

Digit Recognizer

Digit Recognition of Handwritten Digits (Kaggle Competition)

The goal of this project is to correctly identify the digits from a dataset of thousands of handwritten images. I have used a Random Forest classifier to predict the handwritten digits. Accuracy achieved against the cross validation data is: 95.4%

The dataset used from MNIST (Modified National Institute of Standards and Technology) is a classic dataset for handwritten digits

```
library(randomForest)

## randomForest 4.6-12

## Type rfNews() to see new features/changes/bug fixes.

library(caret)

## Loading required package: lattice

## Loading required package: ggplot2

## 
## Attaching package: 'ggplot2'

## The following object is masked from 'package:randomForest':
## 
##     margin
```

Load the training and test dataset.

The train and test dataset contains the grayscale images of handwritten digits from 0 to 9.

The train dataset has 785 columns. One column is “label”, which specifies the digit drawn by the user. The rest of the 784 columns are pixel values of the associated image.

Each image is a 28 x 28 grayscale image with 784 pixels in total. Each pixel value is associated with it indicate the lightness or the darkness of the pixel - higher the pixel value is, darker will be the pixel.

The test dataset has 784 columns, which represent the 784 pixels of the handwritten image. There is no "label" column in the test dataset (Objective is to create a label column for the test dataset using the predicted values.)

```
# Load the train and test dataset
train.df <- read.csv(file="data/train.csv", header = T)
test.df <- read.csv(file="data/test.csv", header = T)

# "Label" in the train dataset need to be a factor
train.df$label <- as.factor(train.df$label)

# Train dataset is split in 80:20 ratio for cross-validation
div <- sample(1:nrow(train.df), 0.8*nrow(train.df))
train <- train.df[div,]
test <- train.df[-div,]
```

The classifier used for this classification is Random Forest classifier. A random forest can be considered as an ensemble technique. This technique fits N number of decision trees on various sub-samples (with replacement) of the dataset and use voting majority (in case of categorical variables) to improve upon the predictive accuracy and control over-fitting.

```
set.seed(1111)
# Use the 80% of the training dataset to train the model.
# Have set ntree to 500, which is the number of trees to grow for this model
# Have set mtry to 5, which is the number of variables randomly sampled as candidates at each split.
model.rf <- randomForest(label~, data=train, ntree=500, mtry=5)

# Use the model generated to predict the label of hand-written digits in the cross-validation data.
pred.rf.cv <- predict(model.rf, test)

# Find the accuracy of the prediction on cross-validation data
mean(test$label == pred.rf.cv)
```

```
## [1] 0.9567857
```

Accuracy of the random forest classifier against the cross validation data stands at 95.5%, which is pretty good. Now let's test it against the test data.

```
# Use the model generated to predict the label of hand-written digits in the test dataset.
pred.test <- predict(model.rf, test.df)

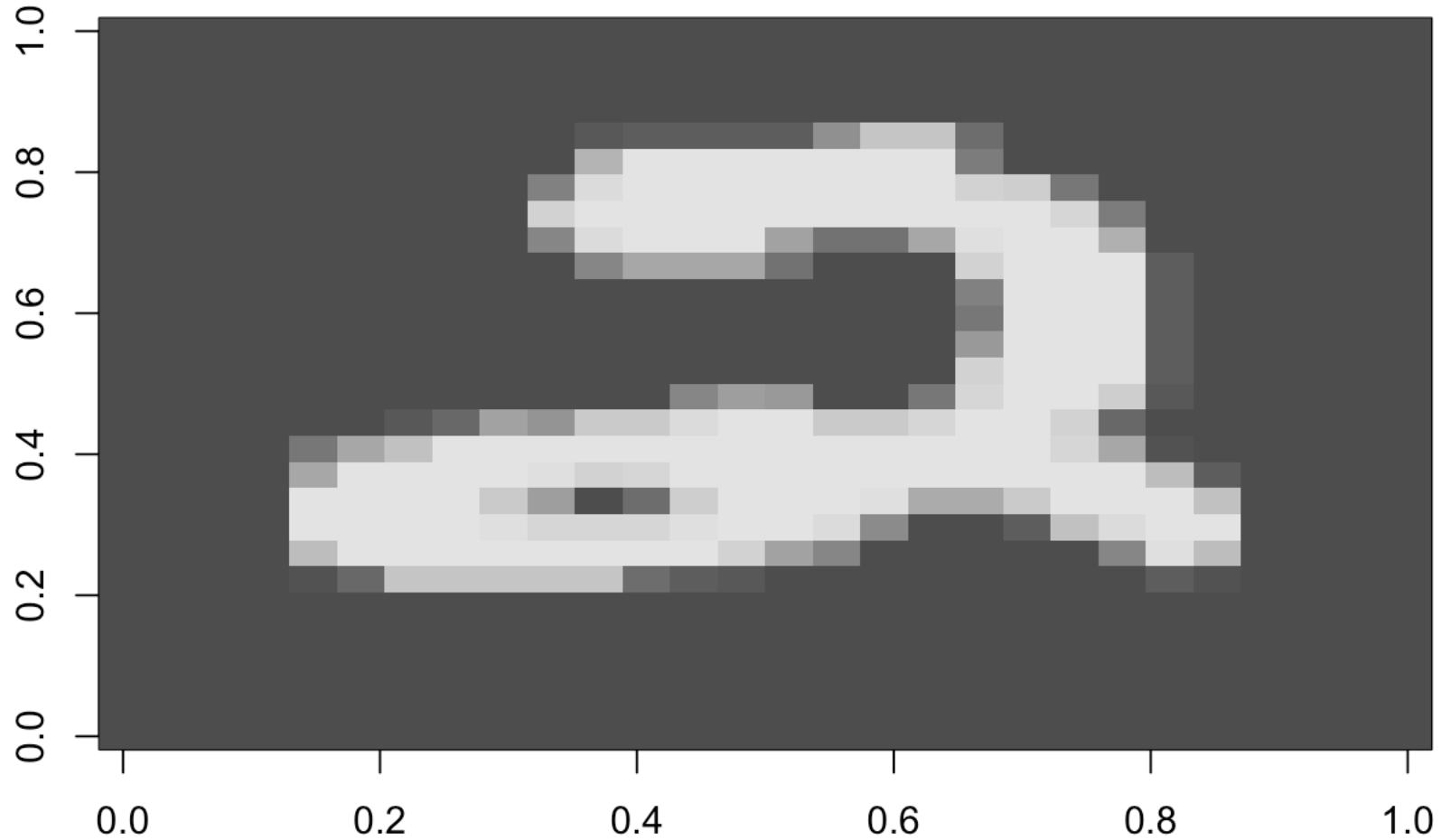
# Append the labels predicted to the test dataset.
test.df$label <- pred.test
```

Test data doesn't include a label field, and hence we will have to manually verify the correctness of the prediction. We can visualize the grayscale images on the test dataset and match them with the predicted values. I have tried to visualize the first 100 samples from the test dataset as below:

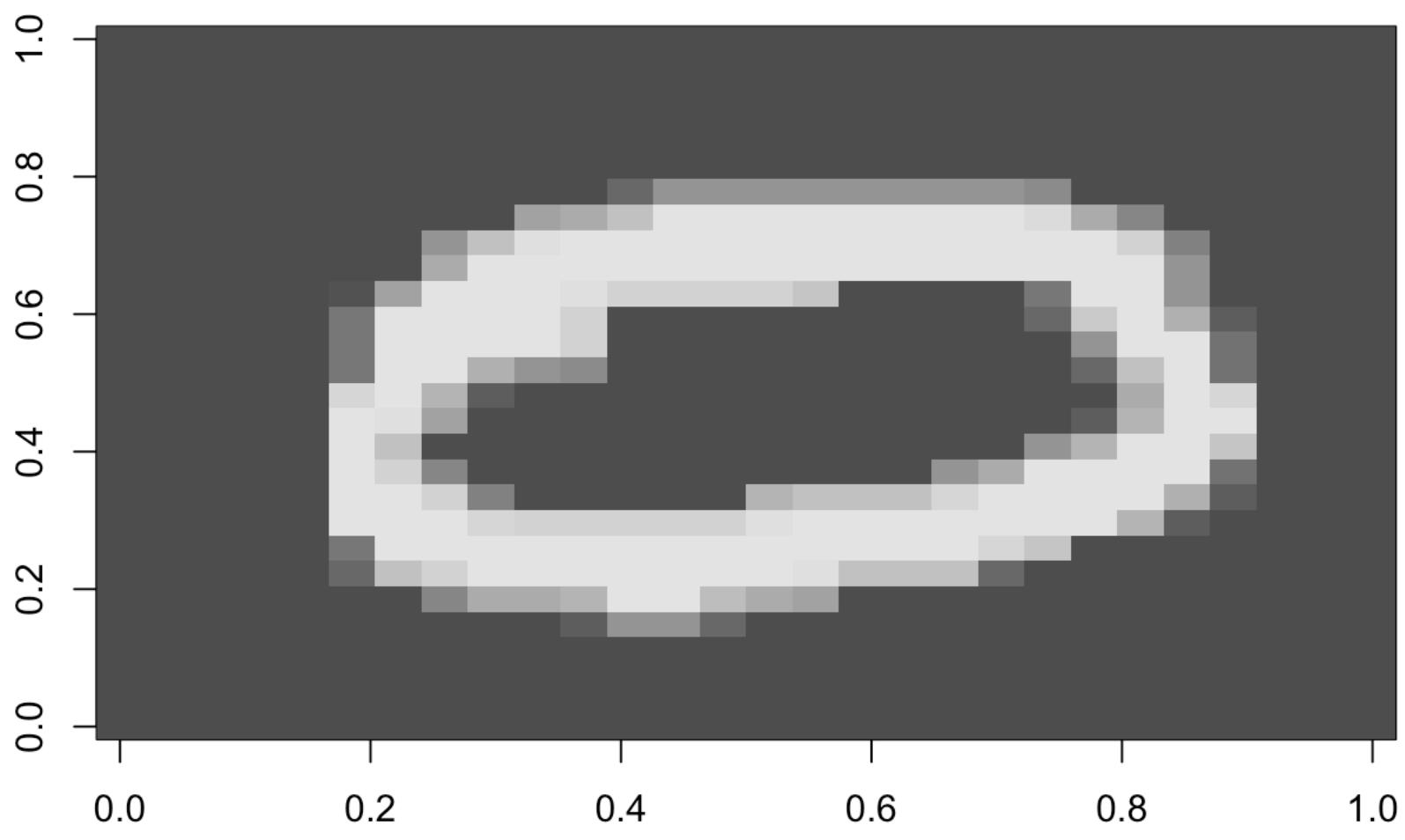
```
rotate <- function(x) t(apply(x, 2, rev))

for (i in 1:100) {
  m = rotate(matrix(unlist(test.df[i, -785]), nrow = 28, byrow = T))
  n <- test.df[i, 785]
  image(m, col = grey.colors(255), main = n)
}
```

