# Raport 3

## Eksploracja danych

Mikołaj Langner, Marcin Kostrzewa
nr albumów: 255716, 255749

2021-04-19

## Spis treści

# 1 Wstęp

# 2 Zadanie 1

## 2.1 Wczytanie danych i podział na zbiór uczący i testowy

Wczytajmy dane o irysach i podzielmy je na zbiór uczący i testowy w proporcji $1 : 2$.

```r
data(iris)
n <- dim(iris)[1]


train.set.index <- sample(1:n, 2/3*n)
train.set <- iris %>% slice(train.set.index) %>% arrange(Species)
test.set <- iris %>% slice(-train.set.index) %>% arrange(Species)

dummies <- dummyVars(" ~ .", data=iris)


train.dummies <- predict(dummies, newdata = train.set)
train.X <- as.matrix(cbind(rep(1, nrow(train.dummies)), train.dummies[, 1:4]))
train.Y <- train.dummies[, 5:7]
```

```r
test.dummies <- predict(dummies, newdata = test.set)
test.X <- as.matrix(cbind(rep(1, nrow(test.dummies)), test.dummies[, 1:4]))
test.Y <- test.dummies[, 5:7]

Y.hat <- solve(t(train.X) %*% train.X) %*% t(train.X) %*% train.Y

train.proba <- train.X %*% Y.hat
train.prediction <- colnames(train.proba)[apply(train.proba, 1, which.max)]

test.proba <- test.X %*% Y.hat
test.prediction <- colnames(test.proba)[apply(test.proba, 1, which.max)]

train.confusion <- table(train.set$Species, train.prediction)
train.confusion
```

```
##              train.prediction
##               Species.setosa Species.versicolor Species.virginica
##    setosa                  32                  0                 0
##    versicolor               0                 23                11
##    virginica                0                  5                29
```

```r
sum(diag(train.confusion)) / length(train.prediction)
```

```
## [1] 0.84
```

```r
test.confusion <- table(test.set$Species, test.prediction)
test.confusion
```

```
