

Homework01

JongCheolLee

2024-09-30

Data

```
library(MASS)
data(Boston)
```

```
head(Boston)
```

```
##      crim zn  indus chas   nox    rm  age    dis rad tax ptratio  black  lstat
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1 296    15.3 396.90  4.98
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2 242    17.8 396.90  9.14
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2 242    17.8 392.83  4.03
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3 222    18.7 394.63  2.94
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3 222    18.7 396.90  5.33
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3 222    18.7 394.12  5.21
##   medv
## 1 24.0
## 2 21.6
## 3 34.7
## 4 33.4
## 5 36.2
## 6 28.7
```

```
y <- Boston[, 1]
x <- Boston[, -c(1,4,9)]
x <- as.matrix(scale(x))
```

- 506 observations with 14 variables
- **crim**: response variable (11 variables are scaled predictors)

```
dim(x)
```

```
## [1] 506  11
```

```
length(y)
```

```
## [1] 506
```

```
class(x)
```

```
## [1] "matrix" "array"
```

```
colnames(x)
```

```
## [1] "zn"      "indus"   "nox"     "rm"      "age"     "dis"     "tax"
## [8] "ptratio" "black"   "lstat"   "medv"
```

Housing Values in Suburbs of Boston

- **crim**: 마을별 1인당 범죄율.
- **zn**: 25,000 평방 피트 이상의 주거용 토지 비율.
- **indus**: 마을별 비소매업 상업 지역의 비율.
- **nox**: 대기 중 산화질소 농도(10백만분율).
- **rm**: 주택당 평균 방 개수.
- **age**: 1940년 이전에 건설된 자가 거주 주택 비율.
- **dis**: 보스턴의 5개 고용 센터까지의 가중 평균 거리.
- **tax**: 10,000달러당 재산세율.
- **ptratio**: 마을별 학생-교사 비율.
- **black**: $1000(Bk \square 0.63)^2$ 여기서 Bk는 마을별 흑인 인구 비율.
- **lstat**: 저소득층 인구 비율(백분율).
- **medv**: 자가 소유 주택의 중위값(단위: 1,000달러).

Goal Boston의 각 마을의 1인당 범죄율 예측

```
apply(x, 2, function(t) sum(is.na(t)))
```

```
##      zn      indus      nox      rm      age      dis      tax ptratio      black      lstat
##      0          0          0          0          0          0          0          0          0          0
##      medv
##      0
```

```
head(x, 3)
```

```
##      zn      indus      nox      rm      age      dis      tax
## 1  0.2845483 -1.2866362 -0.1440749  0.4132629 -0.1198948  0.140075 -0.6659492
## 2 -0.4872402 -0.5927944 -0.7395304  0.1940824  0.3668034  0.556609 -0.9863534
## 3 -0.4872402 -0.5927944 -0.7395304  1.2814456 -0.2655490  0.556609 -0.9863534
##      ptratio      black      lstat      medv
## 1 -1.4575580  0.4406159 -1.0744990  0.1595278
## 2 -0.3027945  0.4406159 -0.4919525 -0.1014239
## 3 -0.3027945  0.3960351 -1.2075324  1.3229375
```

Question 1.

```

test <- x[1,]
train <- x[-1,]

te <- 1
tran <- (1:nrow(x))[-te]

dist_1 <- function(train, test) {
  diff <- train - matrix(rep(test, nrow(train)),
                        nrow=nrow(train), byrow=T)
  dists <- apply(abs(diff), 1, sum)
  dists <- as.numeric(dists)
  return(dists)
}

dist_2 <- function(train, test) {
  diff <- train - matrix(rep(test, nrow(train)),
                        nrow=nrow(train), byrow=T)
  dists <- sqrt(apply(diff^2, 1, sum))
  dists <- as.numeric(dists)
  return(dists)
}

dist_3 <- function(train, test) {
  diff <- train - matrix(rep(test, nrow(train)),
                        nrow=nrow(train), byrow=T)
  numer <- abs(diff)
  denom <- abs(train) + abs(matrix(rep(test, nrow(train)),
                                nrow=nrow(train), byrow=T))
  dists <- apply((numer/denom), 1, sum)
  dists <- as.numeric(dists)
  return(dists)
}

fhat <- function(dist_func, train, test, target, K) {
  dist_vector <- dist_func(train, test)
  closest_K <- order(dist_vector,
                    decreasing = F)[1:K]
  fhat <- mean(target[closest_K])
  return(fhat)
}

fhat_1 <- fhat(dist_func=dist_1,
              train=train,
              test=test,
              target=y[tran],
              K=10)
fhat_2 <- fhat(dist_func=dist_2,
              train=train,
              test=test,
              target=y[tran],
              K=10)
fhat_3 <- fhat(dist_func=dist_3,