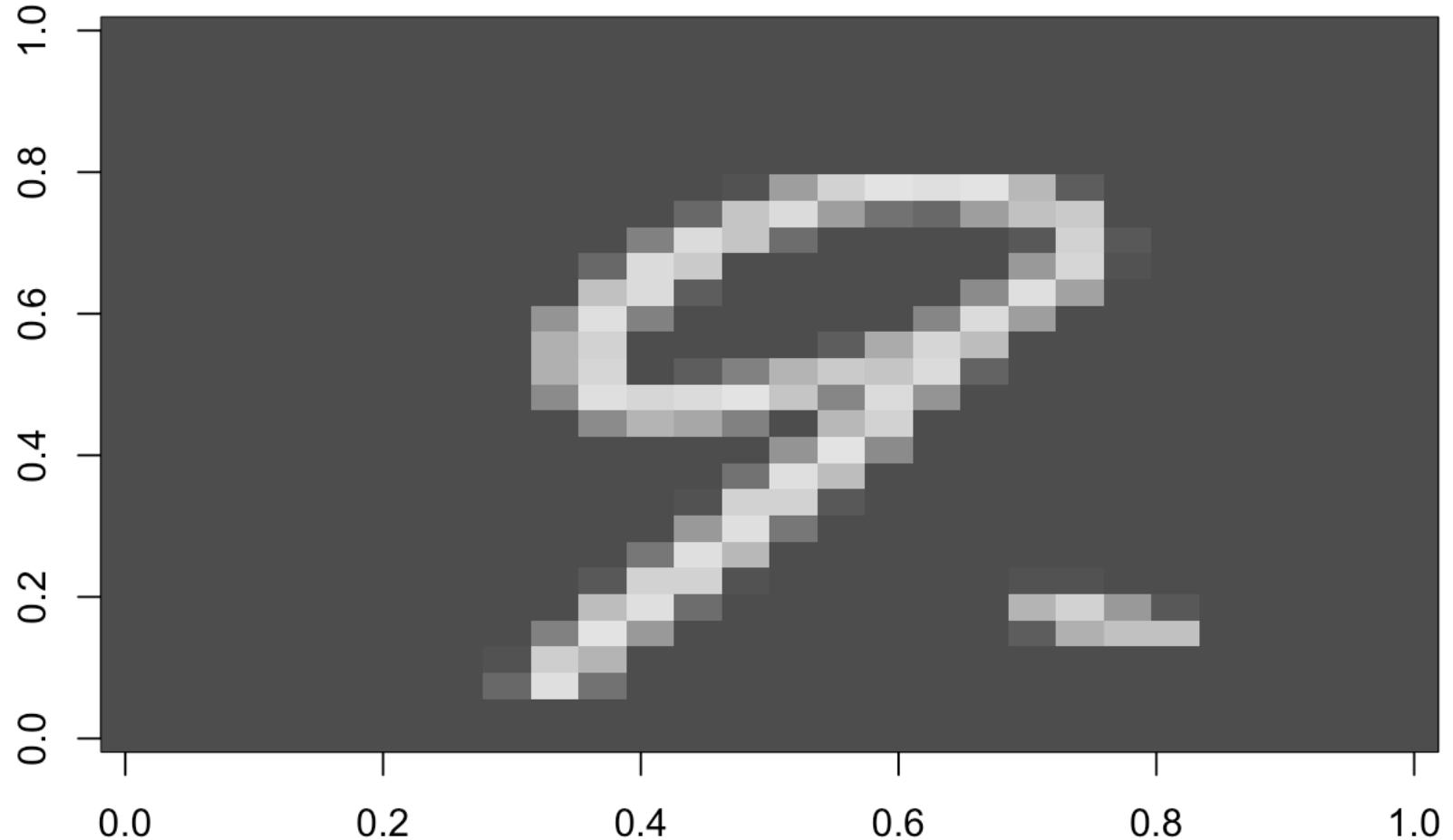
2



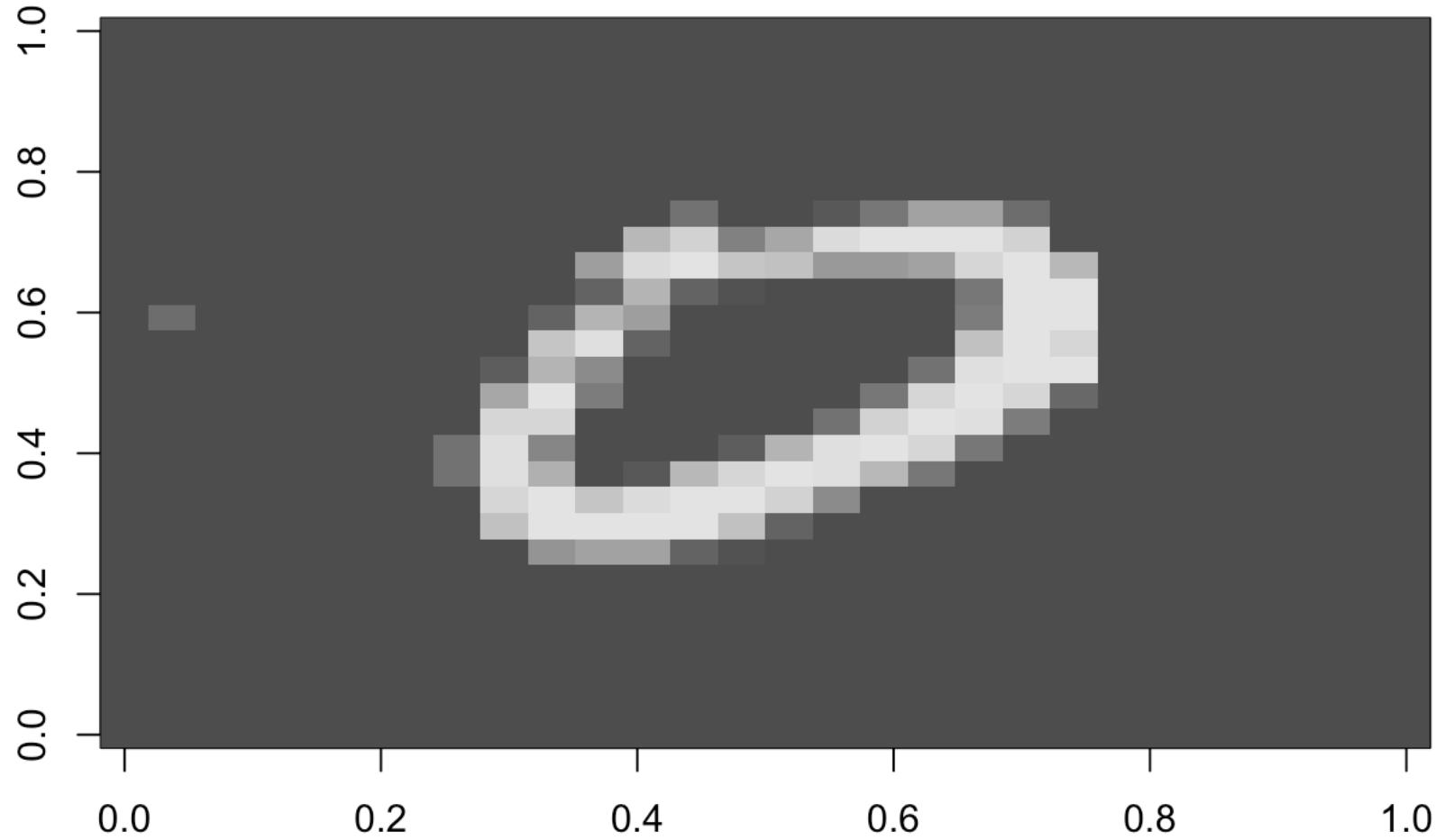
0



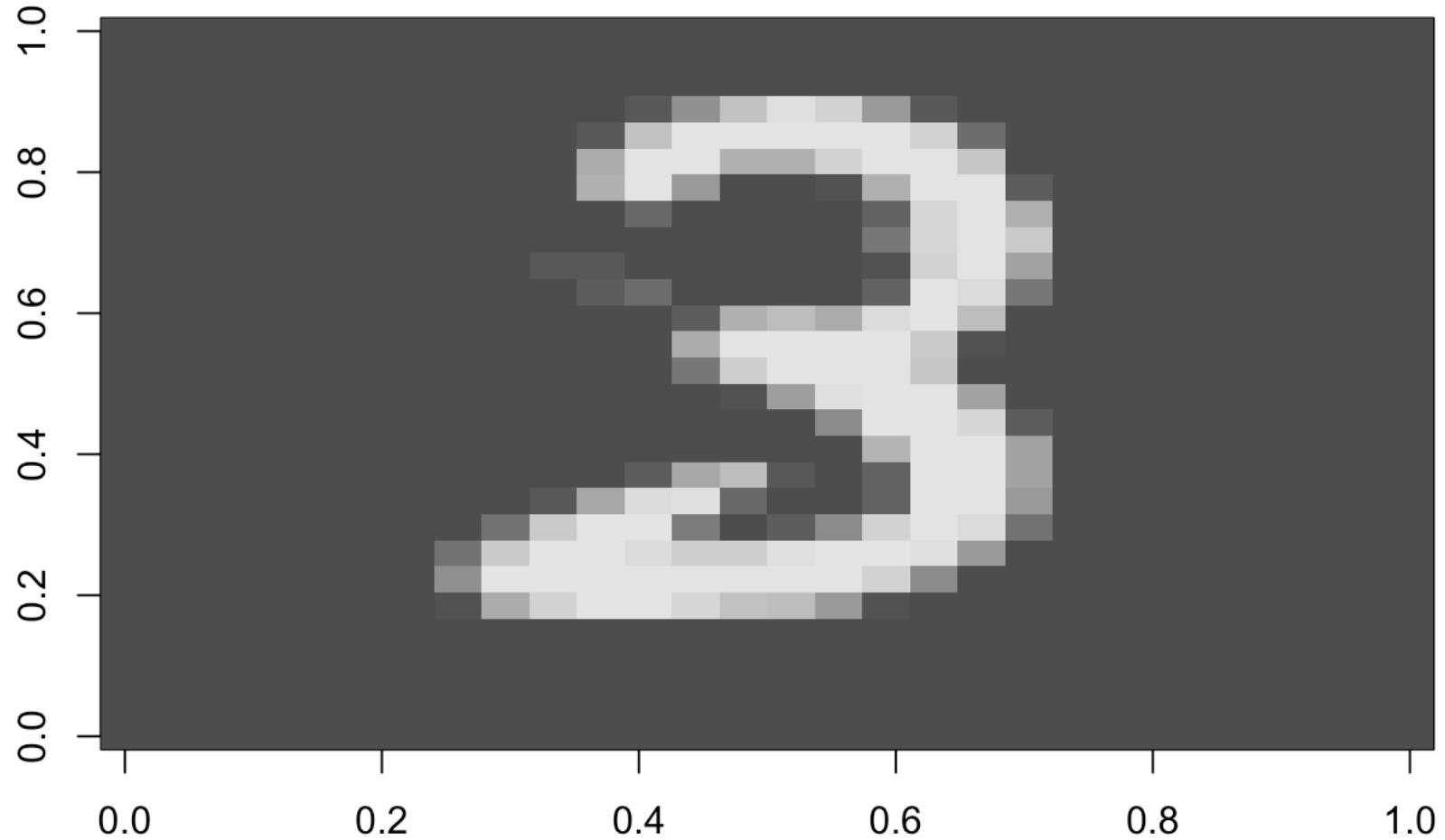
9

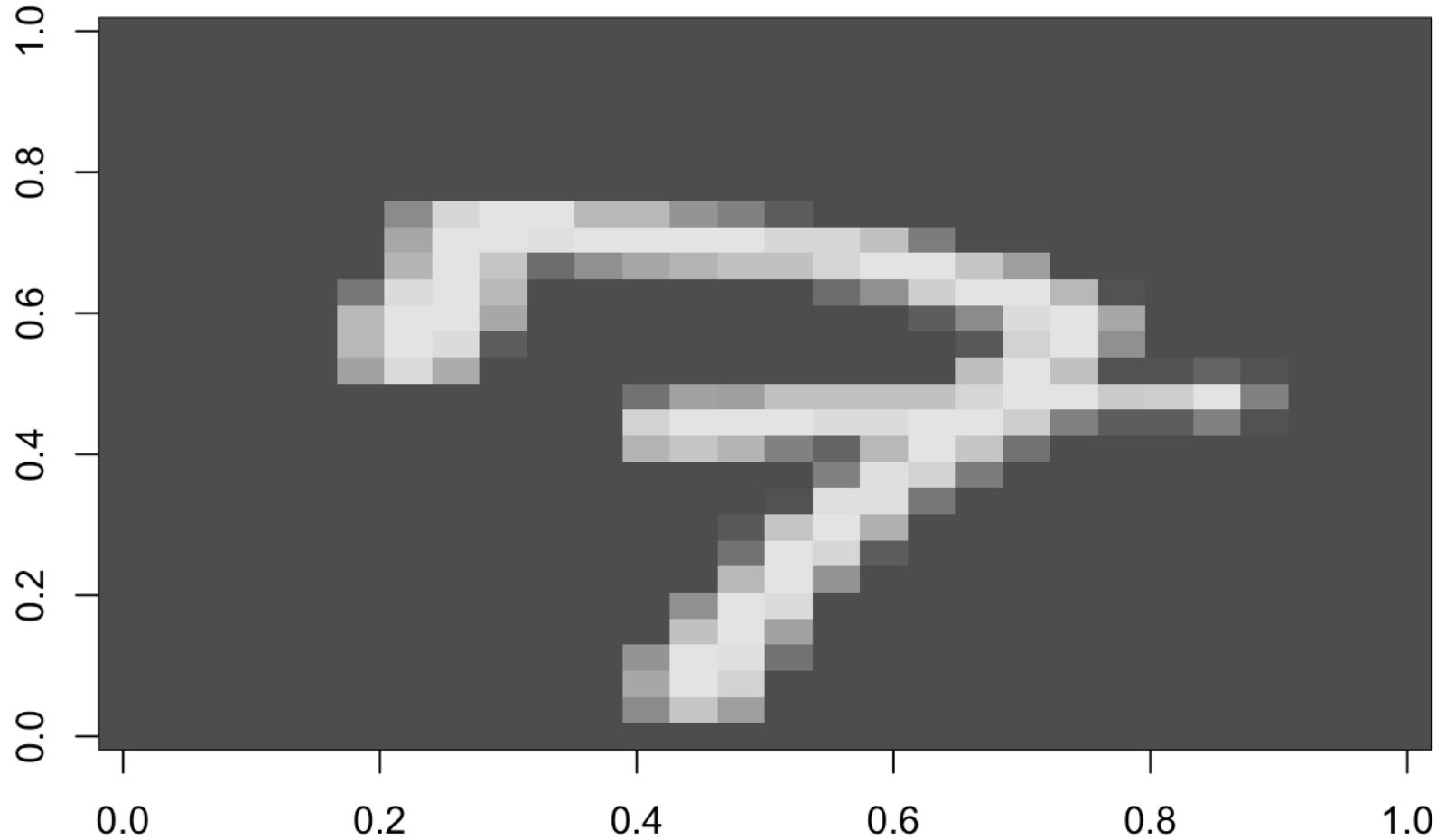


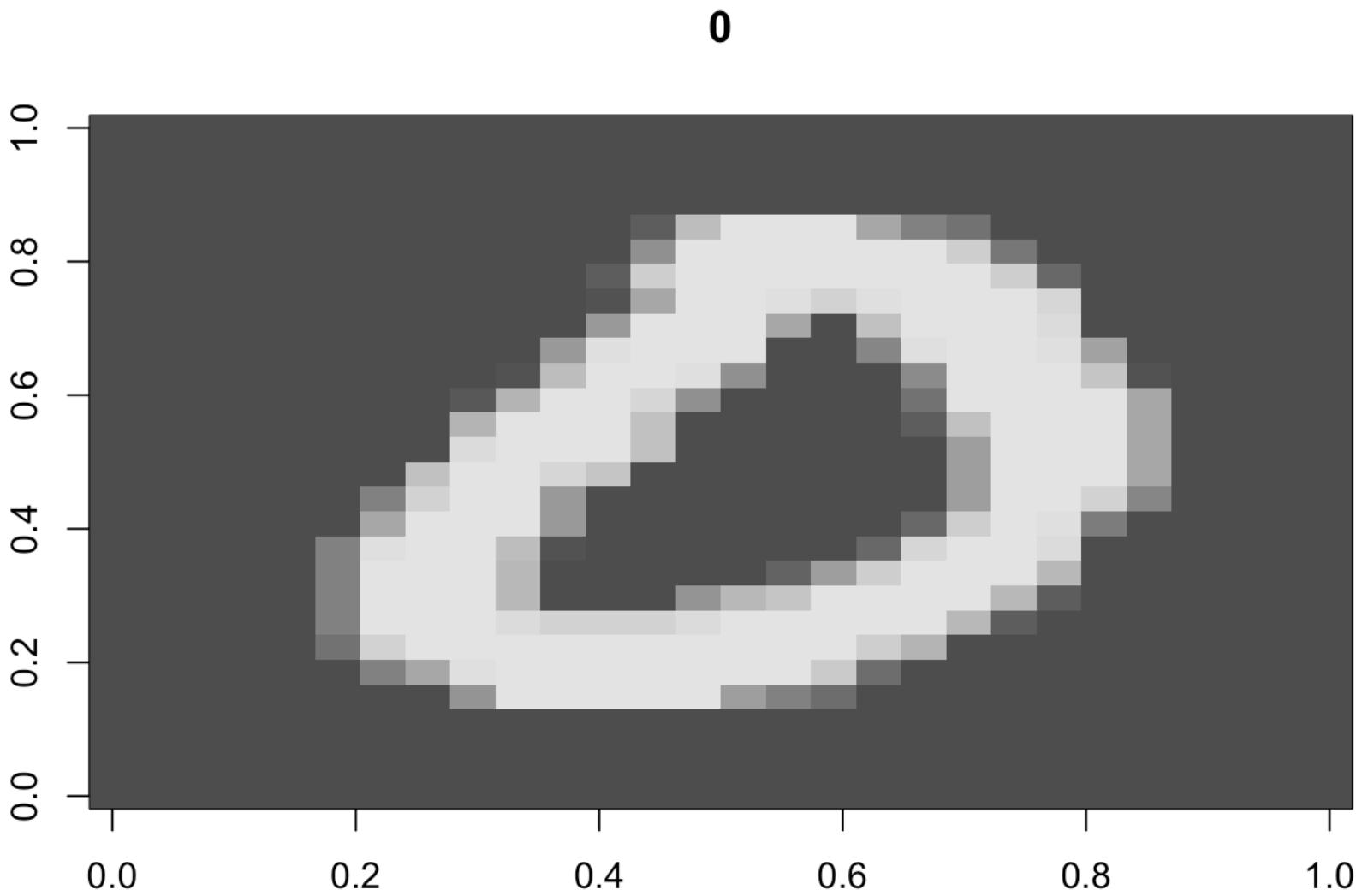
9



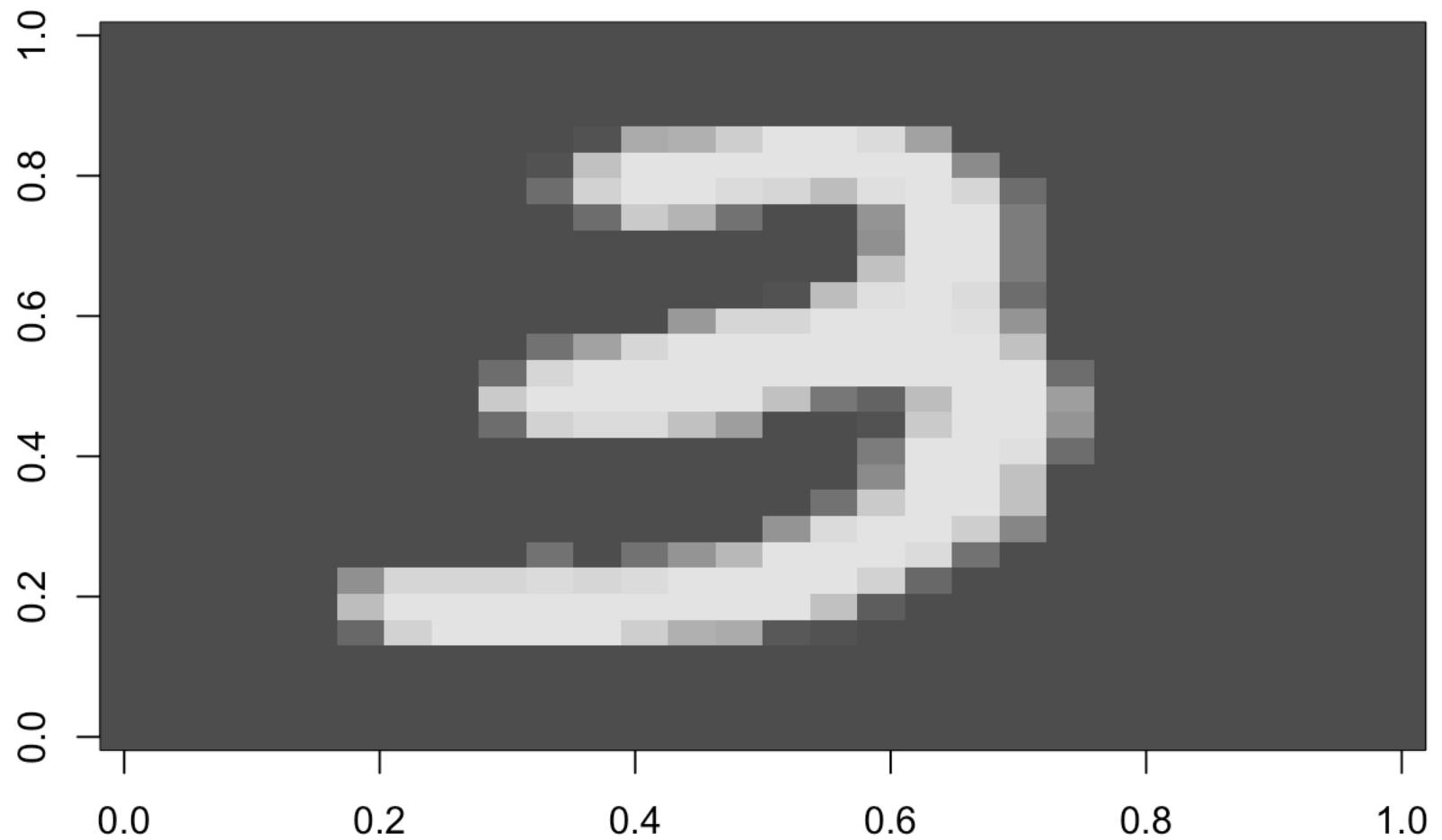
3

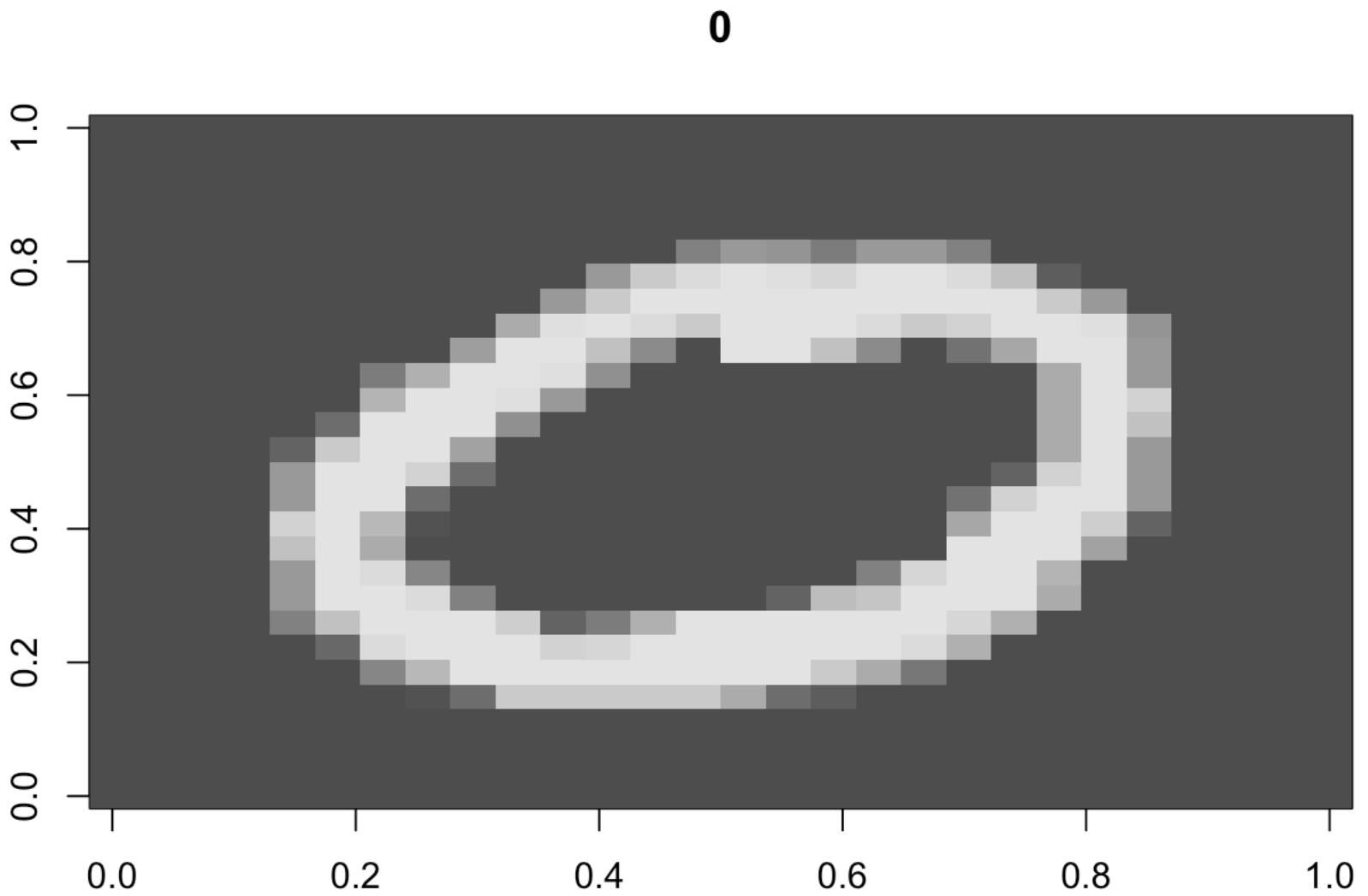




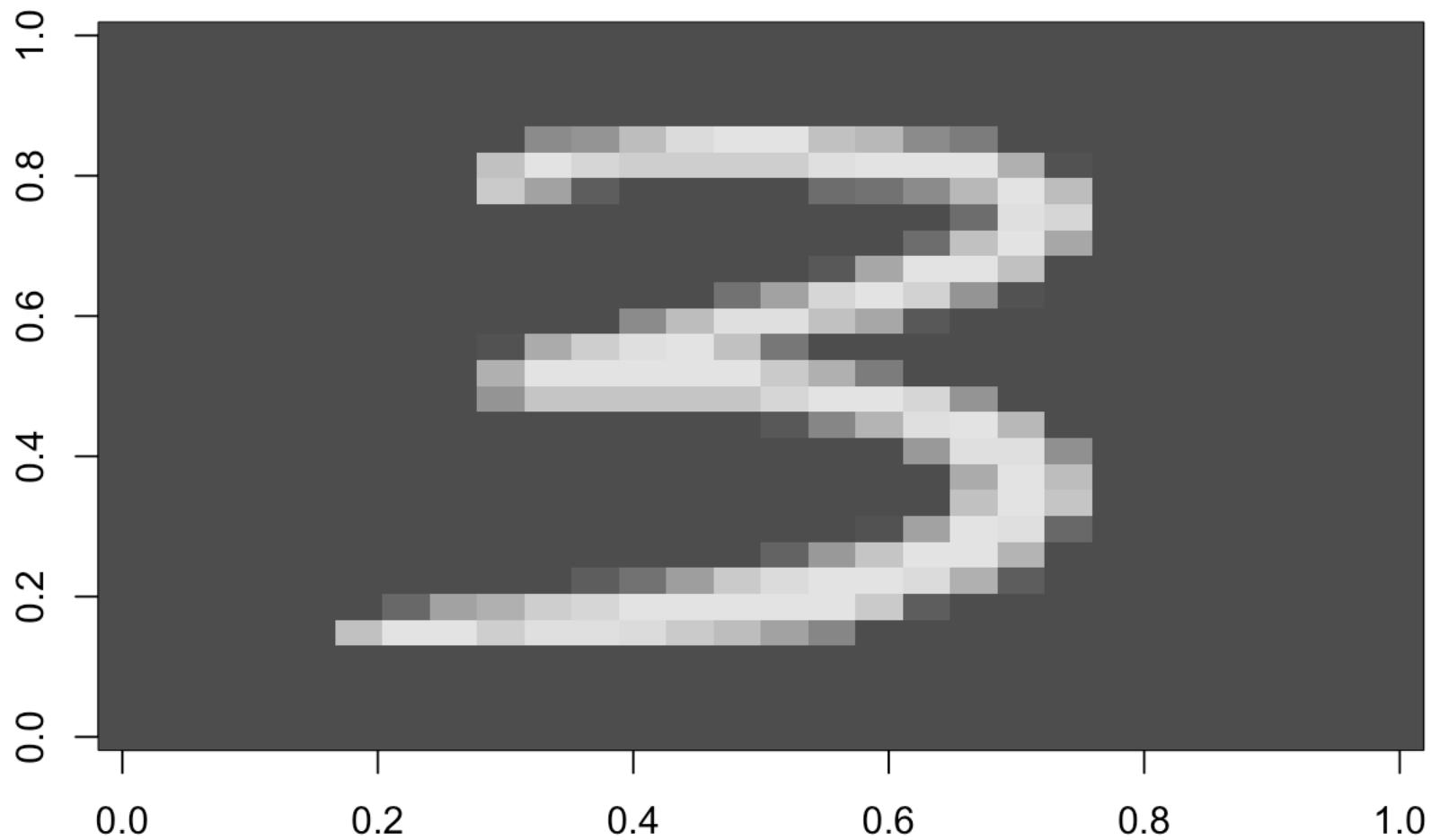


3

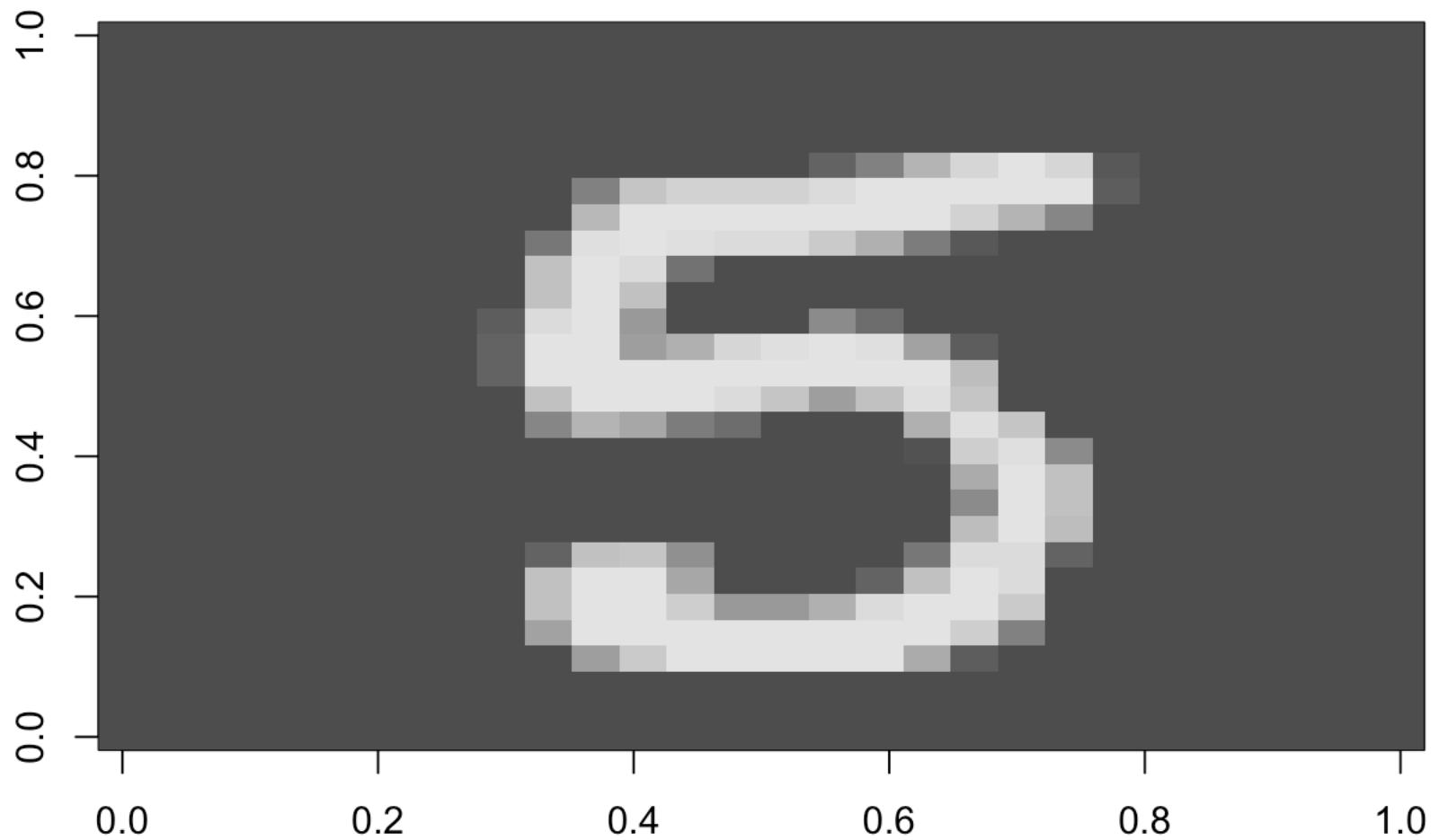


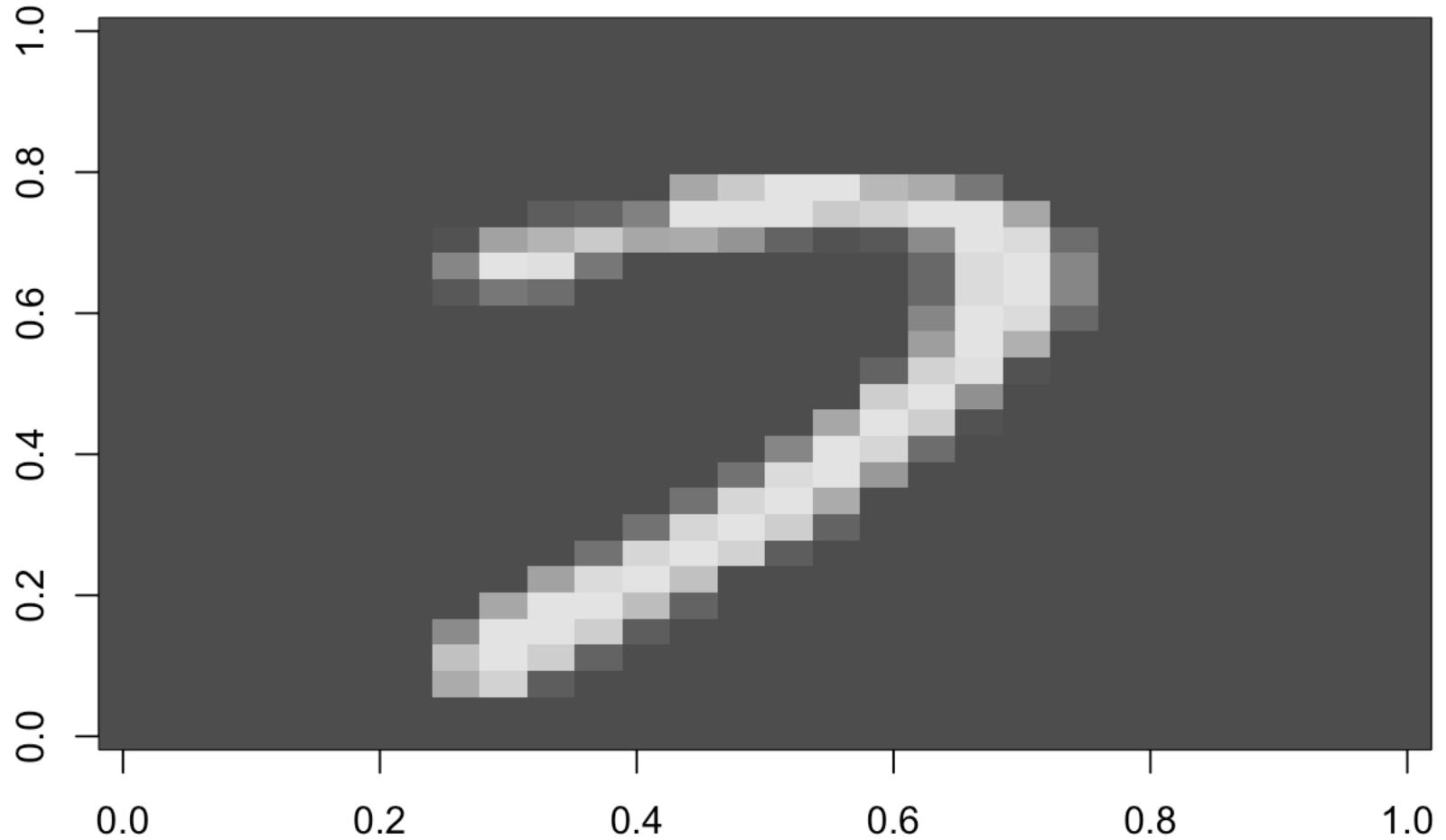


3

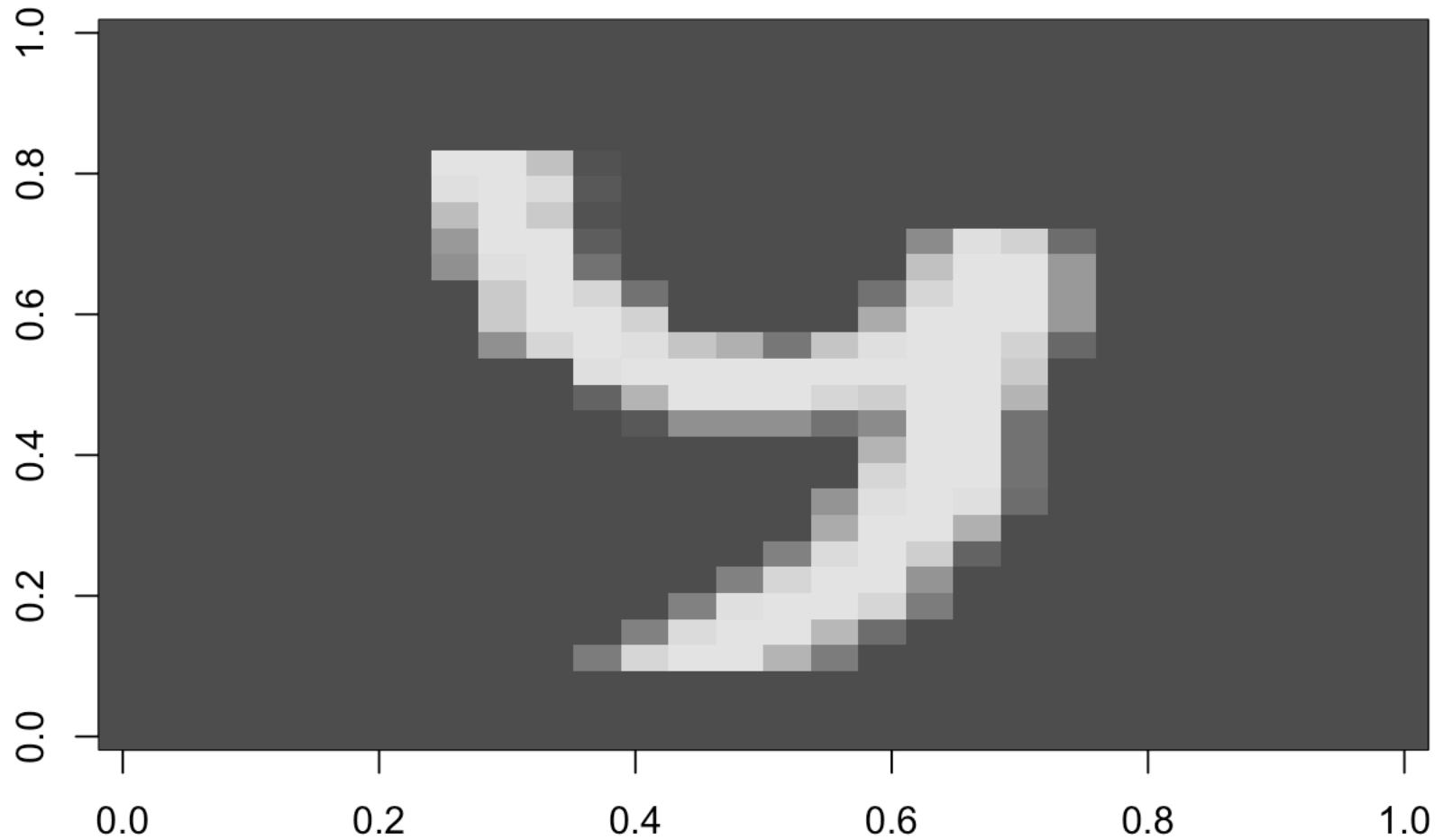


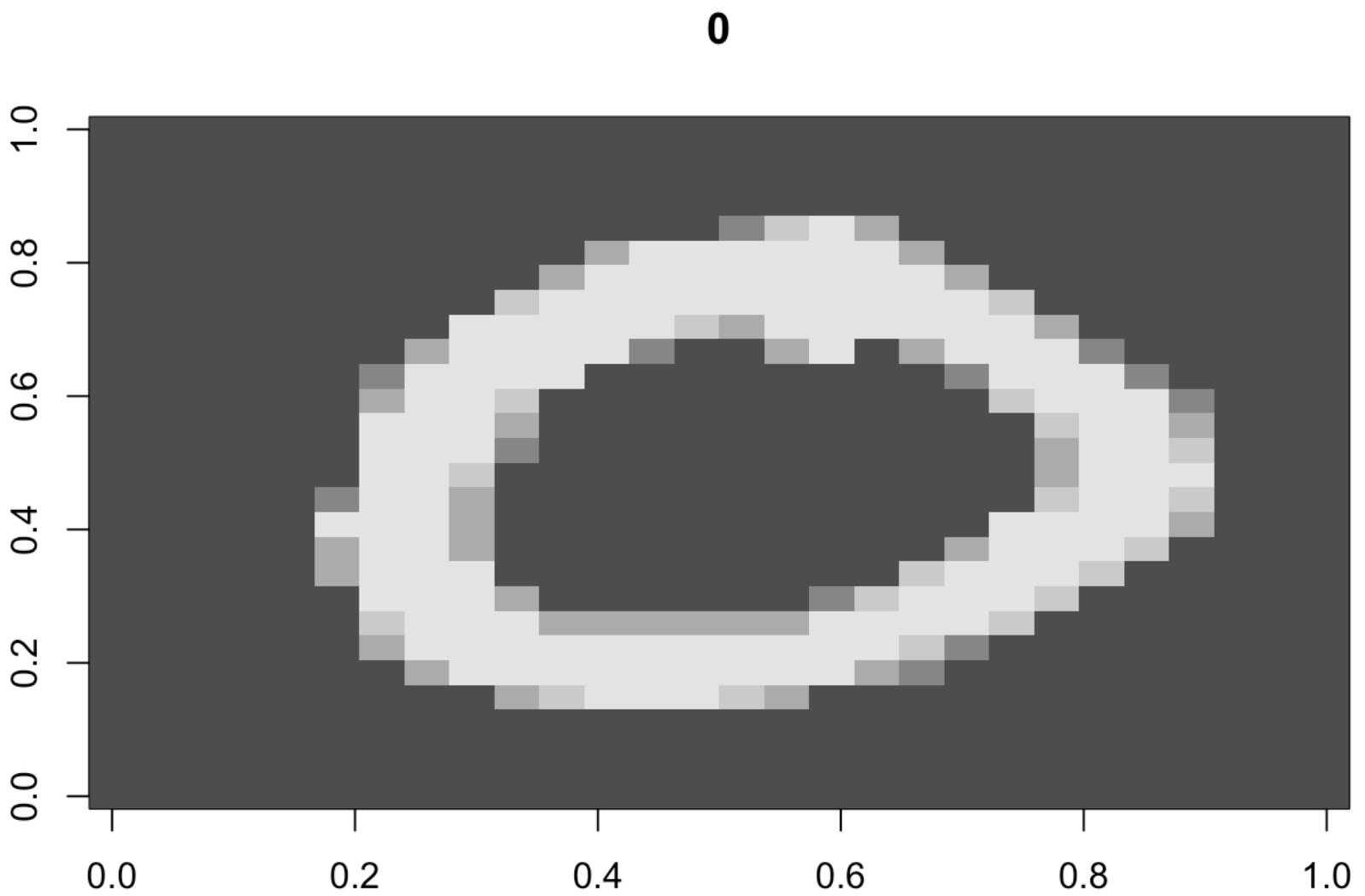
5



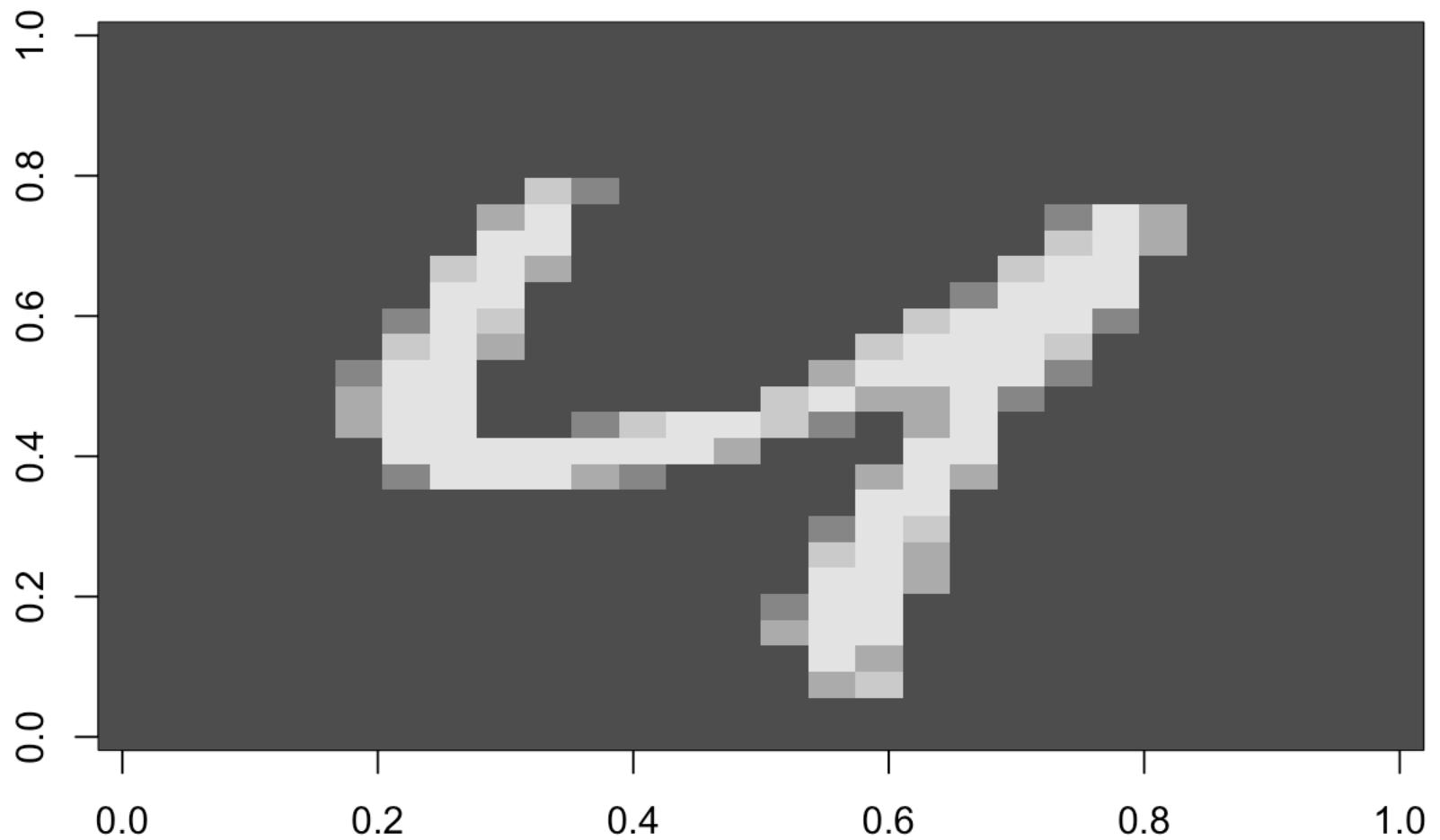


