

SL_EX1

roi hezkiyahu

28 10 2022

```
library(dplyr)
library(ggplot2)
library(purrr)
library(caret)
library(class)
```

Q1

1. Population Optimizer of absolute loss

Prove that for absolute loss: $L_{\text{abs}}(Y, f(X)) = |Y - f(X)|$, EPE is minimized by setting $f^*(x) = \text{Median}(Y|X = x)$

Hint: you may find the following identity useful:

$$\int_{y>c} (y - c) dP(y) = \int_{y>c} Pr(Y > y) dy$$

(a) **Generalization to quantile loss** The τ th quantile loss for $0 < \tau < 1$ is defined as:

$$L_{\tau}(Y, f(X)) = \begin{cases} \tau \times (Y - f(X)) & \text{if } Y - f(X) > 0 \\ -(1 - \tau) \times (Y - f(X)) & \text{otherwise} \end{cases}$$

Prove that the EPE is minimized by setting $f^*(x)$ to be the τ th quantile of $P(Y|X = x)$, i.e., $P(Y \leq f^*(x)|X = x) = \tau$

notice that the median is a specific case for the quantile loss with $\tau = 0.5$ so proving for the general case will cover both questions

1a

$$\begin{aligned} \frac{\partial E_{Y|X}[(L_{\tau}(Y, f(X))|X = x]}{\partial f(x)} &= \frac{\partial}{\partial f(x)} \left[\int_{\min_y}^{f(x)} (1 - \tau) F_Y(y) dy + \int_{f(x)}^{\max_y} -\tau F_Y(y) dy \right] = F_Y(f(x)) - \tau F_Y(f(x)) - \tau + \tau F_Y(f(x)) = \\ &= F_Y(f(x)) - \tau = 0 \iff f(x) = F_Y^{-1}(\tau) \\ &\text{thus the minimizer is the } \tau\text{th quantile of } P(Y|X = x) \end{aligned}$$

Q2

2. **ESL 2.3:** Derive equation (2.24) (expected median distance to origin's nearest neighbor in an ℓ_p ball):

$$d(p, n) = \left(1 - \frac{1}{2}\right)^{1/p}$$

Suggested approach:

- Find the probability that all observations are outside a ball of radius $r < 1$, as a function of r .
- You are looking for r such that this probability is $1/2$.

Plot $d(p, n)$ against p for $n \in \{100, 5000, 100000\}$ and $p \in \{3, 5, 10, 20, 50, 100\}$ (make one curve for every value of n — use the R functions `plot()` and `lines()`) and interpret the graph.

denote $D_c = \min(\|X_1\|, \dots, \|X_n\|)$

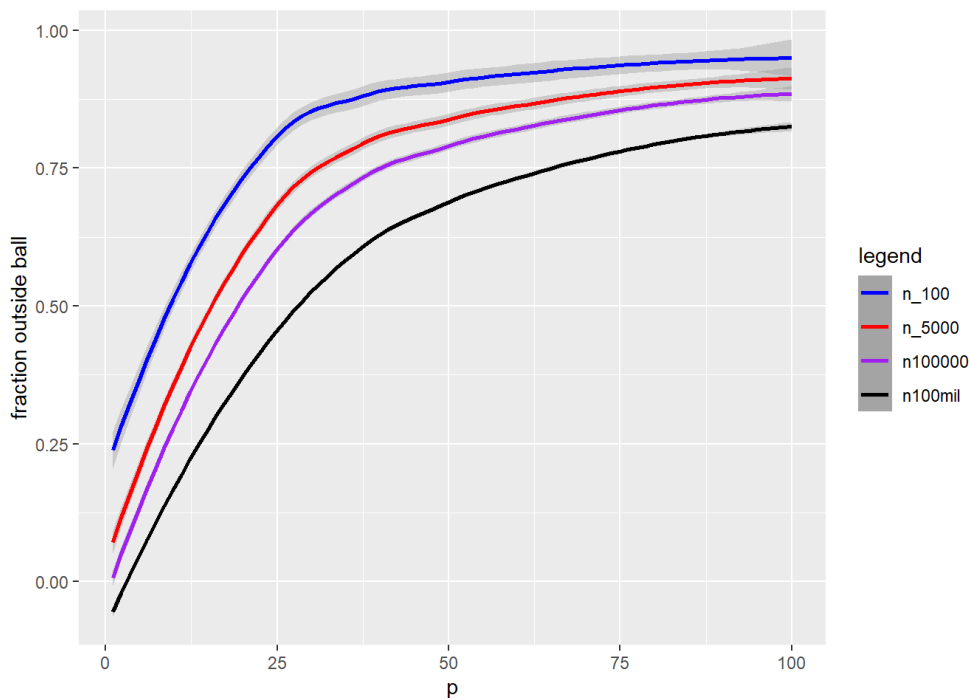
$P(\text{all observations are outside a ball of radius } r) = P(D_c > r)$