```

```

        train=train,
        test=test,
        target=y[tran],
        K=10)

data.frame(l=1:3,fhat=c(fhat_1, fhat_2, fhat_3))

```

```

##    l    fhat
## 1 1 0.115894
## 2 2 0.201866
## 3 3 0.074659

```

```

# [1] 0.115894 0.201866 0.074659

```

Question 2.

```

set.seed(12345)
tran <- sample(nrow(x), 400) # Randomly selected 400 samples
te <- c(1:nrow(x))[-tran]

PE <- function(x, y) {
  train <- x[tran,]
  testset <- x[-tran,]
  y_te <- y[te]

  K_list <- seq(1,100)

  f1_mat <- matrix(0, nrow=nrow(testset), ncol=100)
  f2_mat <- matrix(0, nrow=nrow(testset), ncol=100)
  f3_mat <- matrix(0, nrow=nrow(testset), ncol=100)

  for (i in 1:nrow(testset)) {
    test <- testset[i,]
    y_i <- y_te[i]
    for (k in K_list) {
      fhat_1 <- fhat(dist_1, train, test, y[tran], K=k)
      fhat_2 <- fhat(dist_2, train, test, y[tran], K=k)
      fhat_3 <- fhat(dist_3, train, test, y[tran], K=k)
      f1_mat[i,k] <- (y_i - fhat_1)^2
      f2_mat[i,k] <- (y_i - fhat_2)^2
      f3_mat[i,k] <- (y_i - fhat_3)^2
    }
  }

  PE1 <- apply(f1_mat, 2, function(t) sqrt(mean(t)))
  PE2 <- apply(f2_mat, 2, function(t) sqrt(mean(t)))
  PE3 <- apply(f3_mat, 2, function(t) sqrt(mean(t)))

  # plot(PE1, type="l")
  # lines(PE2, col="blue")