4

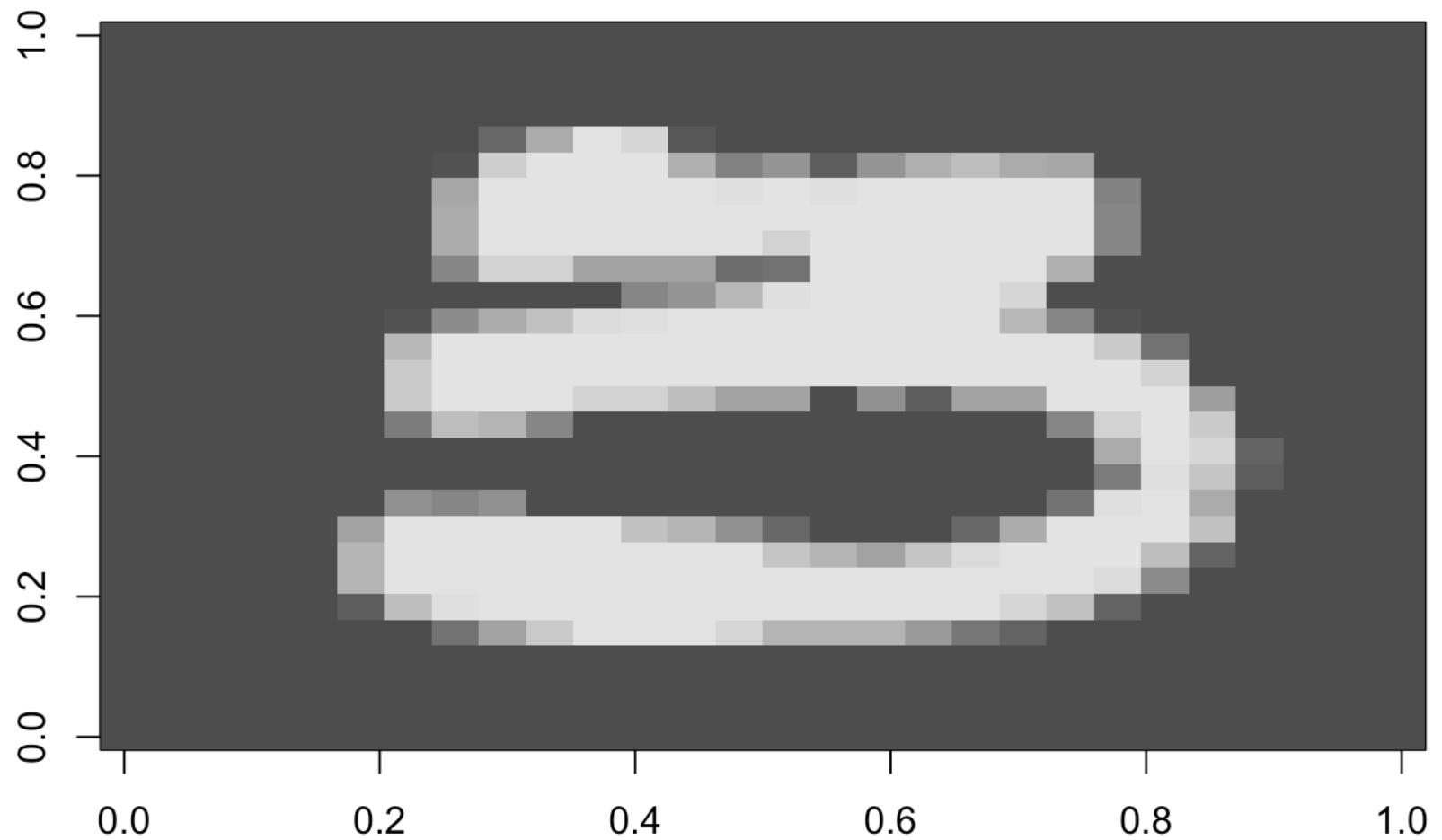




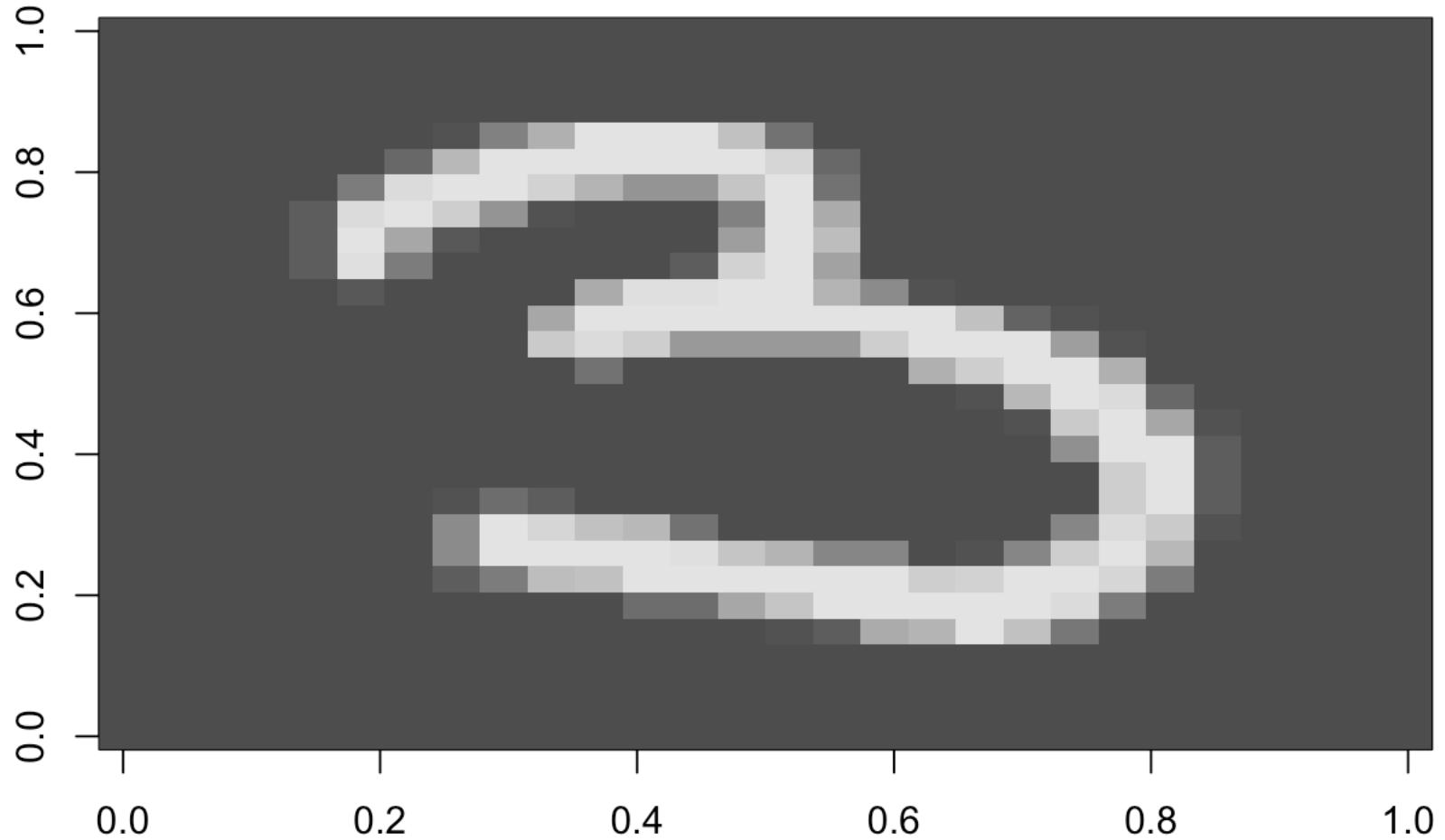
4



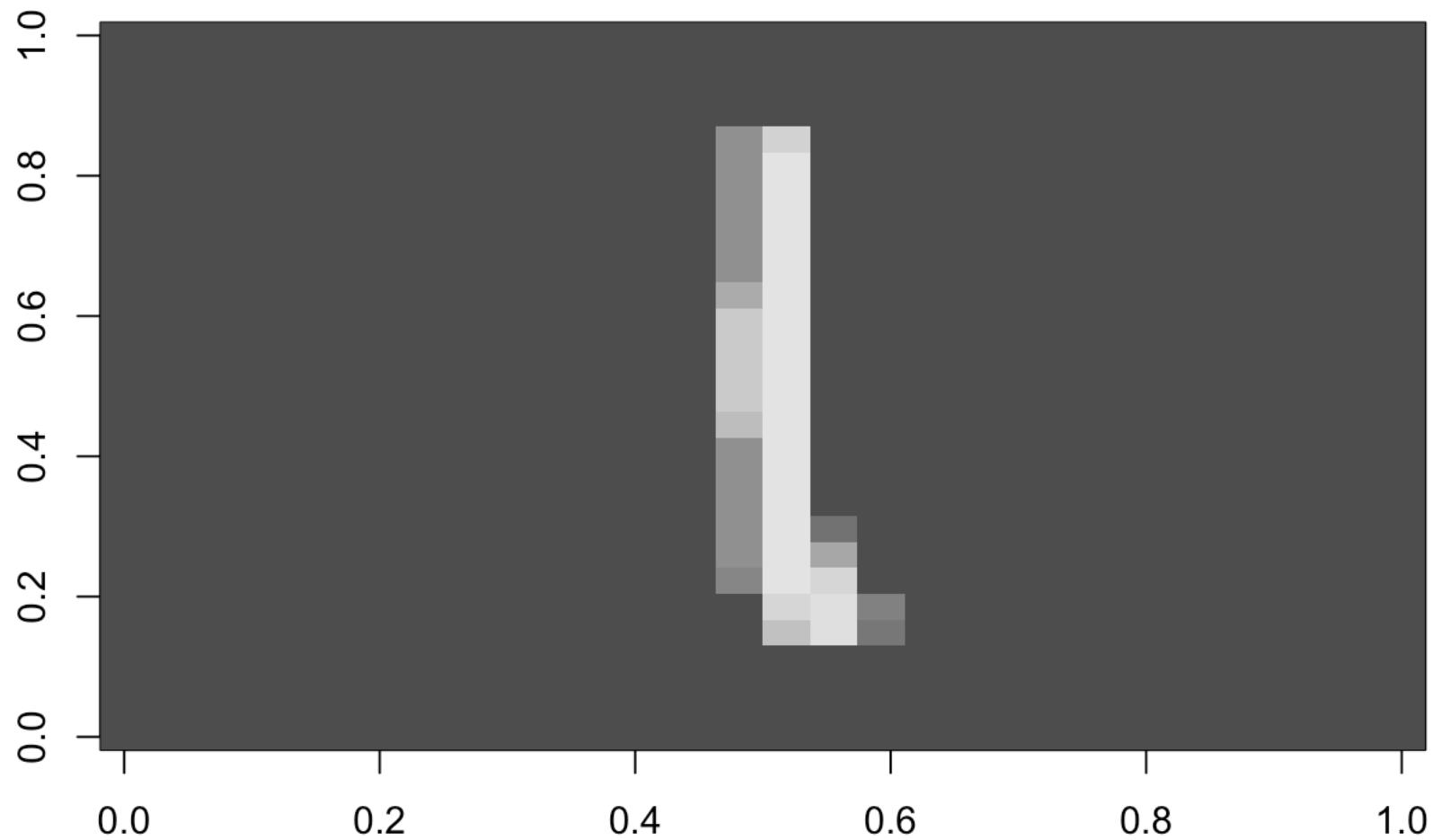
3



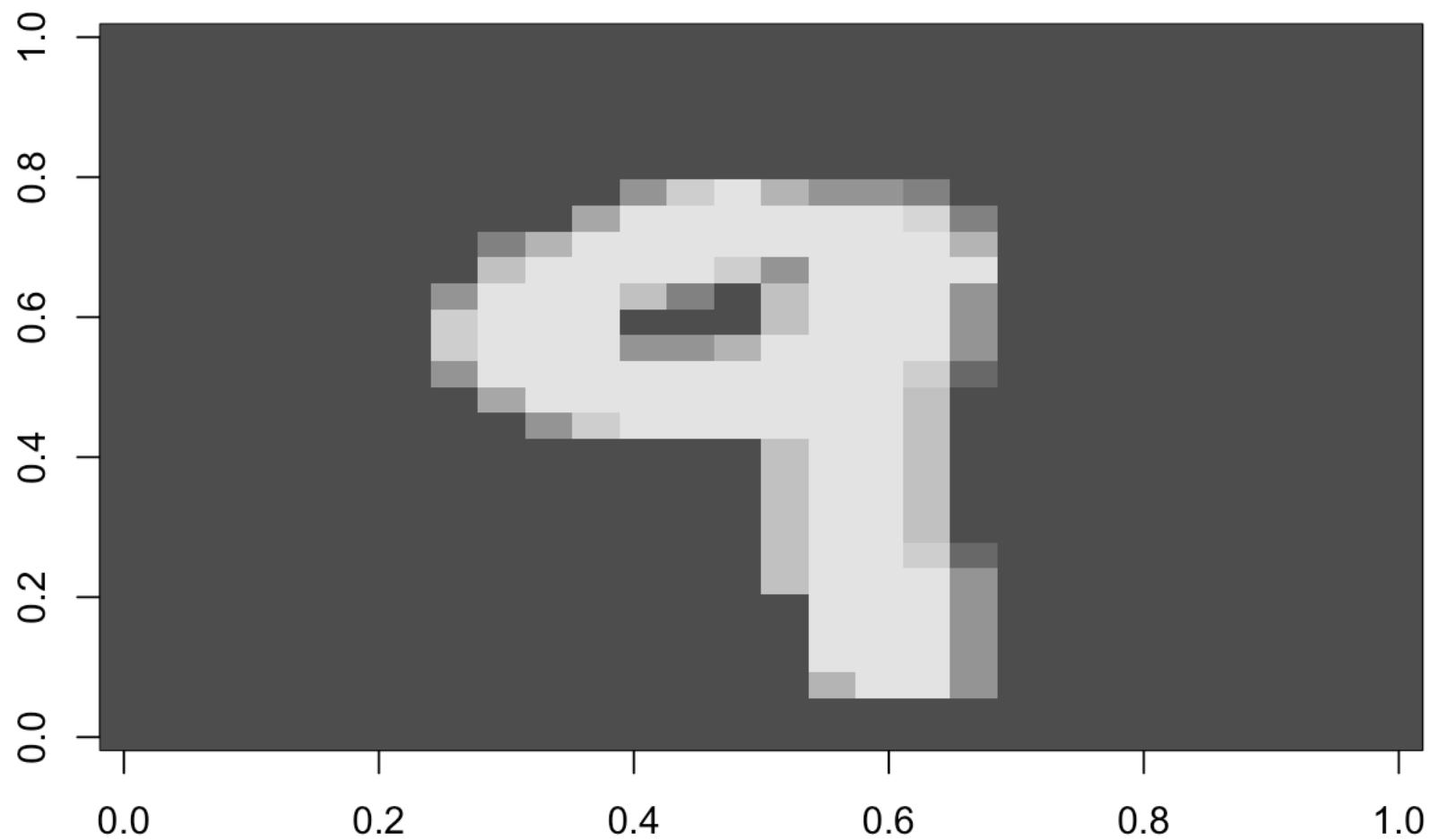
3

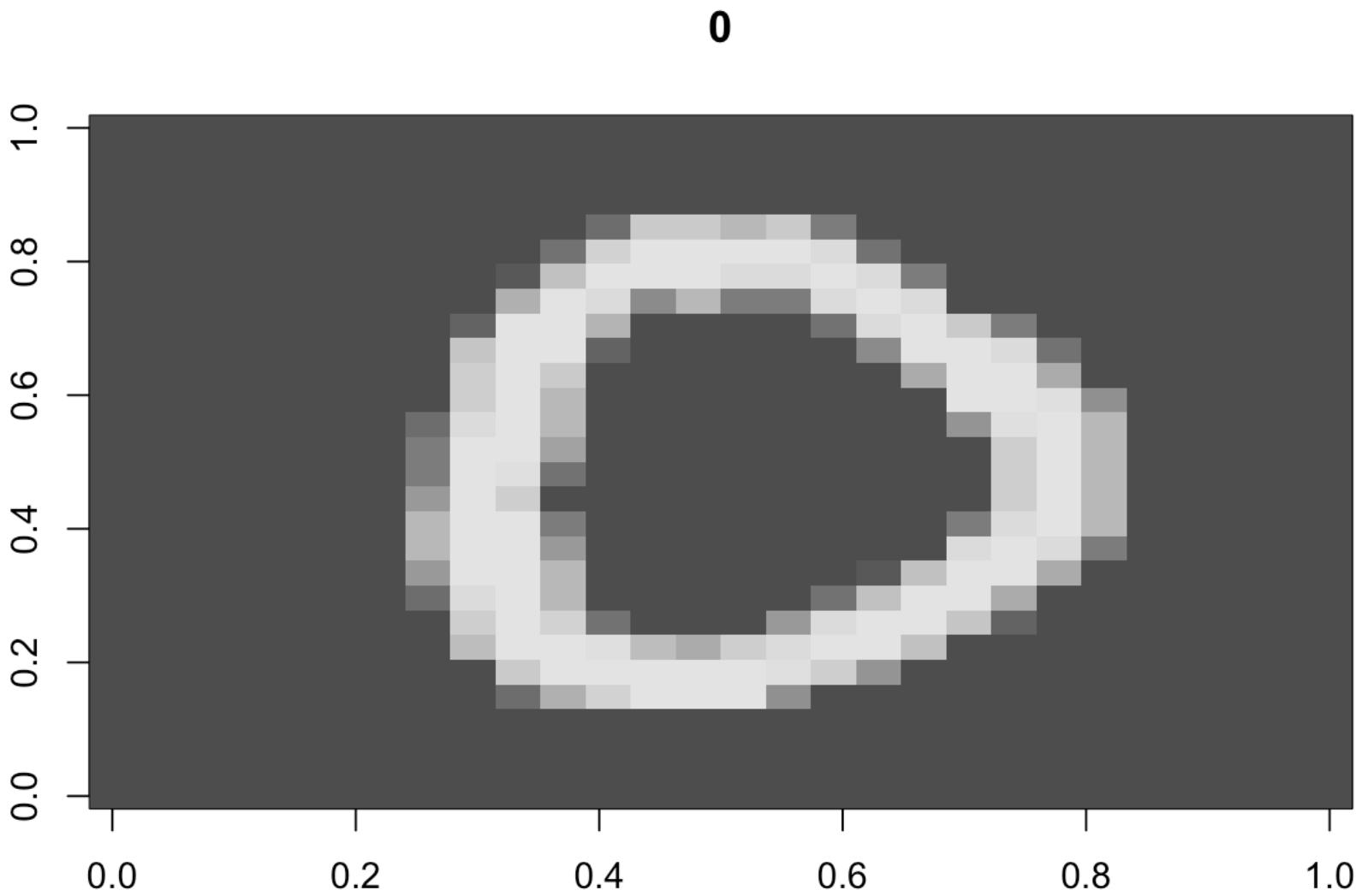


1

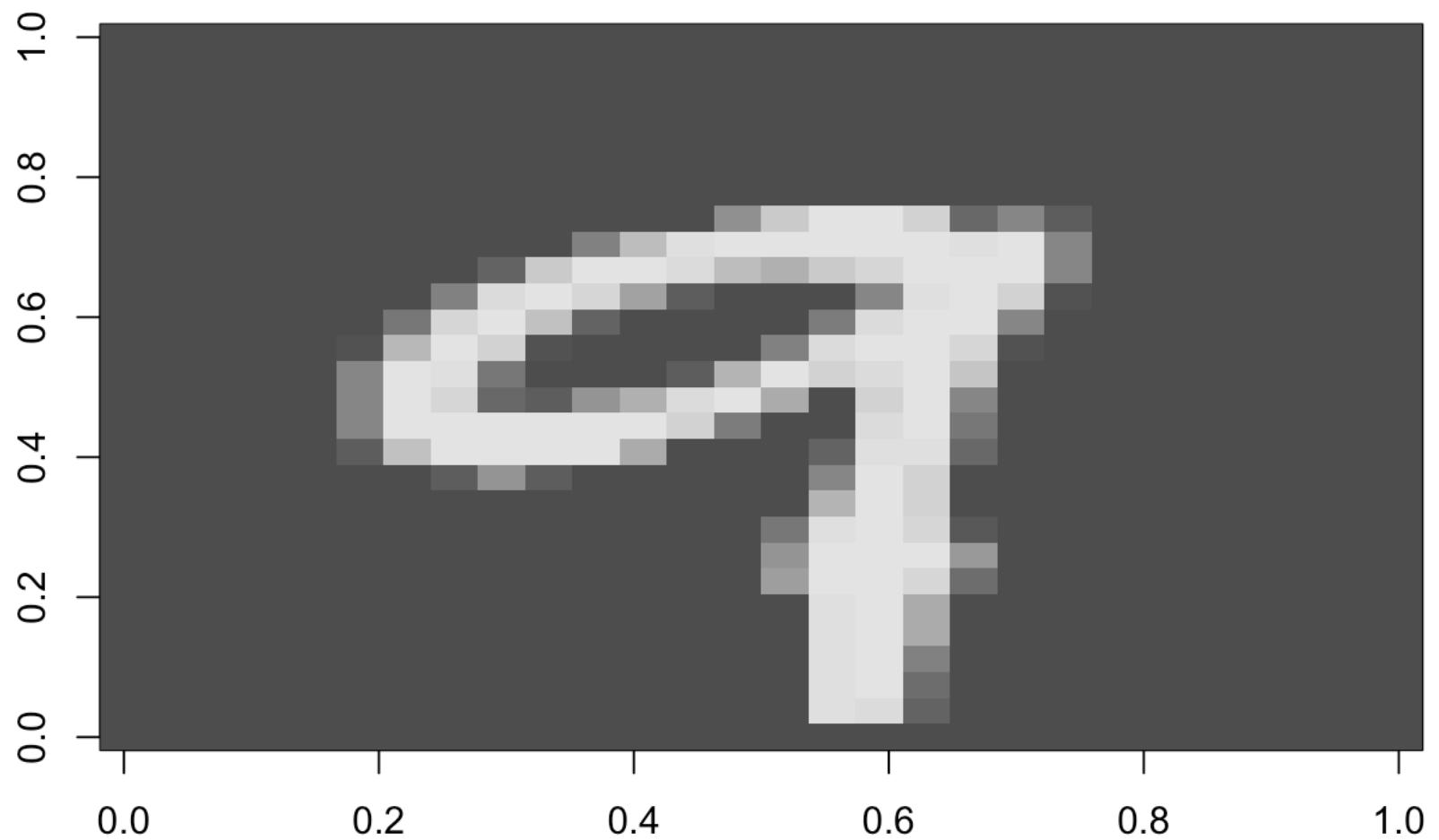


9

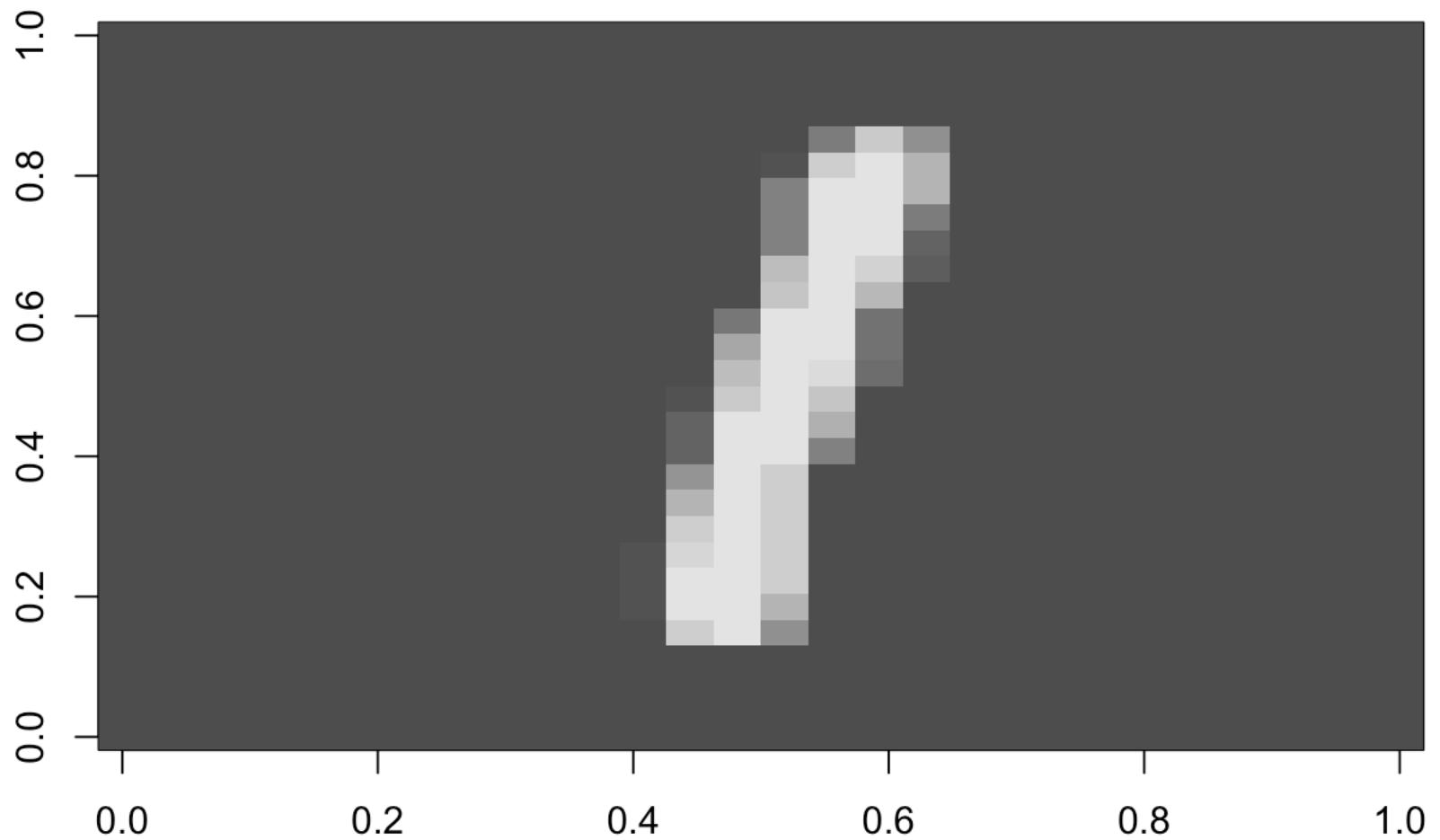




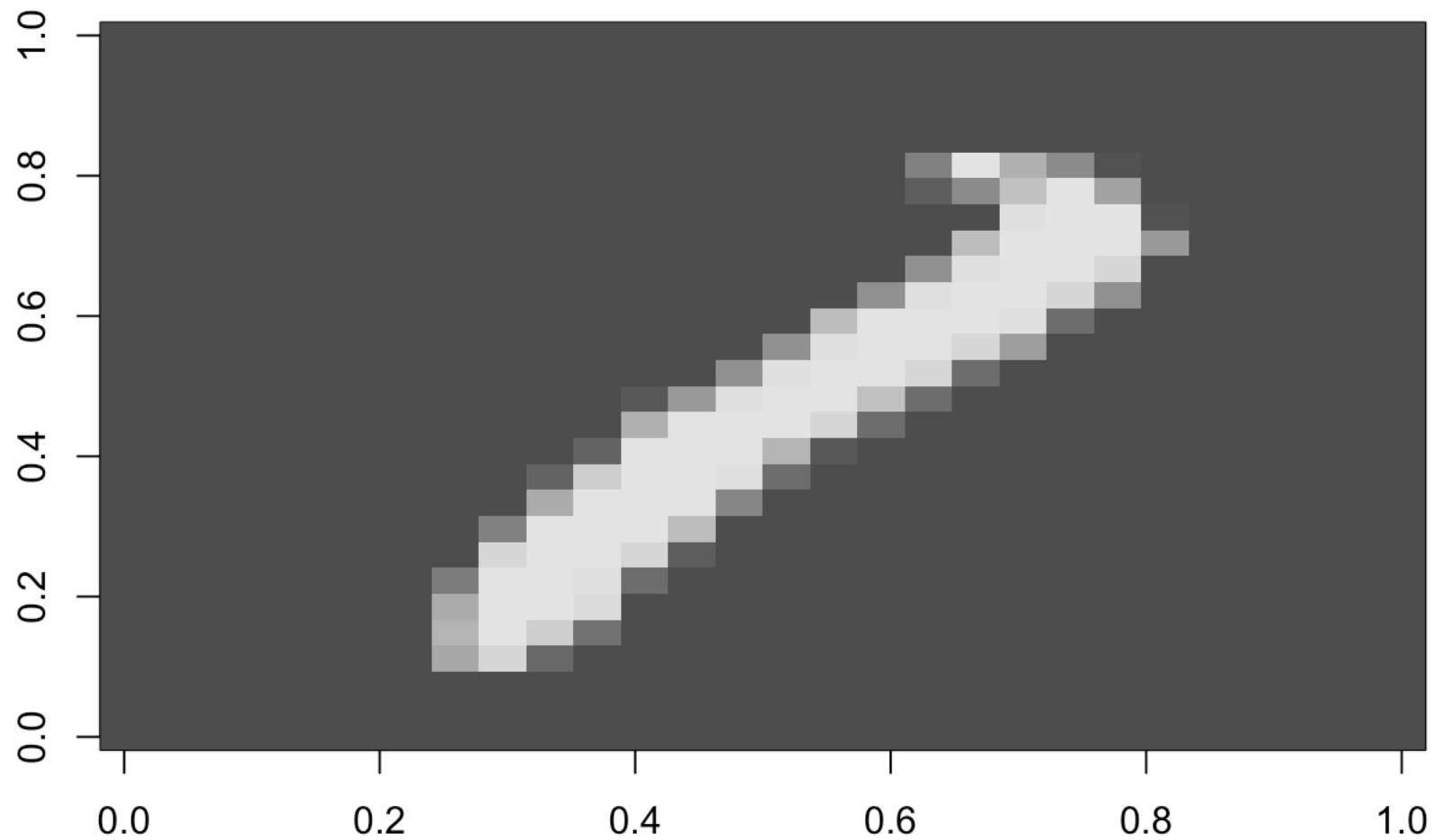
9



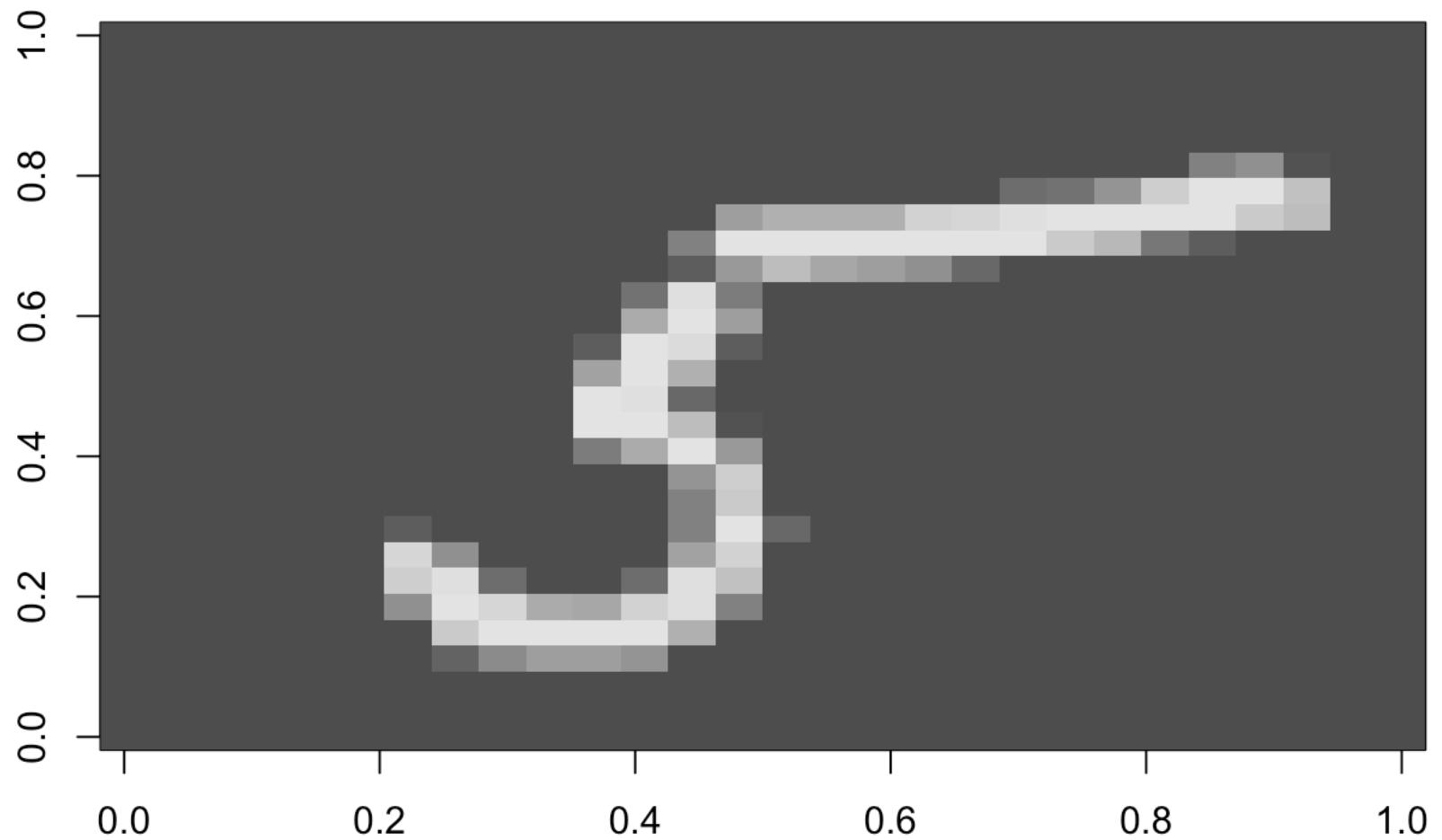
1

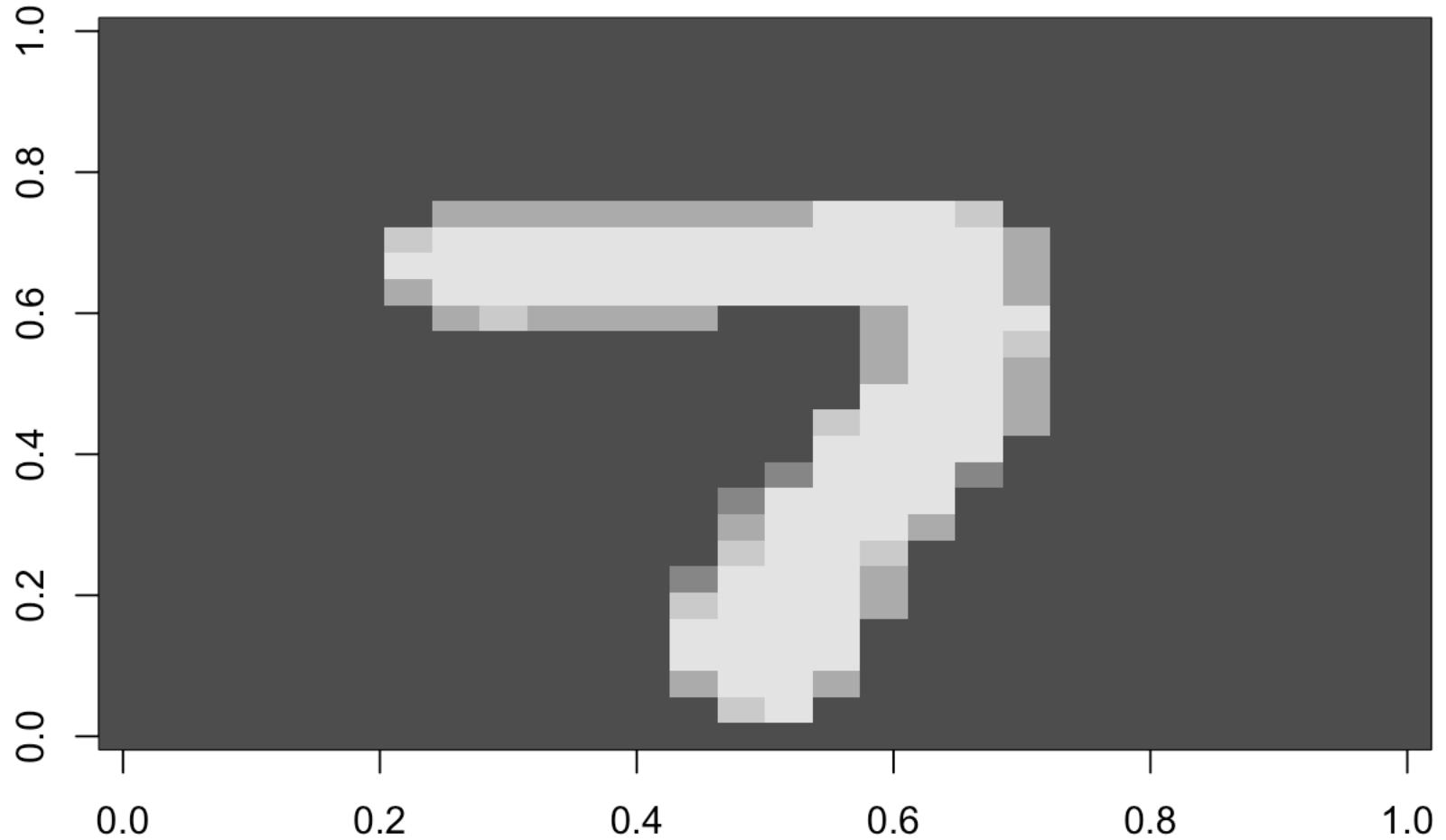


1

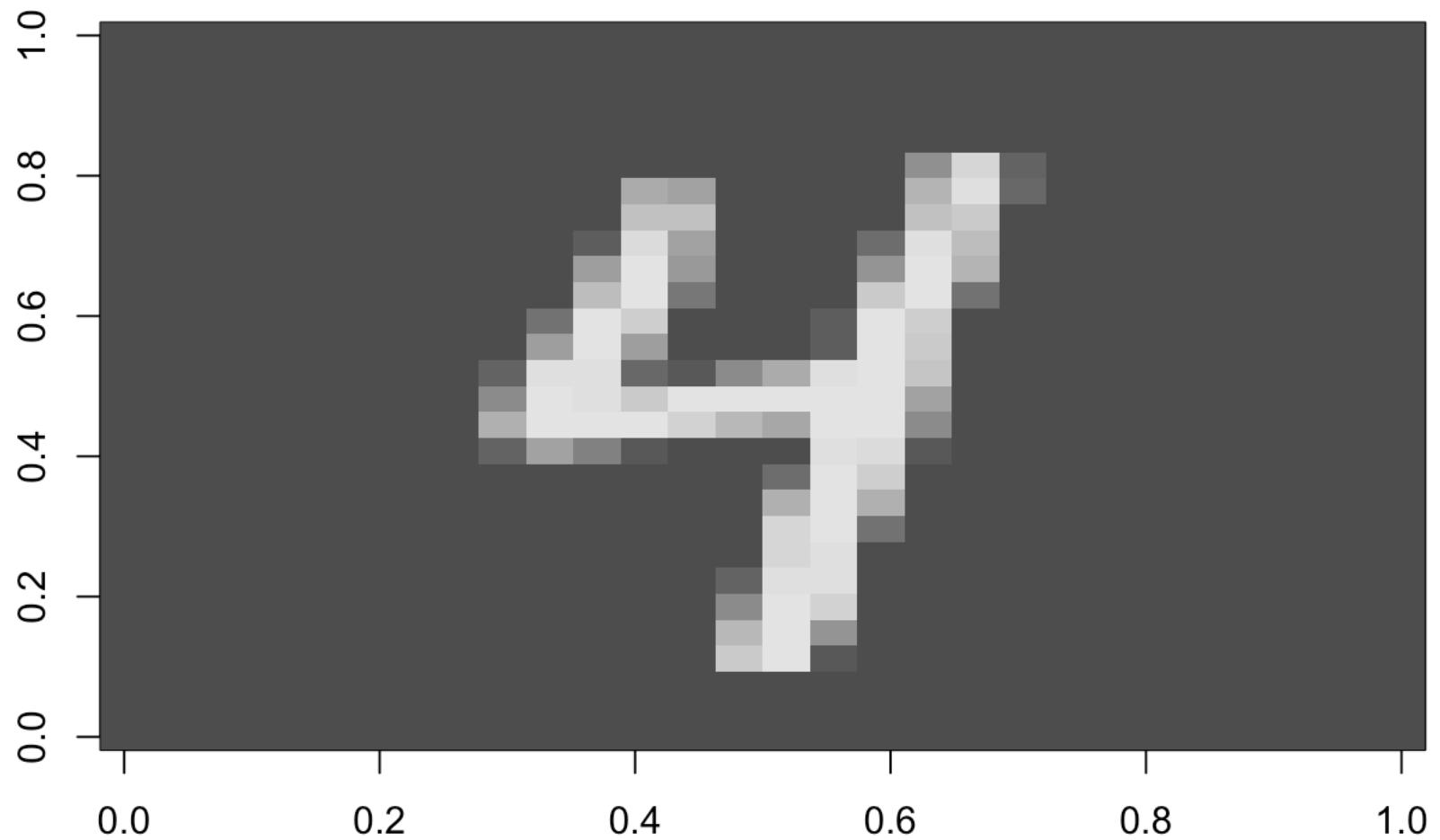


5

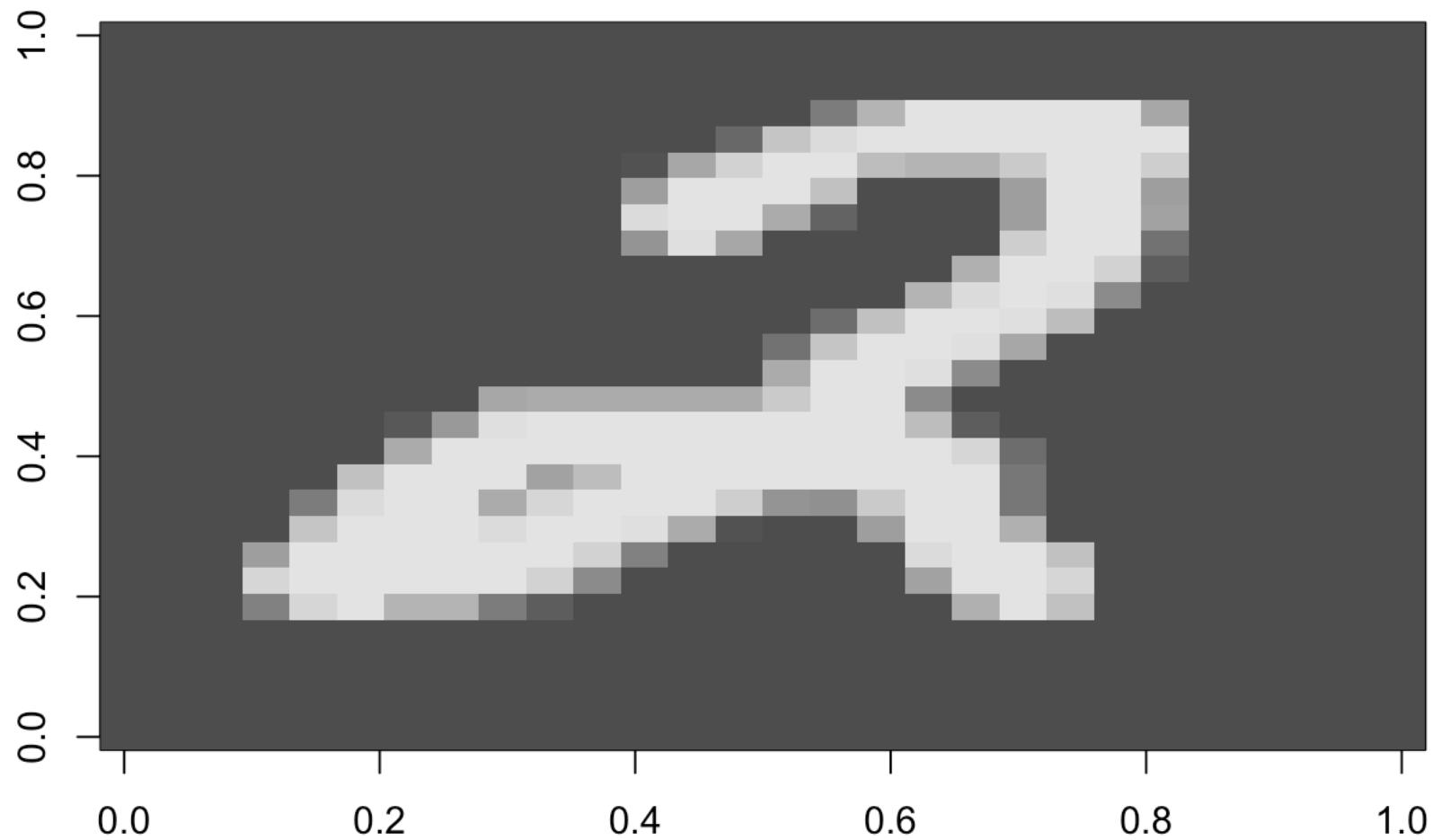


