each...
73              # Calculate the similarity for this observation and our desired response
74              similarities[i, j] <- forecast_kernel(c(location[1], weather[i,]$longitude),
75                                                    c(location[2], weather[i,]$latitude),
76                          c(day, weather[i,]$day), c(clock[j], weather[i,]$hour))
77              total_kernel_sum[j] <- total_kernel_sum[j] + similarities[i, j]
78              weighted_temperature <- similarities[i, j] * weather[i,]$air_temperature
79              weighted_temperatures[j] <- weighted_temperatures[j] + weighted_temperature
80          }
81          # Apply the Nadaraya-Watson kernel regression.
82      } ; return(weighted_temperatures / total_kernel_sum)
83  }
84
85  args <- commandArgs(TRUE)
86  day <- if (!is.na(args[1])) as.Date(args[1]) else as.Date("2013-11-04")
87  latitude <- if (!is.na(args[2])) as.numeric(args[2]) else 58.4274
88  longitude <- if (!is.na(args[3])) as.numeric(args[3]) else 14.826
89  temperatures <- forecast(day, c(longitude, latitude),
90                      weather[indx,]) # Kill me...
91  cat("Forecast (in oC) for the", as.character(day),
92      "at 04:00:00 - 24:00:00 in", longitude, latitude,
93      "(longitude, latitude)\n")
94  cat(temperatures, "\n")
95
96  setEPS()
97  cairo_ps("forecast")
98  plot(clock, temperatures, xlab="Time of Day",
99                          ylab="Temperature Celsius")
100 dev.off()
```

# Introduction to Machine Learning
# Individual Laboration Report –6–

Erik Sven Vasconcelos Jansson
erija578@student.liu.se
Linköping University, Sweden

December 20, 2016

Finally, the last machine learning topic covered are *artificial neural networks*. These estimators are very flexible, such that even a *single layer feedforward neural network* complies with the *universal approximation theorem*, presented by *Csáji* [Csá01]:

**Theorem 1** (Universal Approximation Theorem). *Any artificial feed-forward neural network with a single hidden layer, containing a finite amount of neurons, can approximate any continuous functions on the compact subset $\mathcal{R}^n$ (with restrictions on $\sigma$).*

*Proof. Csáji's* [Csá01] derivation of Theorem 1. □

Basically, the theorem states that even simple neural networks can represent interesting functions, given some suitable subset of *activation functions*. For this assignment, we want to approximate $\sin x$, where we are given *25 observation* for *training set*. Also, we are given a *validation set* of length *25* for checking if our neural network is under/overfitting. We are using the *R package* `neuralnet` for our fit with *10 hidden units* in a *single hidden layer*, also initialized with *random weights* in $[-1, 1]$ interval. See Listing 1 for the entire assignment source code.

For the curious, Equations 1, 2, and 3 are given:

$$\sigma(u) = \frac{1}{1 + e^{-u}} \tag{1}$$

$$\boldsymbol{w}_{(i)} = \boldsymbol{w}_{(i-1)} - \eta_k \nabla E(\boldsymbol{w}_{(i-1)}) \tag{2}$$

$$\hat{y}_j(\boldsymbol{x}) = \sigma(w_0 + \sum_{h=1}^{H} \sigma(w_{0h} + \boldsymbol{w}_h^\intercal \boldsymbol{x})) \tag{3}$$

1. **Sigmoid Activation Function:** "S"-shaped function which converges $\sigma(u) = 1$ as $u \to \infty$ and $\sigma(u) = 0$ as $u \to -\infty$. Used in Equation 3.

2. **Batch Gradient Descent:** finds the "step" in the right direction for *minimizing error $E$*. This is achieved with the *gradient* of $E$ given in respect to the weights $\boldsymbol{w}$; giving a *hyperplane*.

3. **Single-Layer Neural Network Estimator:** uses Equations 1 and 2 to find $\hat{y}_j$ by finding the parameters $\boldsymbol{w}$ in each layer (a linear equation) by means of *gradient descent* and producing a *non-linear result in subsequent layers* by the *activation function*. This is the primary reason why neural networks are so flexible & general.