##              test.prediction
##               Species.setosa Species.versicolor Species.virginica
##    setosa                  18                  0                 0
##    versicolor               0                 12                 4
##    virginica                0                  3                13
```
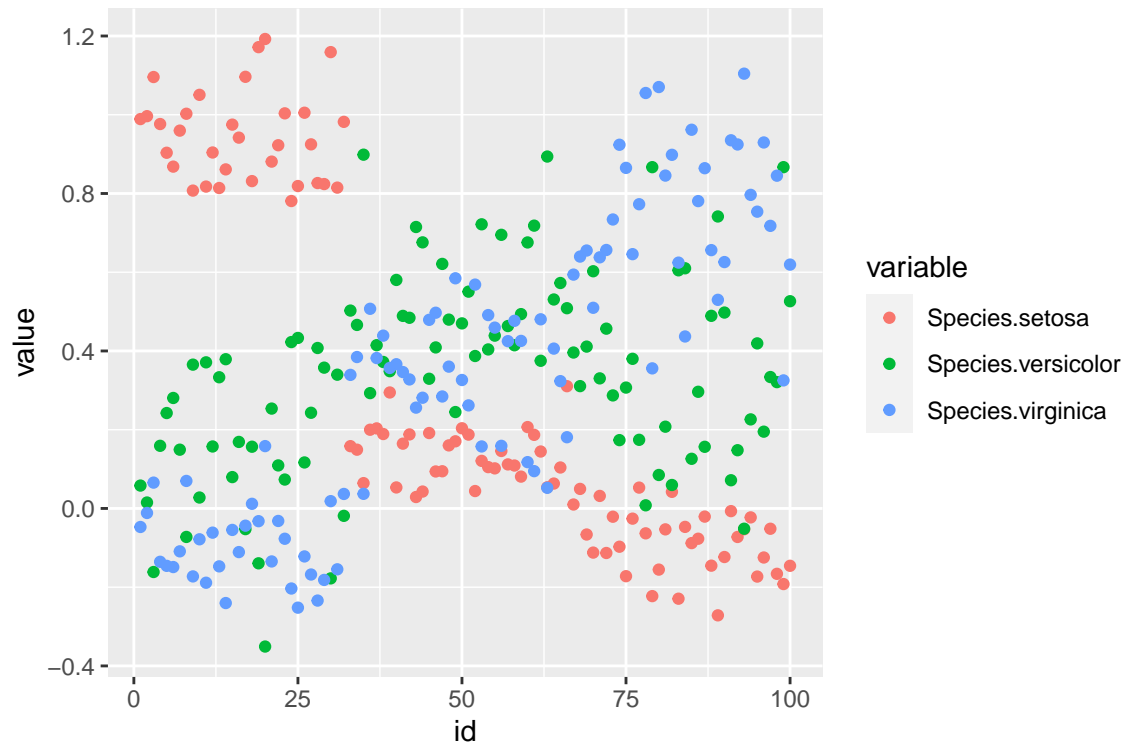
```r
sum(diag(test.confusion)) / length(test.prediction)
```

```
## [1] 0.86
```

```r
train.plot <- melt(as.data.frame(train.proba))
```
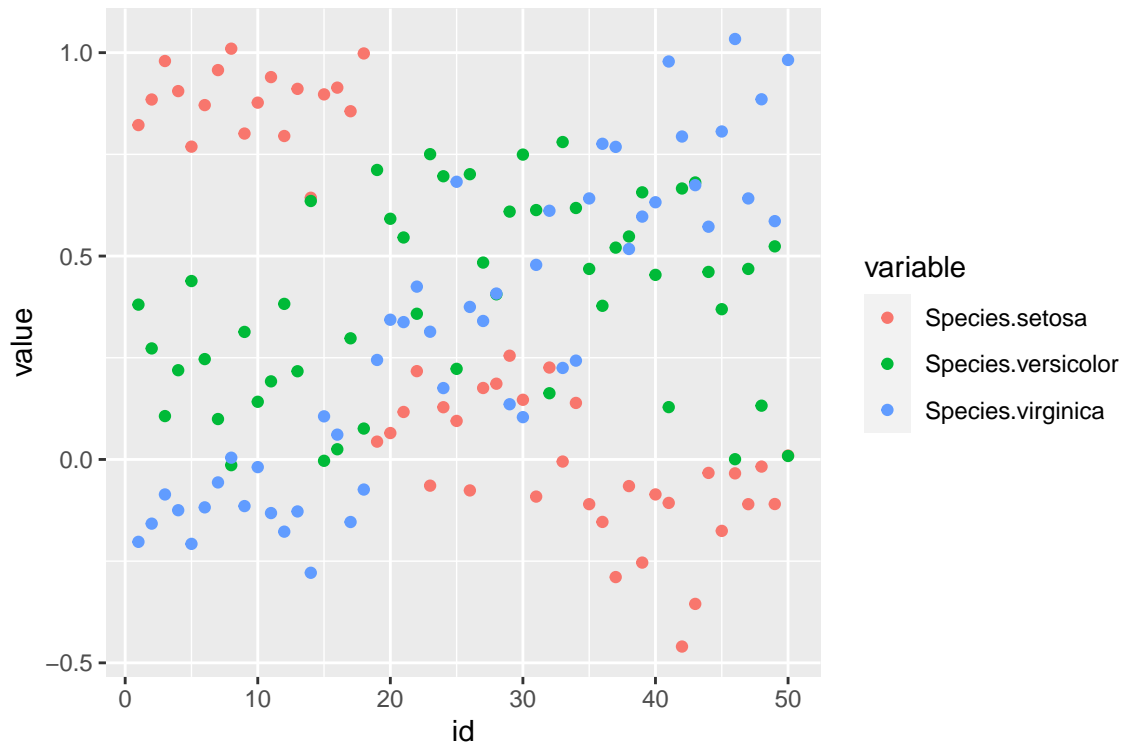
```
## No id variables; using all as measure variables
```

```r
train.plot$id <- as.integer(rownames(train.proba))
ggplot(train.plot, aes(x=id, y=value, color=variable)) +
  geom_point()
```

```
test.plot <- melt(as.data.frame(test.proba))
```

```
## No id variables; using all as measure variables
```

```
test.plot$id <- as.integer(rownames(test.proba))
ggplot(test.plot, aes(x=id, y=value, color=variable)) +
  geom_point()
```

```r
iris.quad <- (iris %>% select(-Species))^2
colnames(iris.quad) <- c("SL", "SW", "PL", "PW")
iris <- cbind(iris, combn(iris %>% select(-Species), 2, FUN = Reduce, f = `*`), iris.qu
```

```r
train.set.index <- sample(1:n, 2/3*n)
train.set <- iris %>% slice(train.set.index) %>% arrange(Species)
test.set <- iris %>% slice(-train.set.index) %>% arrange(Species)
```

```r
dummies <- dummyVars(" ~ .", data=iris)

train.dummies <- predict(dummies, newdata = train.set)
train.X <- as.matrix(cbind(rep(1, nrow(train.dummies)), train.dummies[, -c(5:7)]))
train.Y <- train.dummies[, 5:7]

test.dummies <- predict(dummies, newdata = test.set)
test.X <- as.matrix(cbind(rep(1, nrow(test.dummies)), test.dummies[, -c(5:7)]))
test.Y <- test.dummies[, 5:7]
```

```r
Y.hat <- solve(t(train.X) %*% train.X) %*% t(train.X) %*% train.Y
```

```r
train.proba <- train.X %*% Y.hat
train.prediction <- colnames(train.proba)[apply(train.proba, 1, which.max)]

test.proba <- test.X %*% Y.hat
test.prediction <- colnames(test.proba)[apply(test.proba, 1, which.max)]
```