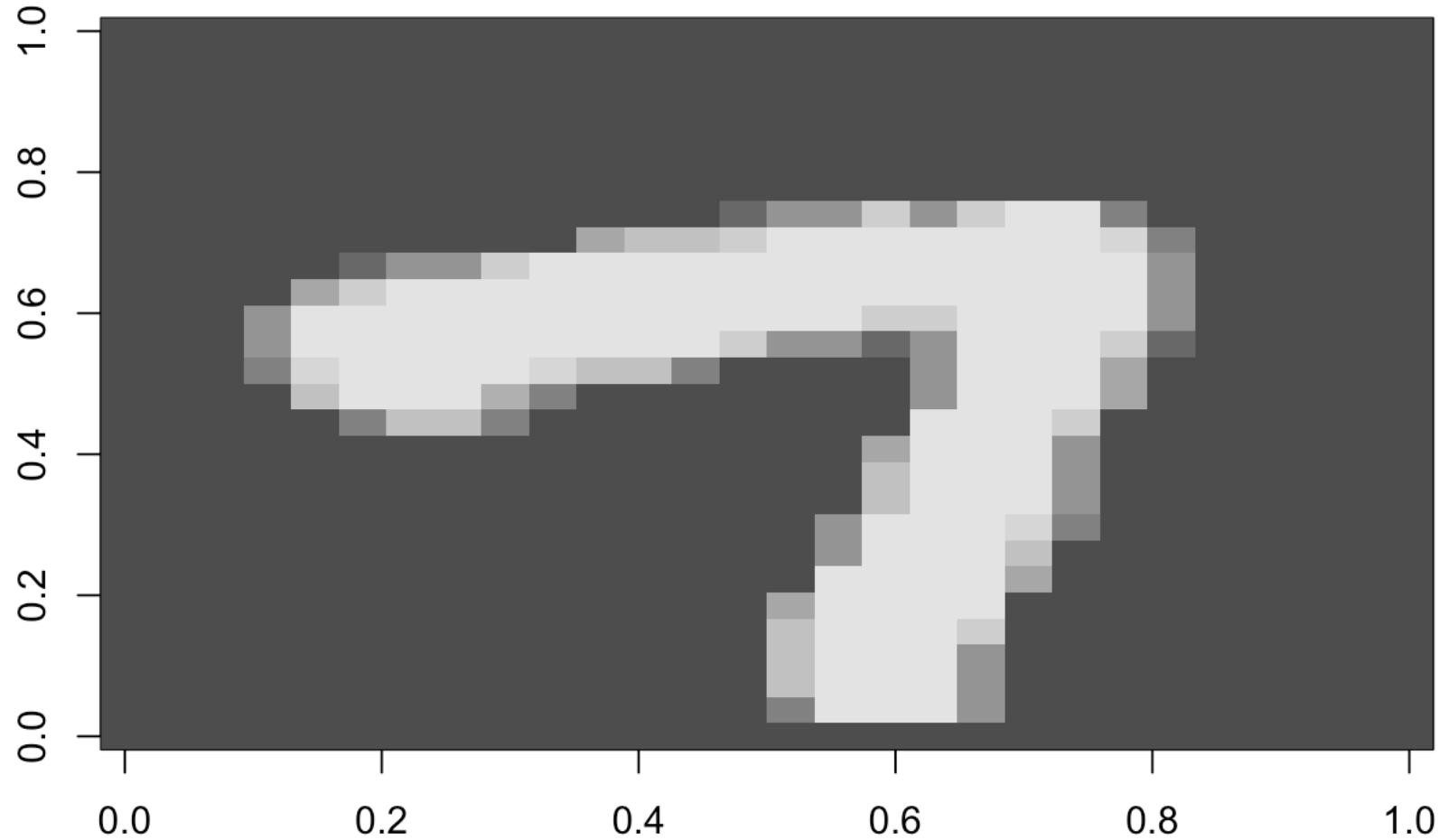


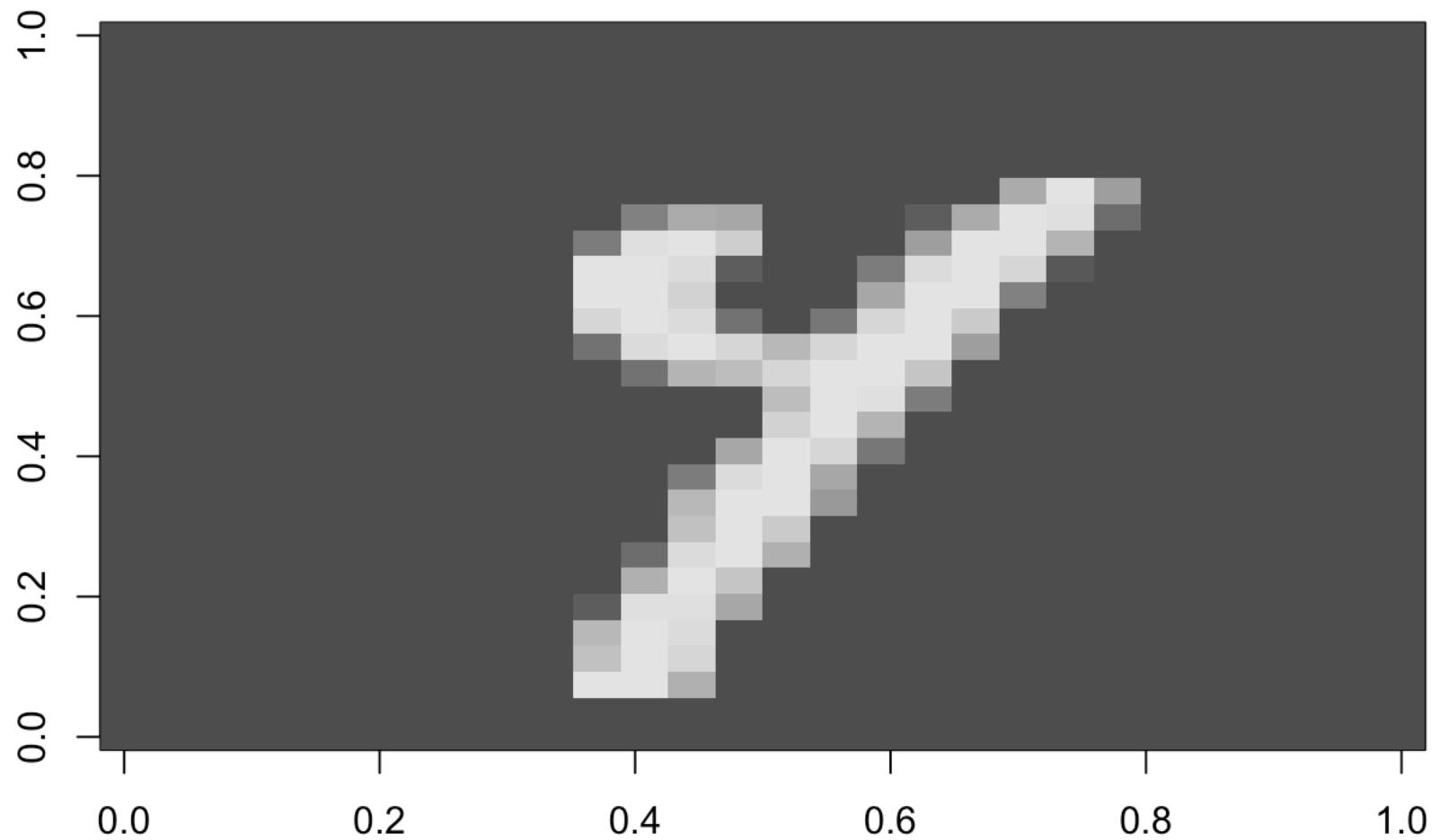
4

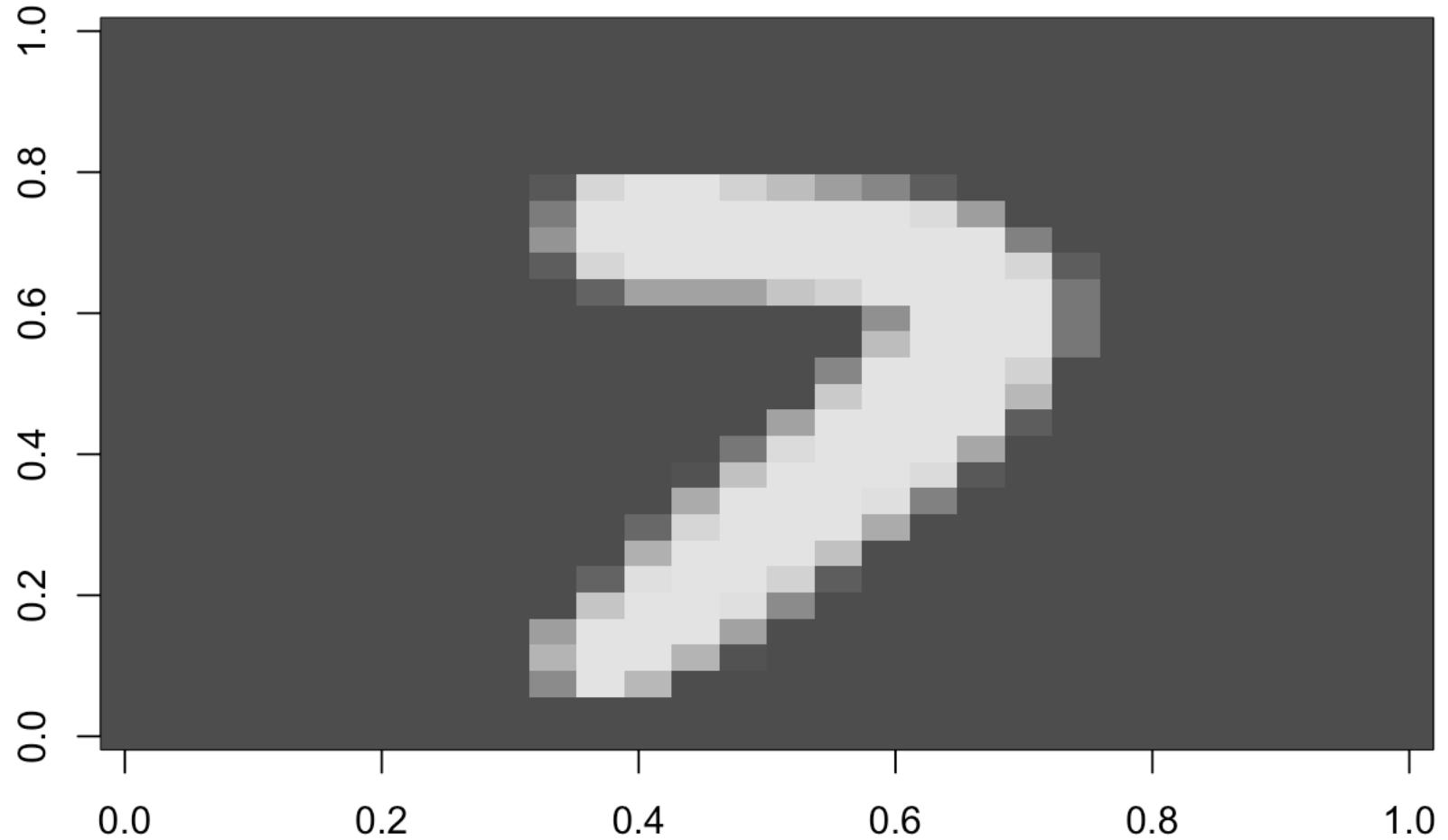


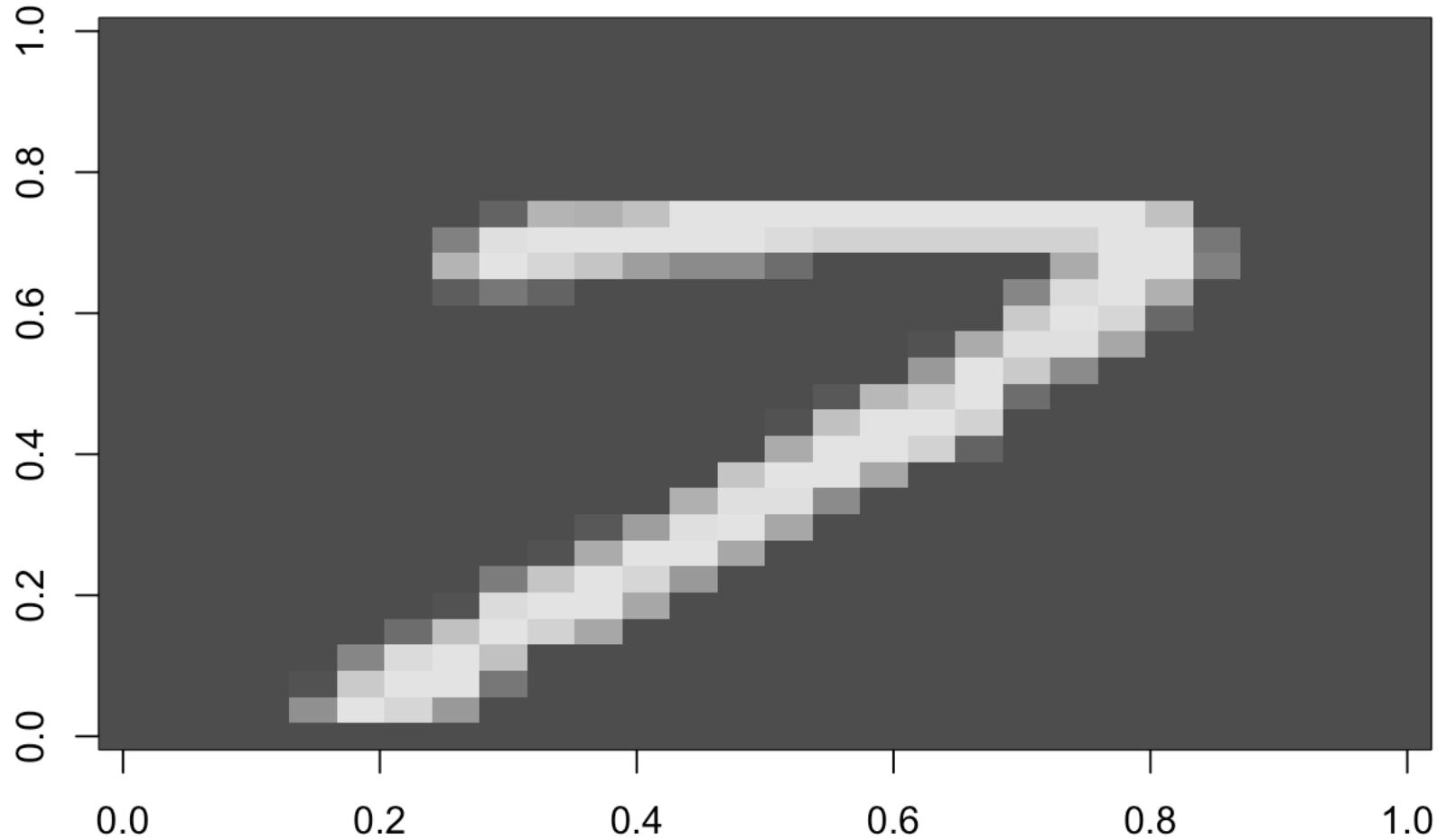
2



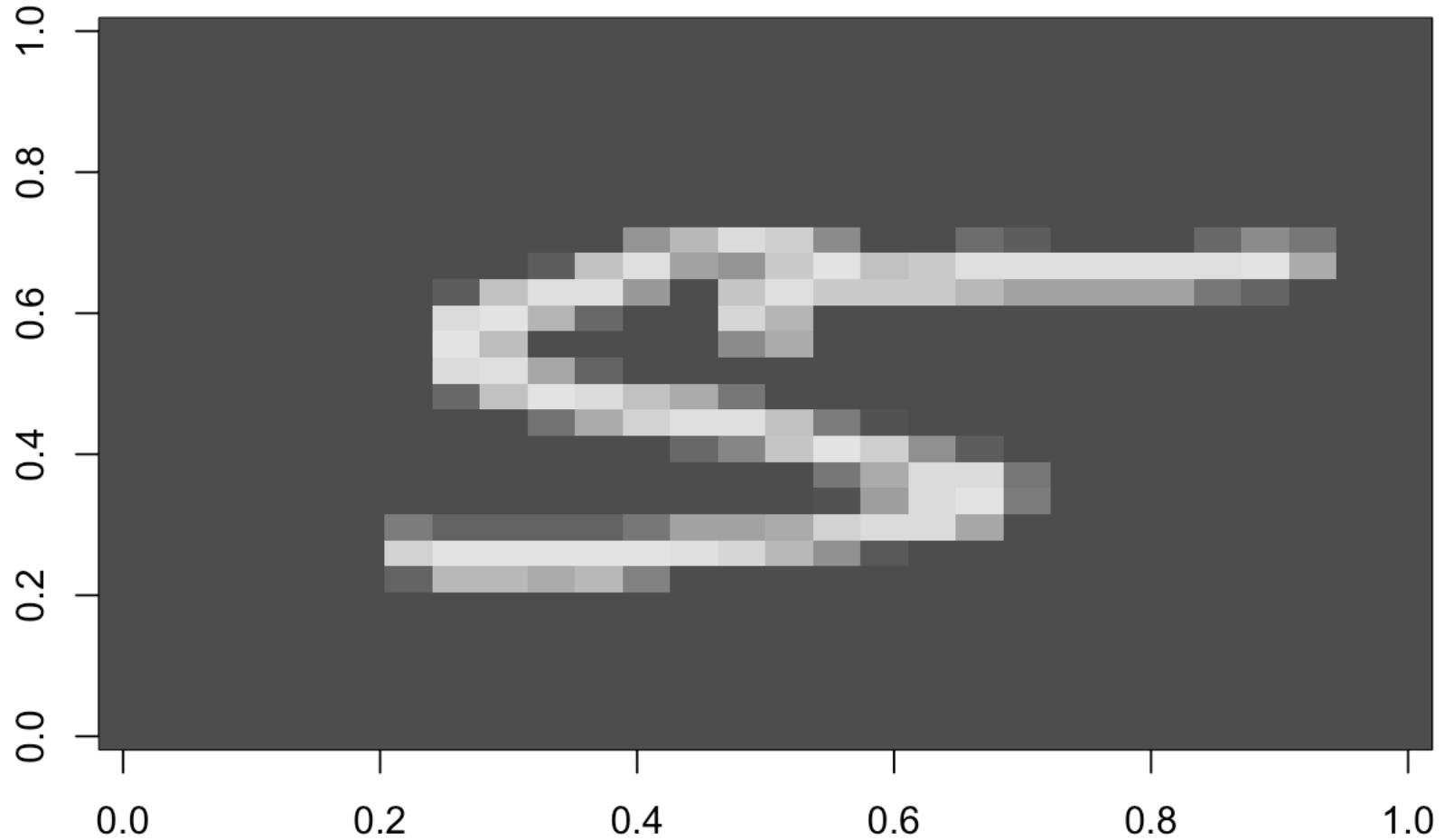




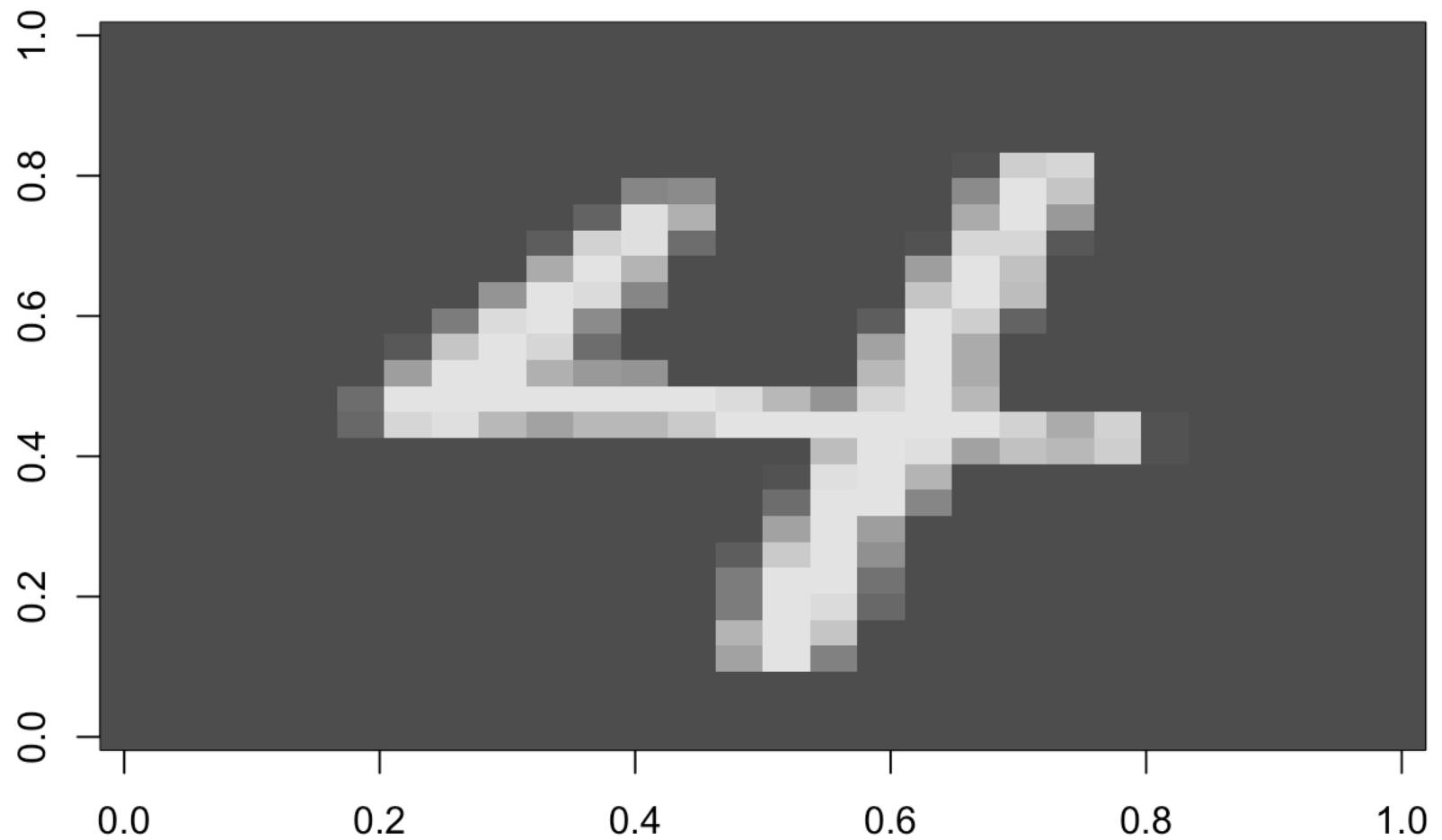




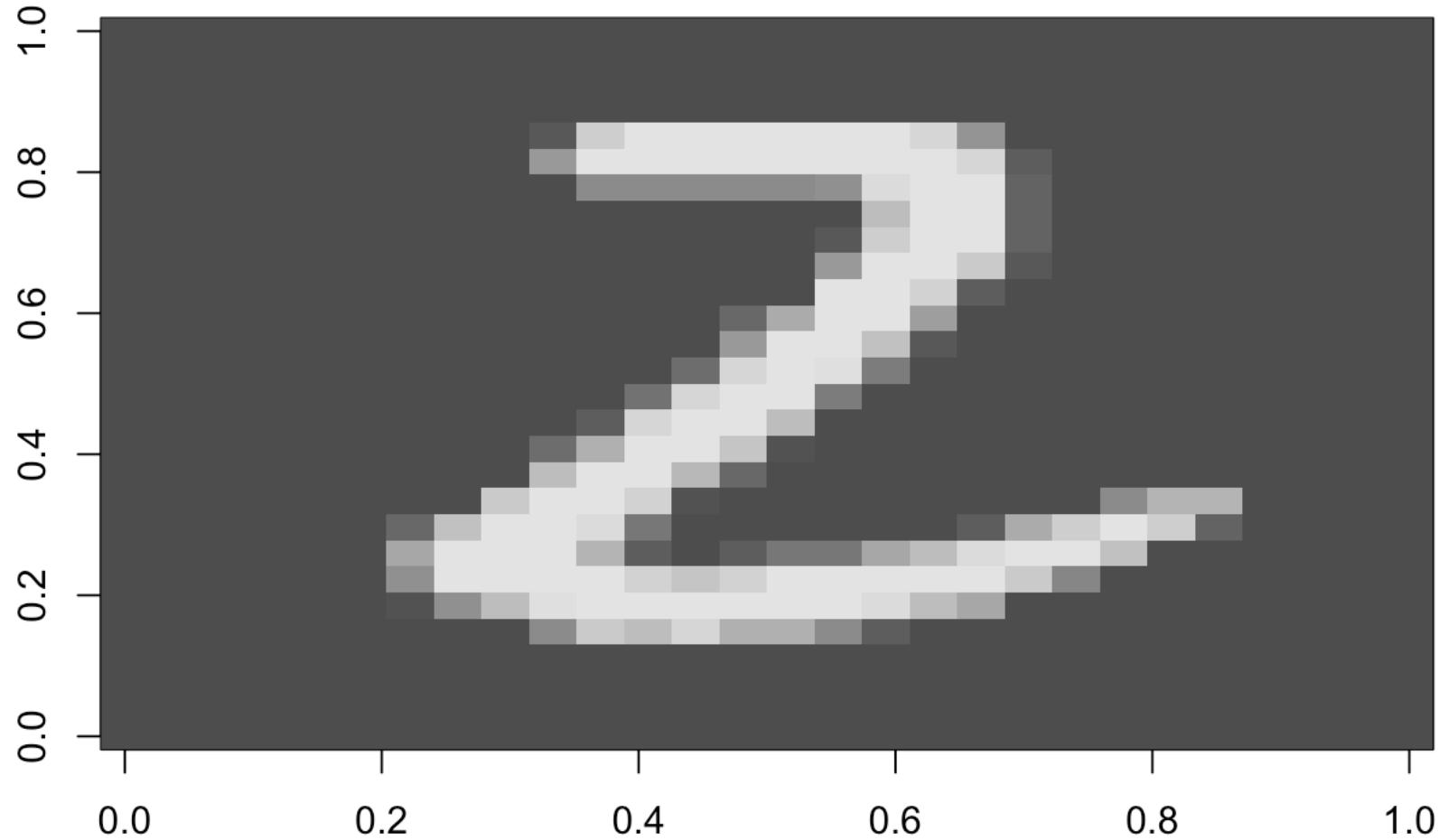
5



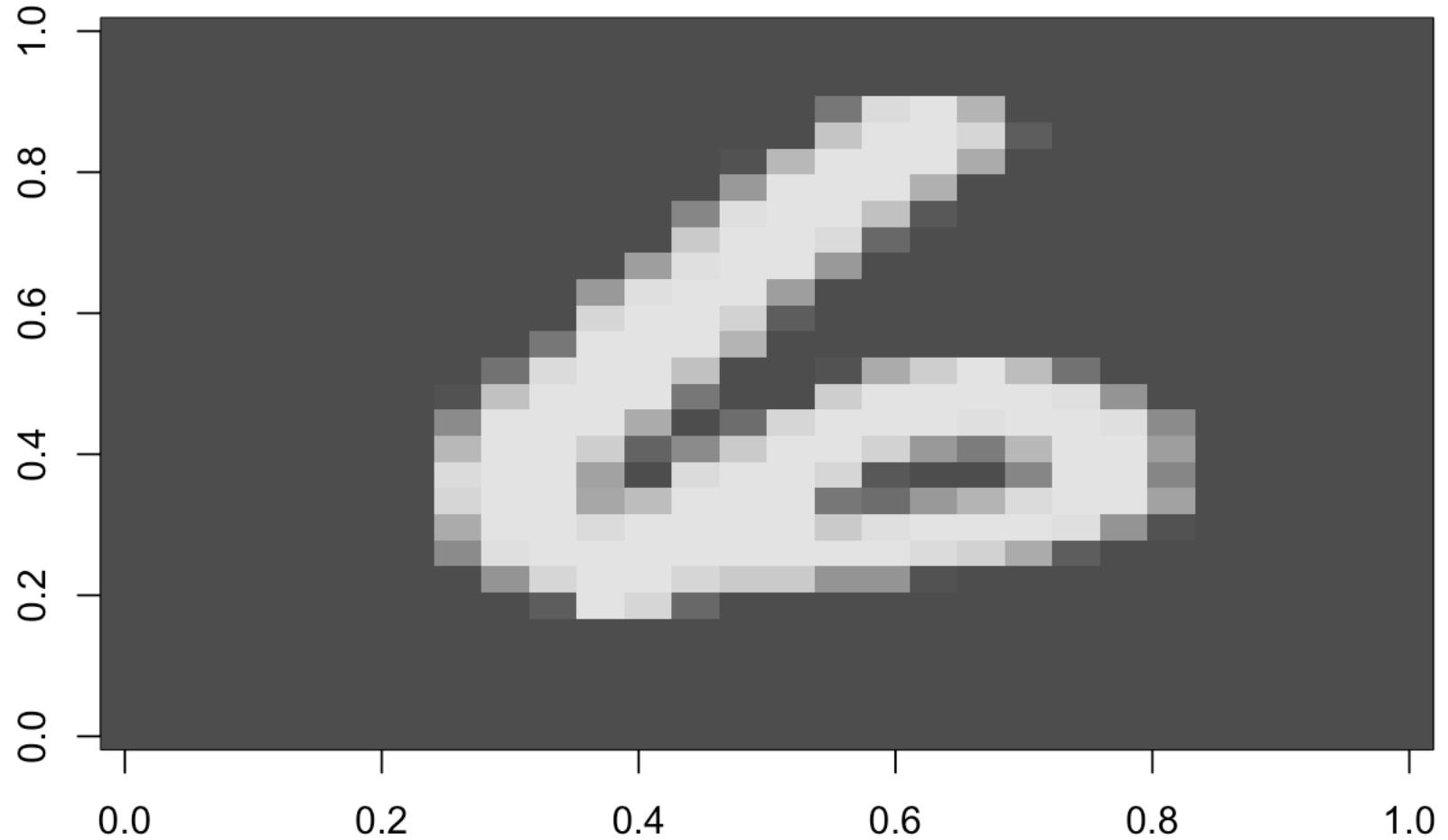
4



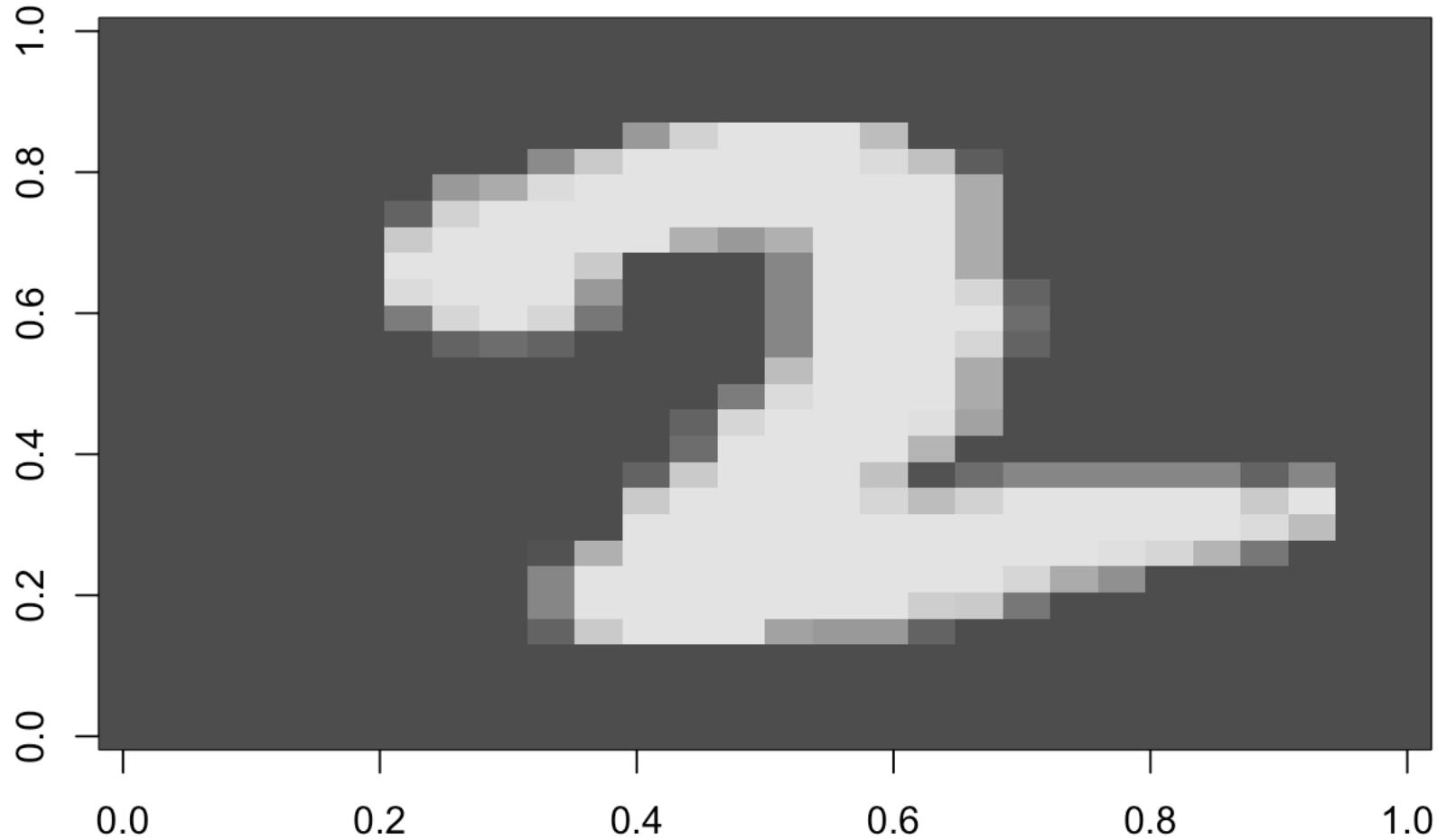
2



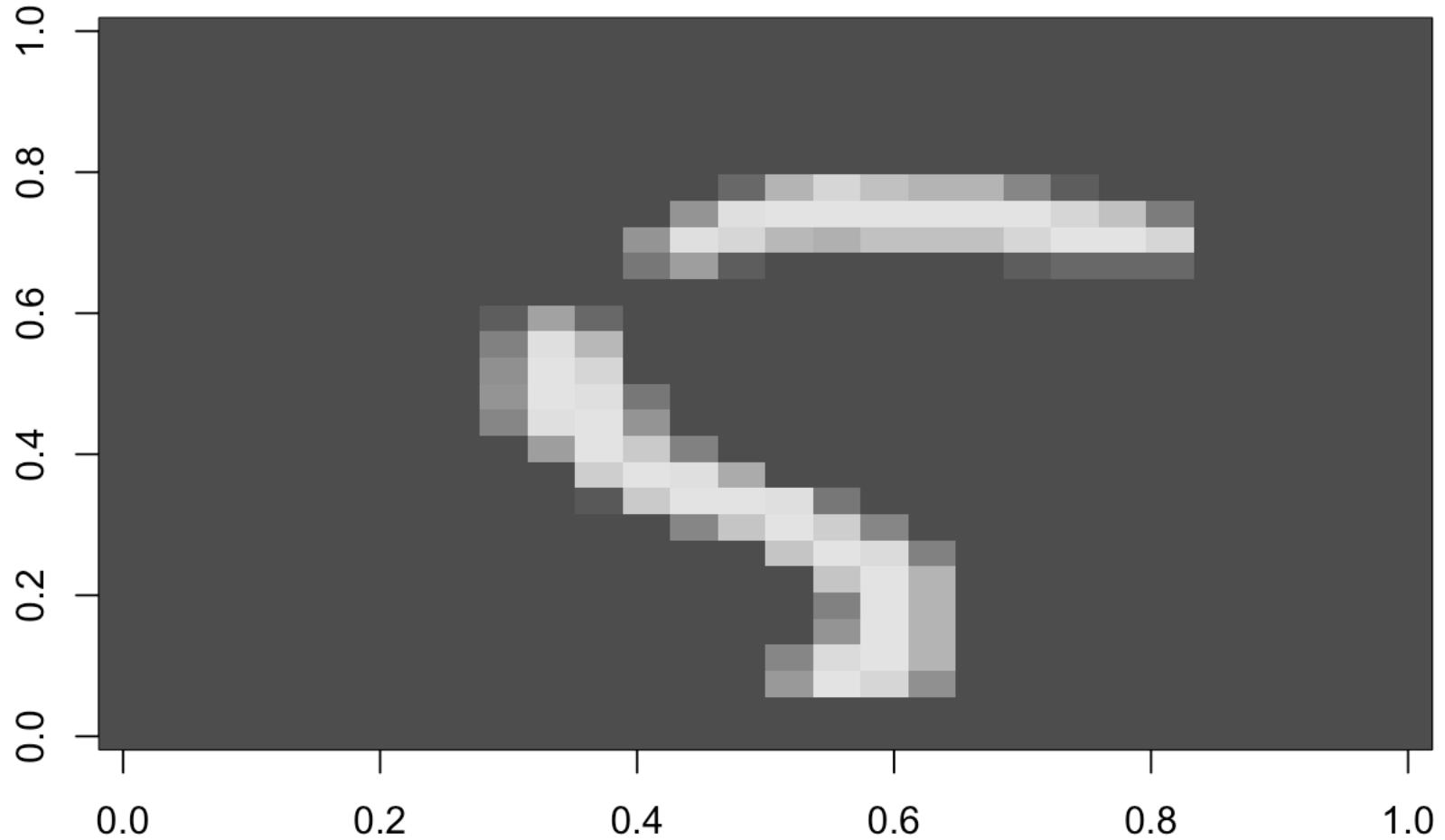
6



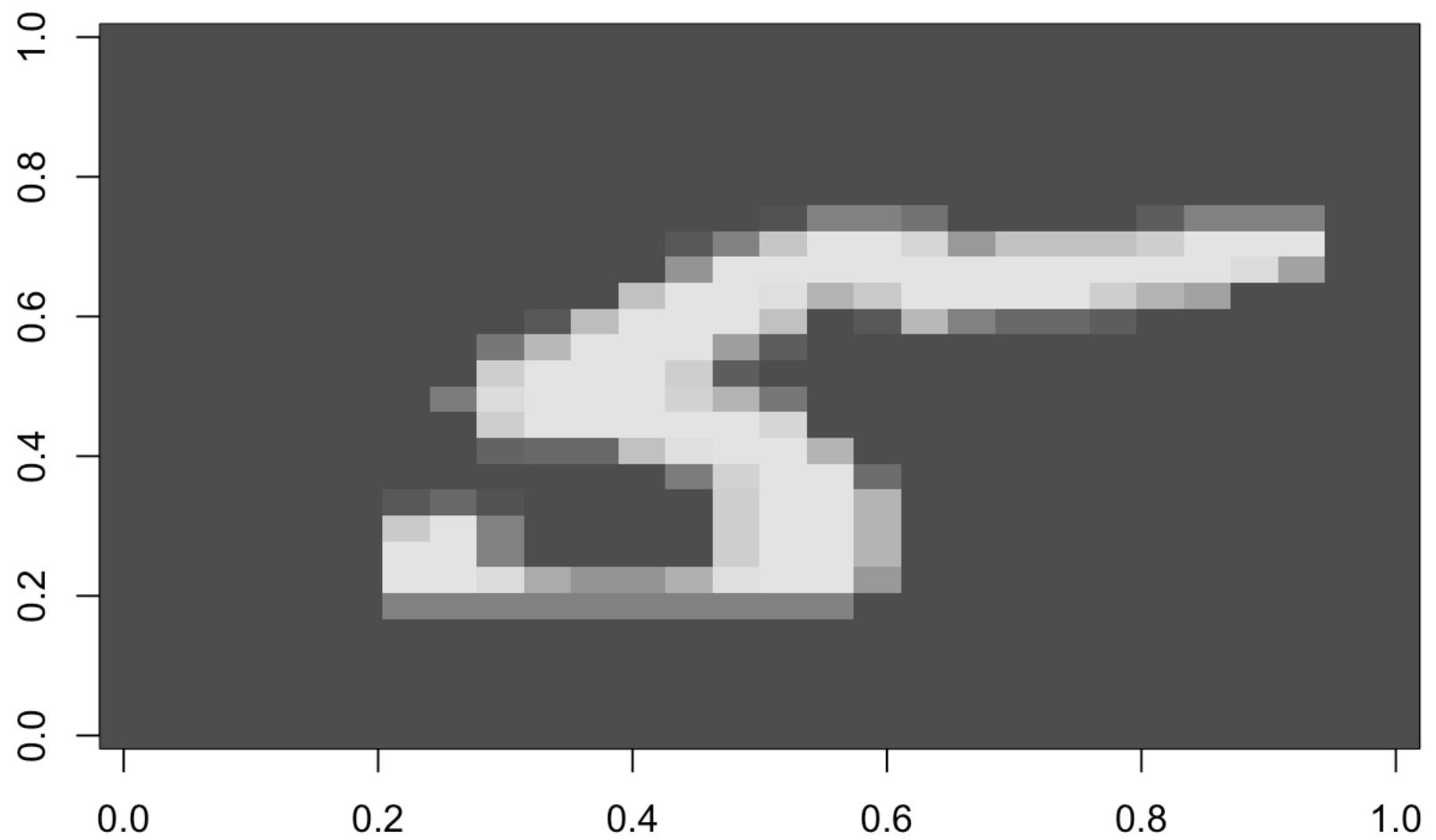
2



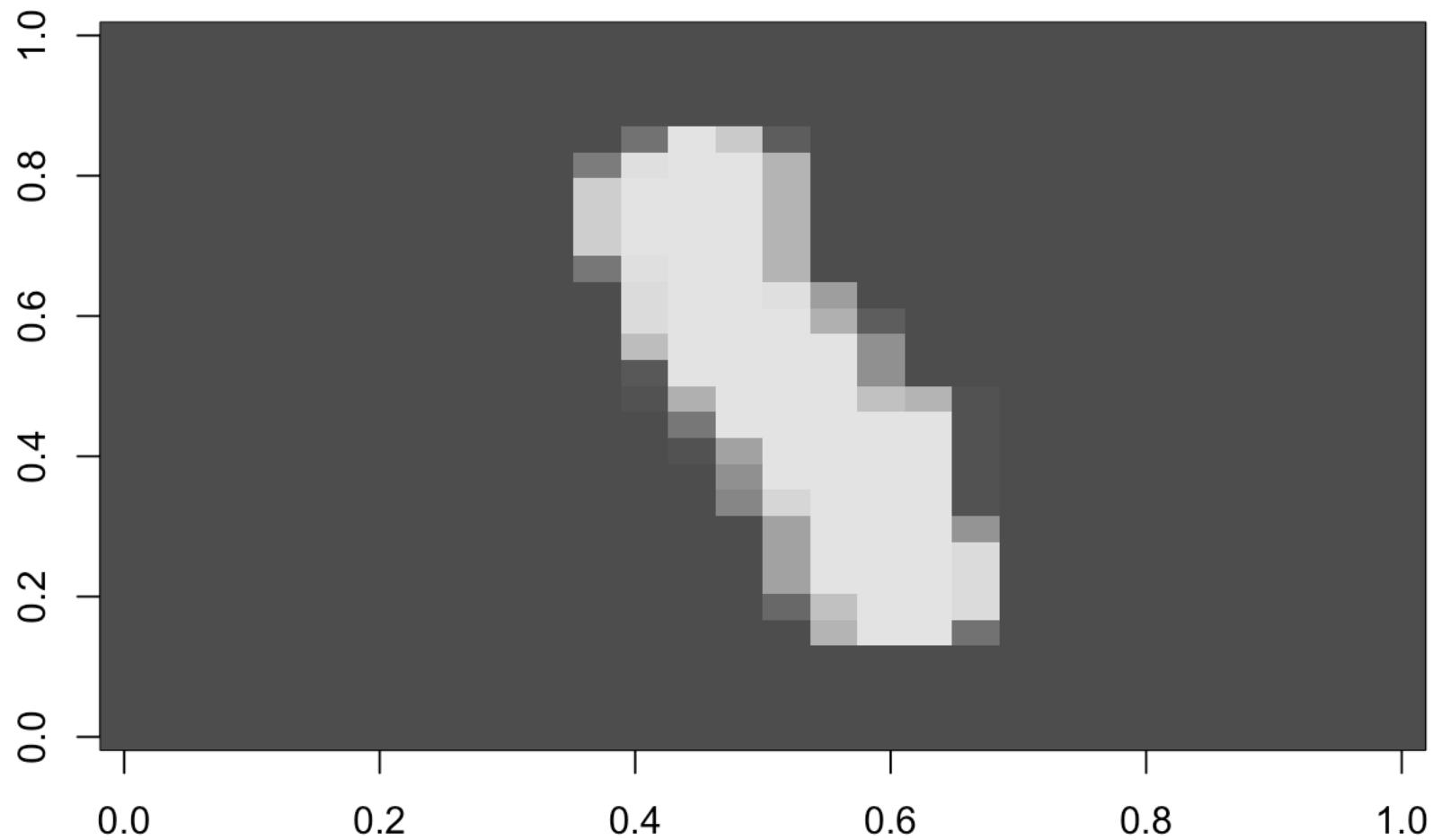
5