By using a *threshold* for the *gradient descent* we can stop the *neural network* from either *overfitting* or *underfitting*. This simply done by increasing the threshold iteratively and taking the validation set's:

| Threshold | S.S.E. |
|-----------|--------|
| 0.001 | 0.01367691527 |
| 0.002 | 0.01262419958 |
| 0.003 | 0.00988418900 |
| 0.004 | 0.00850089424 |
| 0.005 | 0.00955545744 |
| 0.006 | 0.00974372099 |
| 0.007 | 0.01583926857 |
| 0.008 | 0.01649252416 |
| 0.009 | 0.02112490377 |
| 0.010 | 0.02735909554 |

Table 1: Neural Network Values

After finding the "optimal threshold" of *0.004* by picking the $4^{th}$ iteration (where $i = 4$ that is) which gives the least amount of error for a validation set, we plot the best neural network in Figure 1, and also the predictions in Figure 2 for a *sine function*. Notice that the fit is pretty good, and the estimator gives a pretty "spot on" prediction for the function. It seems *neural networks* are incredibly powerful, but take time to train and are harder to reason about (for example, how do we choose the number of hidden layers and units? How long will it take?)

# References

[Csá01]  Balázs Csanád Csáji. Approximation with Artificial Neural Networks. *Faculty of Sciences, Etvs Lornd University, Hungary*, 24:48, 2001.

[FHT09]  Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning*. Springer series in statistics, Berlin, second (11th) edition, 2009.

[GF10]  Frauke Günther and Stefan Fritsch. neuralnet: Training of Neural Networks. *The R Journal*, 2(1):30–38, 2010.
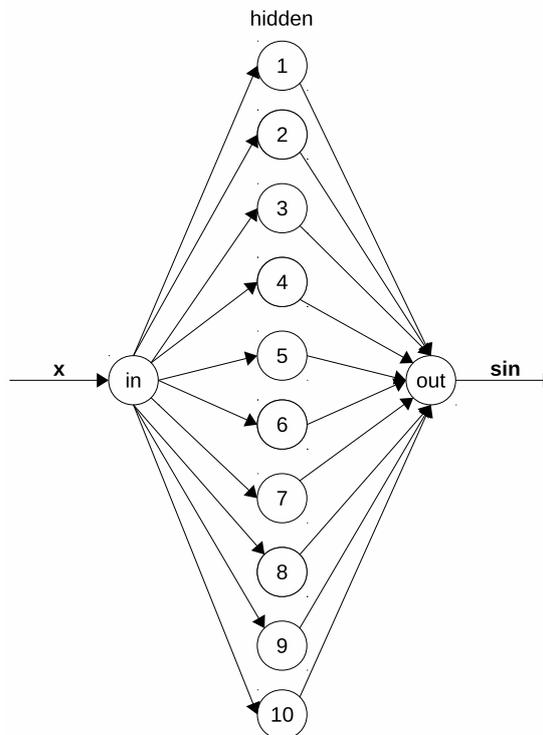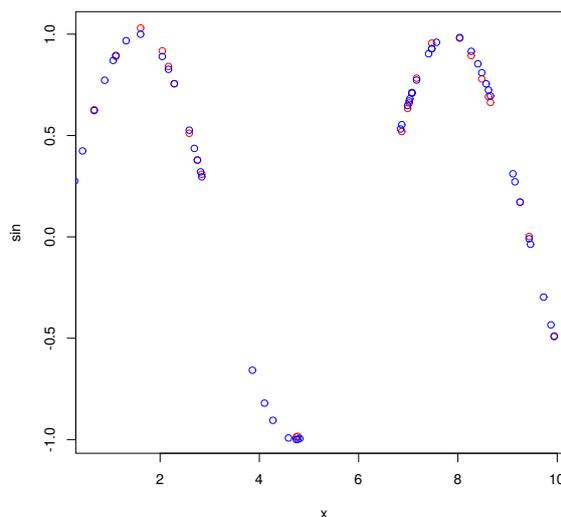
Figure 1: Neural Network



Figure 2: Neural Network's Produced Predictions (in the graph are raw values and predicted values).

# Appendix

## Listing 1: Feed-Forward Backpropagating Neural Network Sine Estimator Script

```
1   library("ggplot2")
2   library("reshape2")
3   library("neuralnet")
4   library("grDevices")
5   set.seed(1234567890)
6
7   variable <- runif(50, 0, 10)
8   sine <- data.frame(x=variable, sin=sin(variable))
9   training <- sine[1:25,] ; testing <- sine[26:50,]
10
11  candidate_error <- Inf
12  units <- 10 # Hidden baby!
