



UNIVERSITY OF  
MICHIGAN

UNIVERSITY OF MICHIGAN  
DATA MINING  
STATS415

---

## ASSIGNMENT 3

*Author:* Shu ZHOU  
*ID:* 19342932  
*Lab Section:* 001

September 28, 2020

1. (a). True.  $E[X] = \int \int xy P_{X,Y}(x,y) dx dy = \int x P_X(x) dx \int y P_Y(y) dy = E[X] E[Y]$

(b). False  $E[\epsilon|x] = 0$ .  $E[E[\epsilon|x]] = E[\epsilon] = E[0] = 0$ . we can just prove that.  
 $\text{cov}(\epsilon, x) = E(\epsilon x) - E(\epsilon) E(x) = E(x) E(\epsilon|x) - E(\epsilon) E(x) = 0$ . However, even though  $\epsilon$  and  $x$  are uncorrelated, they can still be dependent. E.g.  $x \sim U(-1, 1)$ ,  $\epsilon = -x$   
 Then  $E[\epsilon|x] = E(\epsilon) = \int_{-1}^1 \epsilon P(\epsilon) d\epsilon = 0$

(c). True. Already proved in (b)

(d). True.  $P_X = X(X^T X)^{-1} X^T$ , It is obvious that  $X^T X$  is symmetric, i.e.  $[X^T X]^{-1}]^T = (X^T X)^{-1}$   
 $P_X^T = [X(X^T X)^{-1} X^T]^T = (X^T)^T (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = P_X$ . Hence. Q.E.D

(e). True.  $(I - P_X)^2 = (I - P_X)(I - P_X) = I^2 - P_X I - I P_X + P_X^2 = I - P_X - P_X + P_X^2$   
 $P_X^2 = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X[(X^T X)^{-1} X^T X] (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = P_X$   
 Hence  $(I - P_X)^2 = I - P_X$  we can easily conclude that  $(I - P_X)^{\infty} = I - P_X$

(f). False. If  $y_i = x_i^T \beta + \epsilon_i$  where  $\epsilon_i$  is i.i.d., then the fitted values under OLS is  
 $\hat{y} = P_X y = X(X^T X)^{-1} X^T y$ ; Then the residuals are  $y - \hat{y} = (I - P_X) y$

The Covariance matrix of the residuals is.  
 $\text{Cov}((I - P_X) y) = (I - P_X) \text{Cov}(y) (I - P_X)^T = (I - P_X) (\sigma^2 I) (I - P_X)^T = \sigma^2 (I - P_X)$   
 Which is not diagonal, so the residuals might not be identically distributed

(g). True  $\hat{\epsilon} = y - \hat{y} = (I - P_X) y$   
 Then  $\hat{\epsilon}^T X = [(I - P_X) y]^T X = y^T (I - P_X)^T X = y^T (I - P_X) X = y^T (I - X(X^T X)^{-1} X^T) X$   
 $= y^T [X - X(X^T X)^{-1} (X^T X)] = y^T [X - X] = 0$

2 (a) We know that  $\hat{\beta} = (X^T X)^{-1} X^T y = \beta + (X^T X)^{-1} X^T \epsilon$   
 Then  $\text{cov}(\hat{\beta}) = \text{var}(\hat{\beta}|x) = \text{Var}(\beta|x + (X^T X)^{-1} X^T \epsilon|x) = \text{Var}(\beta|x) + \text{Var}((X^T X)^{-1} X^T \epsilon|x)$   
 $= (X^T X)^{-1} X^T \text{var}(\epsilon|x) [X^T X]^{-1}$   
 $= (X^T X)^{-1} X^T \sigma^2 I_n X (X^T X)^{-1}$   
 $= \sigma^2 (X^T X)^{-1} (X^T X) (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$

(b). Consider the vector  $Z = (Z_{1j}, Z_{2j}, \dots, Z_{nj})^T$  for the  $j$ -th feature, Its response is calculated by.  
 Hence,  $\hat{\gamma}_j = Z^T (X^T X)^{-1} X^T y$ , the prediction is unbiased  
 Since  $E(\hat{\gamma}_j) = E(Z^T) E((X^T X)^{-1} X^T y) = E(Z^T) E(\hat{\beta}) = Z^T \beta = E(y_j)$