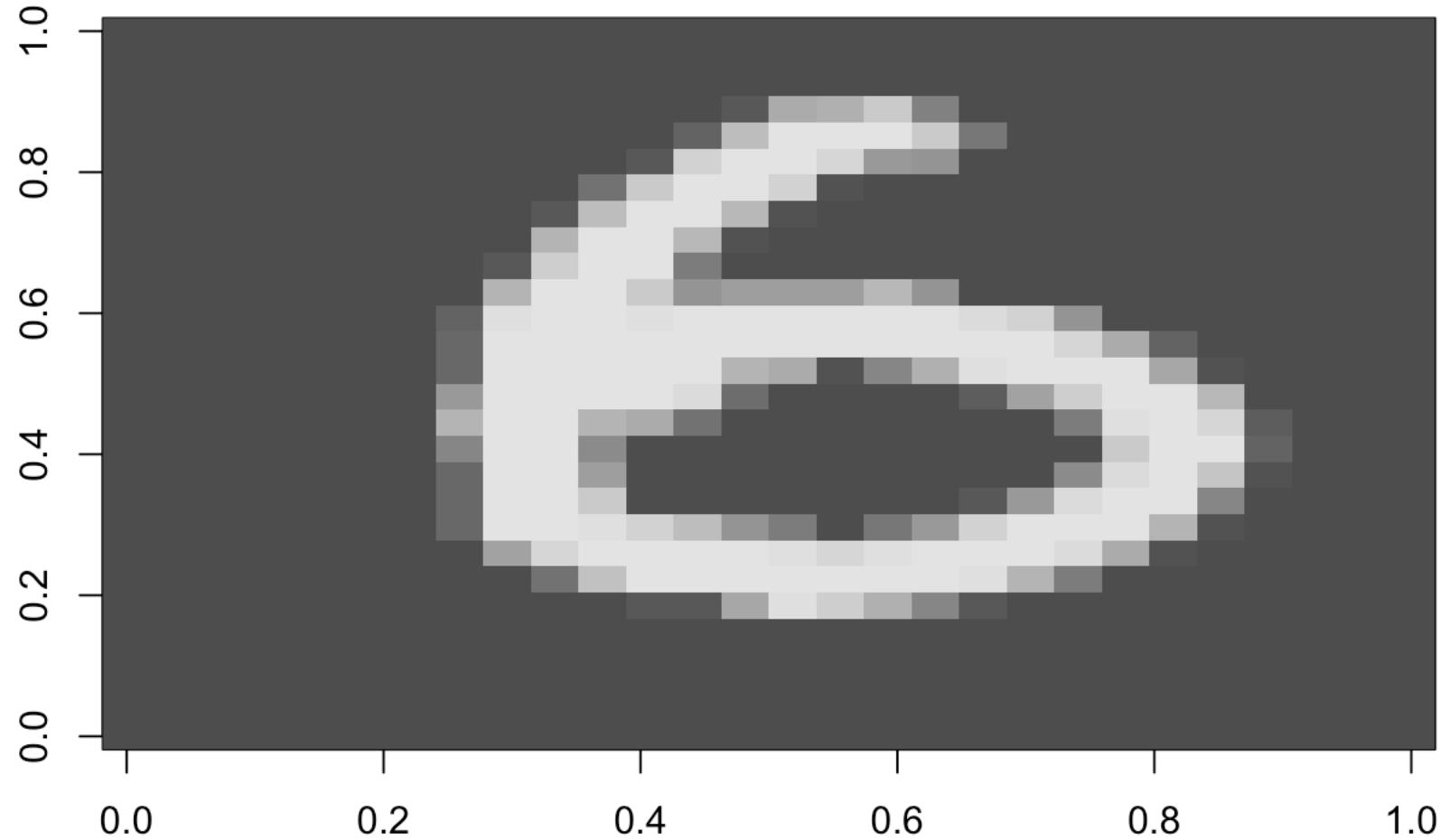


5

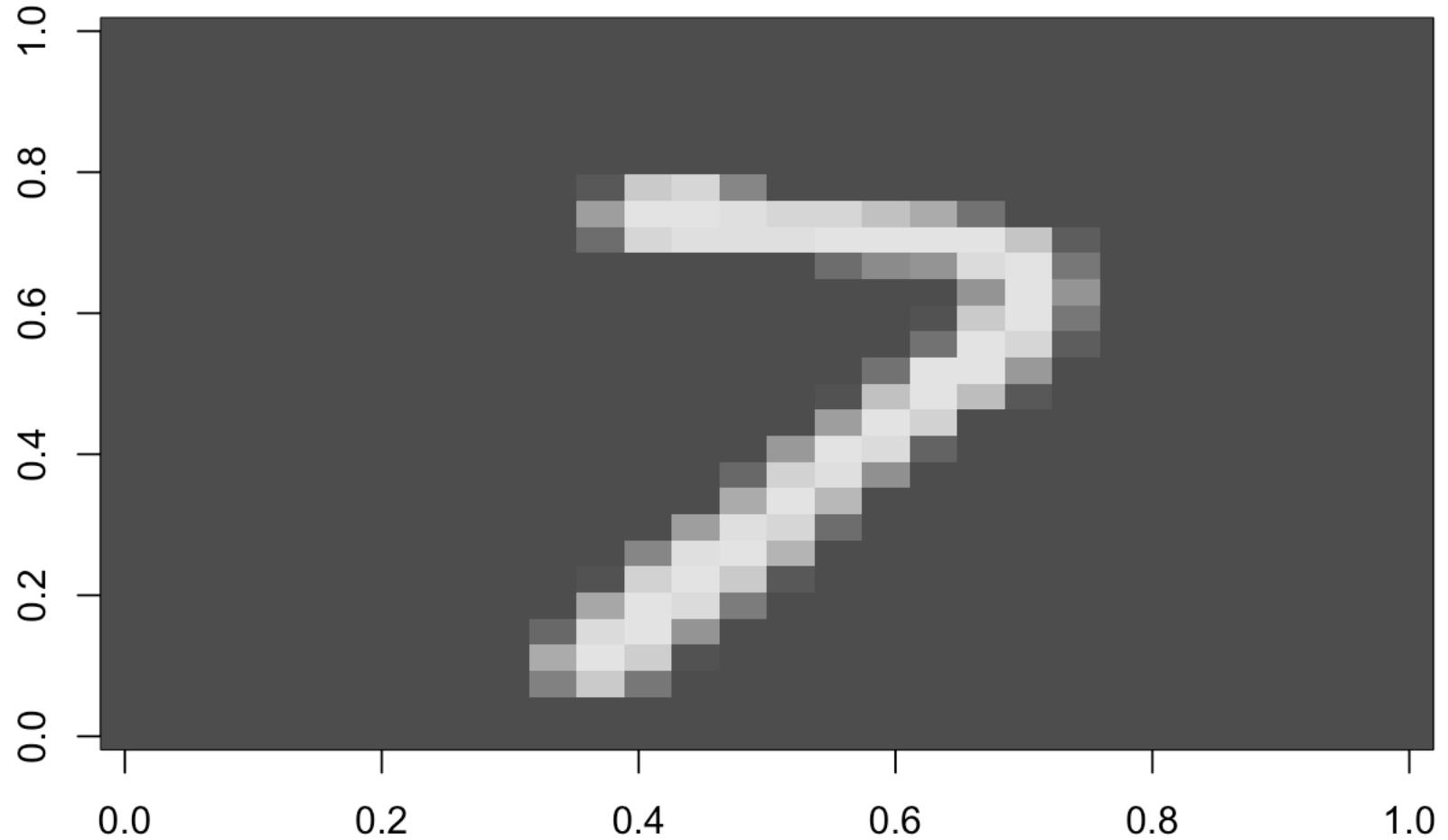


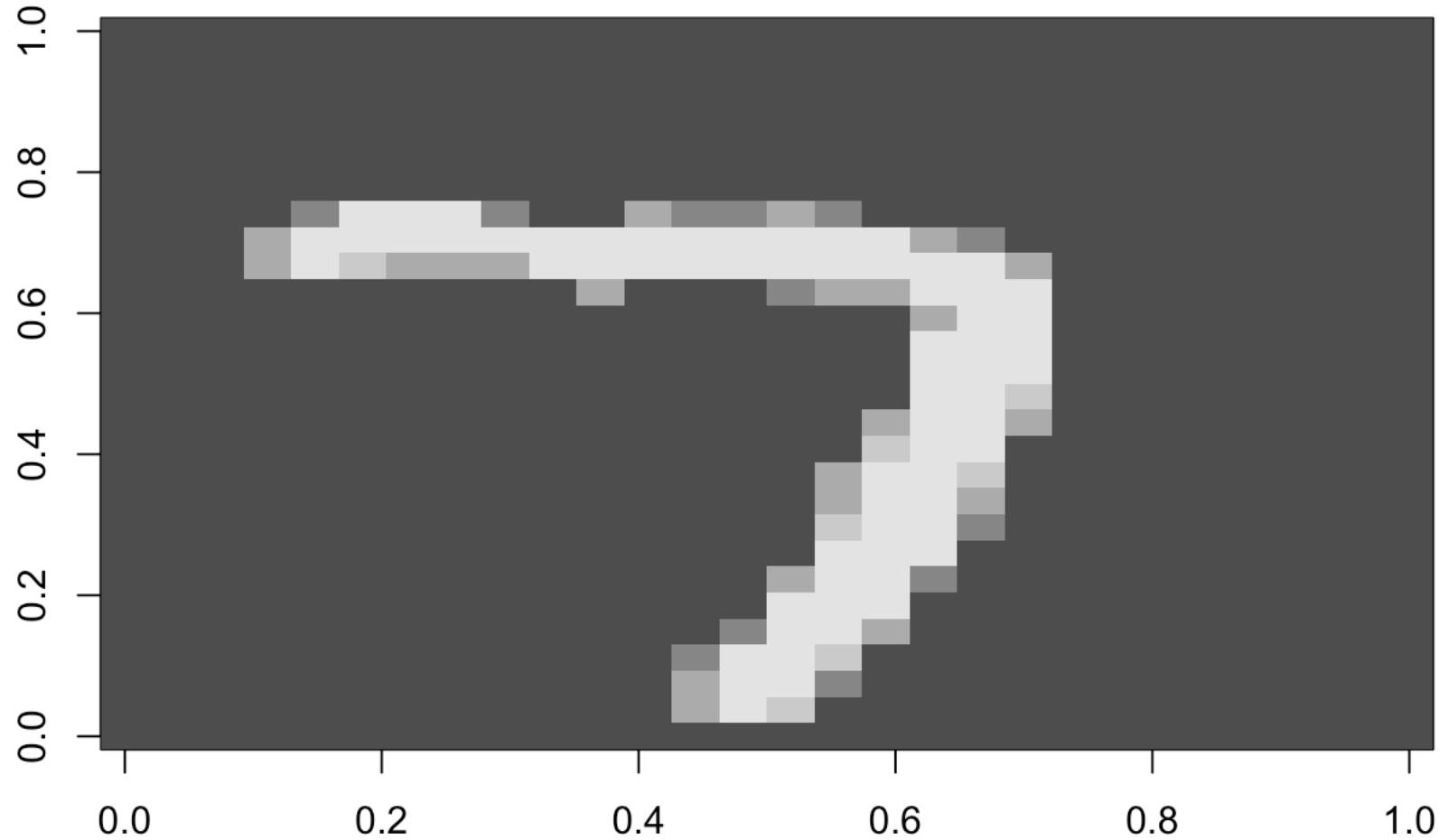
1



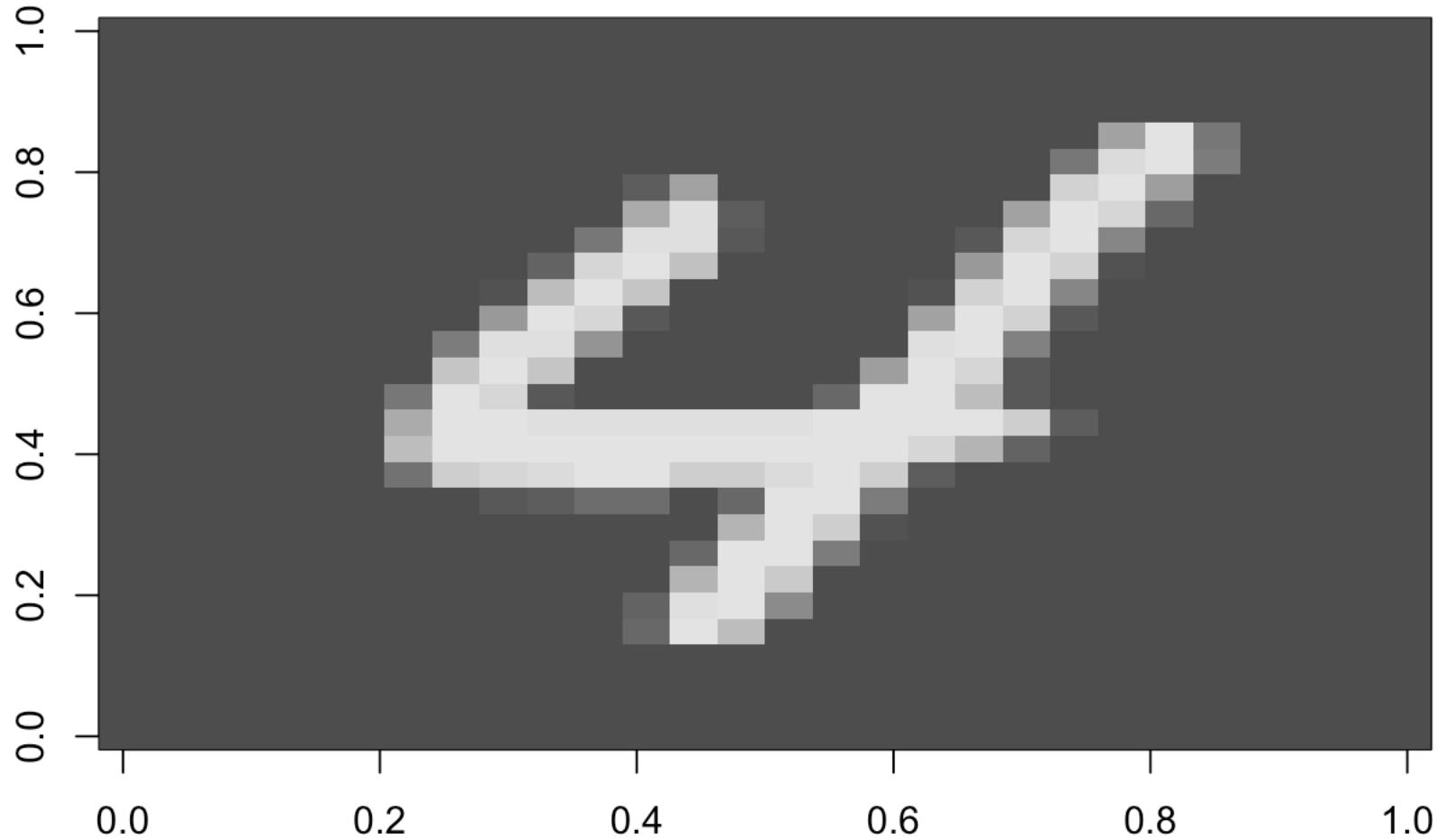


7

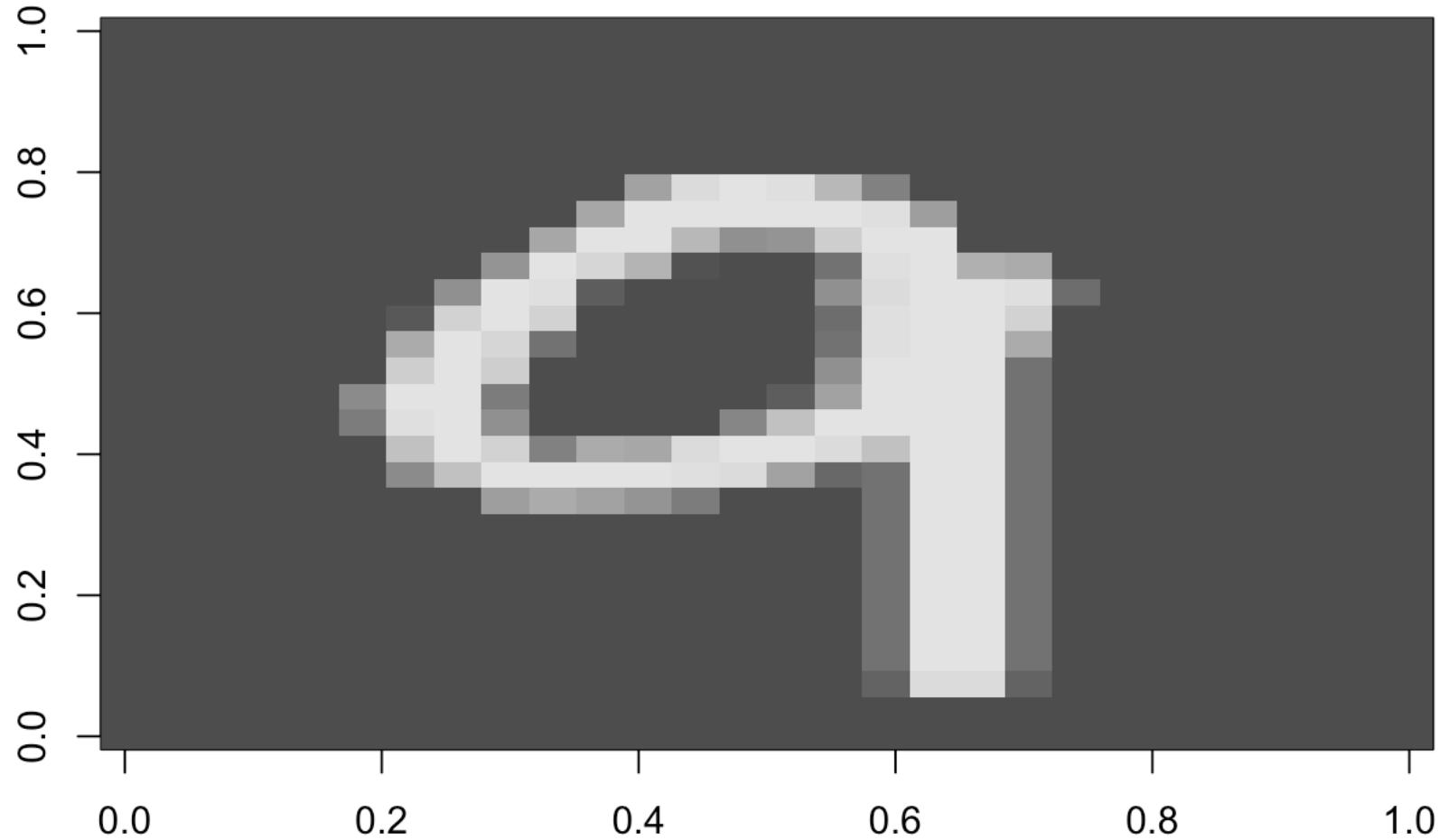




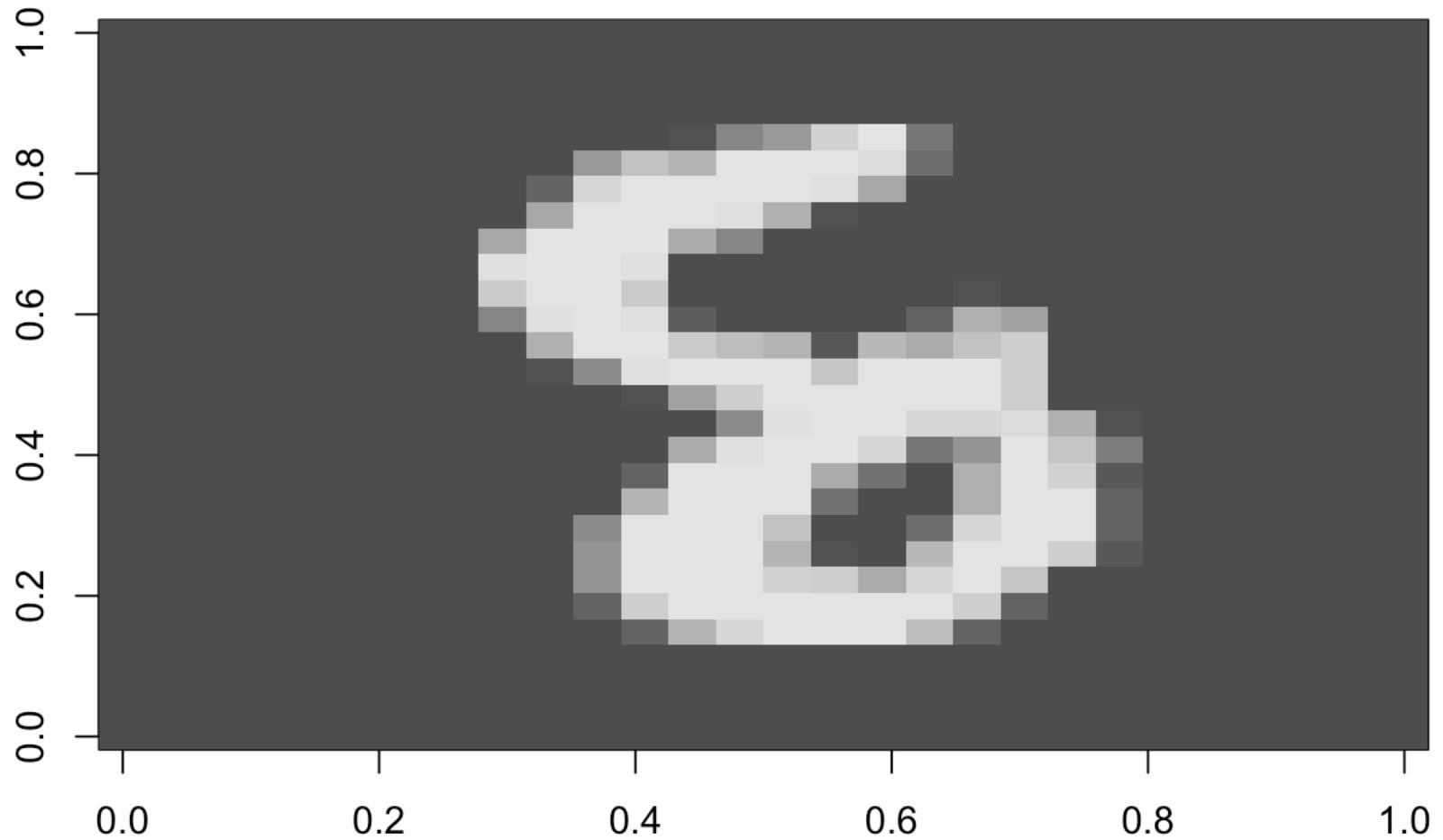
4

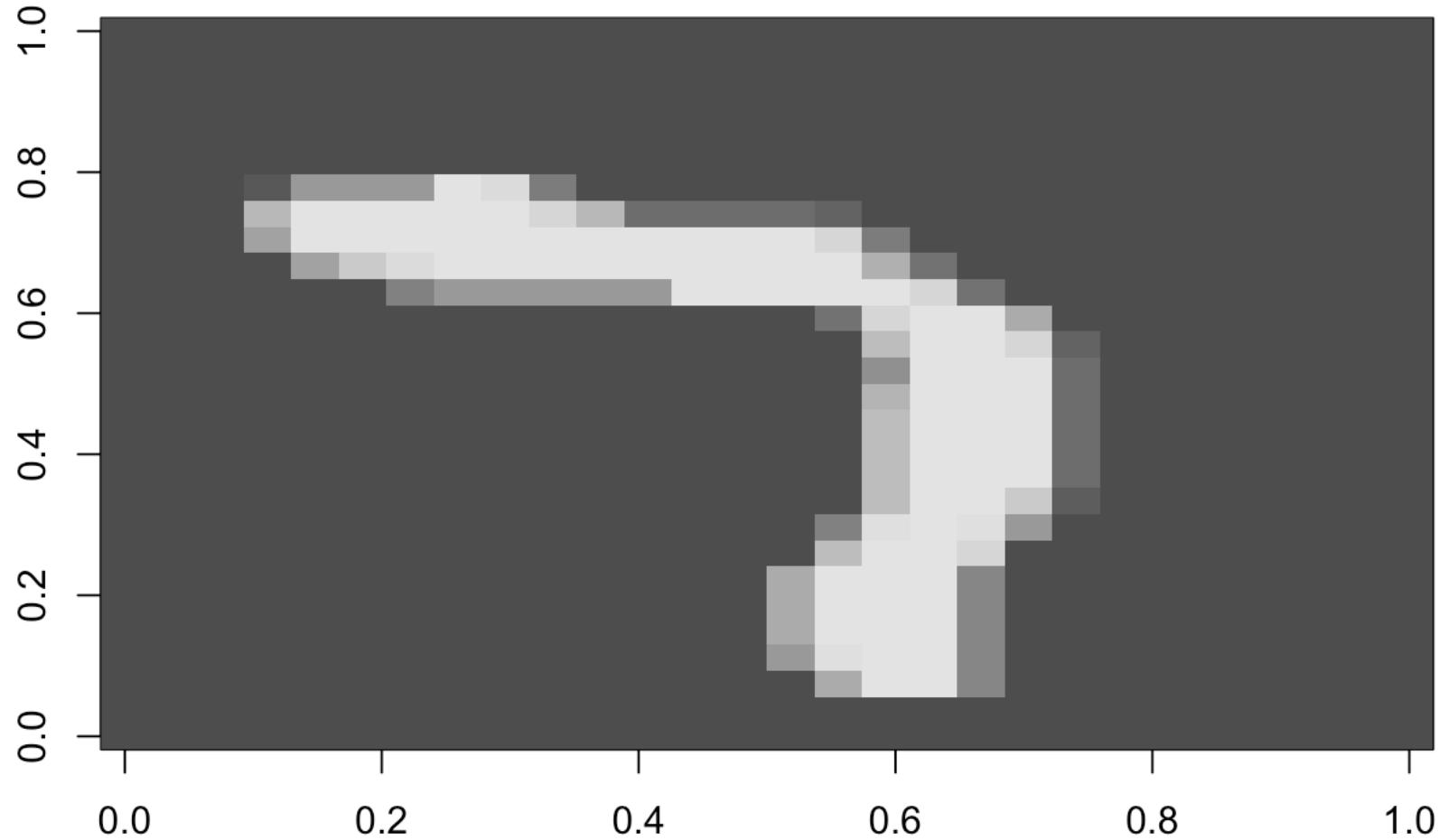


9

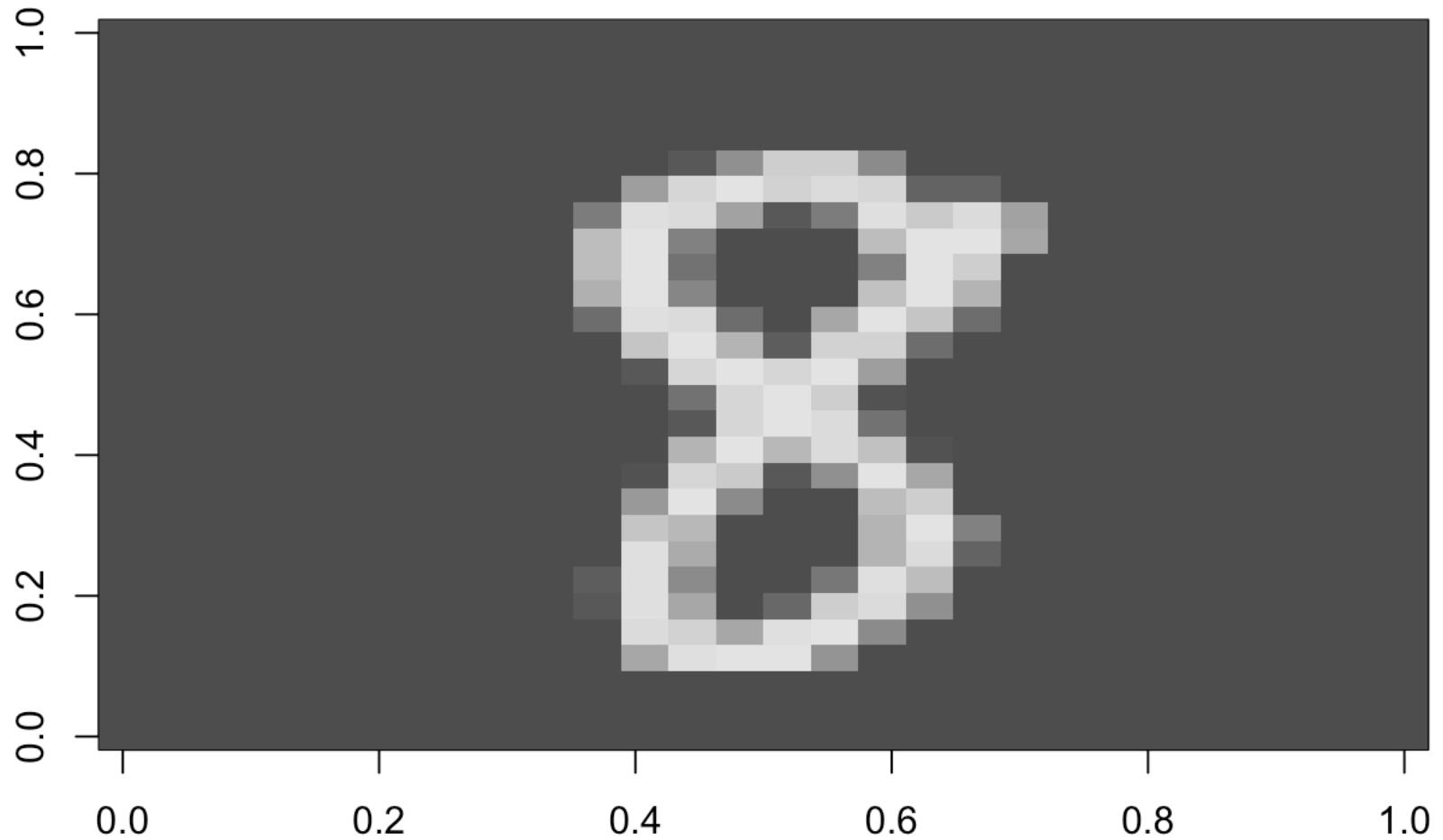


8

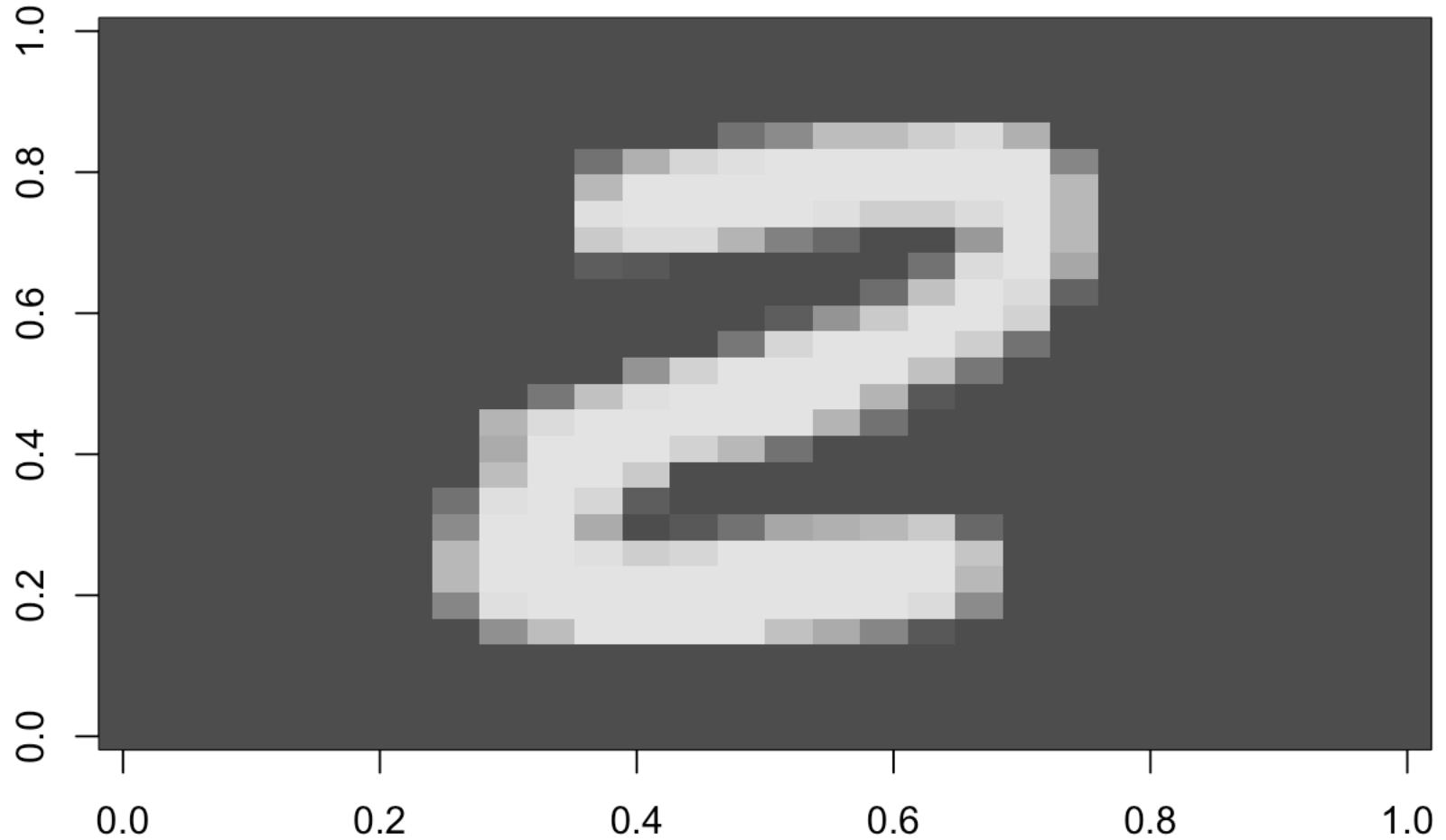




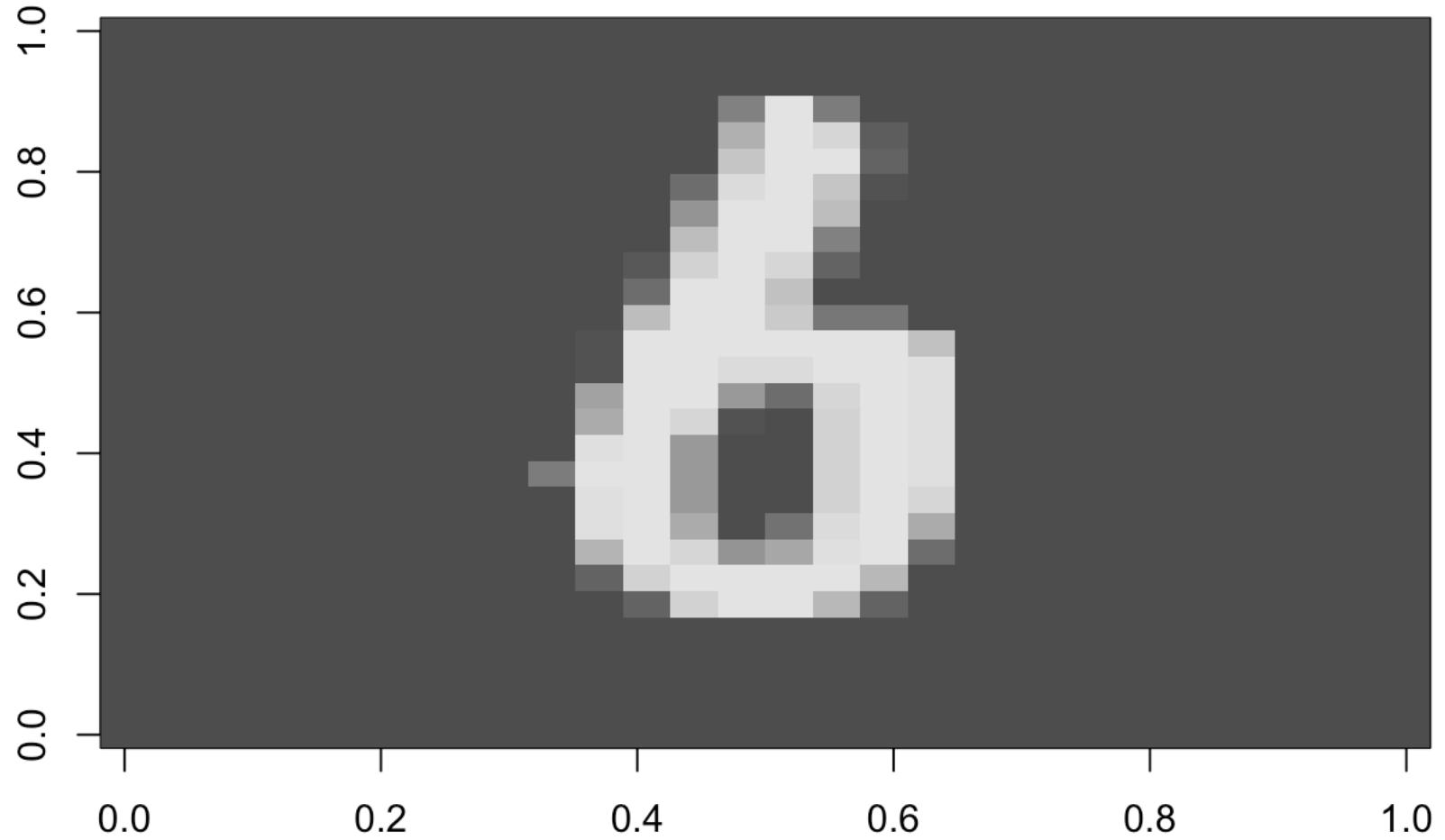
8

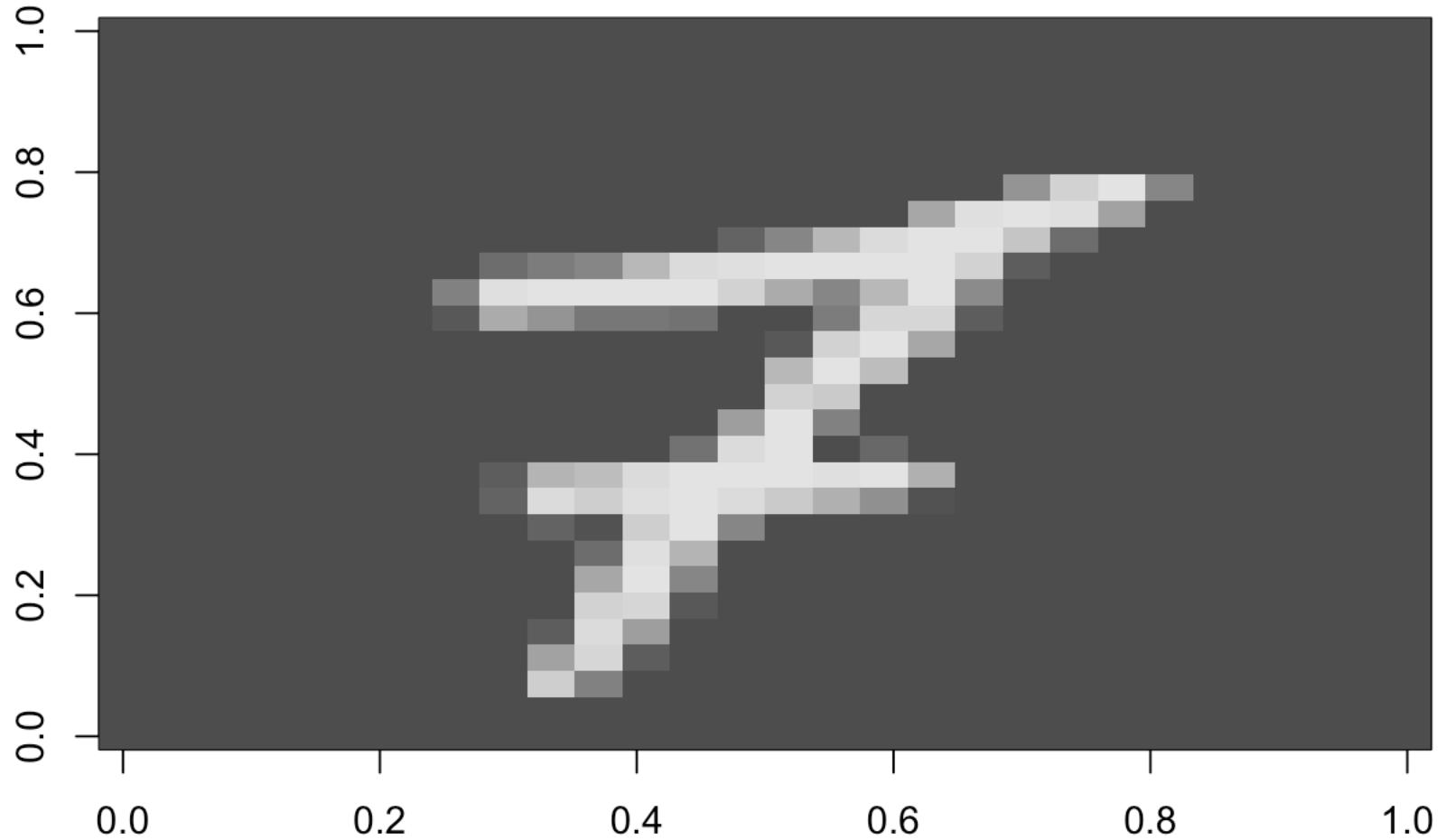


2

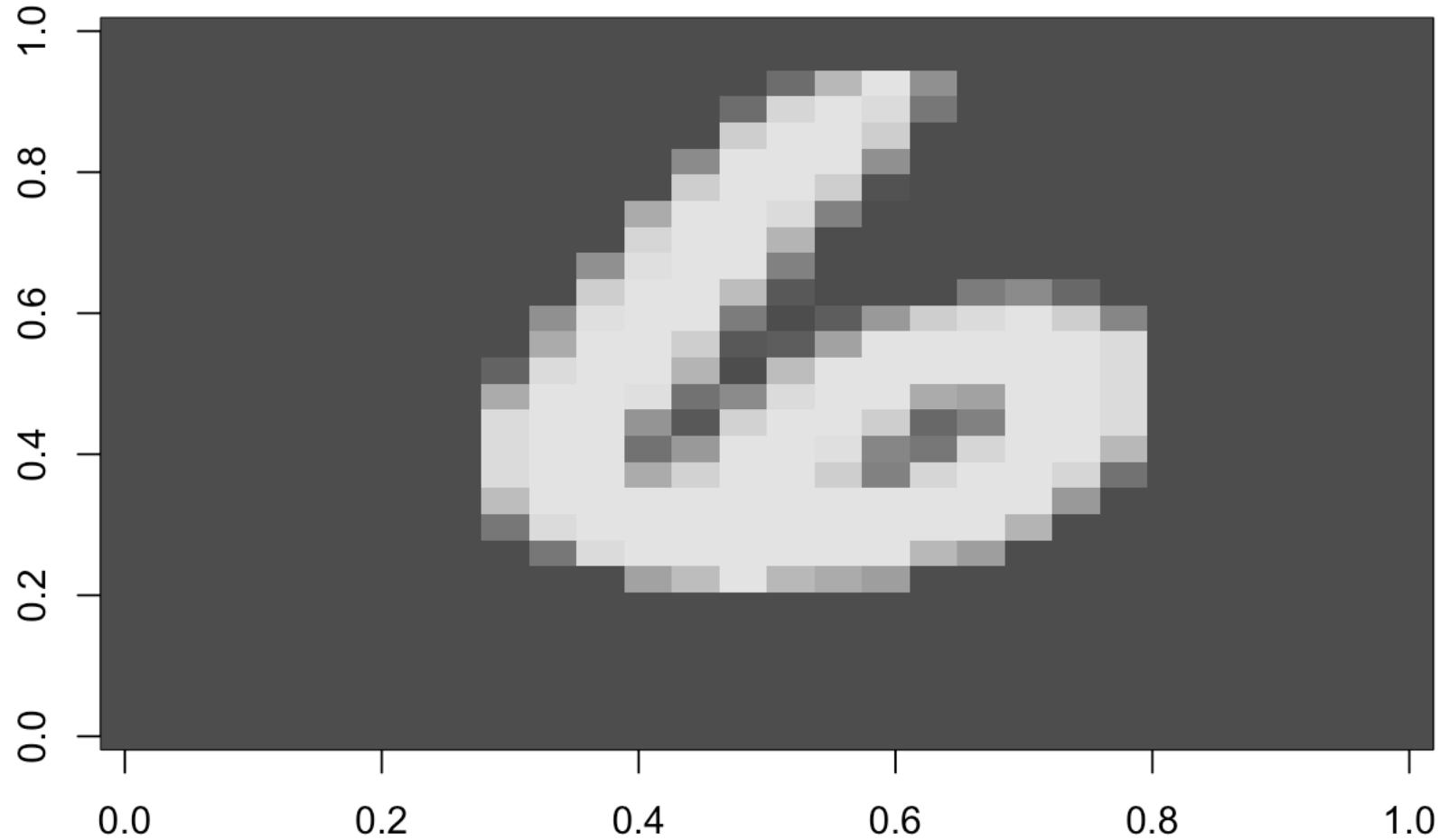


6

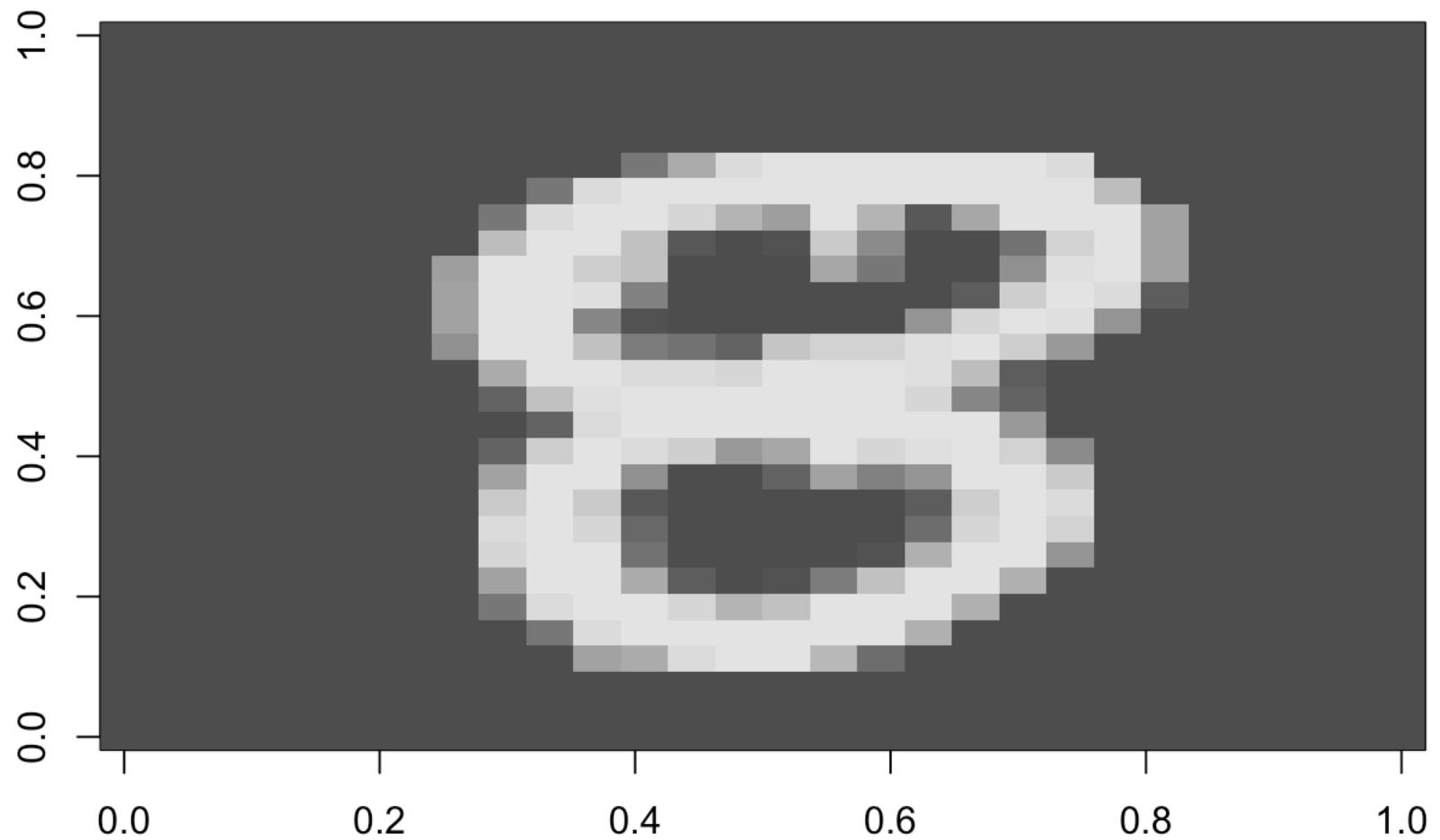