```

```
# lines(PE3, col="red")

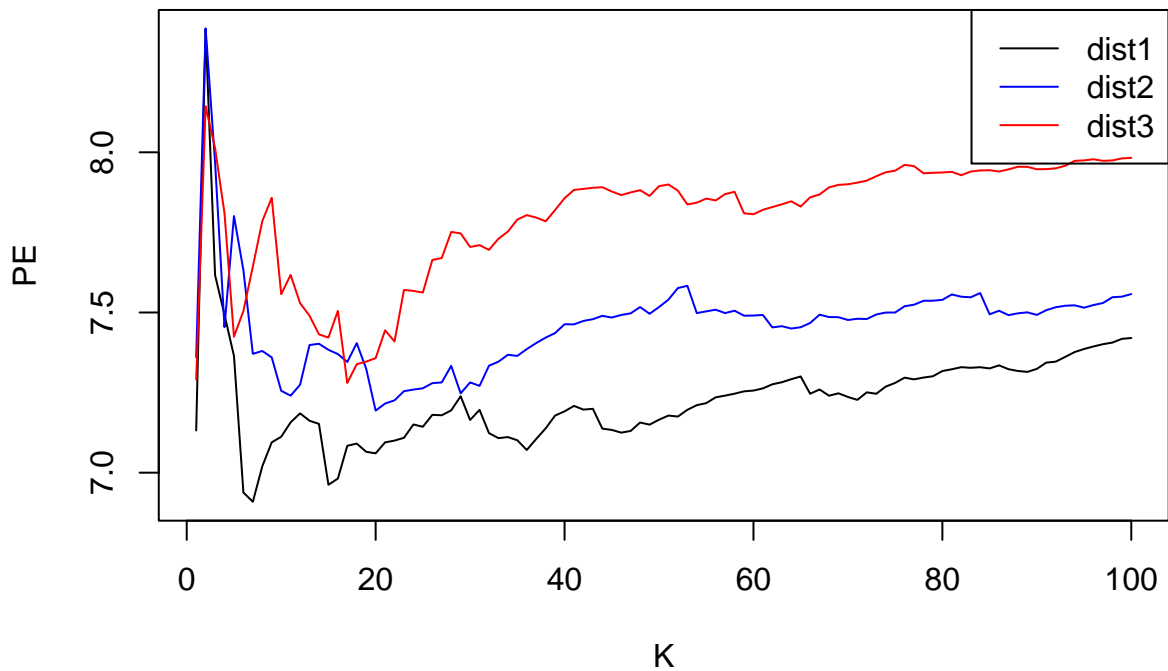
result_df <- data.frame(K_opt=c(which.min(PE1),
                                   which.min(PE2),
                                   which.min(PE3)),
                        PE_value=c(min(PE1),
                                   min(PE2),
                                   min(PE3)))

return(list(result_df, cbind(PE1, PE2, PE3)))
}

res2 <- PE(x,y)

plot(res2[[2]][,1], type="l",
      xlab="K", ylab="PE", main="3 PE(f) lines over K")
lines(res2[[2]][,2], col="blue")
lines(res2[[2]][,3], col="red")
legend("topright", legend=c("dist1", "dist2", "dist3"),
      col=c("black", "blue", "red"), lty=1)
```

3 PE(f) lines over K



```
res2[[1]]
```

```
##   K_opt PE_value
## 1    7 6.909187
## 2   20 7.193706
## 3   17 7.280037
```