(c)  $\text{Var}(\hat{\gamma}_j) = \sigma^2 Z^T (X^T X)^{-1} Z = Z^T \sigma^2 (\hat{\beta}) Z$

$s^2(\hat{\gamma}_j) = \text{MSE}(Z^T (X^T X)^{-1} X) = Z^T s^2(\hat{\beta}) Z$

The 1- $\alpha$  confidence interval of  $E(\hat{\gamma}_j)$  is given by:  $\hat{\gamma}_j \pm t(1-\alpha/2; n-p) s(\hat{\gamma}_j)$

### Q3.

(a)

```
library(ISLR)
data("Carseats")
inTrain <- createDataPartition(Carseats$Sales, p = 0.75, list = FALSE)
training <- Carseats[inTrain,]
testing <- Carseats[-inTrain,]

model1 <- lm(Sales ~ ., data = training)
summary(model1)

##
## Call:
## lm(formula = Sales ~ ., data = training)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.78574 -0.73181 -0.00855  0.66672  3.15417
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.9238208   0.7131565   8.306 3.86e-15 ***
## CompPrice     0.0899659   0.0048764  18.449 < 2e-16 ***
## Income        0.0165904   0.0021691   7.648 3.04e-13 ***
## Advertising    0.1173044   0.0131905   8.893 < 2e-16 ***
## Population     0.0004820   0.0004258   1.132  0.2586
## Price        -0.0957211   0.0031392 -30.492 < 2e-16 ***
## ShelfLocGood   4.9680452   0.1822050  27.266 < 2e-16 ***
## ShelfLocMedium 2.1178758   0.1466664  14.440 < 2e-16 ***
## Age           -0.0480296   0.0037441 -12.828 < 2e-16 ***
## Education     -0.0223281   0.0228958  -0.975  0.3303
## UrbanYes       0.2260941   0.1312760   1.722  0.0861 .
## USYes         -0.1646478   0.1781490  -0.924  0.3561
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.037 on 289 degrees of freedom
## Multiple R-squared:  0.8703, Adjusted R-squared:  0.8653
## F-statistic: 176.3 on 11 and 289 DF,  p-value: < 2.2e-16

pred1 <- predict(model1, testing)
Result1 <- postResample(pred1, testing$Sales)
head(Result1)

##      RMSE Rsquared      MAE
## 0.9882446 0.8803269 0.7875085

##      RMSE Rsquared      MAE
## 1.1487637 0.8415894 0.9375886

model2 <- lm(Sales ~ CompPrice+Income+Advertising+Price+ShelveLoc, data = training)
summary(model2)

##
## Call:
```

```
## lm(formula = Sales ~ CompPrice + Income + Advertising + Price +
##     ShelfLoc, data = training)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.6917 -0.9249 -0.0156  0.9100  4.2619
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.480998   0.670739   3.699 0.000258 ***
## CompPrice      0.094901   0.006060  15.660 < 2e-16 ***
## Income         0.016877   0.002710   6.228 1.63e-09 ***
## Advertising    0.112785   0.011250  10.026 < 2e-16 ***
## Price         -0.094083   0.003934 -23.913 < 2e-16 ***
## ShelfLocGood   4.873567   0.228355  21.342 < 2e-16 ***
## ShelfLocMedium 2.045924   0.183390  11.156 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.302 on 294 degrees of freedom
## Multiple R-squared:  0.7918, Adjusted R-squared:  0.7875
## F-statistic: 186.3 on 6 and 294 DF, p-value: < 2.2e-16

pred2 <- predict(model2, testing)
Result2 <- postResample(pred2, testing$Sales)
head(Result2)
```

```
##      RMSE  Rsquared      MAE
## 1.1519276 0.8340922 0.9325066
```

```
##      RMSE  Rsquared      MAE
##1.3994254 0.7627533 1.1437243
```

The error of the second model is no better than the first model, though the second model abolished all the non-significant variables. It is mainly because the fact that “non-significant” is not equivalent to “no-impact”. If the coefficient is not equal to zero, but we failed to reject the mistake of our null-hypothesis, then it would cause a type-II error.