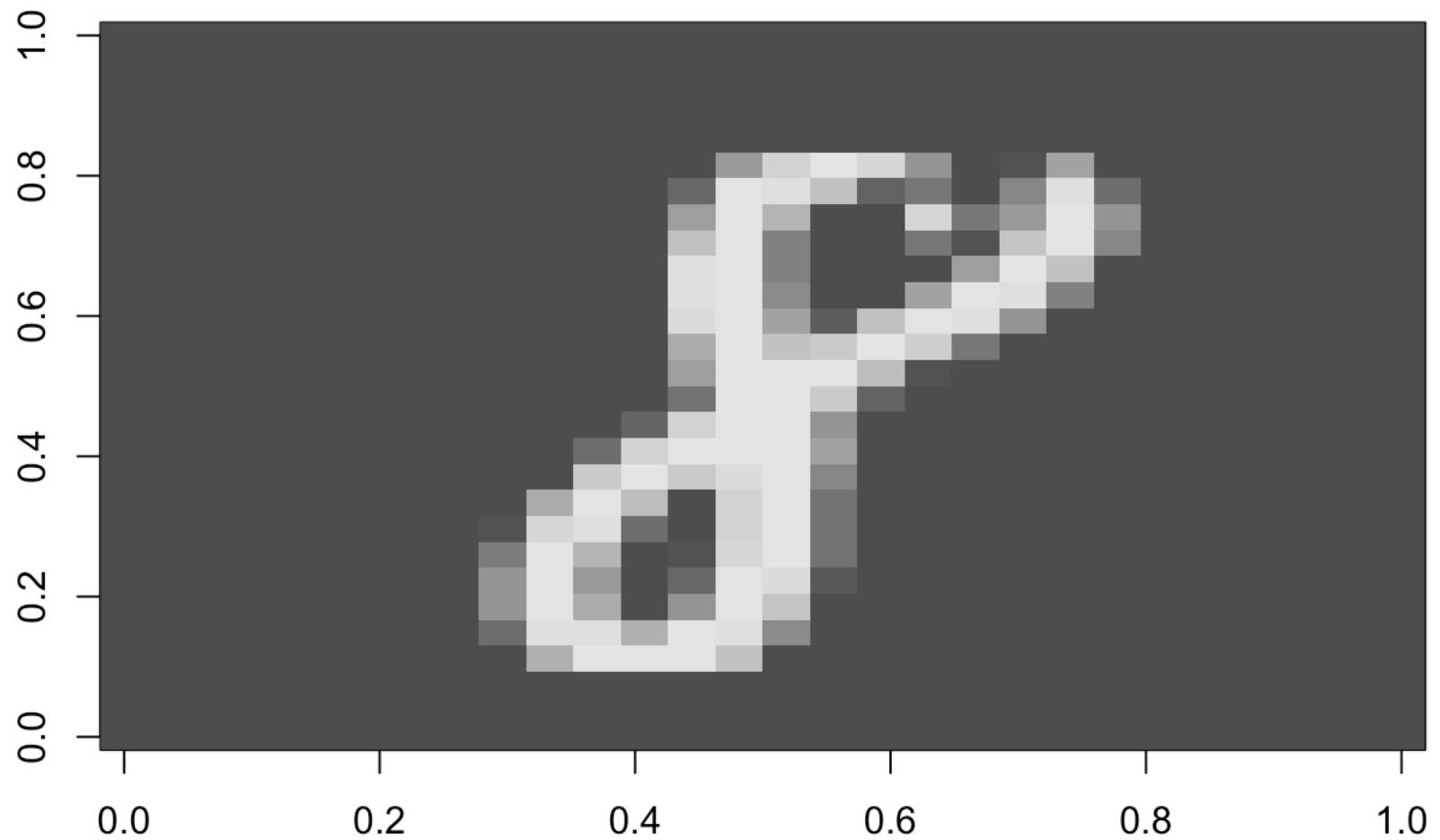
6



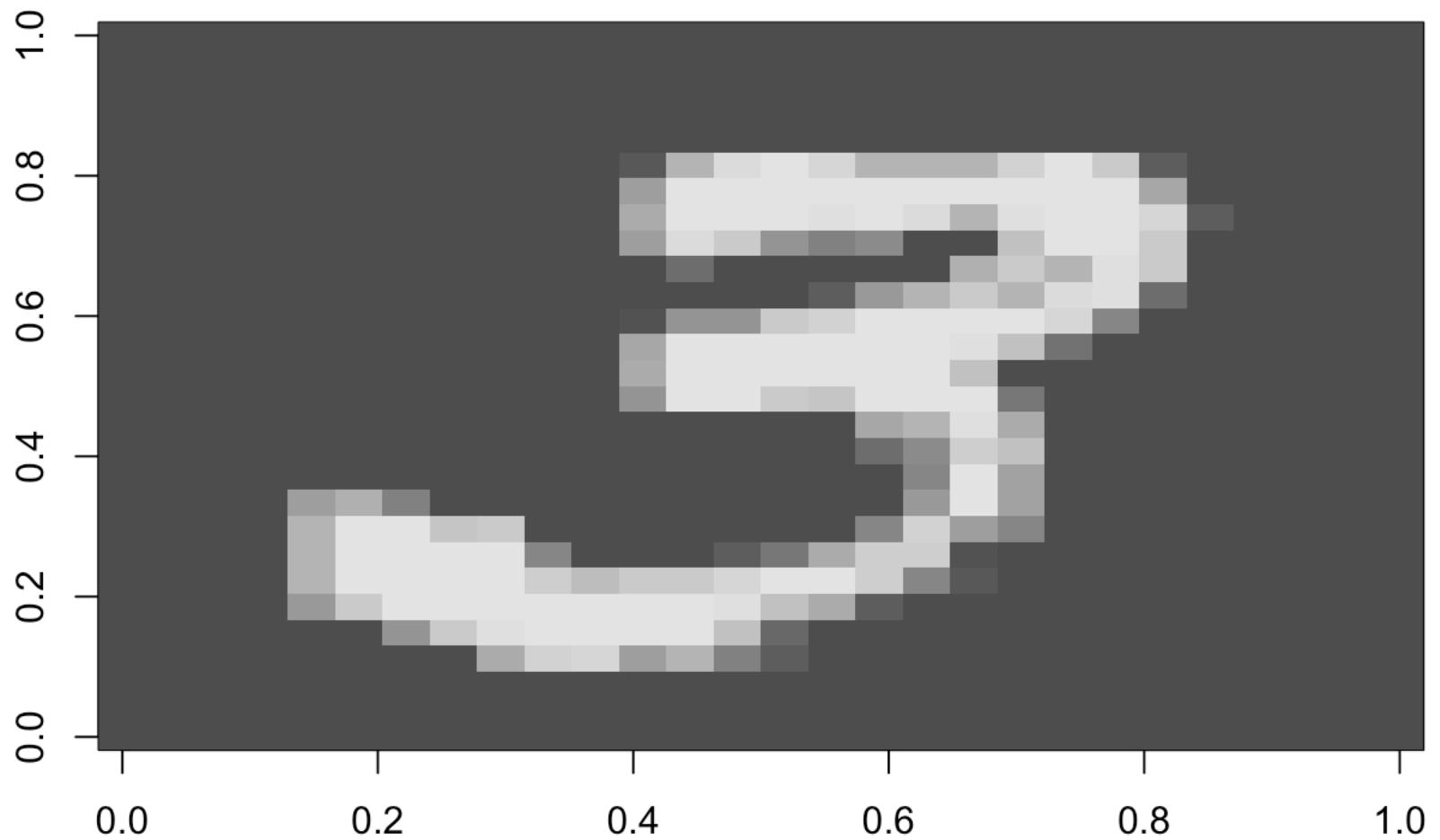
8



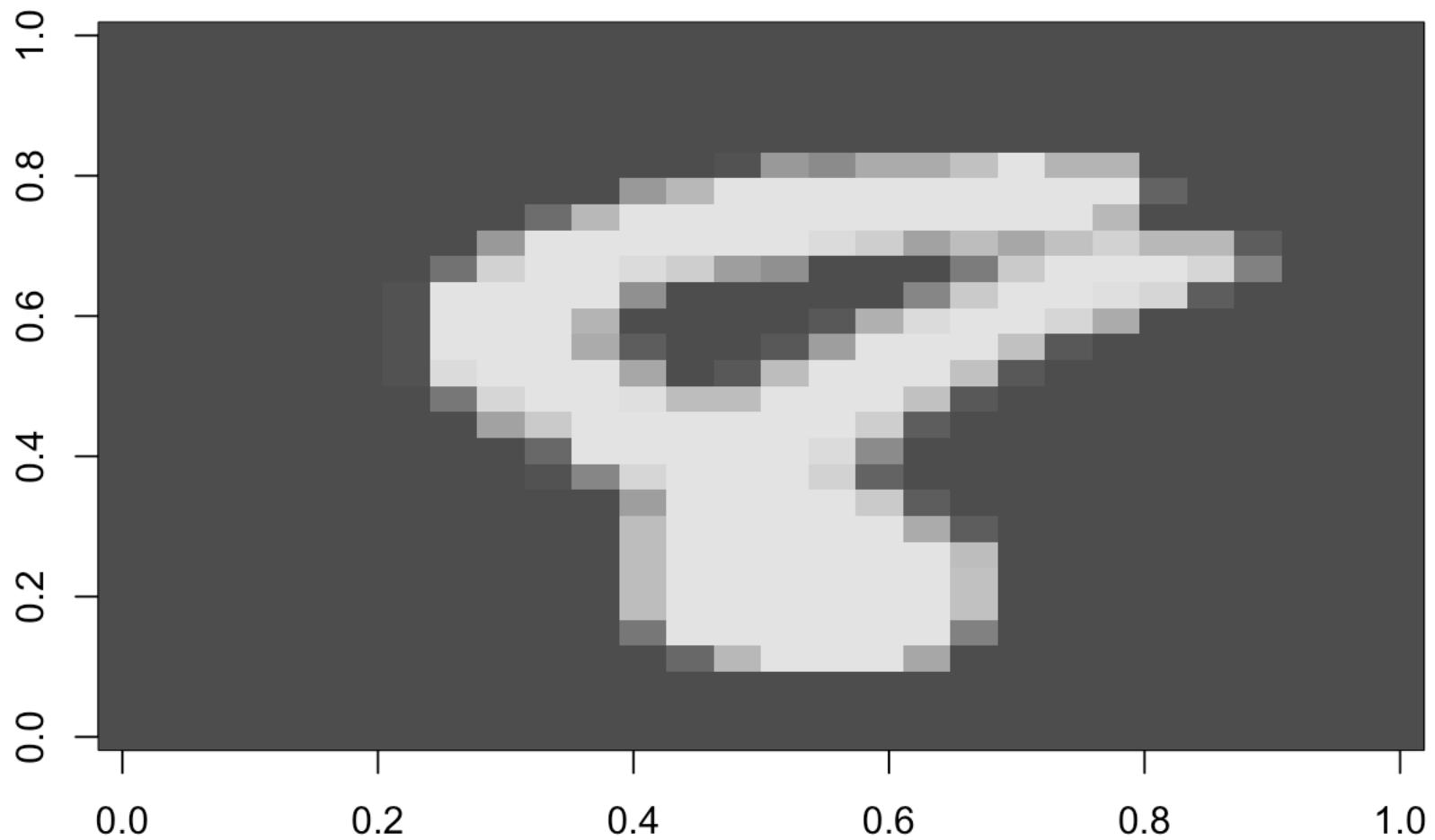
8



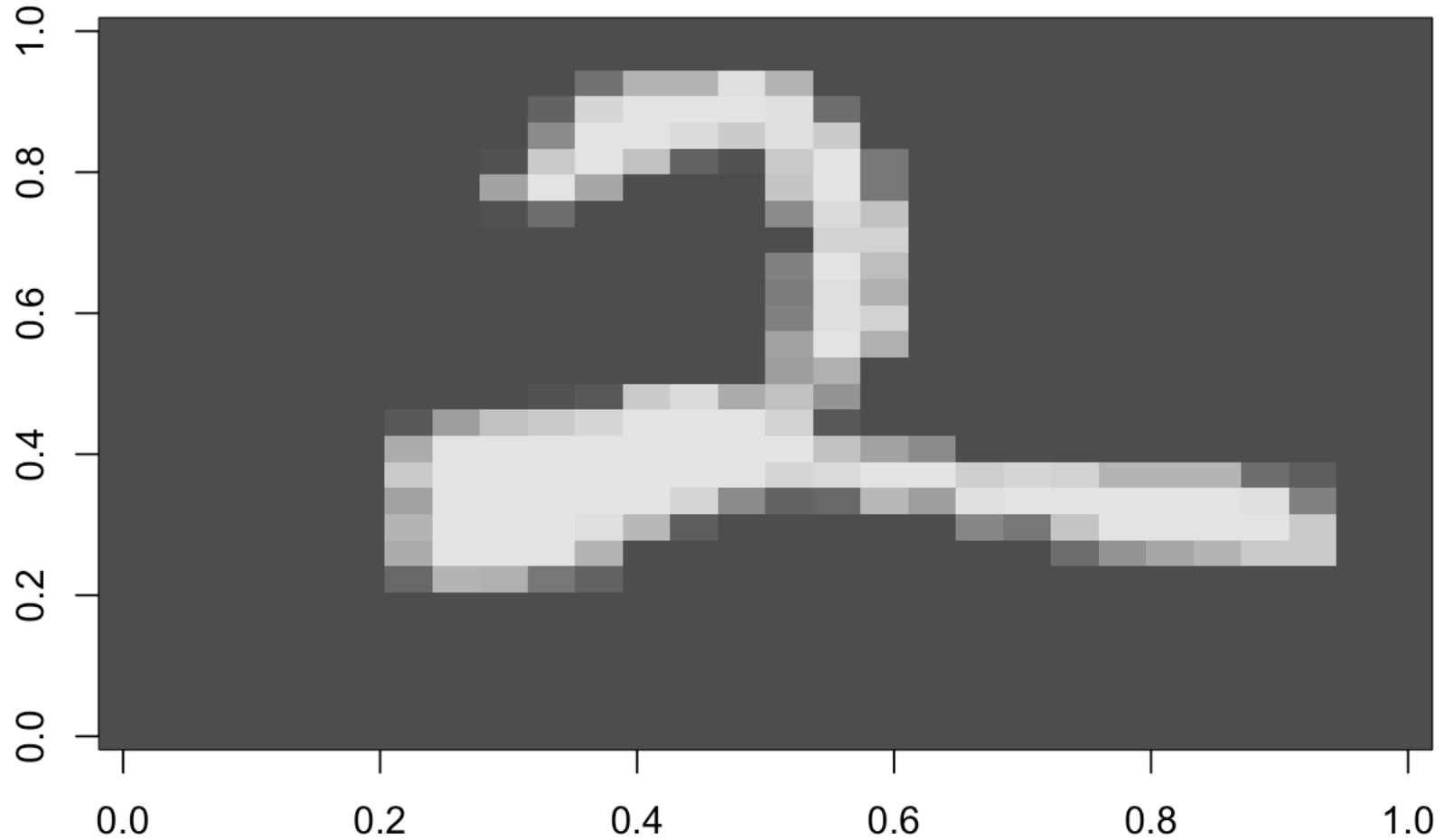
3



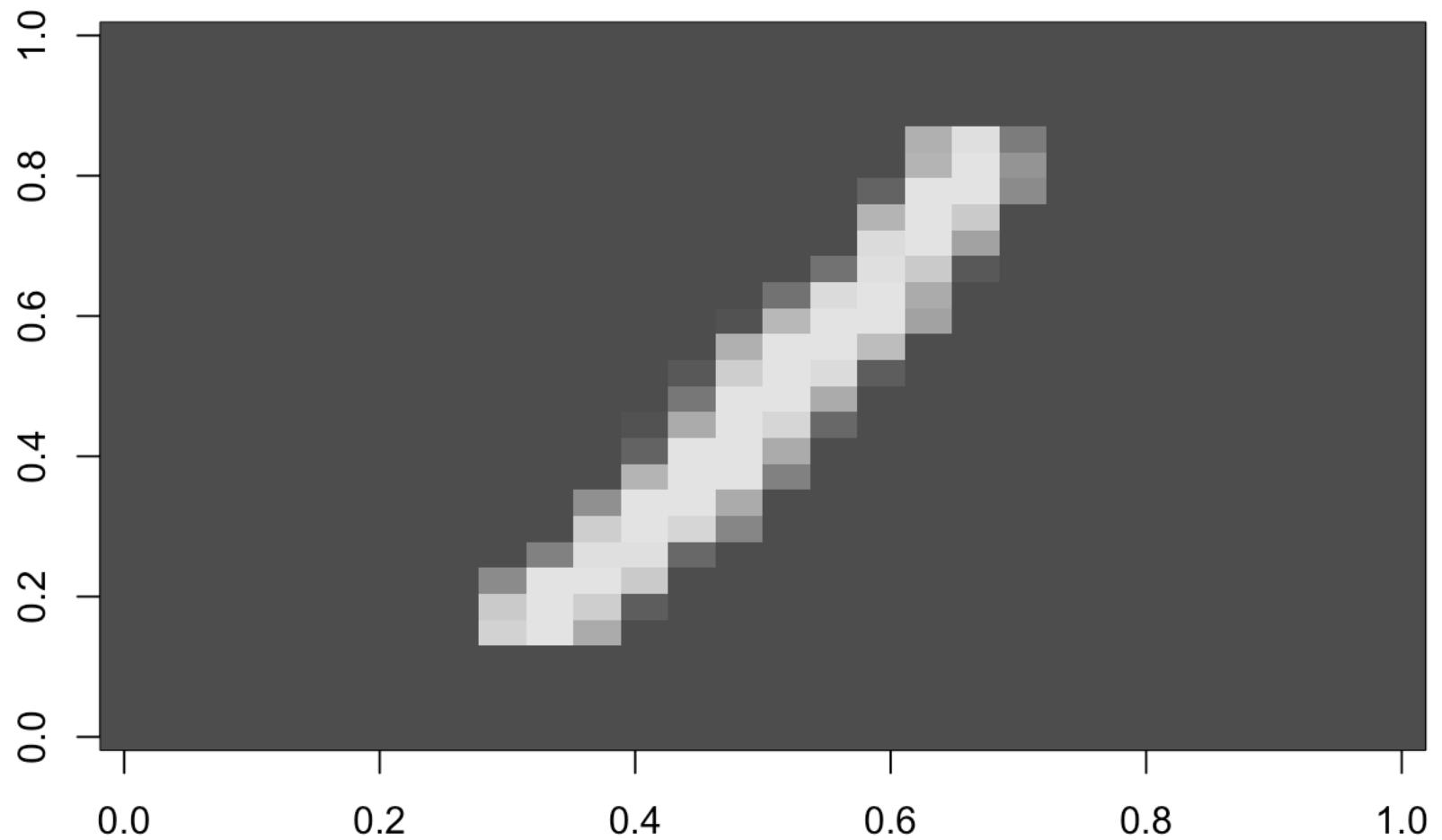
8



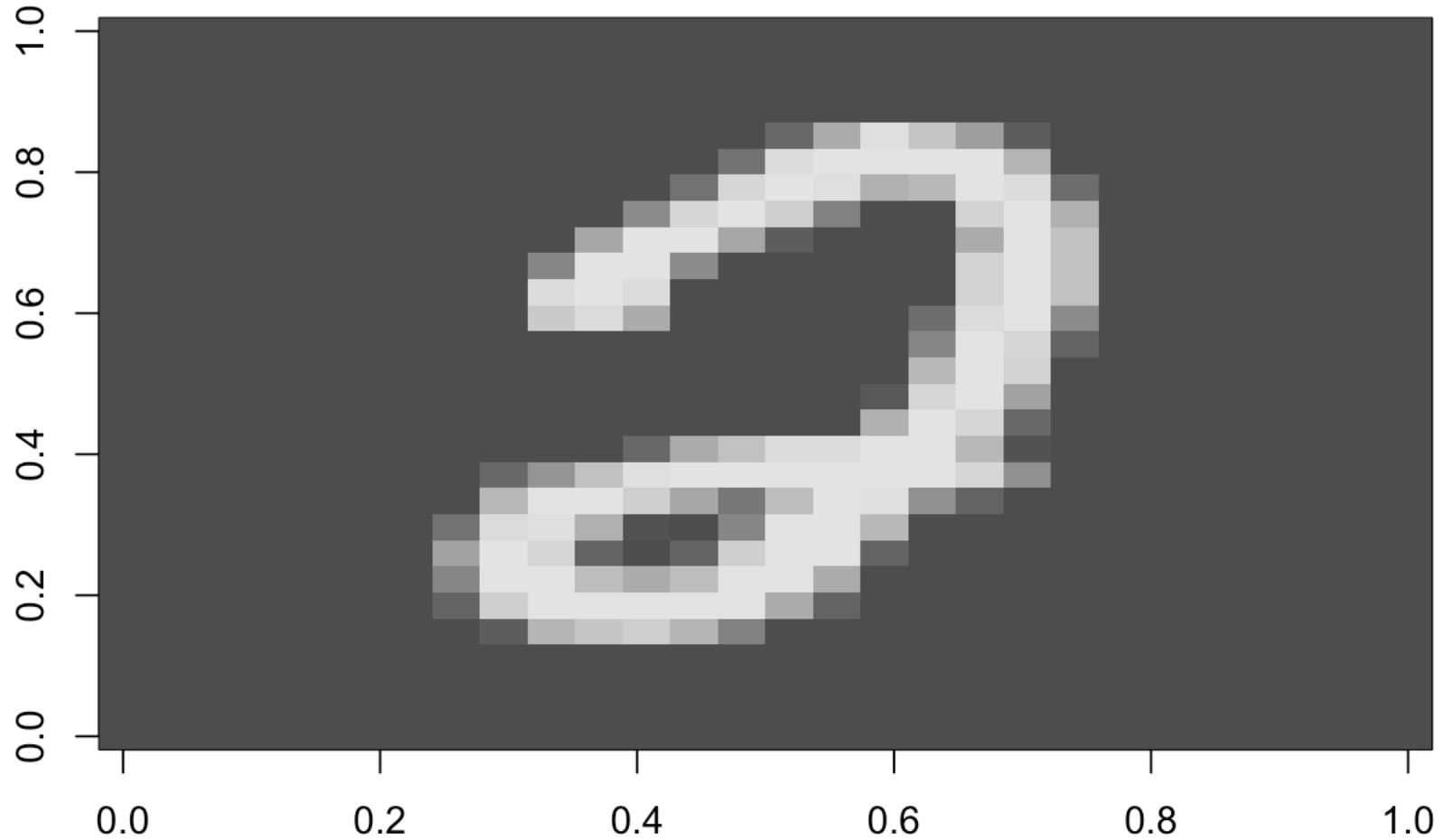
2



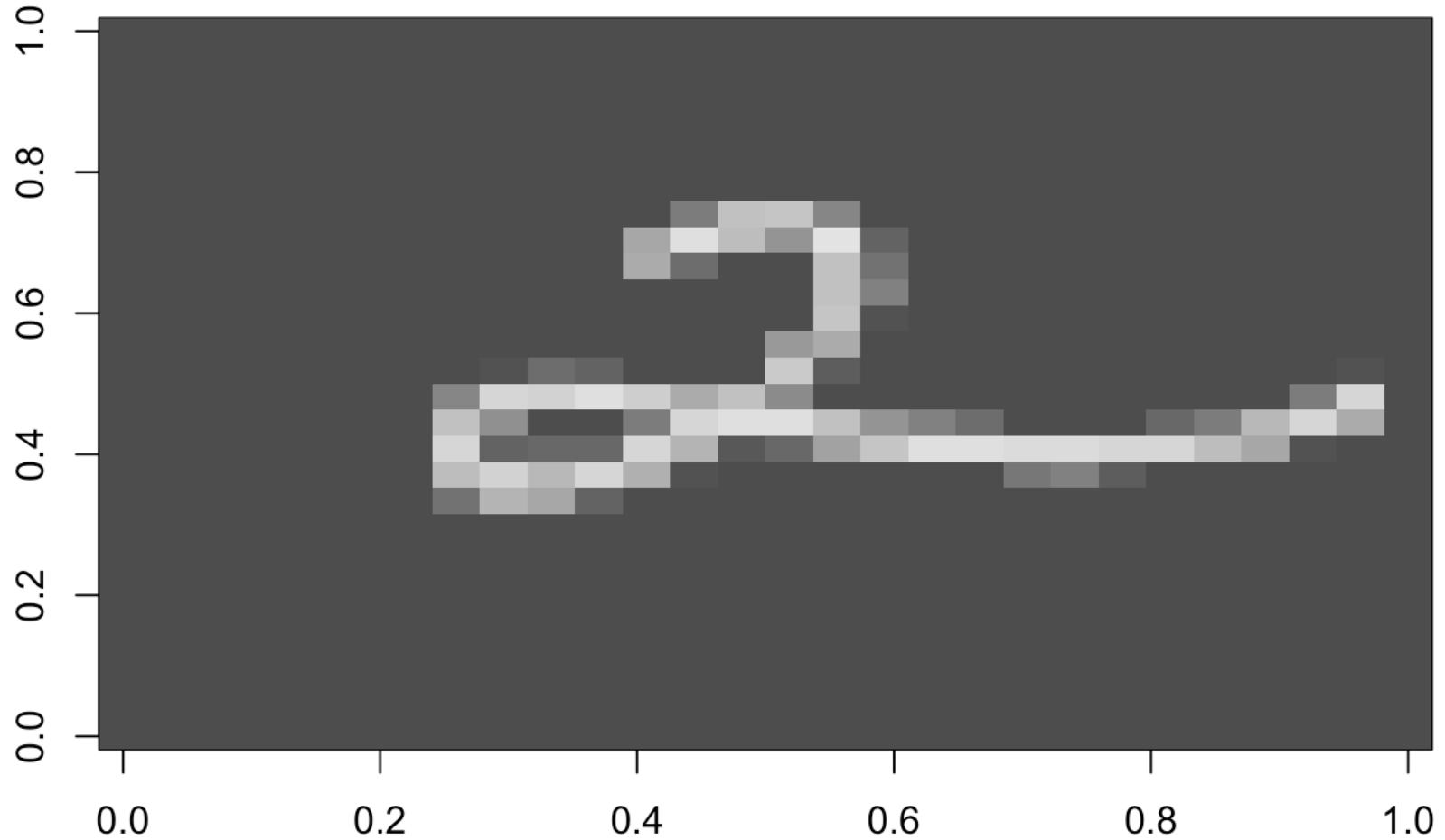
1

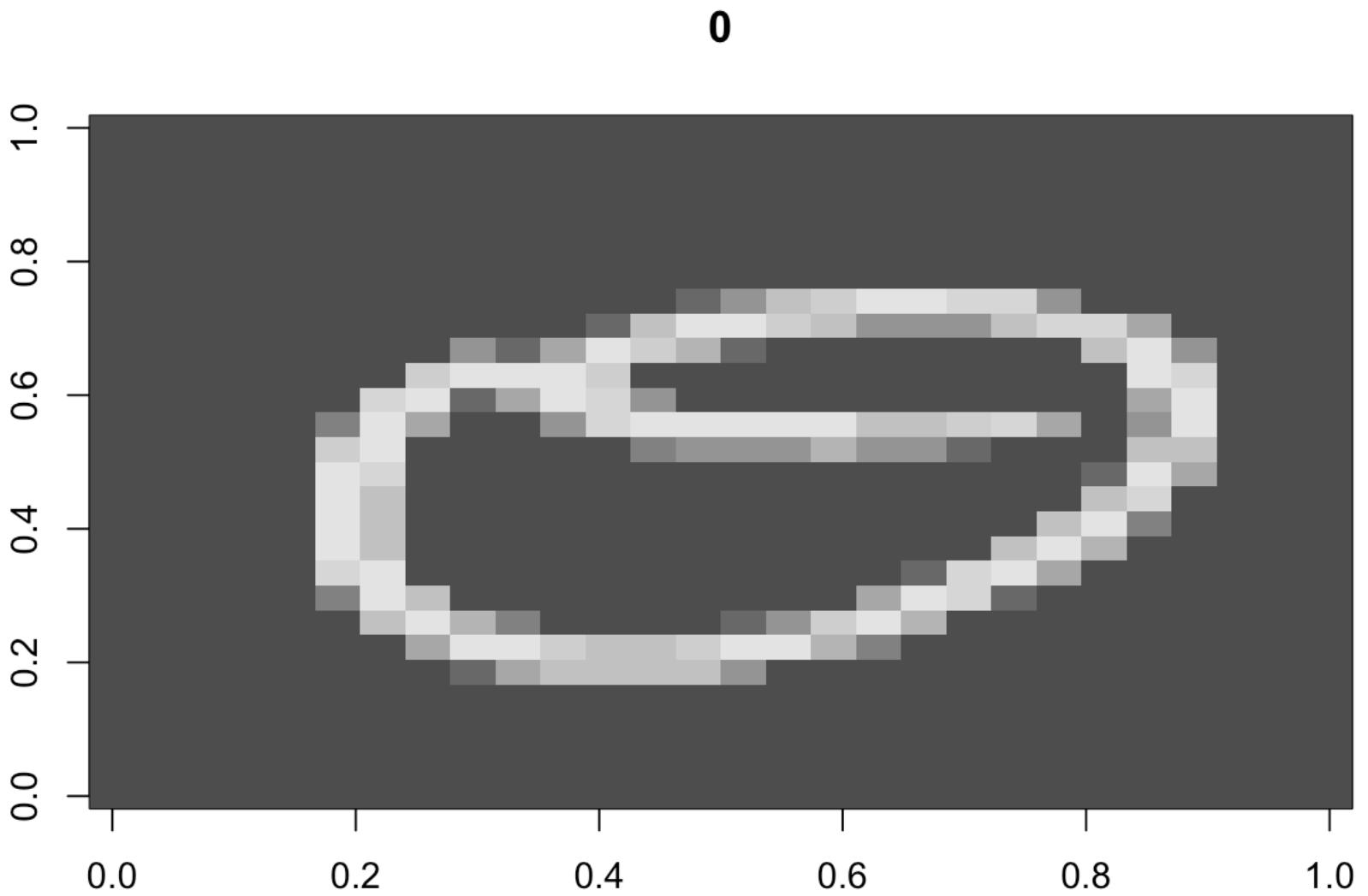


2

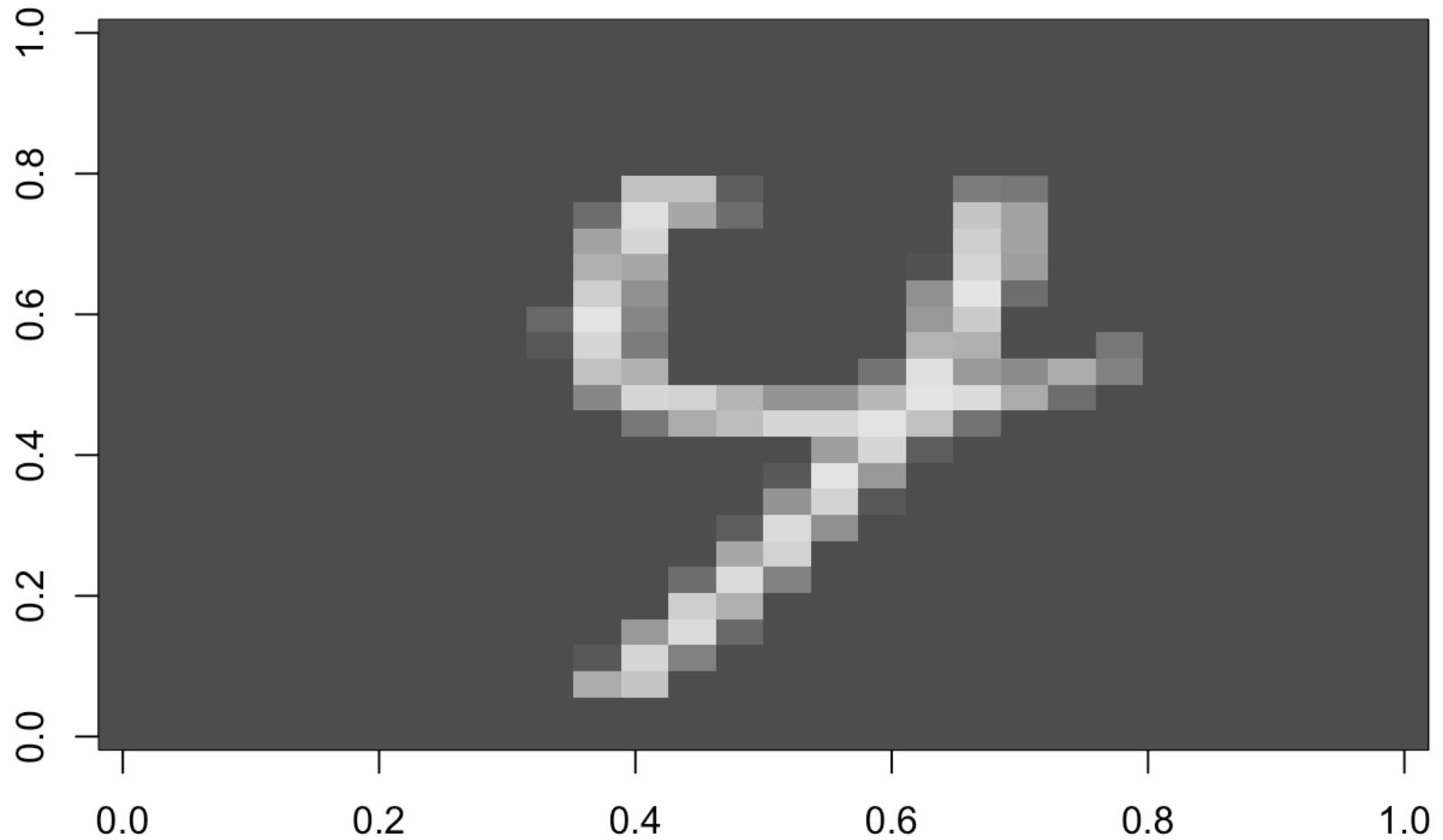


2

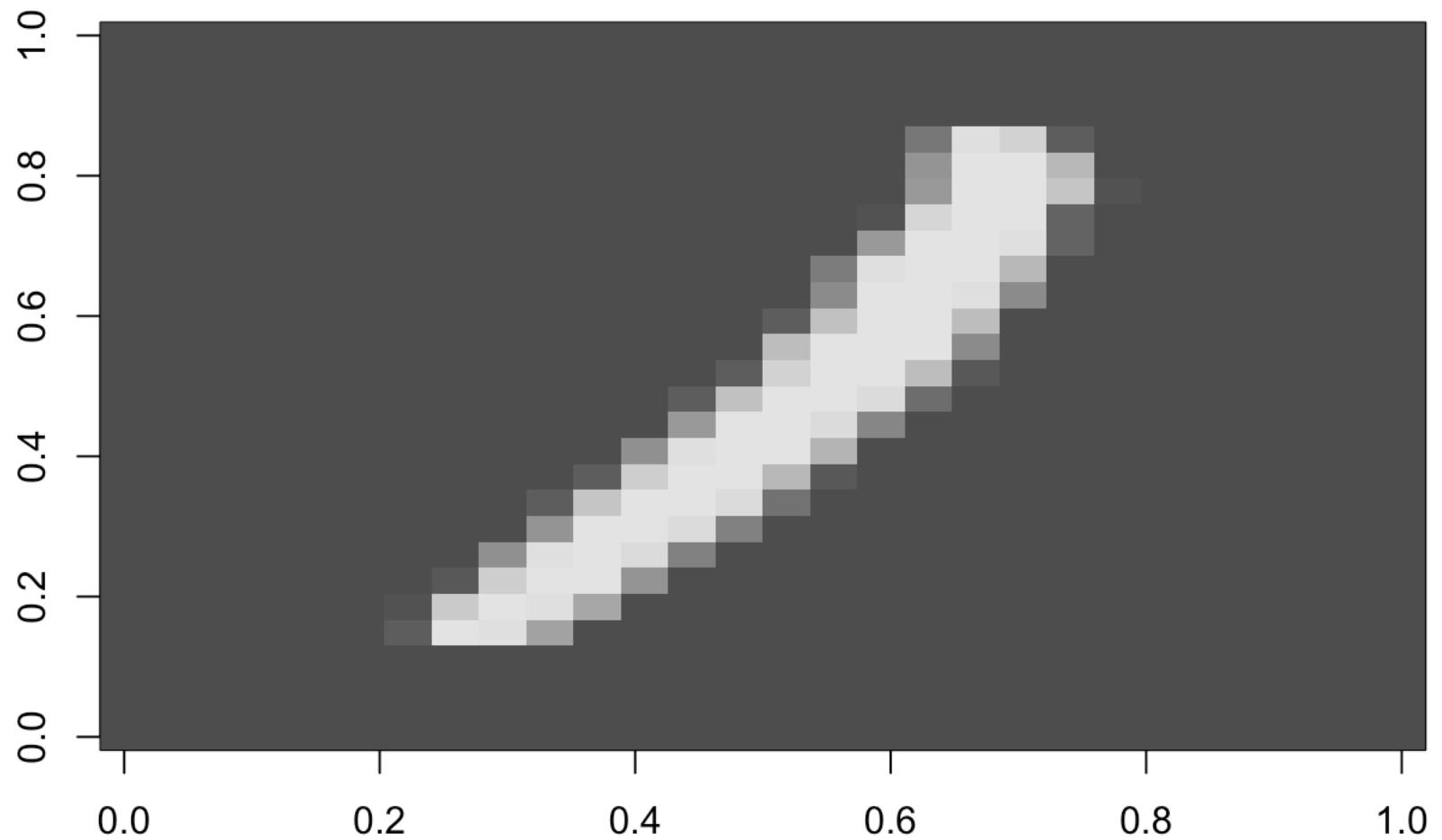


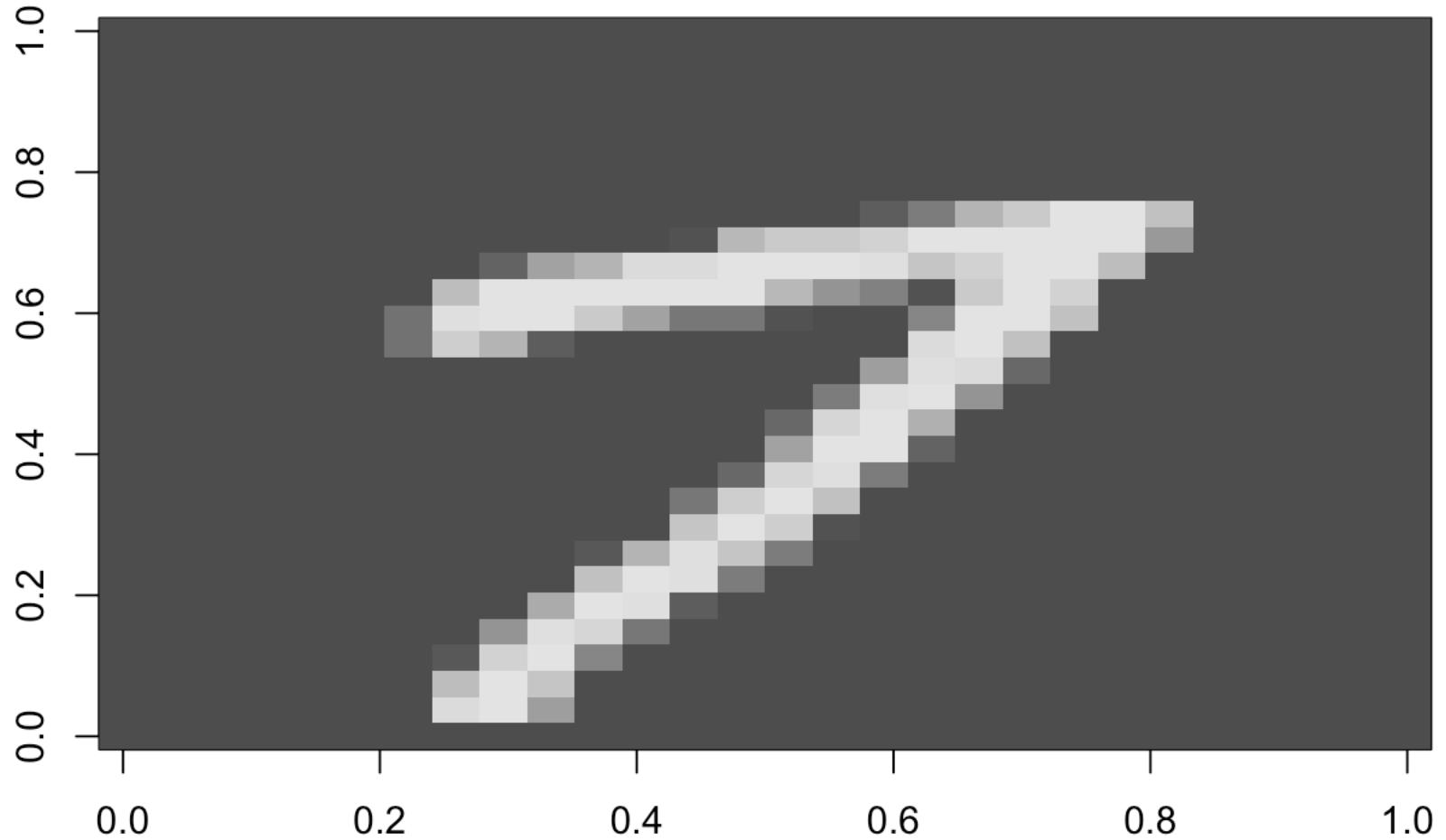


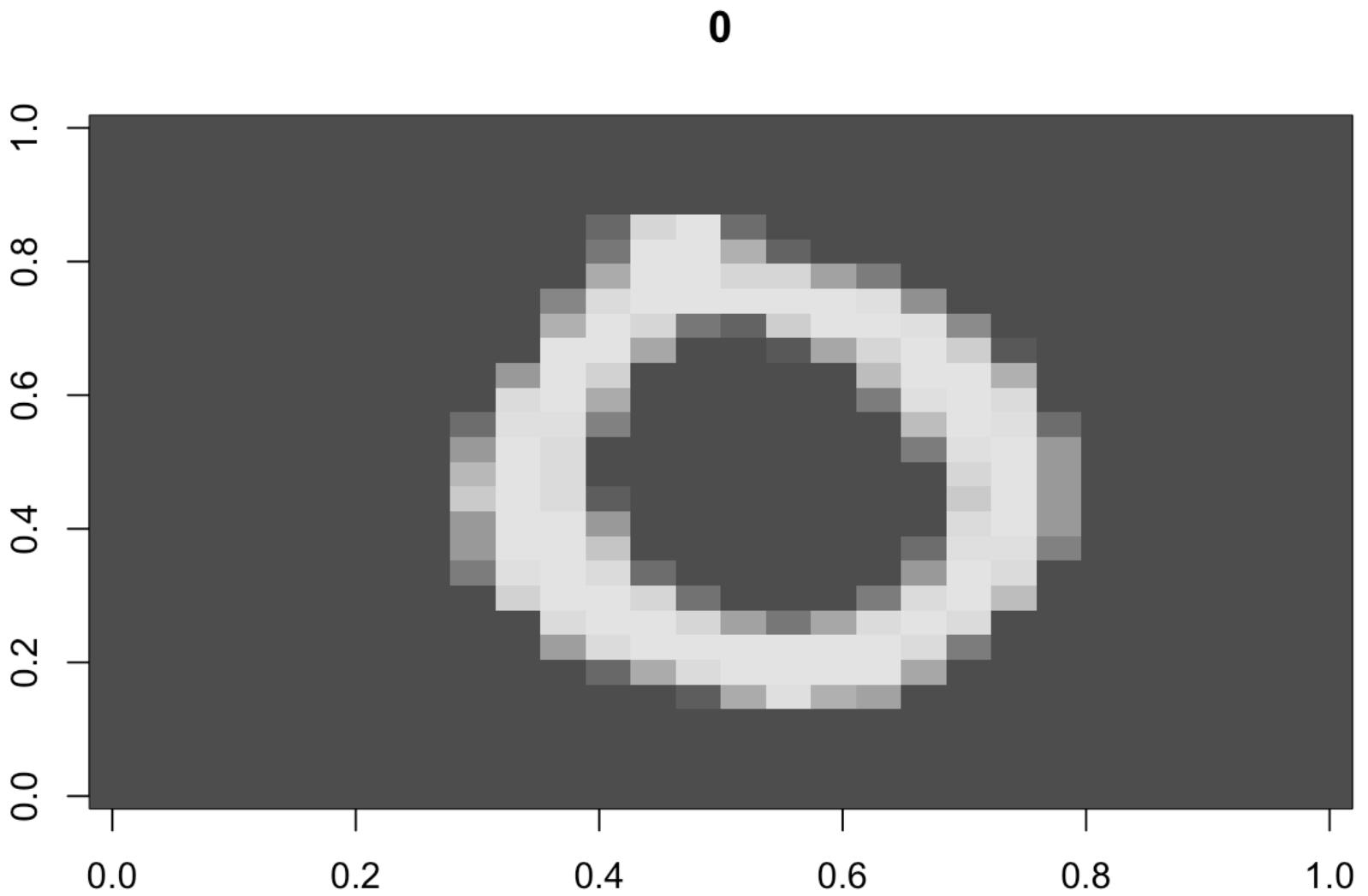
4