$$P(D_c > r) = \prod_{i=1}^n P(\|X_i\| > r) = P(\|X_1\| > r)^n = (1 - P(\|X_1\| < r))^n = (1 - r^p)^n$$

$$\text{we are looking for such } r \text{ such that } P(D_c > r) = \frac{1}{2} \Rightarrow (1 - r^p)^n = \frac{1}{2} \Rightarrow \left(1 - \frac{1}{2^{1/n}}\right)^{1/p} = r$$

```
radius_function <- function(n){
  p <- 1:100
  return( (1 - 1/2^(1/n))^(1/p))
}
n_s <- c(100,5000,100000,100000000)
df <- map(n_s,radius_function)
df <- as.data.frame(df)
colnames(df) <- c("n_100","n_5000","n_100000","n_100000000")
df["p"] = 1:100
colors = c("n_100"="blue", "n_5000" = "red", "n100000"="purple", "n100mil"= "black")
ggplot(data = df) +
  geom_smooth(aes(x = p, y = n_100,color = "n_100")) +
  geom_smooth(aes(x = p, y = n_5000,color = "n_5000")) +
  geom_smooth(aes(x = p, y = n_100000,color = "n100000")) +
  geom_smooth(aes(x = p, y = n_100000000,color = "n100mil")) +
  labs(x="p", y = "fraction outside ball", color = "legend") +
  scale_color_manual(values = colors)
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



we can see that as p grows the fraction of observations outside the ball increases and even for a very high number of observations

Q3

3. **ESL 2.7:** Compare classification performance of k -NN and linear regression on the `zipcode` data, on the task of separating the digits 2 and 3. Use $k \in \{1, 3, 5, 7, 15\}$. Plot training and test error for k -NN choices and linear regression. Comment on the shape of the graph.

```

acc_err <- function(y_true,y_pred){
  return (1-mean(y_true == y_pred))
}

df_train <- read.table("zip.train") %>% rename("y" = "V1") %>%
  filter(y %in% c(2,3)) %>% mutate(y = y-2)
df_test <- read.table("zip.test") %>% rename("y" = "V1") %>%
  filter(y %in% c(2,3)) %>% mutate(y = y-2)

get_knn_res <- function(k,df_train,df_test){
  X_tr <- as.matrix(df_train %>% select(-y))
  X_te <- as.matrix(df_test %>% select(-y))
  y_tr <- df_train %>% select(y)
  train_preds <- knn(train = X_tr,test = X_te,cl = as.matrix(y_tr),k = k)
  test_preds <- knn(train = X_tr,test = X_te,cl = as.matrix(y_tr),k = k)
  train_err <- acc_err(df_train$y,train_preds)
  test_err <- acc_err(df_test$y,test_preds)
  return( c(train_err, test_err))
}