13  candidate_threshold <- Inf
14  weights <- runif(50, -1, +1)
15
16  for (threshold_attempt in 1:10) {
17      thresholdi <- threshold_attempt / 1000
18      nn <- neuralnet(sin~x, training, units,
19                      startweights = weights,
20                      threshold = thresholdi)
21
22      predicted <- compute(nn, testing$x)
23      error <- sum((testing$sin - predicted$net.result)^2)
24      cat("NN Threshold", thresholdi, "->", error, "SSE \n")
25      if (error < candidate_error) {
26          candidate_error = error
27          candidate_threshold = thresholdi
28      }
29  }
30
31  nn <- neuralnet(sin~x, training, units,
32                      candidate_threshold,
33                  startweights = weights)
34  predicted <- compute(nn, testing$x)
35
36  plot(nn)
37  setEPS()
38  cairo_ps("predictions.eps")
39  plot(testing$x, predicted$net.result, col = "red",
40      xlab = "x", ylab = "sin")
41  points(sine, col = "blue")
42  dev.off()
```

## Listing 2: Output About the Produced Neural Network in the Assignment

```
1   $response
2               sin
3   1    0.31115890803
4   2   -0.65787112371
5   3    0.85356988285
6   4    0.92820698816
7   5    0.71194544538
8   6    0.95969186755
9   7    0.27531467859
10  8   -0.03662256168
```

3

```
11   9  -0.29718457265
12  10 -0.43427724087
13  11  0.27176755816
14  12  0.96762993527
15  13  0.87023024548
16  14  0.90319426880
17  15  0.53211225475
18  16 -0.90515370065
19  17 -0.99209419164
20  18  0.75493516282
21  19  0.43639270658
22  20  0.42400734122
23  21  0.77254200174
24  22  0.68138797265
25  23  0.32070401674
26  24 -0.99484612705
27  25 -0.82027249428
28
29  $covariate
30                [,1]
31   [1,] 9.1083657276
32   [2,] 3.8595812093
33   [3,] 8.4019783861
34   [4,] 7.4727497180
35   [5,] 7.0754500036
36   [6,] 7.5690891198
37   [7,] 0.2789170155
38   [8,] 9.4614087138
39   [9,] 9.7265206091
40  [10,] 9.8740137112
41  [11,] 9.1495487187
42  [12,] 1.3156641065
43  [13,] 1.0556694935
44  [14,] 7.4103393406
45  [15,] 6.8442786904
46  [16,] 4.2733338126
47  [17,] 4.5865617390
48  [18,] 8.5692227143
49  [19,] 2.6900070859
50  [20,] 0.4378655413
51  [21,] 0.8828348410
52  [22,] 7.0328426105
53  [23,] 2.8151199827
54  [24,] 4.8139597056
55  [25,] 4.1034799209
56
57  $err.fct
58  function (x, y)
59  {
60      1/2 * (y - x)^2
61  }
62  <environment: 0x339a758>
63  attr(,"type")
64  [1] "sse"
65
66  $act.fct
67  function (x)
68  {
69      1/(1 + exp(-x))
70  }
71  <environment: 0x339a758>
72  attr(,"type")
```

```
73   [1] "logistic"
74
75   $linear.output
76   [1] TRUE
77
78   $data
79                 x              sin
80   1   9.1083657276   0.31115890803
81   2   3.8595812093  -0.65787112371
82   3   8.4019783861   0.85356988285
83   4   7.4727497180   0.92820698816
84   5   7.0754500036   0.71194544538
85   6   7.5690891198   0.95969186755
86   7   0.2789170155   0.27531467859
87   8   9.4614087138  -0.03662256168
88   9   9.7265206091  -0.29718457265
89   10  9.8740137112  -0.43427724087
90   11  9.1495487187   0.27176755816
91   12  1.3156641065   0.96762993527
92   13  1.0556694935   0.87023024548
93   14  7.4103393406   0.90319426880
94   15  6.8442786904   0.53211225475
95   16  4.2733338126  -0.90515370065
96   17  4.5865617390  -0.99209419164
97   18  8.5692227143   0.75493516282
98   19  2.6900070859   0.43639270658
99   20  0.4378655413   0.42400734122
100  21  0.8828348410   0.77254200174
101  22  7.0328426105   0.68138797265
102  23  2.8151199827   0.32070401674
103  24  4.8139597056  -0.99484612705
104  25  4.1034799209  -0.82027249428
105
106  $net.result
107  $net.result[[1]]
108              [,1]
109  1    0.30018554412
110  2   -0.64466078637
111  3    0.82577426828
112  4    0.95531707389
113  5    0.71041413678
114  6    0.98604140527
115  7    0.27220011825
116  8   -0.02393016651
117  9   -0.27977135370
118  10  -0.42452441939
119  11   0.26376340300
120  12   0.97423750701
121  13   0.86489547875
122  14   0.92926095080
123  15   0.49438343511
124  16  -0.91926885456
125  17  -0.98908629285
126  18   0.72304892189
127  19   0.42816650637
128  20   0.43013044493
129  21   0.76891173652
130  22   0.67393025023