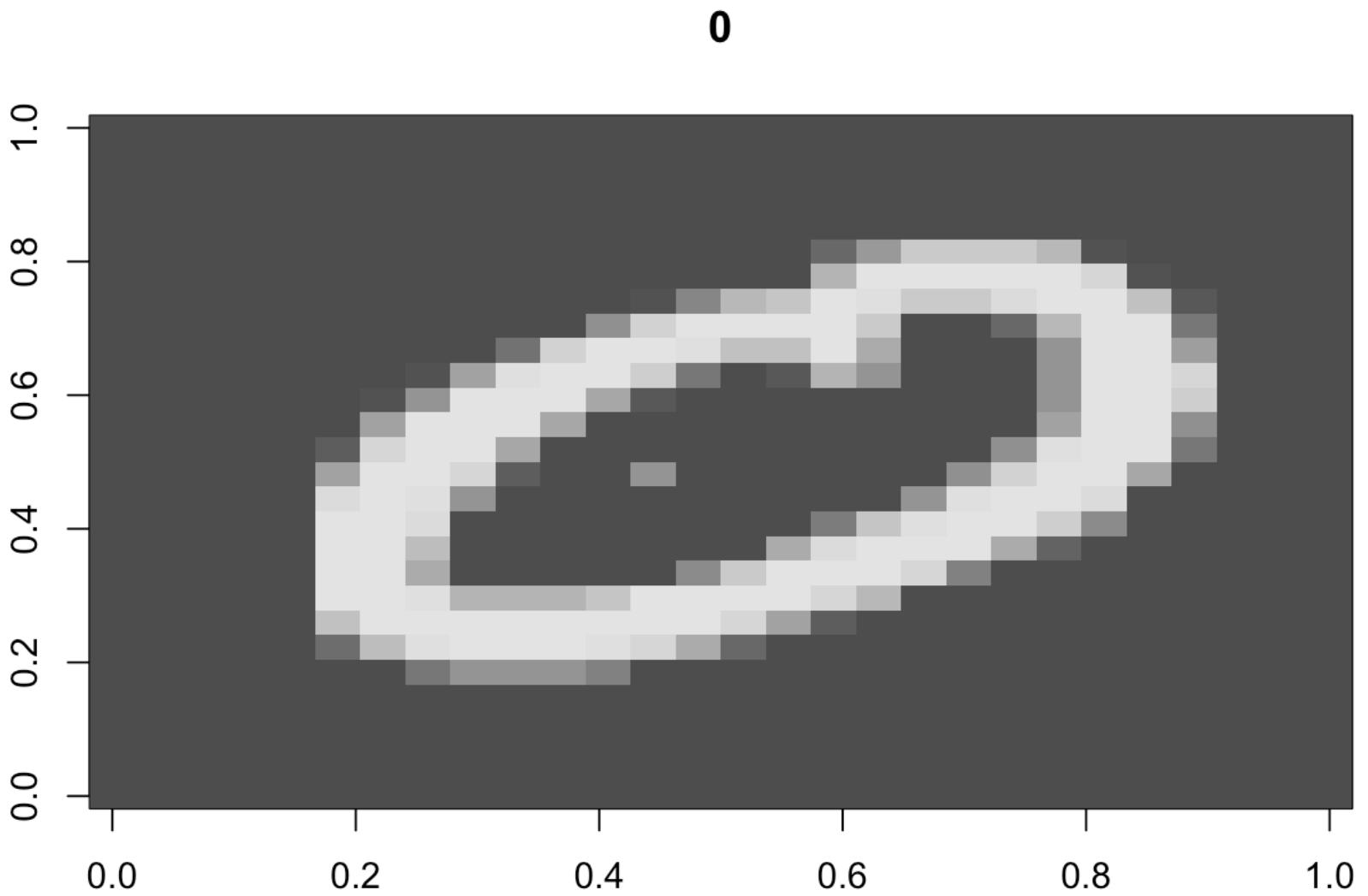


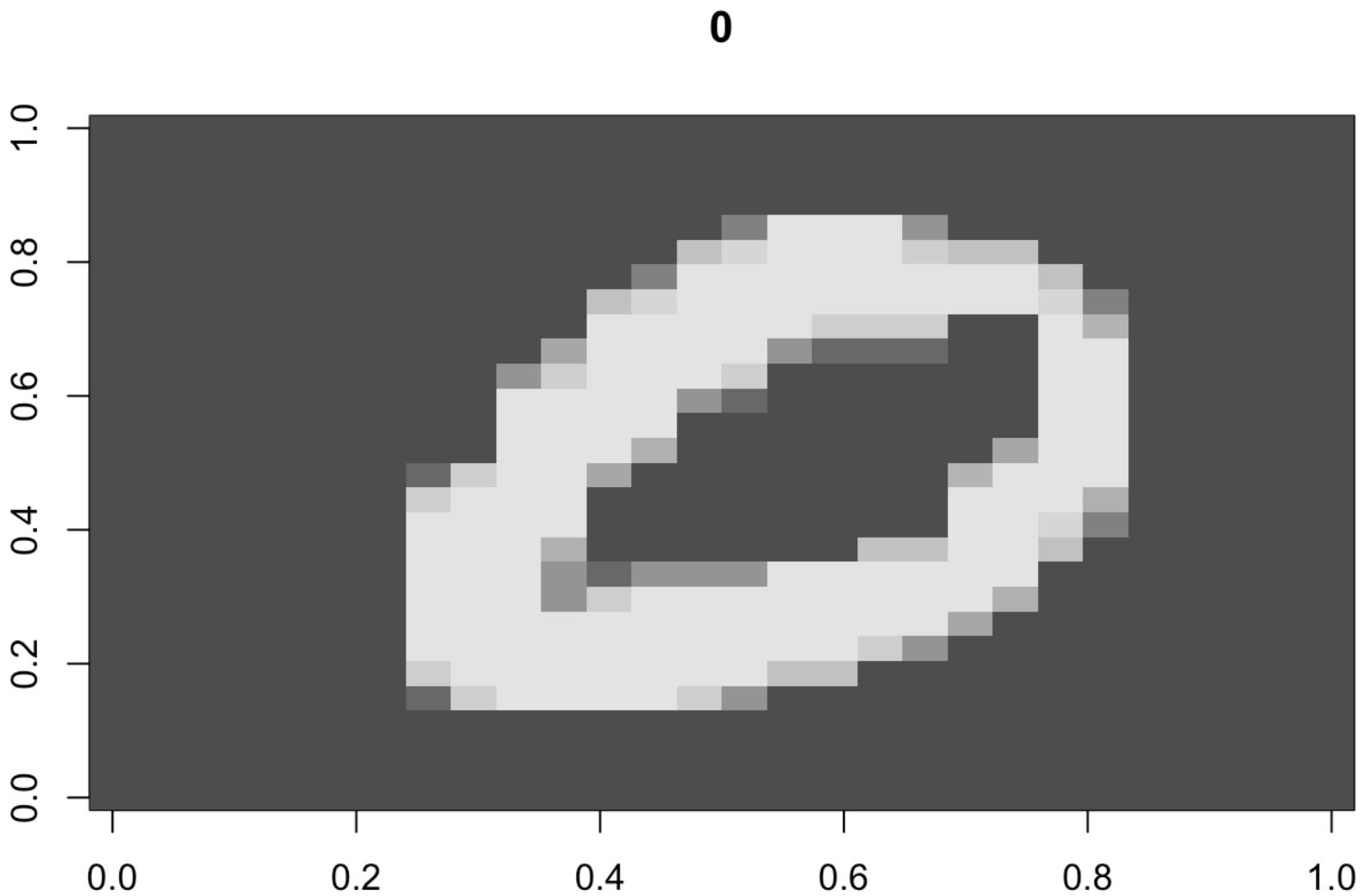
1



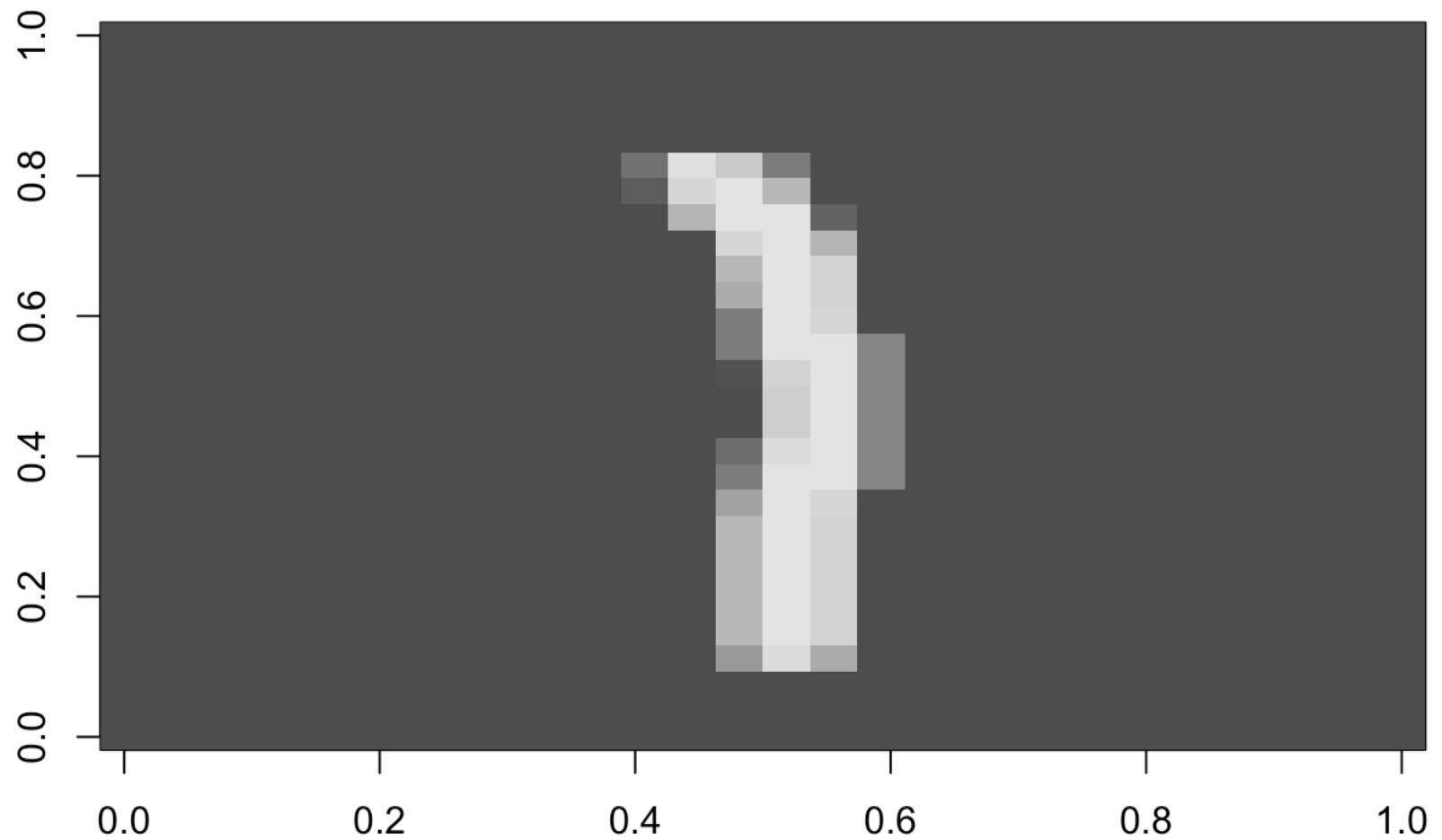




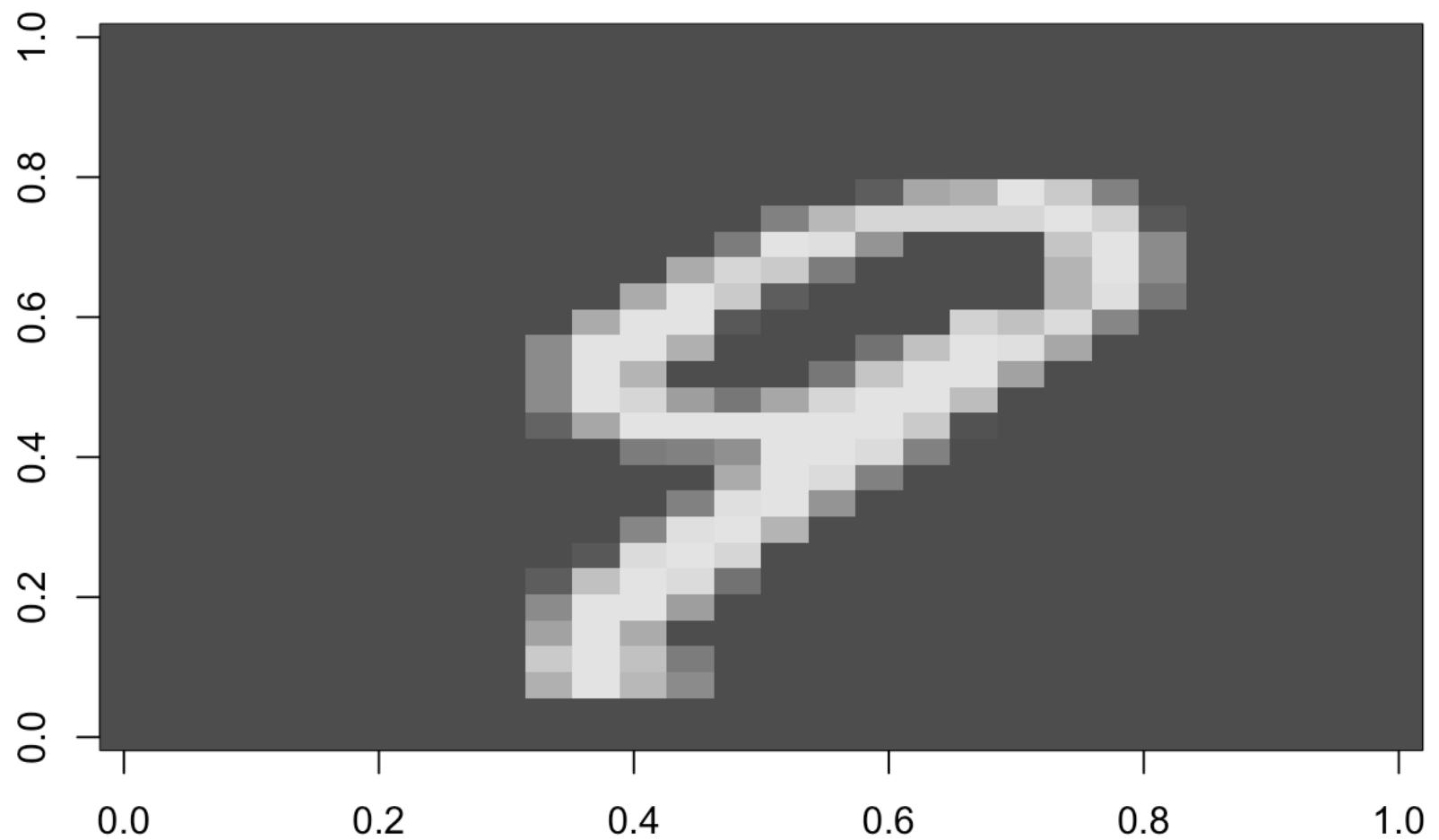


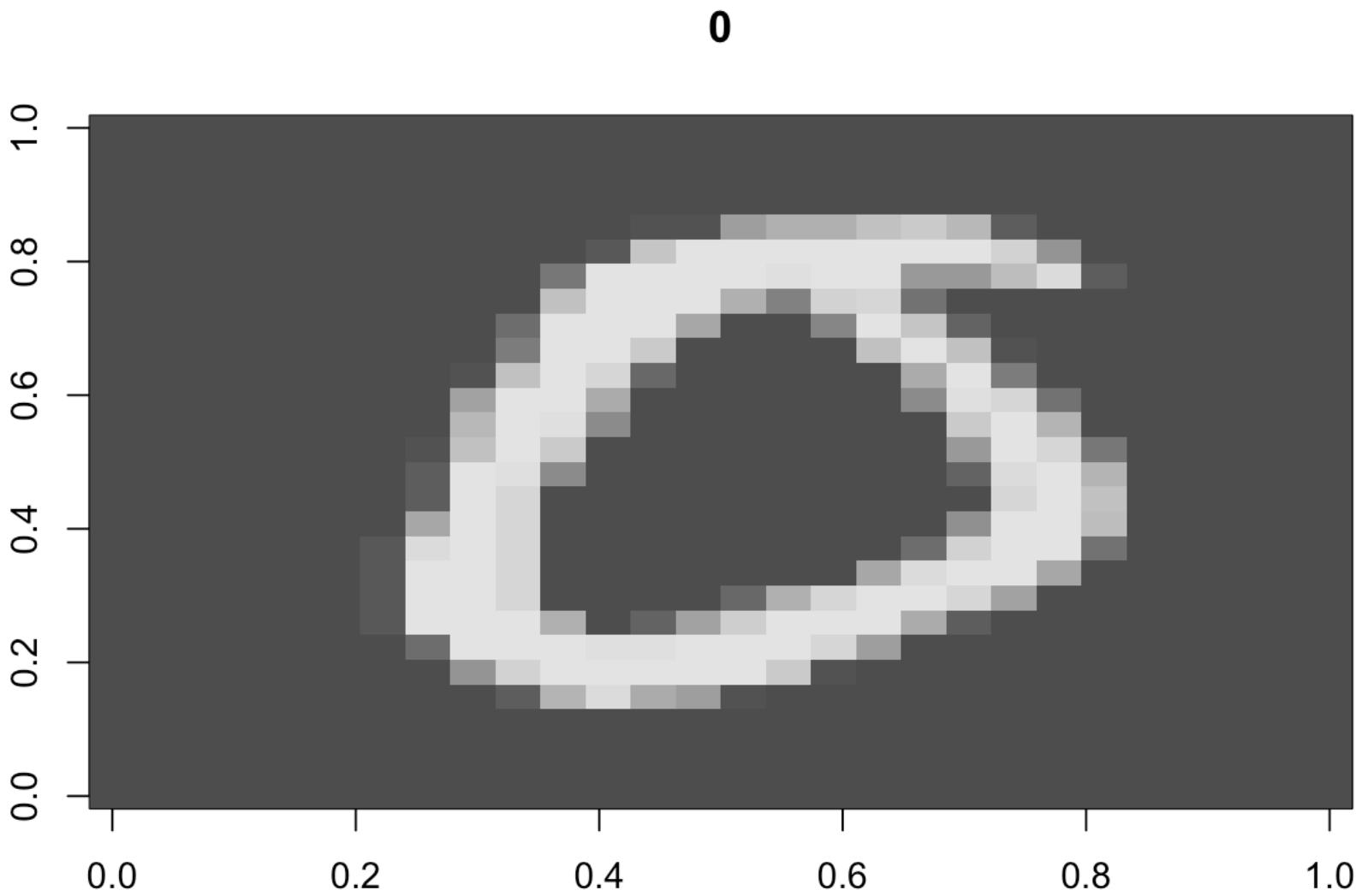


1

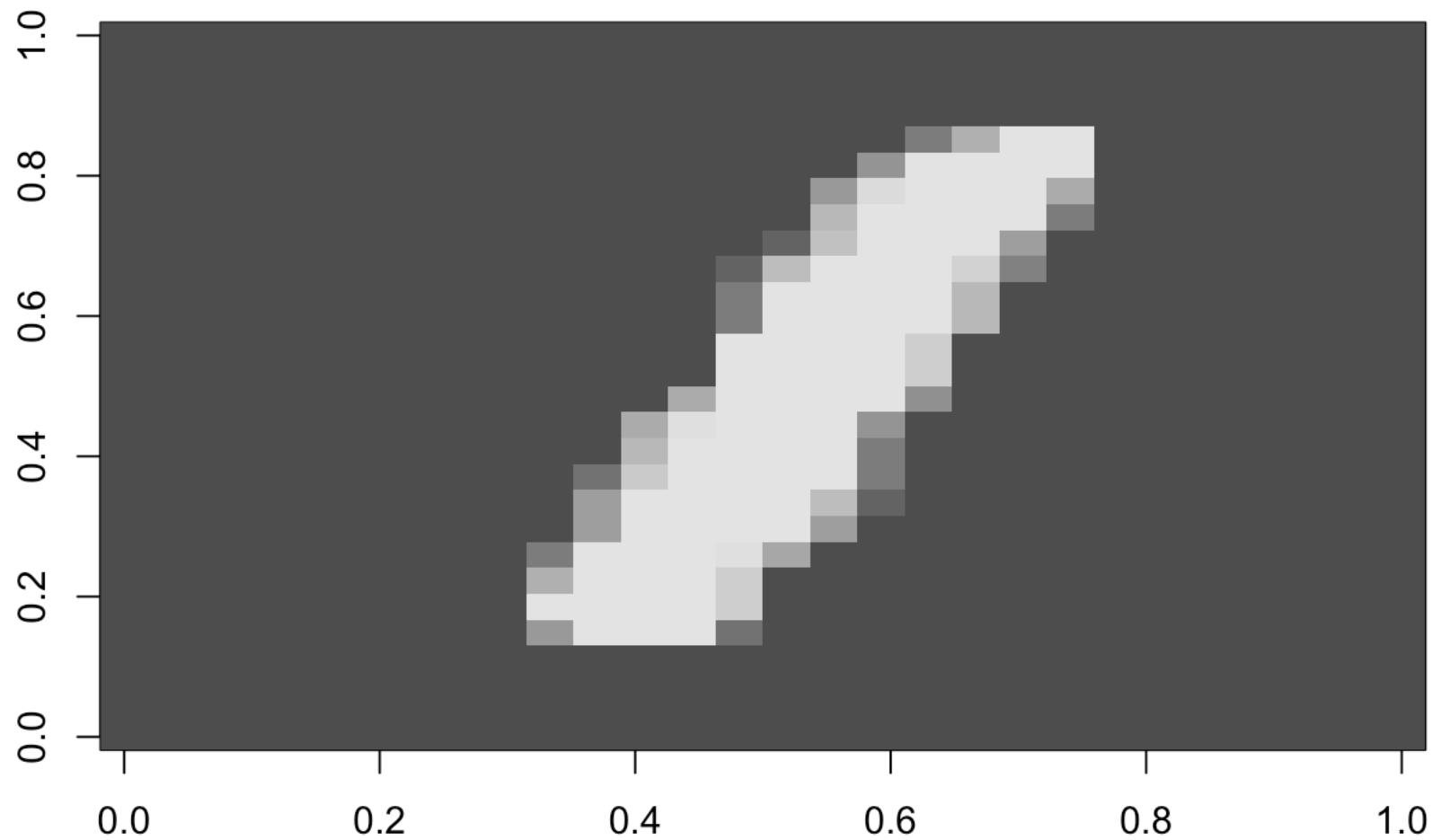


9

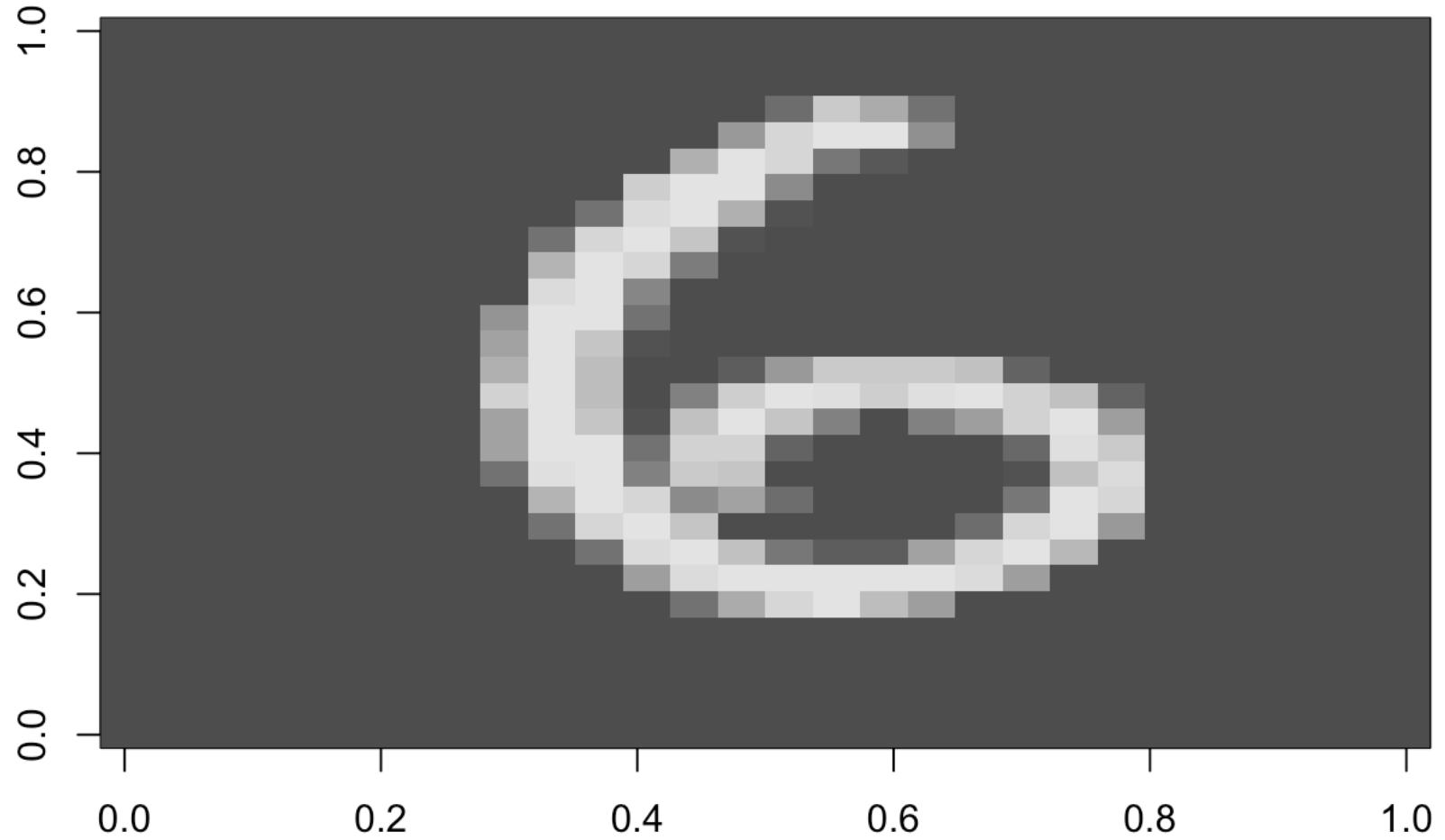




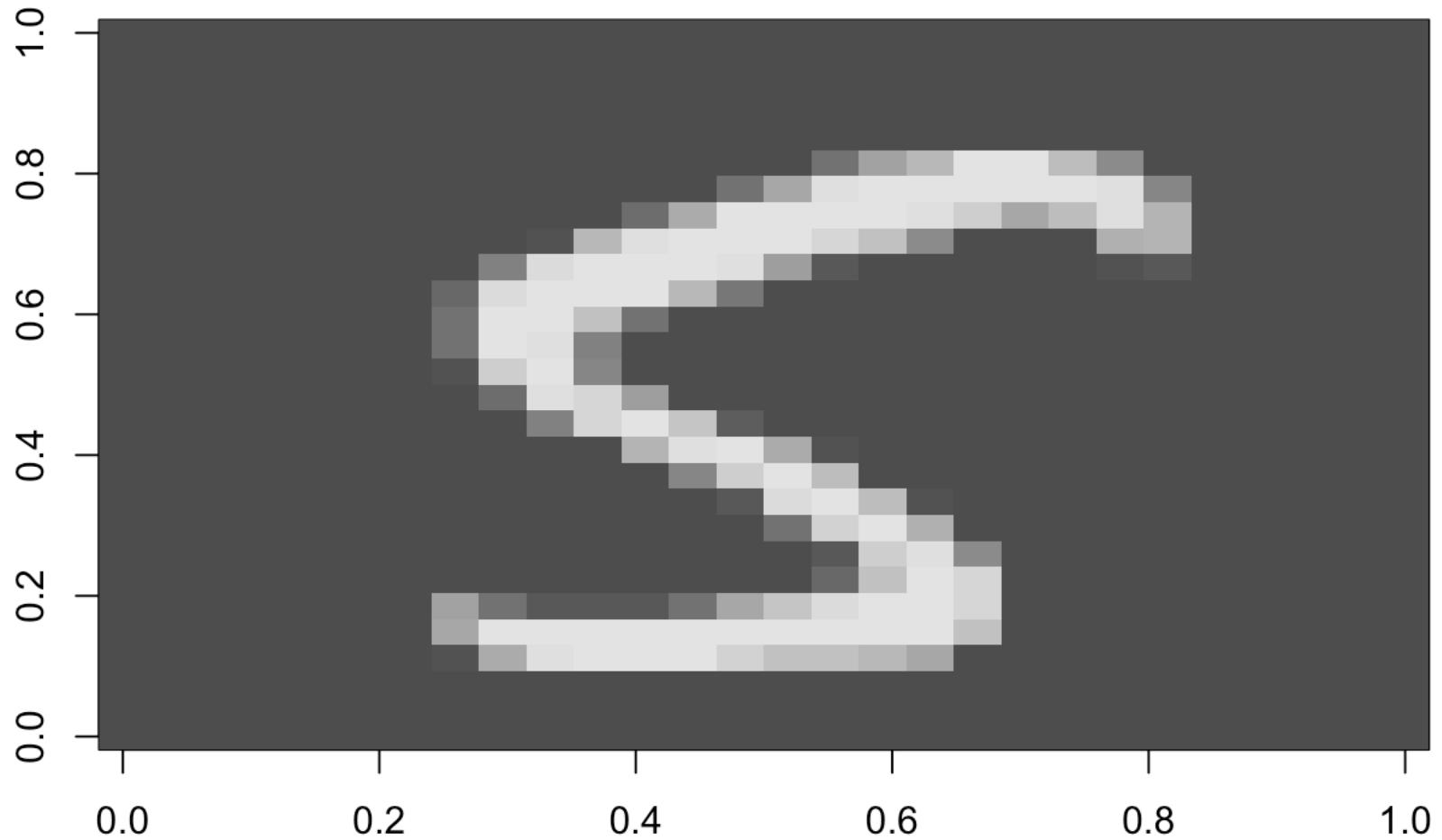
1



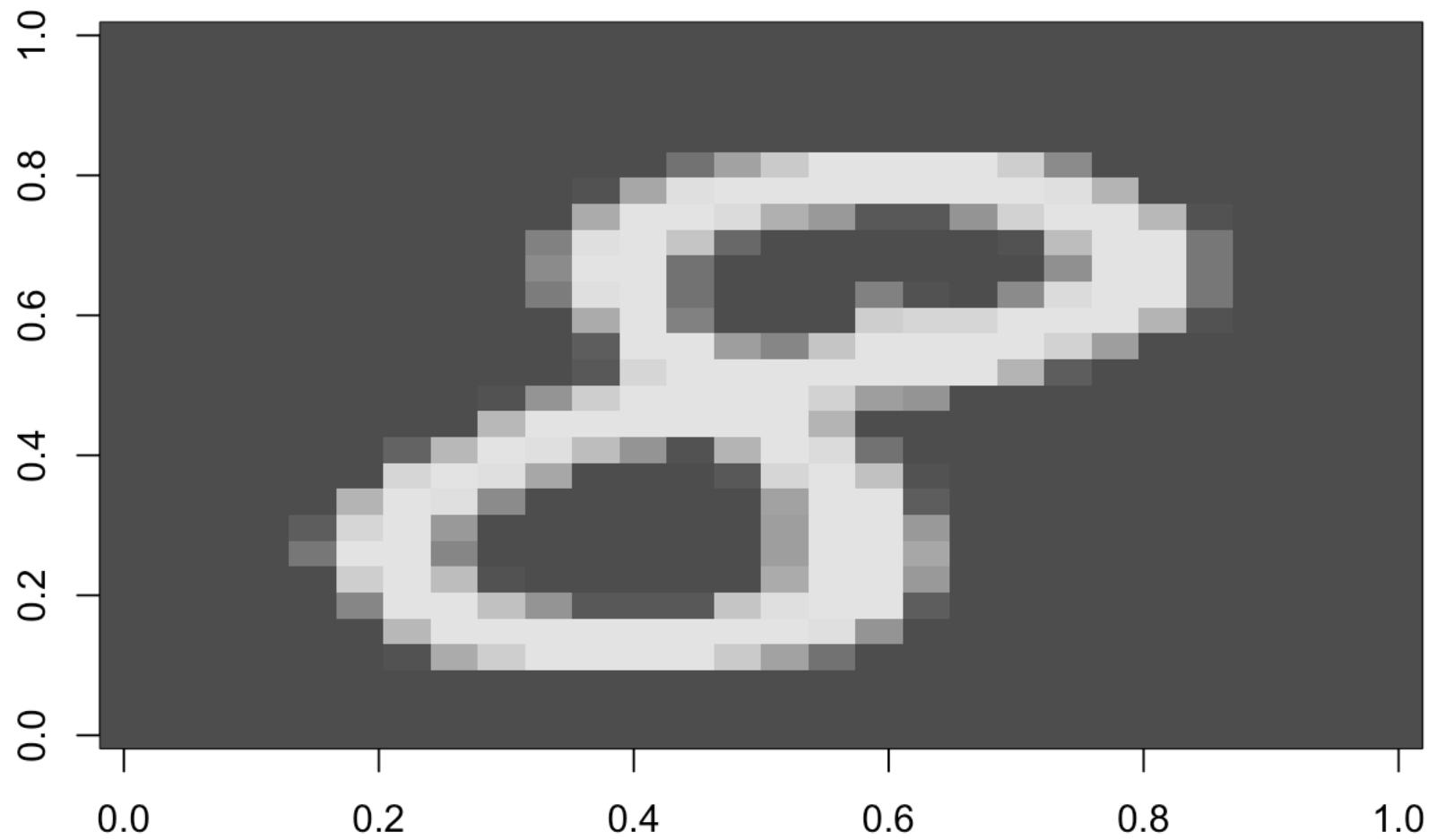
6



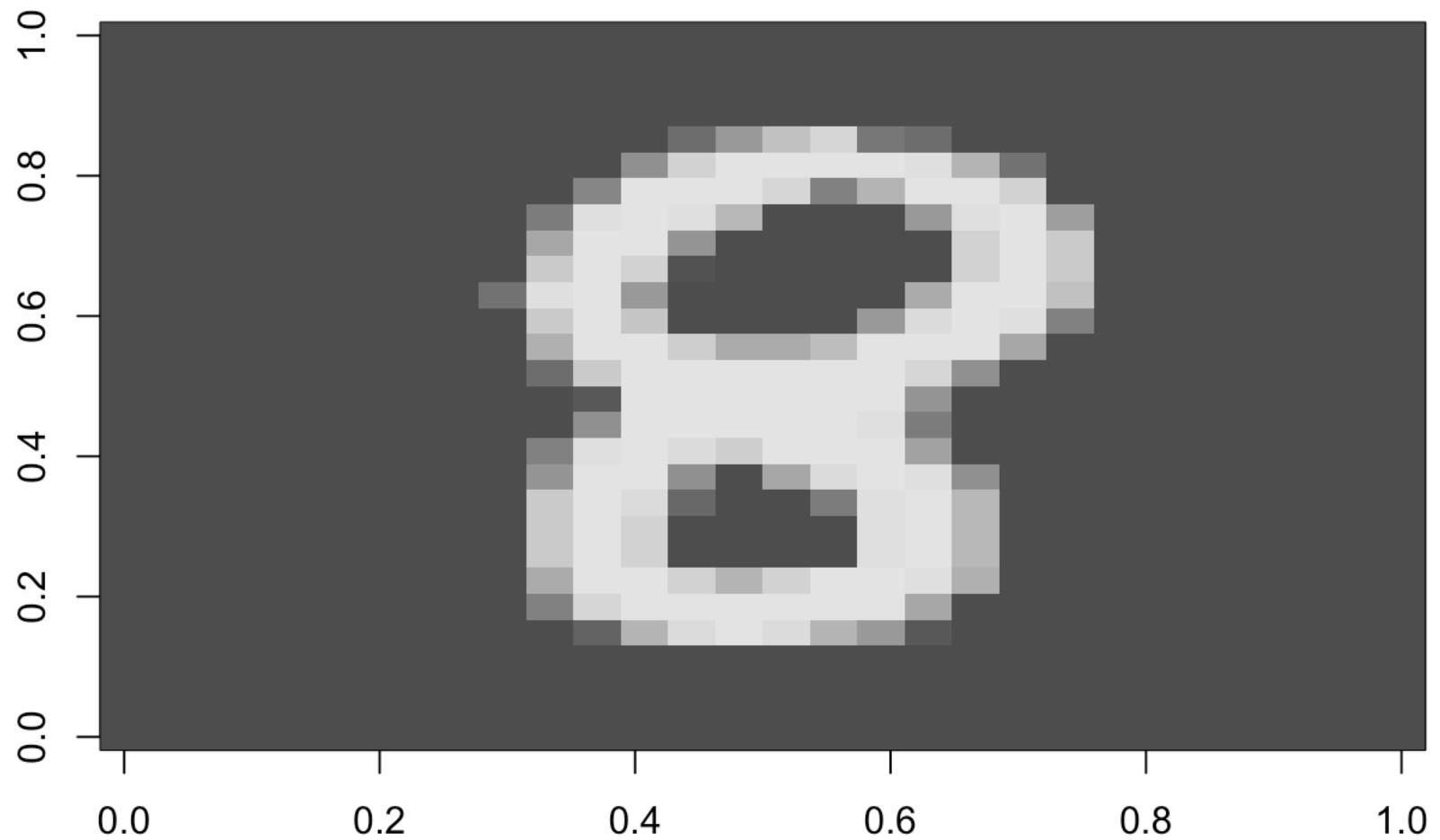
5



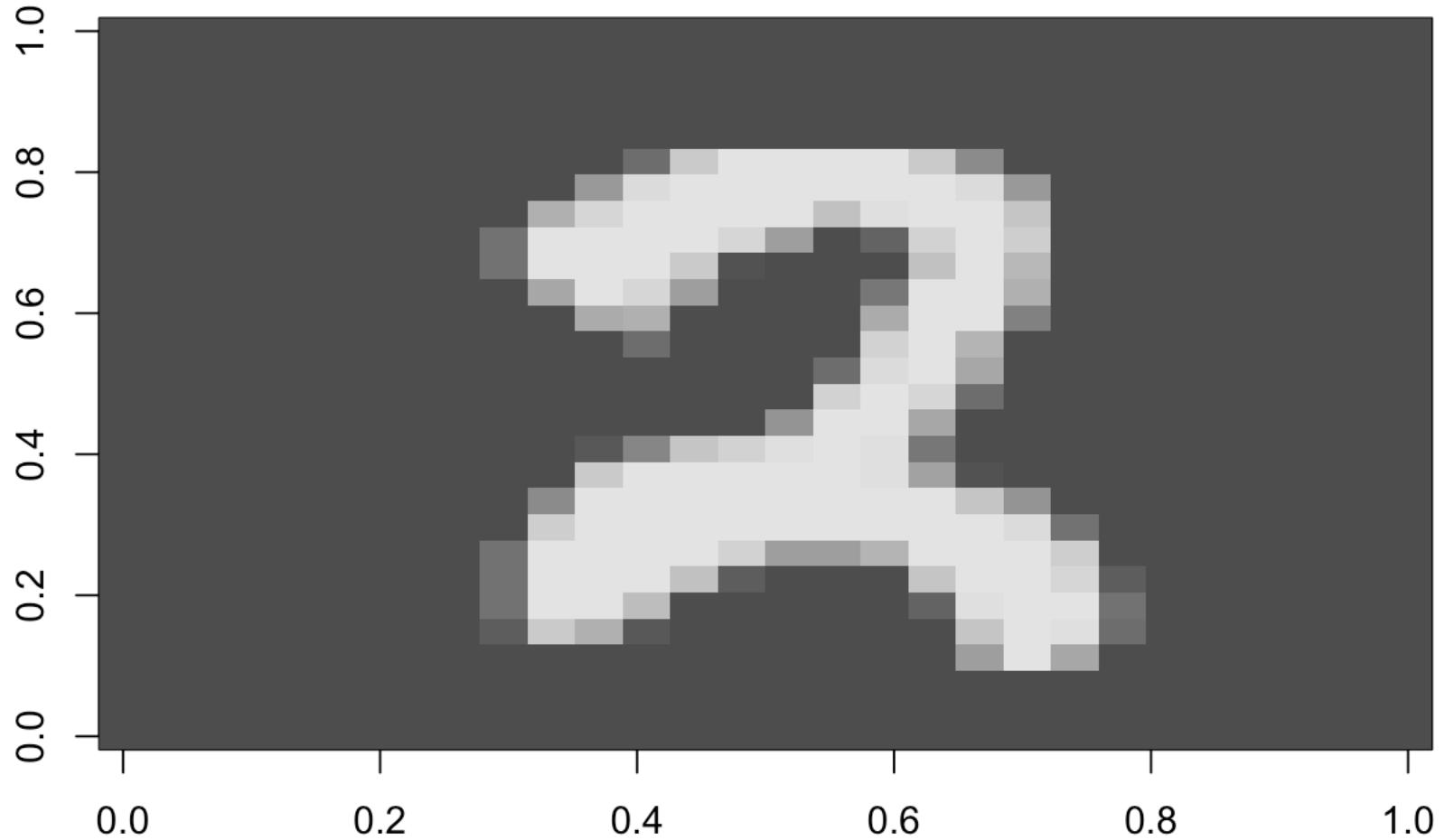
8



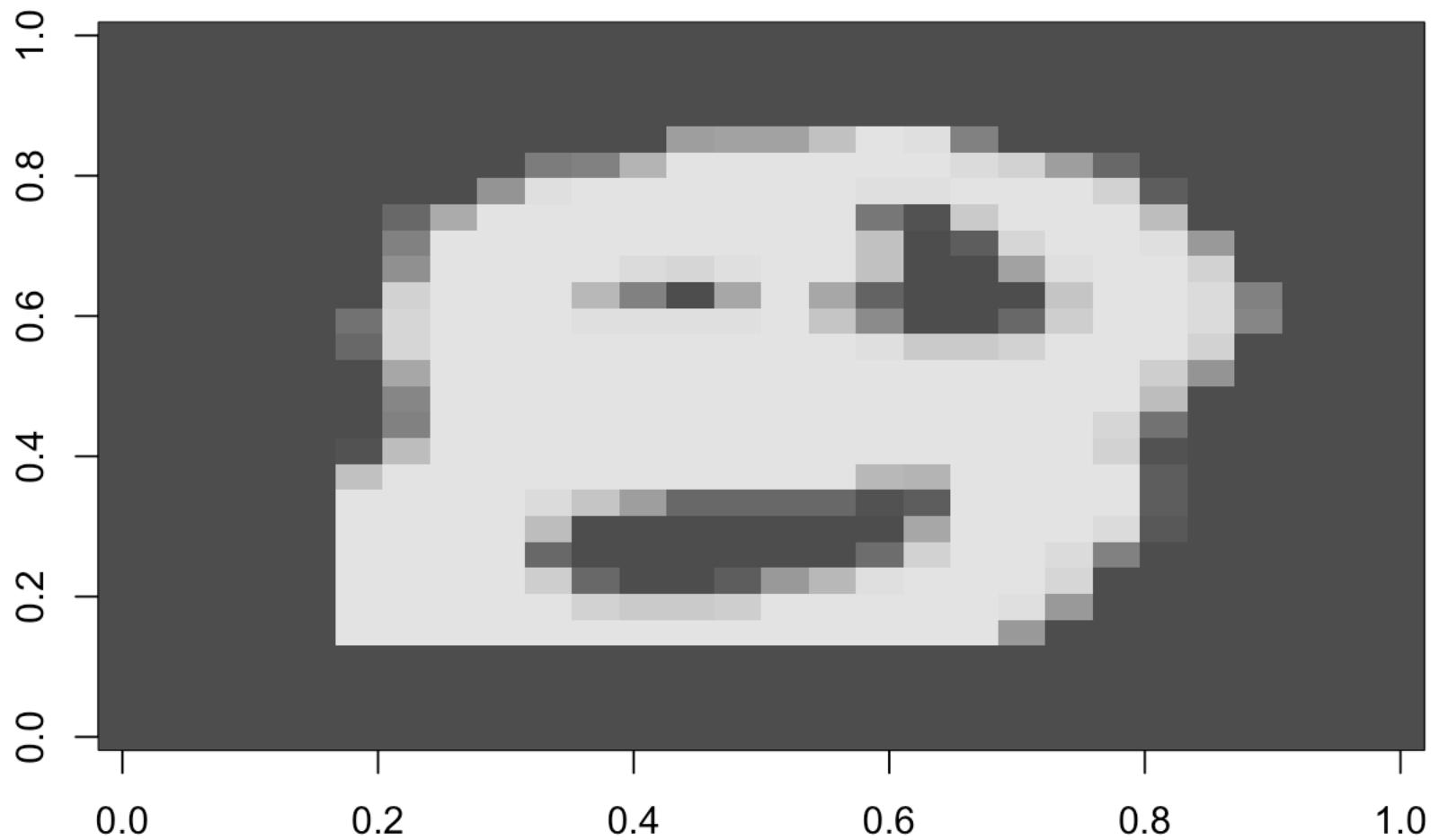
8



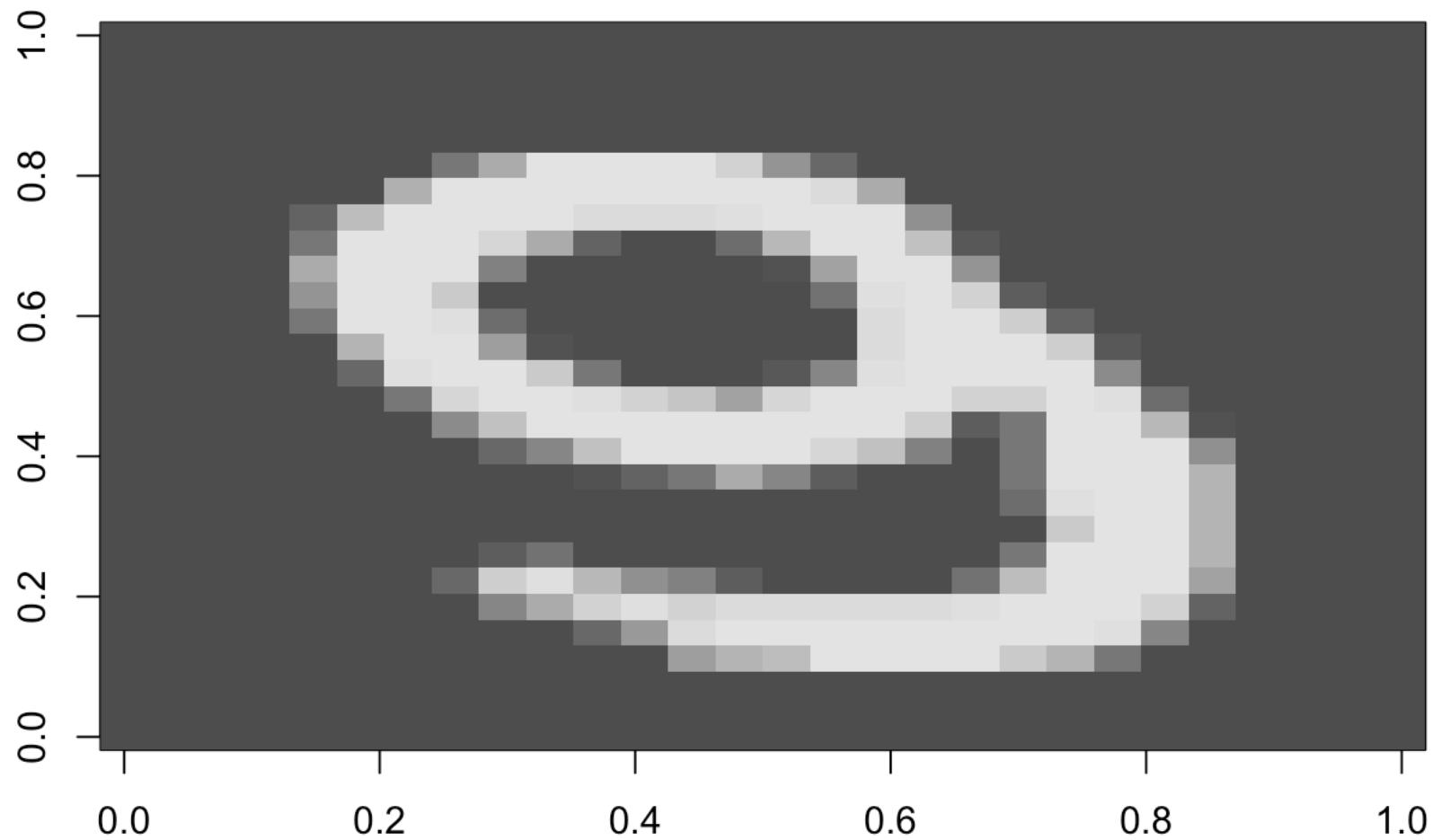