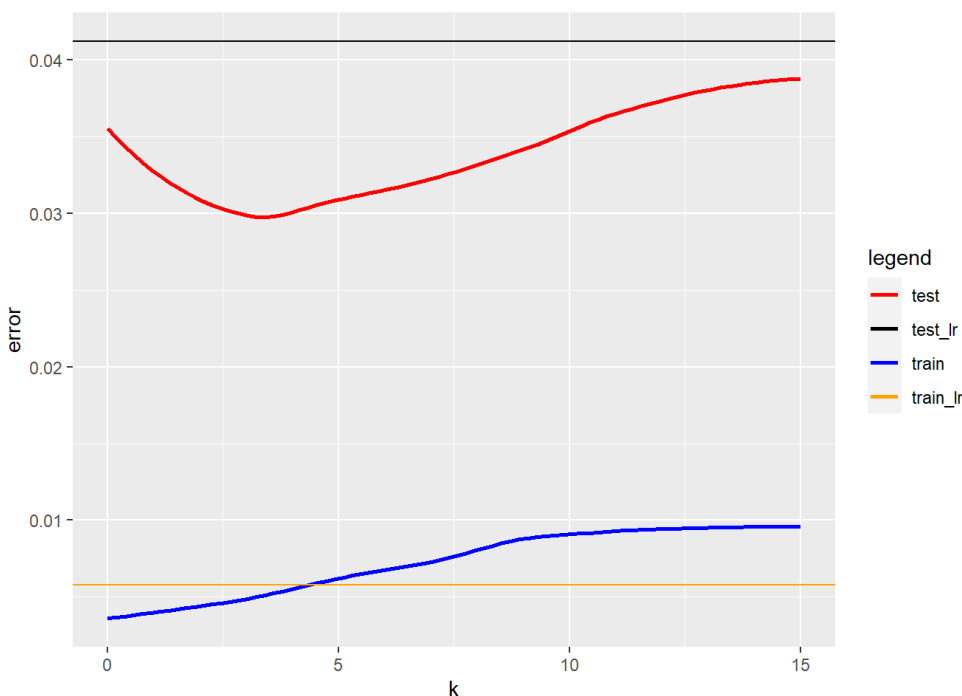
zip_lm <- lm(y~.,data = df_train)
lr_train_err <- acc_err(df_train$y,(predict(zip_lm,df_train) > 0.5)*1)
lr_test_err <- acc_err(df_test$y,(predict(zip_lm,df_test) > 0.5)*1)
k_s = 1:15
knn_train_res <- c()
knn_test_res <- c()
for (i in 1:length(k_s)){
  knn_res <- get_knn_res(k=k_s[i],df_train,df_test)
  knn_train_res <- c(knn_train_res,knn_res[1])
  knn_test_res <- c(knn_test_res,knn_res[2])
}
result_df <- tibble(k = c(0,k_s),train_res = c(lr_train_err,knn_train_res) ,test_res = c(lr_test_err,knn_test_res))
colors = c("train"="blue", "test" = "red","train_lr" = "orange", "test_lr" = "black")
ggplot(data = result_df) +
  geom_smooth(aes(x = k, y = train_res,color = "train"),se=FALSE) +
  geom_smooth(aes(x = k, y = test_res,color = "test"),se=FALSE) +
  geom_hline(aes(yintercept = lr_train_err, color = "train_lr")) +
  geom_hline(aes(yintercept = lr_test_err, color = "test_lr")) +
  labs(x="k", y = "error", color = "legend") +
  scale_color_manual(values = colors)

```