4

```
train.confusion <- table(train.set$Species, train.prediction)
train.confusion
```

```
##             train.prediction
##              Species.setosa Species.versicolor Species.virginica
##   setosa                 33                  0                 0
##   versicolor              0                 35                 1
##   virginica               0                  0                31
```

```
sum(diag(train.confusion)) / length(train.prediction)
```

```
## [1] 0.99
```

```
test.confusion <- table(test.set$Species, test.prediction)
test.confusion
```

```
##             test.prediction
##              Species.setosa Species.versicolor Species.virginica
##   setosa                 17                  0                 0
##   versicolor              0                 14                 0
##   virginica               0                  5                14
```
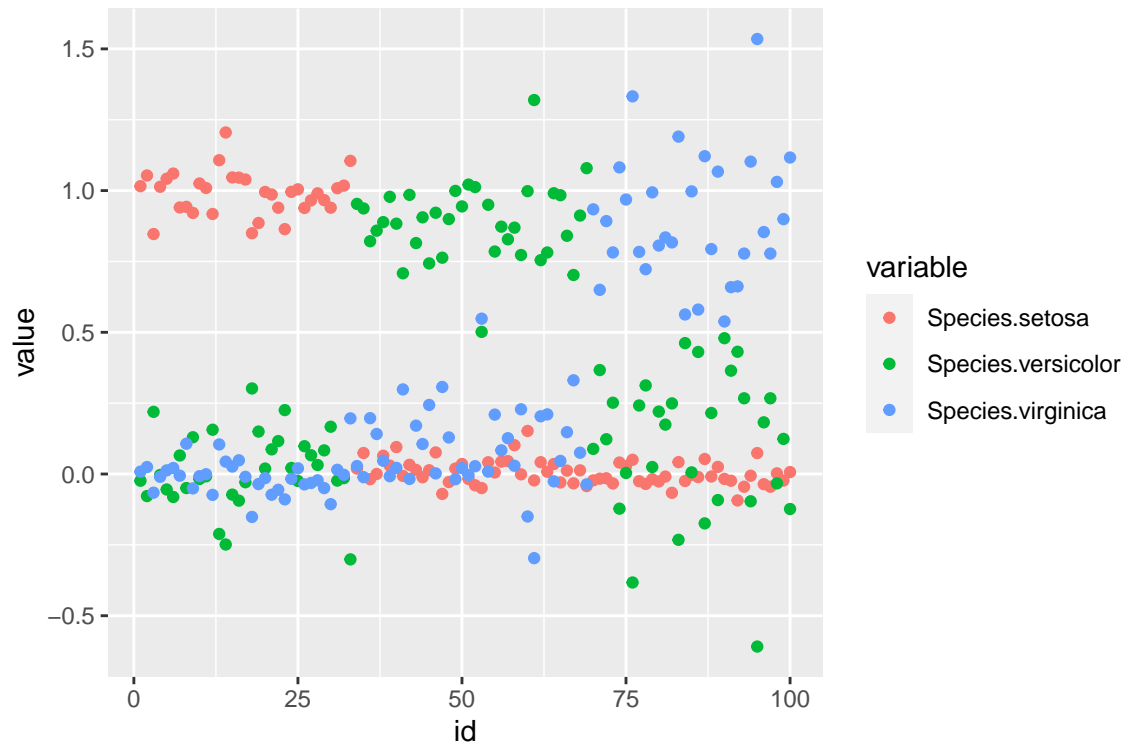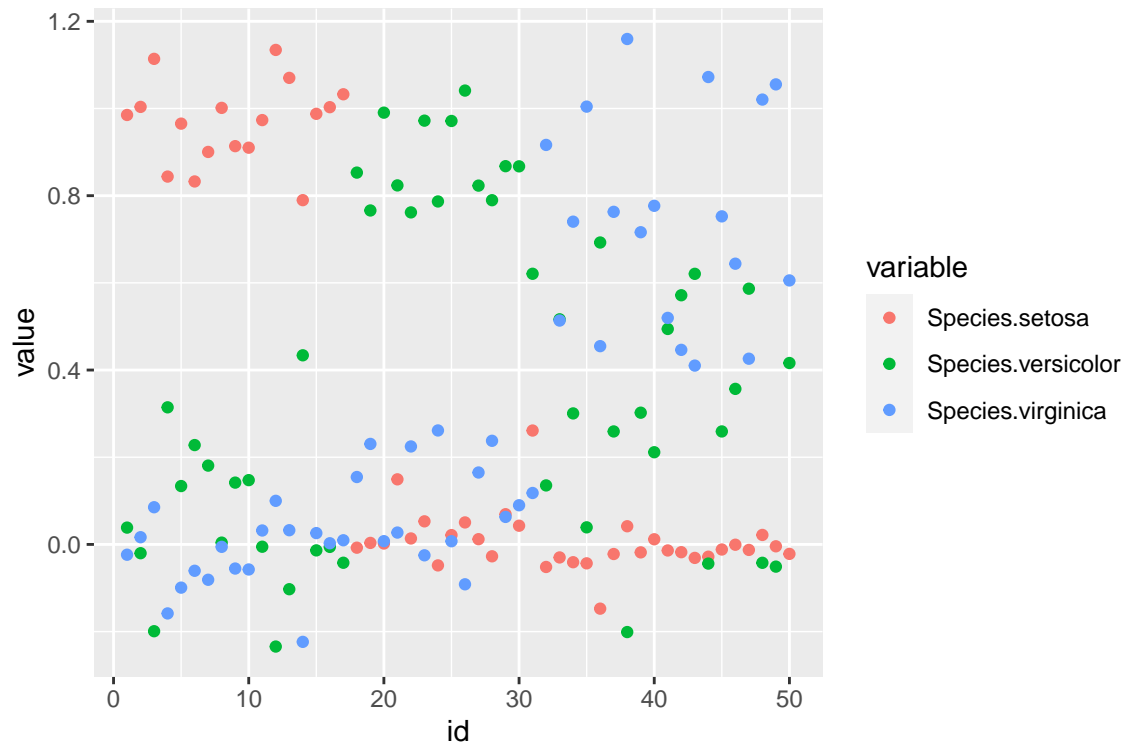
```
sum(diag(test.confusion)) / length(test.prediction)
```

```
## [1] 0.9
```

```
train.plot <- melt(as.data.frame(train.proba))
```

```
## No id variables; using all as measure variables
```

```
train.plot$id <- as.integer(rownames(train.proba))
ggplot(train.plot, aes(x=id, y=value, color=variable)) +
  geom_point()
```

```
test.plot <- melt(as.data.frame(test.proba))
```

```
## No id variables; using all as measure variables
```

```
test.plot$id <- as.integer(rownames(test.proba))
ggplot(test.plot, aes(x=id, y=value, color=variable)) +
  geom_point()
```

# 3 Zadanie 2