Question 3.

```
set.seed(1234)
foldID <- sample(rep(1:10, length=nrow(x)))

CVE <- function(x, y){
  res_mat <- matrix(0, nrow=100, ncol=3)
  for (id in 1:10) {

    train <- x[foldID!=id,]
    testset <- x[foldID==id,]

    y_te <- y[foldID==id]

    K_list <- 1:100

    f1_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    f2_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    f3_mat <- matrix(0, nrow=nrow(testset), ncol=100)

    for (i in 1:nrow(testset)) {
      test <- testset[i,]
      y_i <- y_te[i]
      for (k in K_list) {
        fhat_1 <- fhat(dist_1, train, test, y[foldID!=id], K=k)
        fhat_2 <- fhat(dist_2, train, test, y[foldID!=id], K=k)
        fhat_3 <- fhat(dist_3, train, test, y[foldID!=id], K=k)
        f1_mat[i,k] <- (y_i - fhat_1)^2
        f2_mat[i,k] <- (y_i - fhat_2)^2
        f3_mat[i,k] <- (y_i - fhat_3)^2
      }
    }
    mkPE1 <- apply(f1_mat, 2, sum)
    mkPE2 <- apply(f2_mat, 2, sum)
    mkPE3 <- apply(f3_mat, 2, sum)

    res_mat <- res_mat + cbind(mkPE1, mkPE2, mkPE3)
  }

  CVE_mat <- sqrt(res_mat/506)

  # plot(CVE_mat[,1], type="l")
  # lines(CVE_mat[,2], col="blue")
  # lines(CVE_mat[,3], col="red")

  K_opt <- apply(CVE_mat, 2, which.min)
  CVE_value <- apply(CVE_mat, 2, min)

  result_df <- data.frame(K_opt=c(K_opt[1],
                                  K_opt[2],
                                  K_opt[3]),
                          CVE_value=c(CVE_value[1],
                                       CVE_value[2],
```

```

                                CVE_value[3]))
  rownames(result_df) = c("1","2","3")

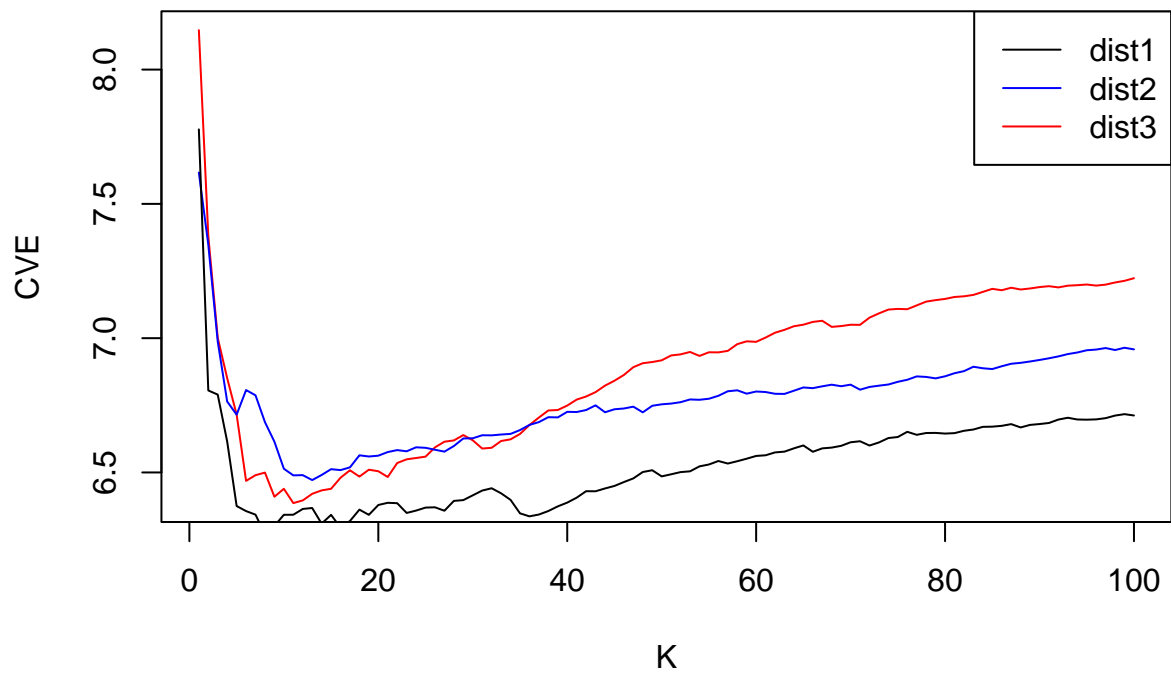
  return(list(result_df, CVE_mat))
}

res3 <- CVE(x,y)

plot(res3[[2]][,3], type="l", col="red",
      xlab="K", ylab="CVE", main="3 CVE(f) lines over K")
lines(res3[[2]][,2], col="blue")
lines(res3[[2]][,1])
legend("topright", legend=c("dist1", "dist2", "dist3"),
      col=c("black", "blue", "red"), lty=1)

```

3 CVE(f) lines over K



```
res3[[1]]
```

```

##   K_opt CVE_value
## 1     8  6.269773
## 2    13  6.471639
## 3    11  6.385981

```

Question 4.

```

ghat <- function(dist_func, train, test, target, K) {
  dist_vector <- dist_func(train, test)