```

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'

```



we can see that as k increases the knn model error rate increase as well in the train and test set.
we can also see that linear regression has the worst test results

Q4

4. **ESL 2.9 (second edition only)** Consider a linear regression model, fit by least squares to a set of training examples $T = \{(X_1, Y_1), \dots, (X_N, Y_N)\}$, drawn i.i.d from some population. Let $\hat{\beta}$ be the least squares estimate. Suppose we also have some other ("test") data drawn independently from the same distribution $\{(\tilde{X}_1, \tilde{Y}_1), \dots, (\tilde{X}_M, \tilde{Y}_M)\}$. Prove that:

$$\frac{1}{N} \mathbb{E} \left(\sum_{i=1}^N (Y_i - X_i^T \hat{\beta})^2 \right) \leq \frac{1}{M} \mathbb{E} \left(\sum_{i=1}^M (\tilde{Y}_i - \tilde{X}_i^T \hat{\beta})^2 \right),$$

that is, the expected squared error in-sample is always bigger than out of sample in least squares fitting. Note that the values X are also random variables here, and the expectation is over everything that is random, including X, Y and $\hat{\beta}$.

Hint: There are several ways to prove this. One starts from considering the best possible linear model we derived in class:

$$\beta^* = (E(XX^T))^{-1}E(XY),$$

and comparing both sides to it.

Note: Students who find more than one valid way to prove the result will get a bonus grade.

* **Extra credit problem: Optimality of k-NN in fixed dimension**

Assume $X \sim \text{Unif}([0, 1]^p)$, and $Y = f(X) + \epsilon$ with $\epsilon \sim (0, \sigma^2)$ (that is, $f(x) = E(Y|X = x)$).

Assume f is Lipschitz: $\|x_1 - x_2\| < \delta \Rightarrow |f(x_1) - f(x_2)| < c\delta$, $\forall x_1, x_2 \in [0, 1]^p$. Choose any sequence $k(n)$ such that:

$$\begin{aligned} k(n) &\xrightarrow{n \rightarrow \infty} \infty \\ k(n)/n &\xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

Then:

$$\text{EPE}(\text{k-NN using } k(n)) \xrightarrow{n \rightarrow \infty} \text{EPE}(f) = \sigma^2$$

(The proof does not have to be completely formal, for example you can replace a binomial with its normal approximation without proof of the relevant asymptotics).

the expected error is the same for all observations so we can assume $M=N=1$

denote E_{tr}, E_{te} the expected train and test error

$$E((\tilde{Y} - \tilde{X}\hat{\beta})^2) \geq E((\tilde{Y} - \tilde{X}\tilde{\beta})^2) \quad (\tilde{\beta} \text{ being LSE for the test set}) \Rightarrow$$

$$E_{te} \geq E((\tilde{Y} - \tilde{X}\tilde{\beta})^2)$$

$$E((\tilde{Y} - \tilde{X}\tilde{\beta})^2) = E((Y - X\hat{\beta})^2) \quad (\text{its an expected value thus the point of estimate doesn't matter})$$

to conclude we get:

$$E_{te} = E((\tilde{Y} - \tilde{X}\hat{\beta})^2) \geq E((\tilde{Y} - \tilde{X}\tilde{\beta})^2) = E((Y - X\hat{\beta})^2) = E_{tr}$$

try #2, i am not completely sure about the variance claim here

the expected error is the same for all observations so we can assume $M=N=1$

denote E_{tr}, E_{te} the expected train and test error

$$E_{tr} = \sigma^2 + \text{bias}(f(x)) + \text{var}(f(x)) = \sigma^2 + \text{var}(f(x)) \quad (\hat{\beta} \text{ is unbiased})$$

$$E_{te} = \sigma^2 + \text{bias}(f(\tilde{x})) + \text{var}(f(\tilde{x}))$$

$$\text{var}(f(\tilde{x})) = \text{var}(f(x)), \text{ and } \text{bias}(f(\tilde{x})) \geq 0$$

$$\text{thus we get: } E_{te} \geq E_{tr}$$

Extra

let $x_0 \in X$ be a fixed point

$$\forall \epsilon \forall k \exists m_k; \forall n > m_k, \quad \|x_0 - x_j\| < \epsilon_k \quad \forall x_j \in N_k(x_0)$$

$$\text{thus from } f \text{ being Lipschitz: } |f(x_0) - f(x_1)| < \epsilon_k \delta$$

$$\begin{aligned} \text{thus } |\hat{f}(x_0) - f(x_0)| &= \left| \sum_{x_j \in N_k(x_0)} \frac{f(x_j)}{k} - f(x_0) \right| = \left| \sum_{x_j \in N_k(x_0)} \frac{f(x_j) - f(x_0) + f(x_0)}{k} - f(x_0) \right| \leq \\ &\leq \sum_{x_j \in N_k(x_0)} \frac{|f(x_j) - f(x_0)|}{k} + \left| \sum_{x_j \in N_k(x_0)} \frac{f(x_0)}{k} - f(x_0) \right| = \sum_{x_j \in N_k(x_0)} \frac{|f(x_j) - f(x_0)|}{k} \leq \epsilon_k \delta \end{aligned}$$

thus for $\varepsilon_k \xrightarrow{n \rightarrow \infty} 0 : bias(\hat{f}(x_0)) \xrightarrow{n \rightarrow \infty} 0$

$$V(E(\hat{f}(x_0))) = V\left(\sum_{x_j \in N_k(x_0)} \frac{E(f(x_j))}{k}\right) \xrightarrow{n \rightarrow \infty} V\left(\sum_{x_j \in N_k(x_0)} \frac{f(x_0)}{k}\right) = \sum_{x_j \in N_k(x_0)} \frac{V(f(x_0))}{k} = 0$$

we get that the bias and variance terms both converge to zero as n goes to ∞ so from the prediction error decomposition we get that

$$EPE(k\text{-nn using } k(n)) \xrightarrow{n \rightarrow \infty} \sigma^2$$