131  23   0.32904062832
132  24  -0.97934291894
133  25  -0.83193245662
134
```

```
135
136  $weights
137  $weights[[1]]
138  $weights[[1]][[1]]
139             [,1]           [,2]         [,3]         [,4]          [,5]
140  [1,] 0.3718763846 -10.931812117 8.275060257 7.8216380497  1.551228805
141  [2,] 0.5081317757   1.628735531 -2.289303563 0.1166254588 -0.594052483
142             [,6]          [,7]         [,8]          [,9]         [,10]
143  [1,]   4.7453831203 -0.6070429987 9.38110372186 -0.1377222063  5.958245683
144  [2,] -0.4968713811  0.1924524647 0.09270470267  3.1685937697 -2.602280174
145
146  $weights[[1]][[2]]
147               [,1]
148   [1,] -0.06529714022
149   [2,] -0.70033511720
150   [3,]  3.84669442716
151   [4,]  2.74586539382
152   [5,] -0.75978348453
153   [6,] -9.09090264177
154   [7,]  6.65416477397
155   [8,] -8.24869453812
156   [9,] -0.20335986240
157  [10,]  1.00884789102
158  [11,]  2.01492912933
159
160
161
162  $startweights
163  $startweights[[1]]
164  $startweights[[1]][[1]]
165             [,1]          [,2]          [,3]         [,4]          [,5]
166  [1,] 0.4591262657 -0.5031667114  0.9014065554 0.9589186665  0.3389983536
167  [2,] 0.5853618421  0.4307389595 -0.8788680169 0.4978831881 -0.5358421607
168             [,6]          [,7]         [,8]          [,9]         [,10]
169  [1,]  0.6293357364 -0.4028865024 0.5968314419 -0.07648990629  0.7421438890
170  [2,] -0.4464168334 -0.3094406570 0.9041418806 -0.17025629710 -0.8588890089
171
172  $startweights[[1]][[2]]
173               [,1]
174   [1,] -0.2698646844
175   [2,] -0.9254528601
176   [3,]  0.3498687637
177   [4,] -0.7230960354
178   [5,] -0.9836924160
179   [6,] -0.6452815034
180   [7,]  0.8491520174
181   [8,]  0.6410233583
182   [9,] -0.4020887311
183  [10,]  0.9406797229
184  [11,]  0.2142777084
185
186
187
188  $generalized.weights
189  $generalized.weights[[1]]
190             [,1]
191  1  -4.18437409317
192  2   0.84720544205
193  3  -3.90202008113
194  4   8.87306916010
195  5   4.06386829427
196  6  18.83336456829
```

6

```
197  7   5.28857232424
198  8   38.75962317972
199  9   2.72921588482
200  10  1.62852618574
201  11 -4.58128912031
202  12  13.07325336492
203  13  4.31497701909
204  14  6.93919826366
205  15  4.07243773606
206  16  0.23010347764
207  17  0.02853630559
208  18 -3.31477743585
209  19 -3.26159533671
210  20  3.80902182260
211  21  3.41516125086
212  22  3.98677345674
213  23 -3.57280260613
214  24 -0.06878886818
215  25  0.40865817424
216
217
218  $result.matrix
219                                    1
220  error                 0.003576080337
221  reached.threshold     0.003929680826
222  steps             23174.000000000000
223  Intercept.to.1layhid1     0.371876384634
224  x.to.1layhid1             0.508131775660
225  Intercept.to.1layhid2   -10.931812117300
226  x.to.1layhid2             1.628735530657
227  Intercept.to.1layhid3     8.275060257474
228  x.to.1layhid3            -2.289303563425
229  Intercept.to.1layhid4     7.821638049688
230  x.to.1layhid4             0.116625458773
231  Intercept.to.1layhid5     1.551228805249
232  x.to.1layhid5            -0.594052483023
233  Intercept.to.1layhid6     4.745383120333
234  x.to.1layhid6            -0.496871381057
235  Intercept.to.1layhid7    -0.607042998673
236  x.to.1layhid7             0.192452464714
237  Intercept.to.1layhid8     9.381103721859
238  x.to.1layhid8             0.092704702673
239  Intercept.to.1layhid9    -0.137722206303
240  x.to.1layhid9             3.168593769723
241  Intercept.to.1layhid10    5.958245682591
242  x.to.1layhid10           -2.602280174422
243  Intercept.to.sin         -0.065297140222
244  1layhid.1.to.sin         -0.700335117198
245  1layhid.2.to.sin          3.846694427157
246  1layhid.3.to.sin          2.745865393824
247  1layhid.4.to.sin         -0.759783484527
248  1layhid.5.to.sin         -9.090902641770
249  1layhid.6.to.sin          6.654164773966
250  1layhid.7.to.sin         -8.248694538124
251  1layhid.8.to.sin         -0.203359862396
252  1layhid.9.to.sin          1.008847891021
253  1layhid.10.to.sin         2.014929129330
254
255  attr(,"class")
256  [1] "nn"
```