(b) The training error will be lower when  $k = 1$ , however the testing error will be lower when  $k = 20$ . When  $k = 1$ , the closest training sample will be chosen to the test sample (which is itself). Hence, the training error of  $k = 1$  is zero(unbiased). However the test error is maximized since the model is the most noisy.

When  $K=20$ , we choose around a test sample based on that its category and the category of 19 of its closest neighbors, which makes the test error smaller, but the bias will increase and also increase the training error.

```
knn_training <- training[,c(1,2,3,4,6)]
knn_testing <- testing[,c(1,2,3,4,6)]
train_category <- knn_training[,1]
test_category <- knn_testing[,1]

pred3 <- knn(knn_training,knn_testing,cl=train_category,k=1,prob= "True")
pred4 <- knn(knn_training,knn_testing,cl=train_category,k=20,prob= "True")
```

(c) I would standardize variables in the data set first. Standardization removes scale effects caused by use of features with different measurement scales. Once standardization is performed on a set of features, the range and scale of the z-scores should be similar, providing the distributions of raw feature values are alike.

For example, if we want to use a knn model to predict a person's weight based on his height and average

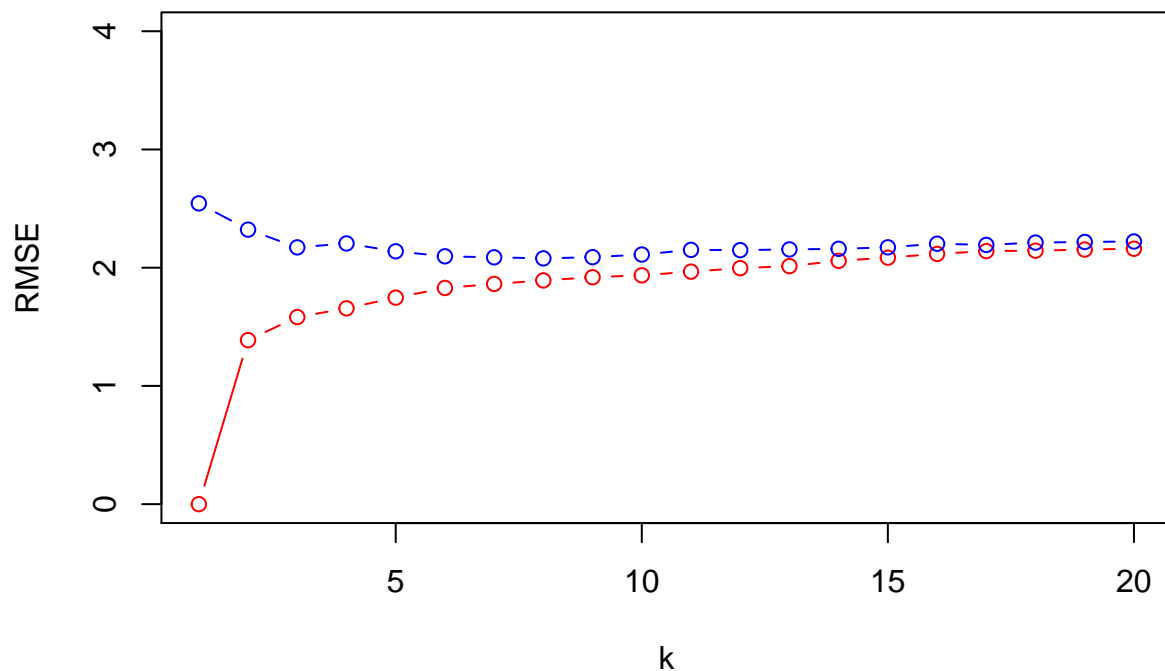
calories taken a day. The first variable is in the range among [150,200], however the second variable is in the range of [1000,3000]. Then second variable will have a much greater influence on the distance between samples and may bias the performance of the classifier.