```

```

closest_K <- order(dist_vector,
                   decreasing = F)[1:K]
ghat <- median(target[closest_K])
return(ghat)
}

set.seed(1234)
foldID <- sample(rep(1:10, length=nrow(x)))

CVE_g <- function(x, y){
  res_mat <- matrix(0, nrow=100, ncol=3)

  for (id in 1:10) {

    train <- x[foldID!=id,]
    testset <- x[foldID==id,]
    y_te <- y[foldID==id]

    K_list <- 1:100

    g1_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    g2_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    g3_mat <- matrix(0, nrow=nrow(testset), ncol=100)

    for (i in 1:nrow(testset)) {
      test <- testset[i,]
      y_i <- y_te[i]
      for (k in K_list) {
        ghat_1 <- ghat(dist_1, train, test, y[foldID!=id], K=k)
        ghat_2 <- ghat(dist_2, train, test, y[foldID!=id], K=k)
        ghat_3 <- ghat(dist_3, train, test, y[foldID!=id], K=k)
        g1_mat[i,k] <- (y_i - ghat_1)^2
        g2_mat[i,k] <- (y_i - ghat_2)^2
        g3_mat[i,k] <- (y_i - ghat_3)^2
      }
    }
    mkPE1 <- apply(g1_mat, 2, sum)
    mkPE2 <- apply(g2_mat, 2, sum)
    mkPE3 <- apply(g3_mat, 2, sum)

    res_mat <- res_mat + cbind(mkPE1, mkPE2, mkPE3)
  }

  CVE_mat <- sqrt(res_mat/506)

  # plot(CVE_mat[,1], type="l")
  # lines(CVE_mat[,2], col="blue")
  # lines(CVE_mat[,3], col="red")

  K_opt <- apply(CVE_mat, 2, which.min)
  CVE_value <- apply(CVE_mat, 2, min)

  result_df <- data.frame(K_opt=c(K_opt[1],

```



```

        K_opt[2],
        K_opt[3]),
        CVE_value=c(CVE_value[1],
                    CVE_value[2],
                    CVE_value[3]))

rownames(result_df) <- 1:3

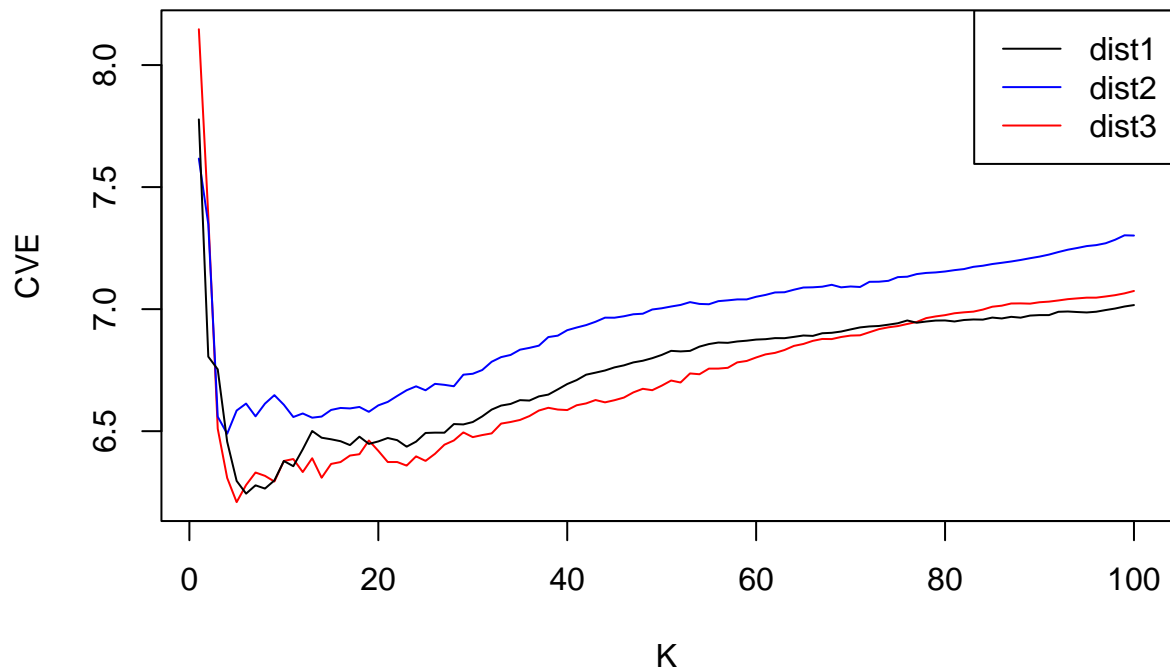
return(list(result_df, CVE_mat))
}

res4 <- CVE_g(x,y)

plot(res4[[2]][,3], type="l", col="red",
      xlab="K", ylab="CVE", main="3 CVE(g) lines over K")
lines(res4[[2]][,2], col="blue")
lines(res4[[2]][,1])
legend("topright", legend=c("dist1", "dist2", "dist3"),
      col=c("black", "blue", "red"), lty=1)