(d)

```
train.RMSE = rep(0,20) ##Derive the error by calculating the RMSE of the regression
test.RMSE = rep(0,20)

for(k in 1 : 20){
fit <- knnreg(knn_training,train_category,k=k)
train.RMSE[k] <- sqrt(sum(abs(predict(fit,knn_training)-train_category)^2)/length(train_category))
test.RMSE[k] <- sqrt(sum(abs(predict(fit,knn_testing)-test_category)^2)/length(test_category))
}

plot(1:20, train.RMSE, xlab = "k", ylab = "RMSE",col='red', type = 'b', ylim = c(0,4))+points(1:20, tes
```

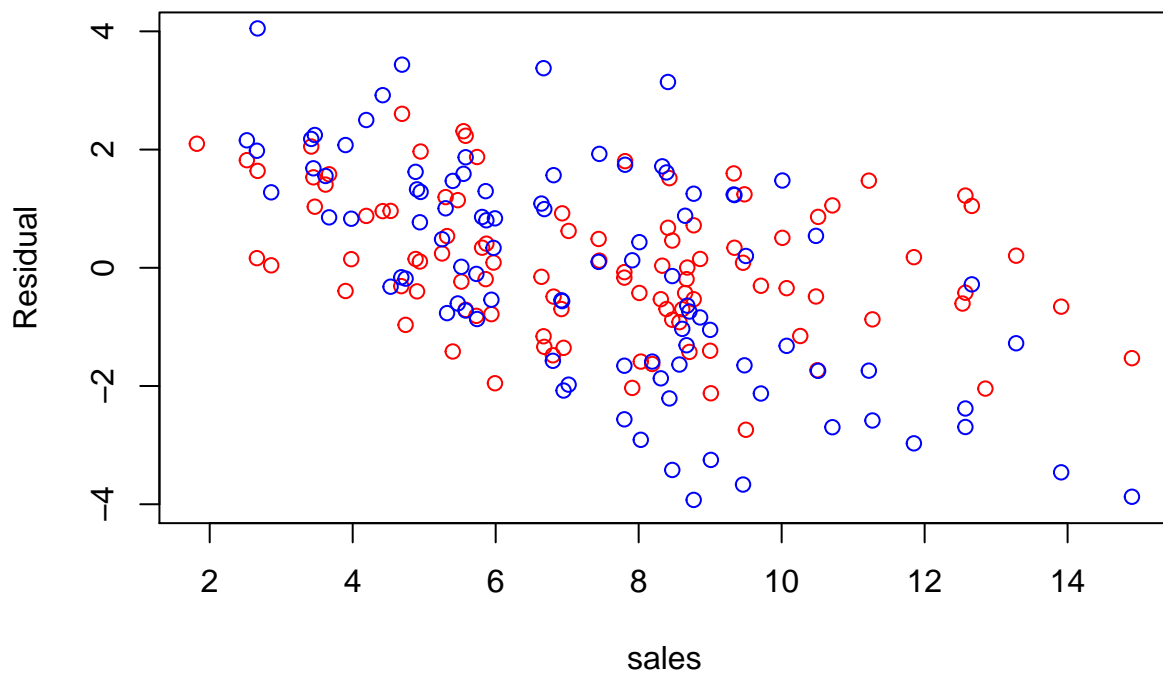


```
## integer(0)
```

Hence, from the result graph, we can see that both train error is minimized when  $k = 1$ , and test error is minimized when  $k = 3$ . (e) We choose  $k = 3$

```
test.residual = rep(0,99)
fit <- knnreg(knn_training,train_category,k=3)
test.residual<- predict(fit,knn_testing)-test_category
model2.res = pred2-test_category

plot(test_category,model2.res,xlab = "sales", ylab = "Residual",col='red', ylim = c(-4,4))+points(test_
```



```
## integer(0)
```

```
##The residual of linear regression is red and the residual of knn regression is blue
```

We can observe that the residual is independent of sales, and the residual of linear regression have less variance than the the residual of knn regression.