```

3 CVE(g) lines over K



```
res4[[1]]
```

```

##   K_opt CVE_value
## 1     6  6.244169
## 2     4  6.489489
## 3     5  6.209326

```

Question 5.

```
hhat <- function(dist_func, train, test, target, K) {
  dist_vector <- dist_func(train, test)
  closest_K <- order(dist_vector,
                     decreasing = F)[1:K]

  min_dist <- min(dist_vector)

  delta <- exp(-(dist_vector[closest_K] - min_dist)^2)
  D_lk <- sum(delta)

  hhat <- sum(delta*target[closest_K])/(D_lk)
  return(hhat)
}

set.seed(1234)
foldID <- sample(rep(1:10, length=nrow(x)))

CVE_h <- function(x, y){
  res_mat <- matrix(0, nrow=100, ncol=3)
  for (id in 1:10) {
    train <- x[foldID!=id,]
    testset <- x[foldID==id,]
    y_te <- y[foldID==id]

    K_list <- 1:100

    h1_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    h2_mat <- matrix(0, nrow=nrow(testset), ncol=100)
    h3_mat <- matrix(0, nrow=nrow(testset), ncol=100)

    for (i in 1:nrow(testset)) {
      test <- testset[i,]
      y_i <- y_te[i]
      for (k in K_list) {
        hhat_1 <- hhat(dist_1, train, test, y[foldID!=id], K=k)
        hhat_2 <- hhat(dist_2, train, test, y[foldID!=id], K=k)
        hhat_3 <- hhat(dist_3, train, test, y[foldID!=id], K=k)
        h1_mat[i,k] <- (y_i - hhat_1)^2
        h2_mat[i,k] <- (y_i - hhat_2)^2
        h3_mat[i,k] <- (y_i - hhat_3)^2
      }
    }
    mkPE1 <- apply(h1_mat, 2, sum)
    mkPE2 <- apply(h2_mat, 2, sum)
    mkPE3 <- apply(h3_mat, 2, sum)

    res_mat <- res_mat + cbind(mkPE1, mkPE2, mkPE3)
  }

  CVE_mat <- sqrt(res_mat/506)
```

```

# plot(CVE_mat[,1], type="l")
# lines(CVE_mat[,2], col="blue")
# lines(CVE_mat[,3], col="red")

K_opt <- apply(CVE_mat, 2, which.min)
CVE_value <- apply(CVE_mat, 2, min)

result_df <- data.frame(K_opt=c(K_opt[1],
                                K_opt[2],
                                K_opt[3]),
                        CVE_value=c(CVE_value[1],
                                    CVE_value[2],
                                    CVE_value[3]))

rownames(result_df) <- 1:3

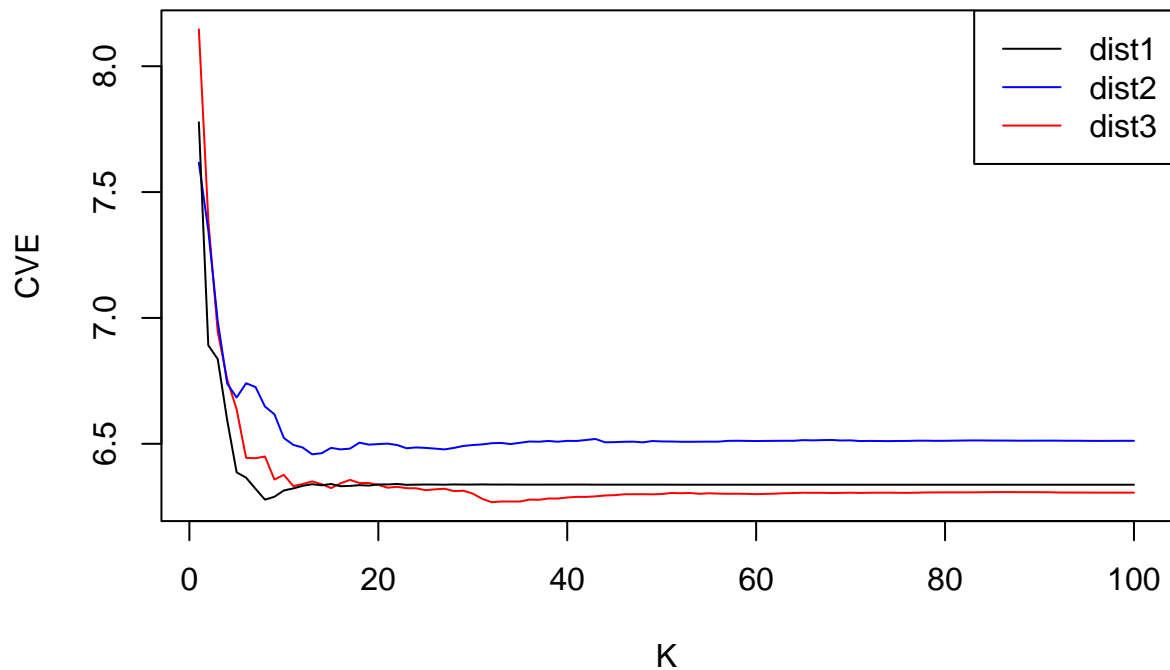
return(list(result_df, CVE_mat))
}

res5 <- CVE_h(x,y)

plot(res5[[2]][,3], type="l", col="red",
      xlab="K", ylab="CVE", main="3 CVE(h) lines over K")
lines(res5[[2]][,2], col="blue")
lines(res5[[2]][,1])
legend("topright", legend=c("dist1", "dist2", "dist3"),
      col=c("black", "blue", "red"), lty=1)

```

3 CVE(h) lines over K



```
res5[[1]]
```

```
##    K_opt CVE_value
## 1      8  6.278009
## 2     13  6.458178
## 3     32  6.268069
```

Question 6.

Result Table

	\hat{f}		\hat{g}		\hat{h}	
	K	CVE	K	CVE	K	CVE
d_1	8	6.270	6	6.244	8	6.278
d_2	13	6.472	4	6.489	13	6.458
d_3	11	6.386	5	6.209	32	6.268

Figure 1: Summary

결과 해석

- 본 교차검증 결과의 최소 **CVE** 관점에서는 예측함수 **g**와 거리함수 **d3**의 조합이 가장 좋았다.
- 하지만 **f**와 **g**는 **K**가 변함에 따라 **CVE**가 다소 불안정하여, **K**에 민감하다고 볼 수 있다. 이는 새로운 데이터에 대해 모델의 성능이 불안정해지는 문제를 초래할 수 있다.
- 반면 예측함수 **h**는 test sample과 특징이 더 가까울 수록 큰 가중치를 부여하여, 주변 샘플들의 target(crim) 가중평균을 예측값으로 사용하므로 그래프를 비교하였을 때 **K**가 증가해도 가장 안정적인 예측 성능을 보였다. 즉, **CVE**가 극도로 낮진 않더라도 새로운 sample의 예측을 다른 함수보다 더 안정적으로 수행할 수 있다.
- 또한 **h**는 거리함수 **d3**를 사용했을 때, 전반적으로 낮은 **CVE**값을 나타냈다.
- 결론적으로 예측의 안정성(일반화 성능)까지 고려했을 때는 예측함수 **h**와 거리함수 **d3**를 조합한 모델이 가장 좋다.