2



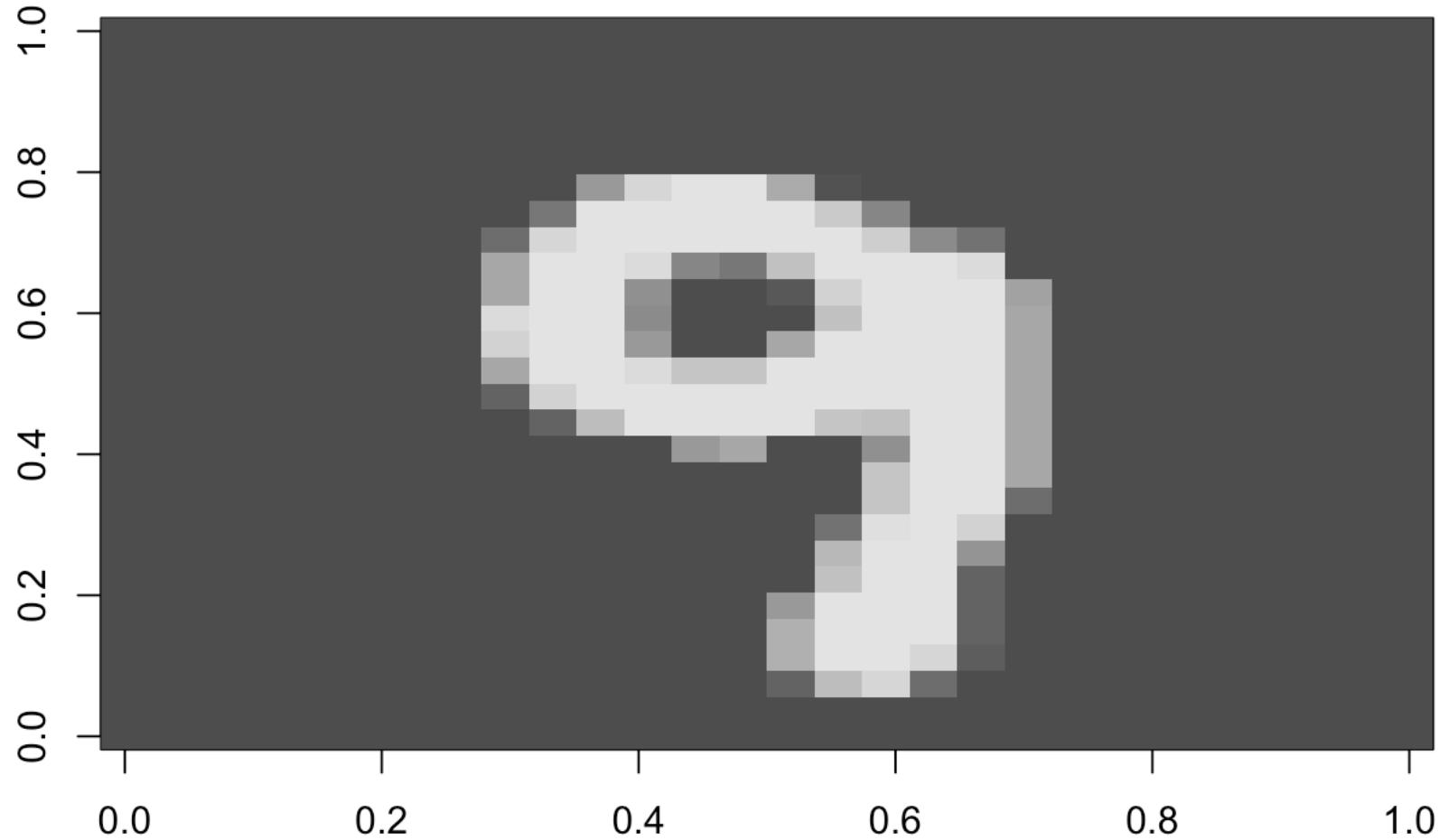
8



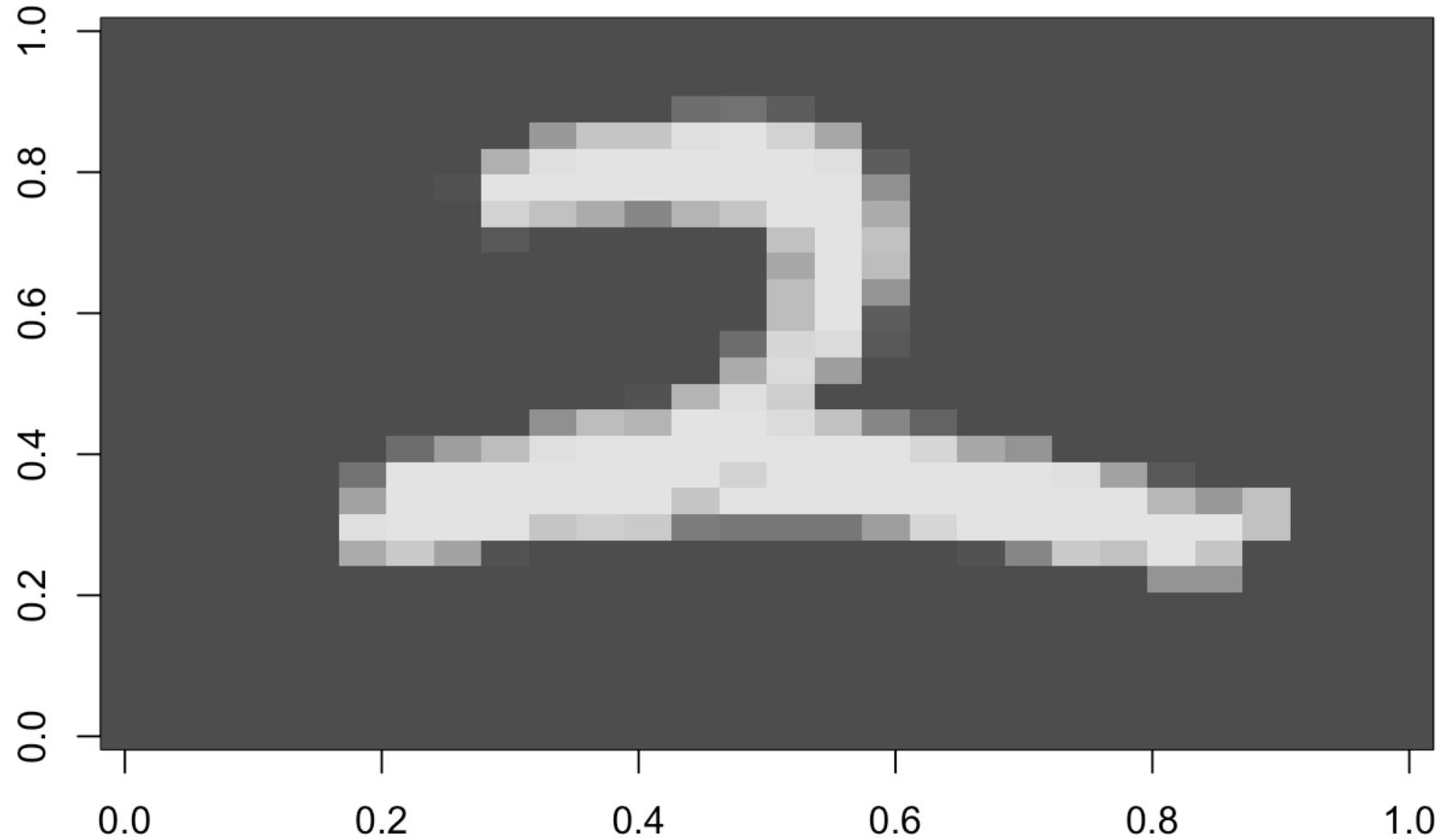
3



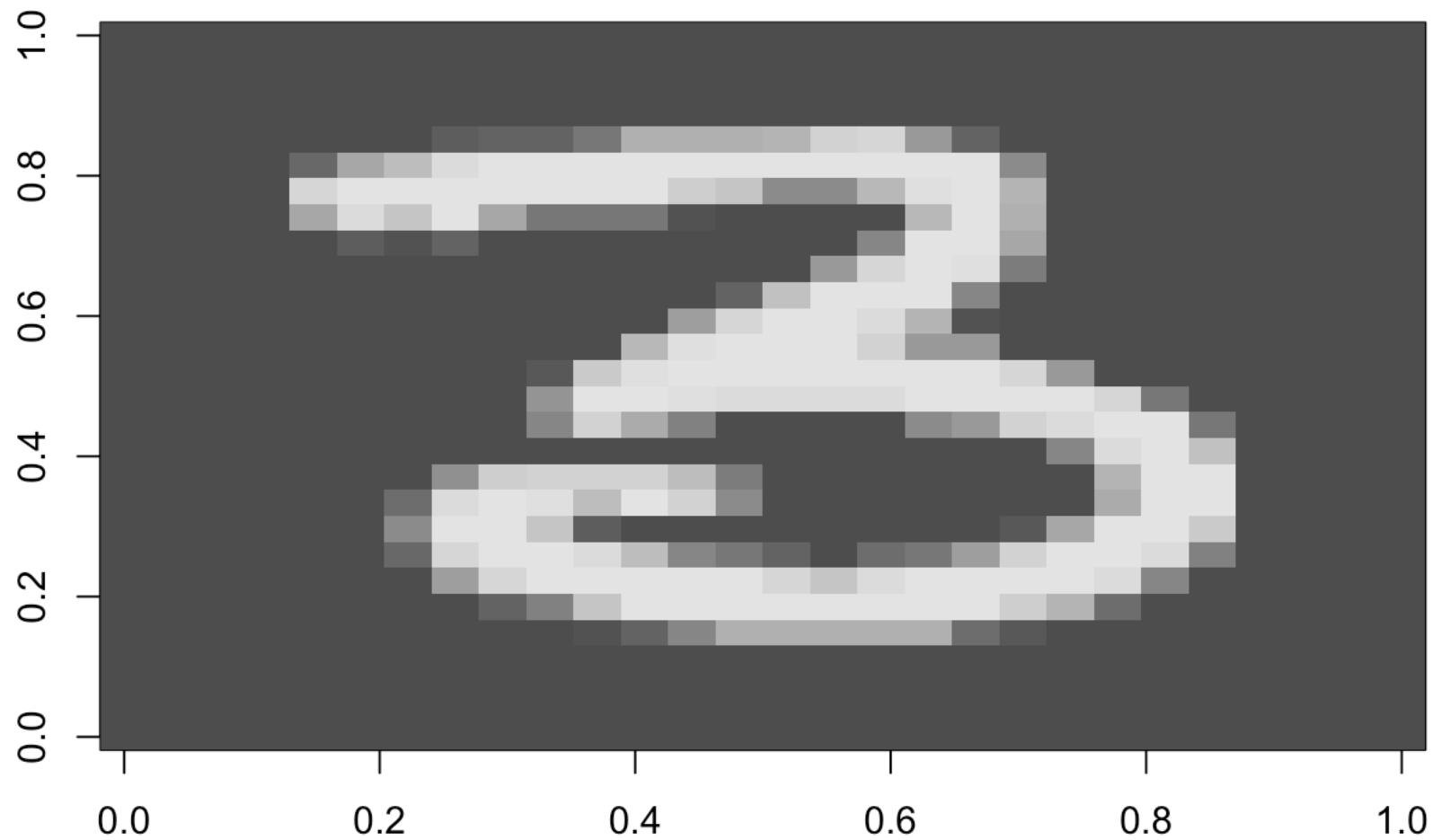
9



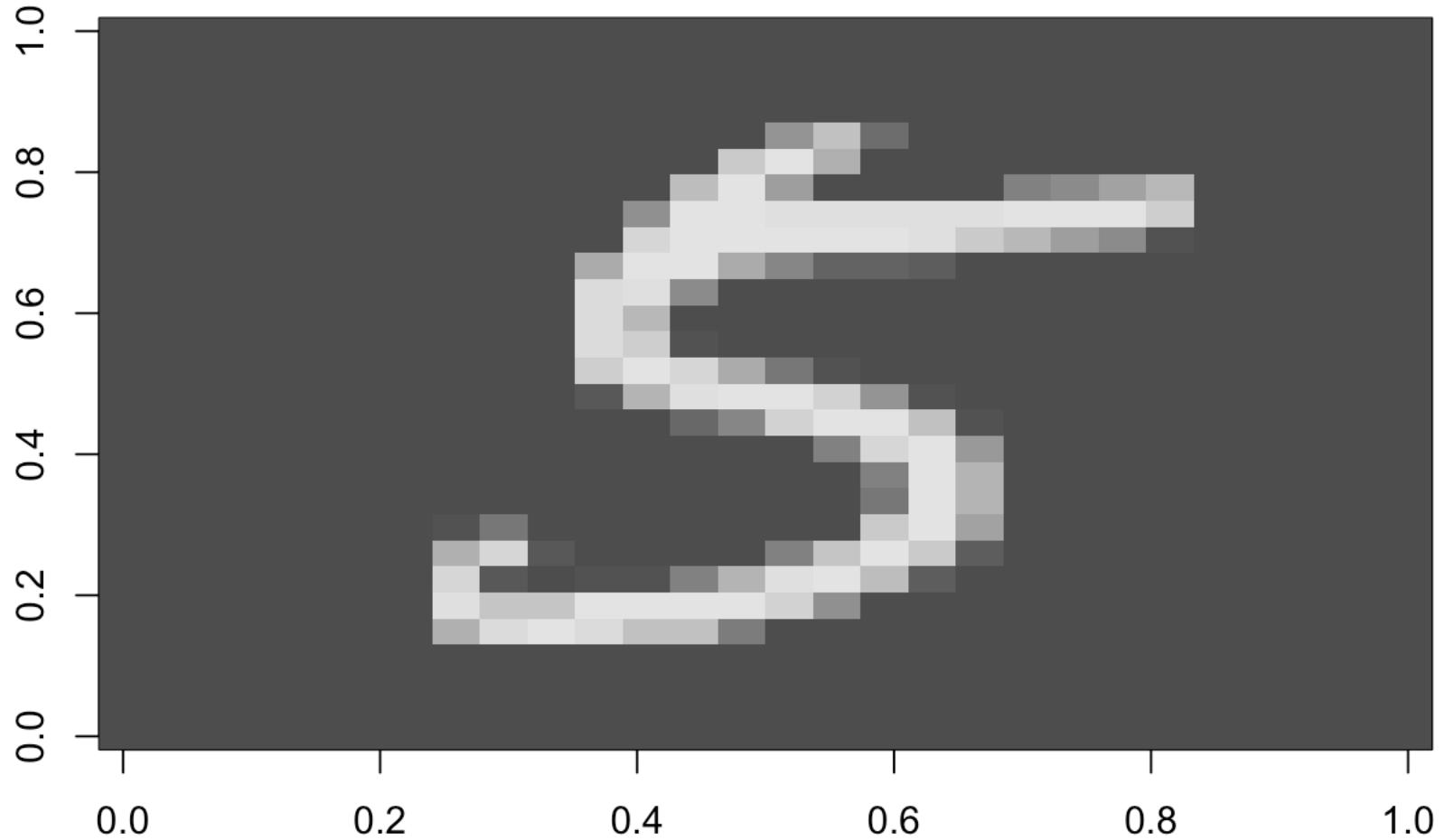
2



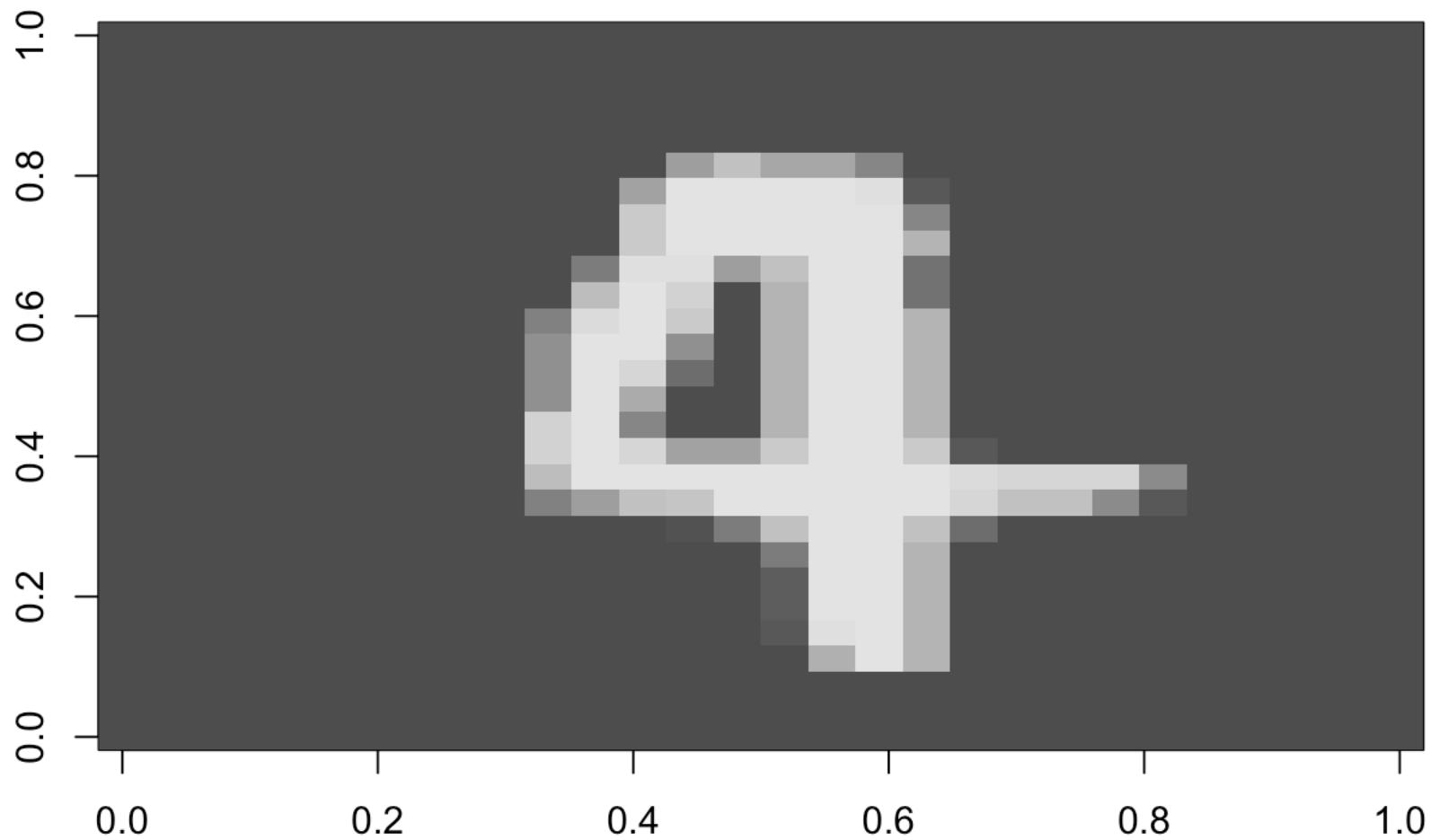
3



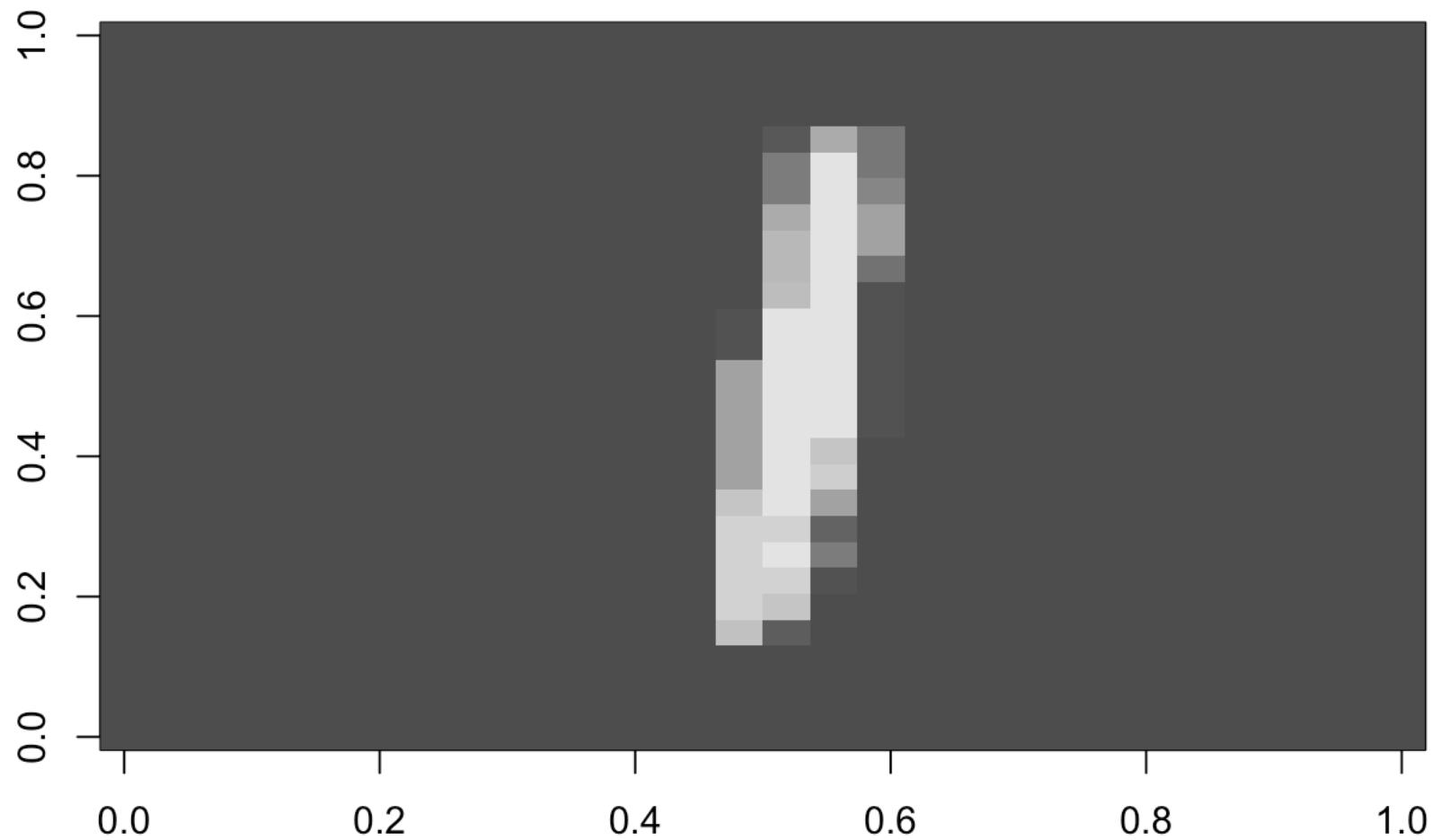
5

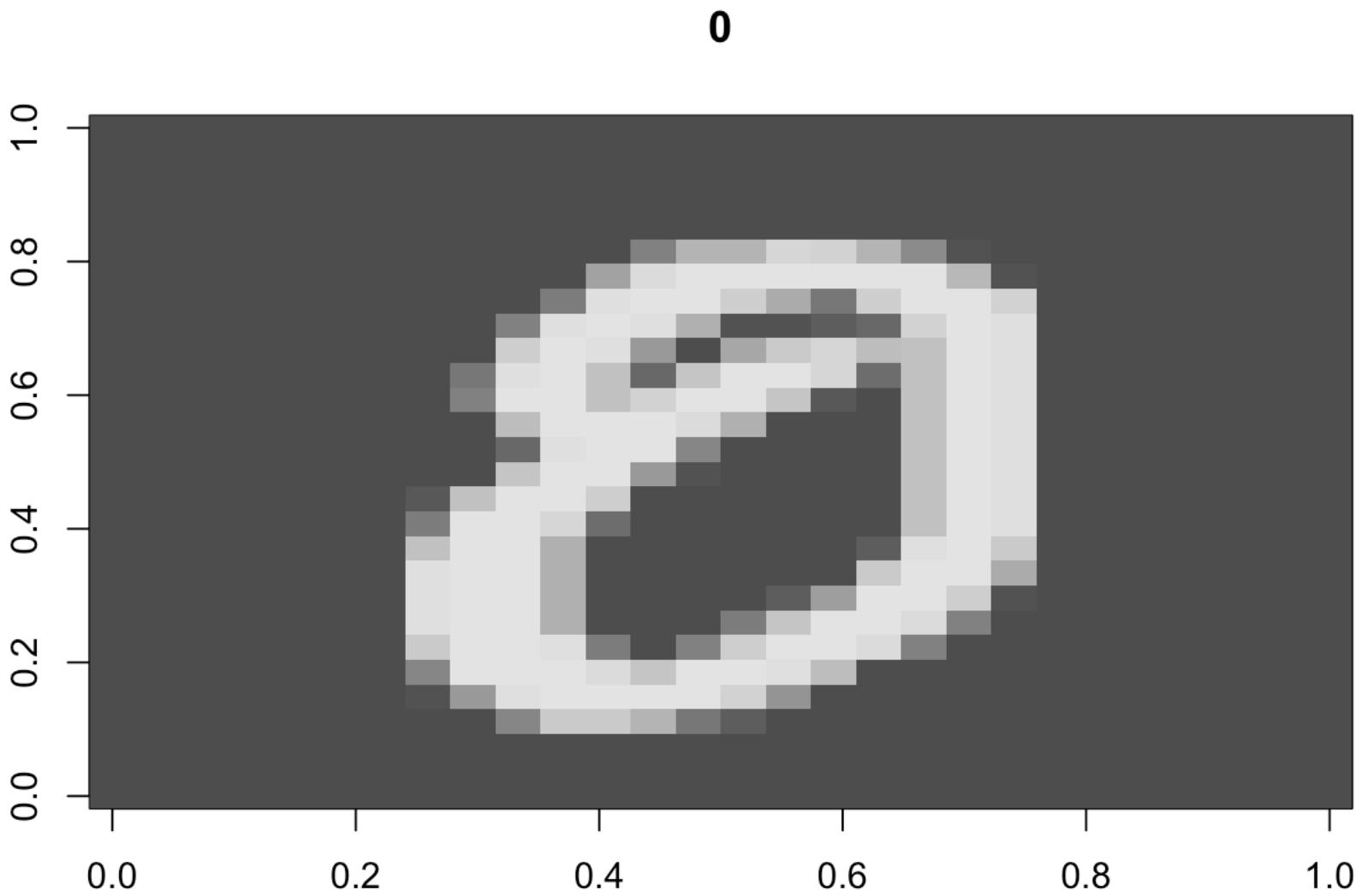


4

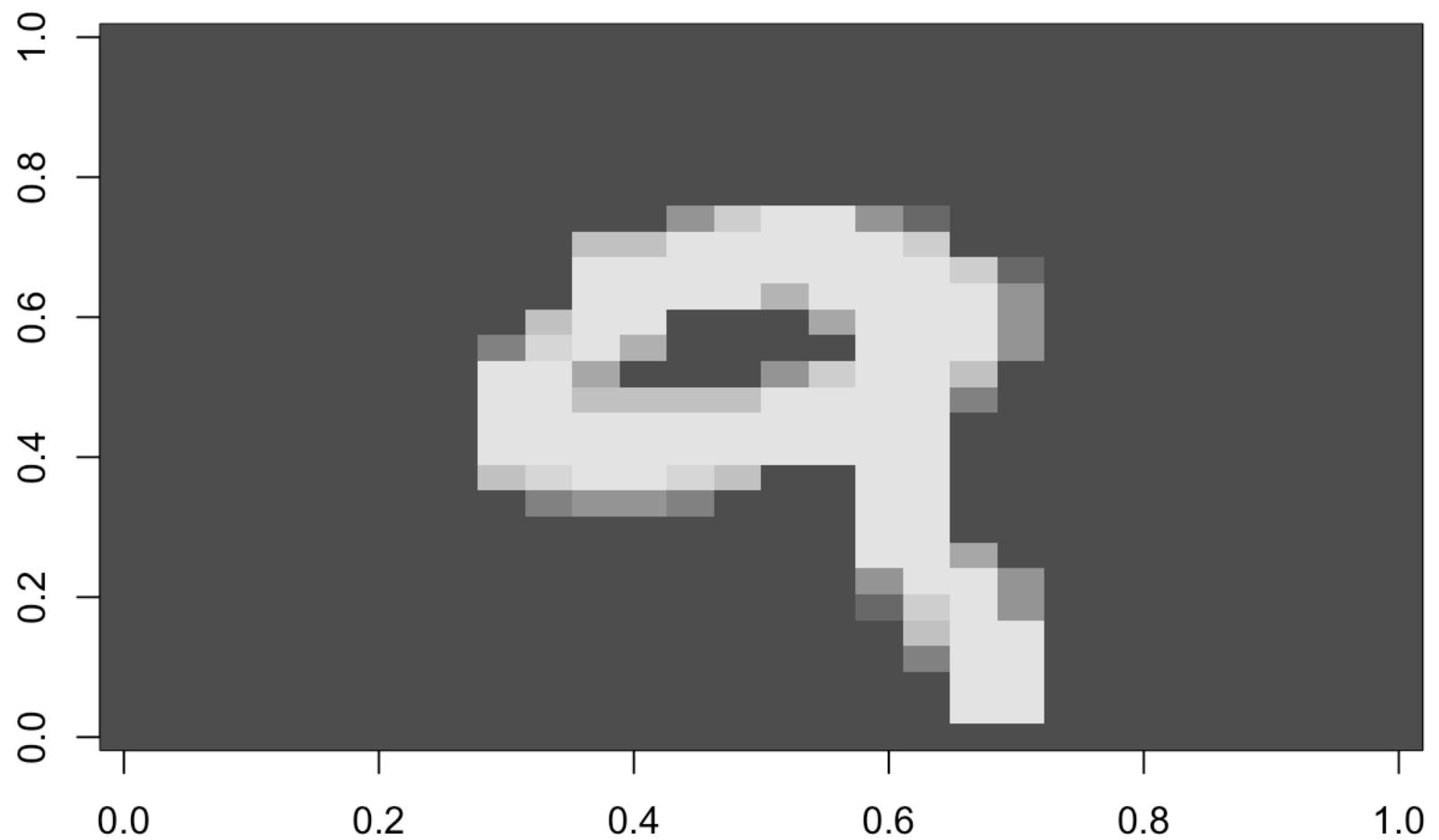


1

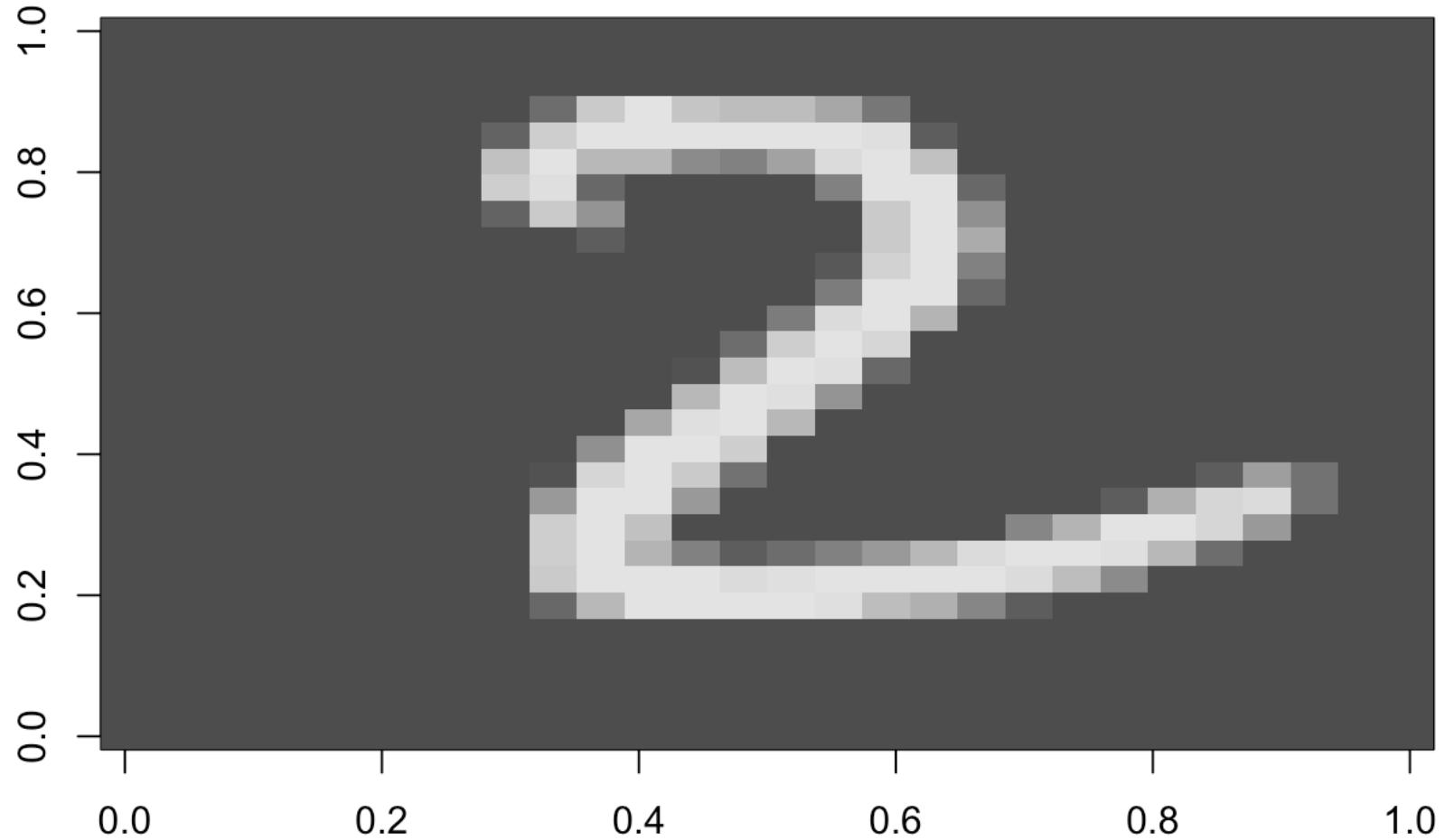




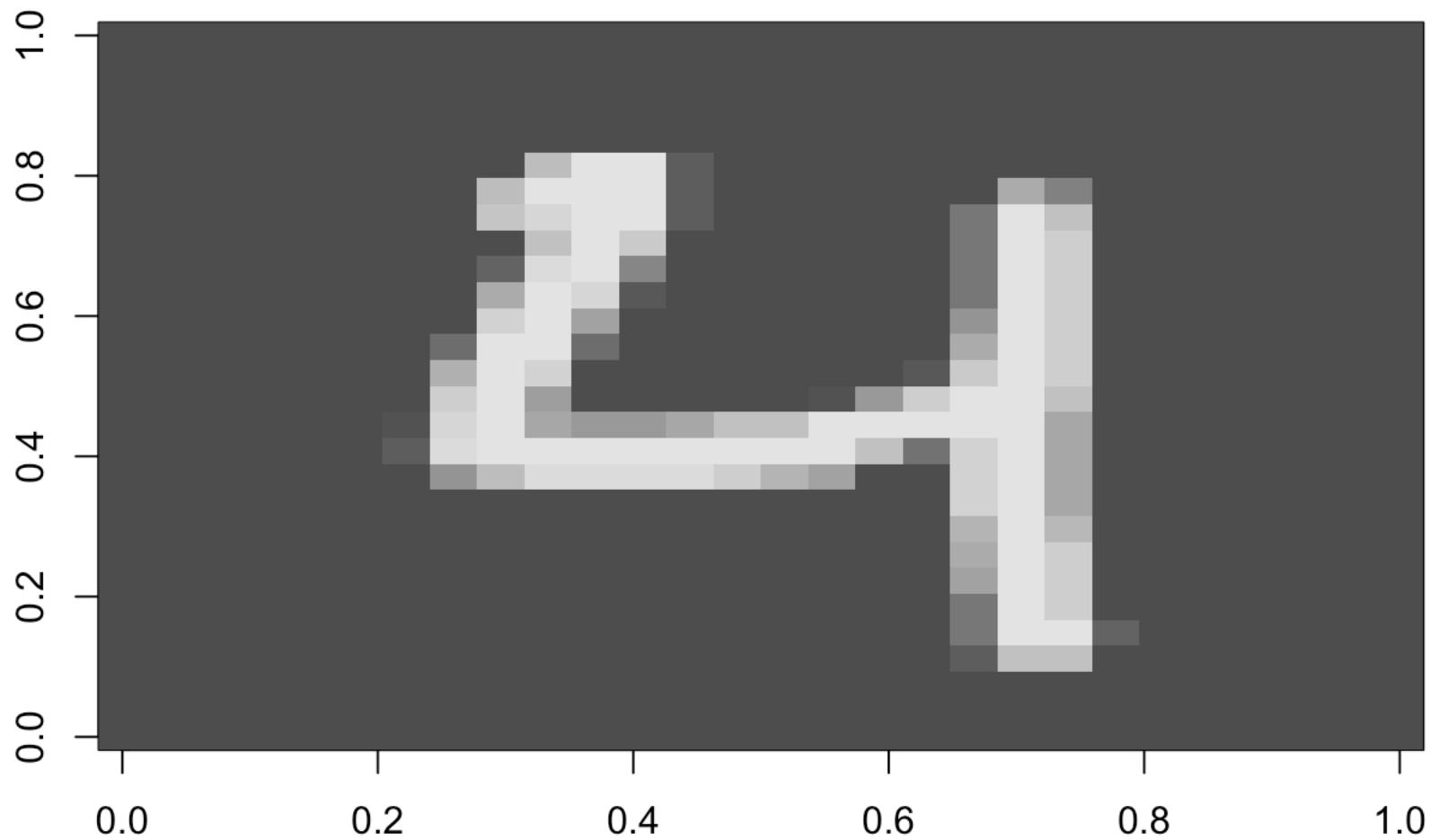
9



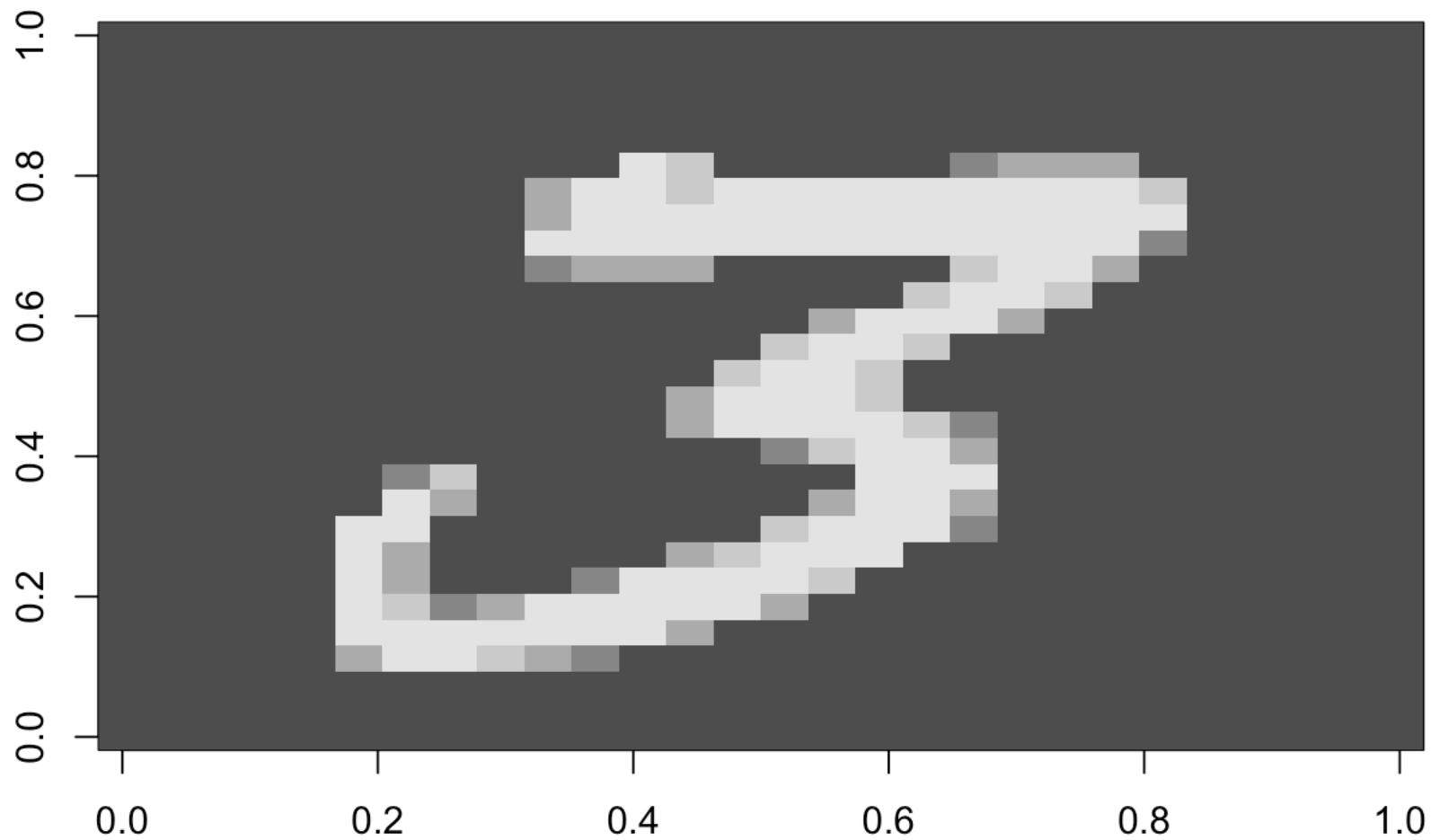
2



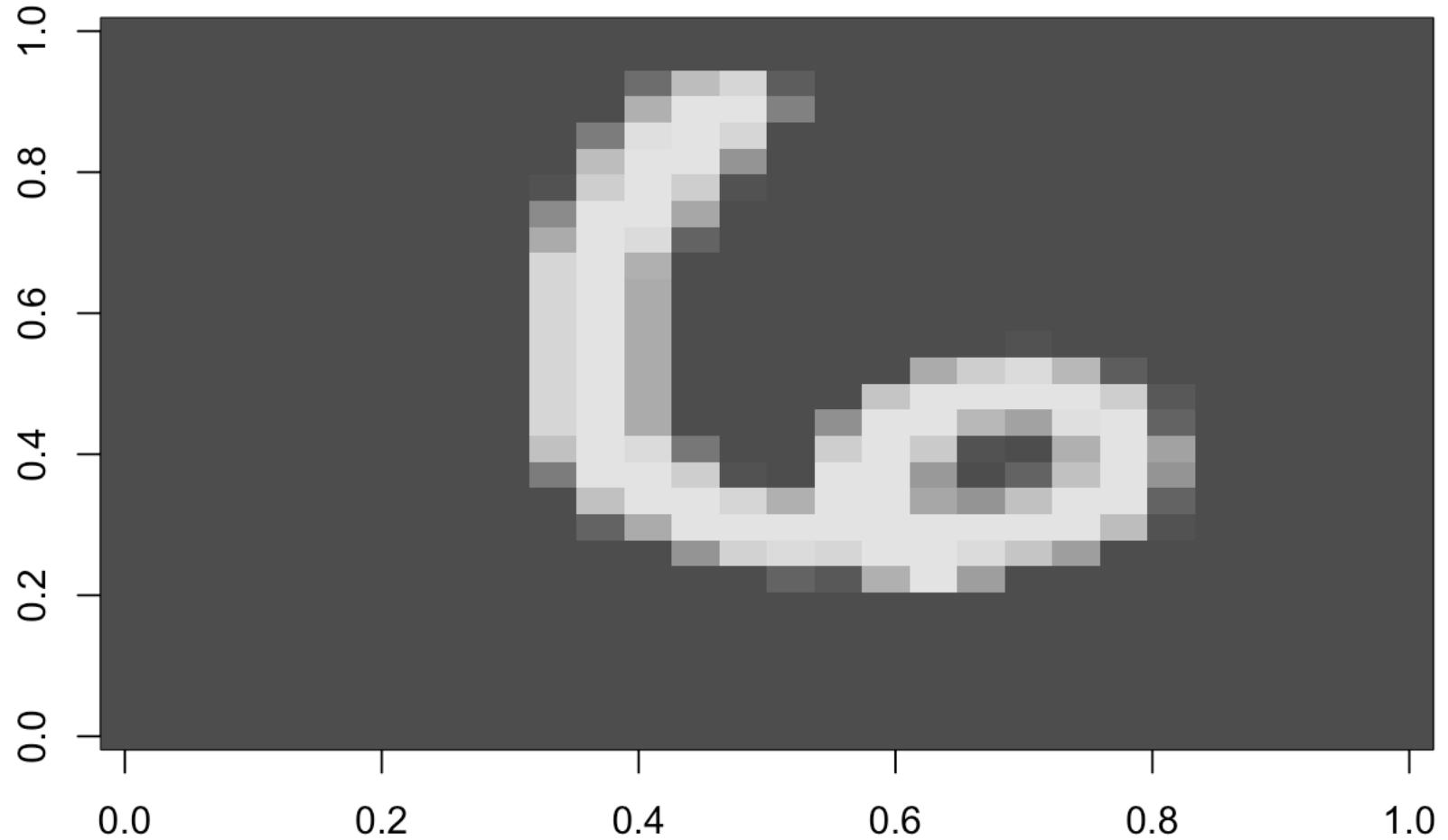
4

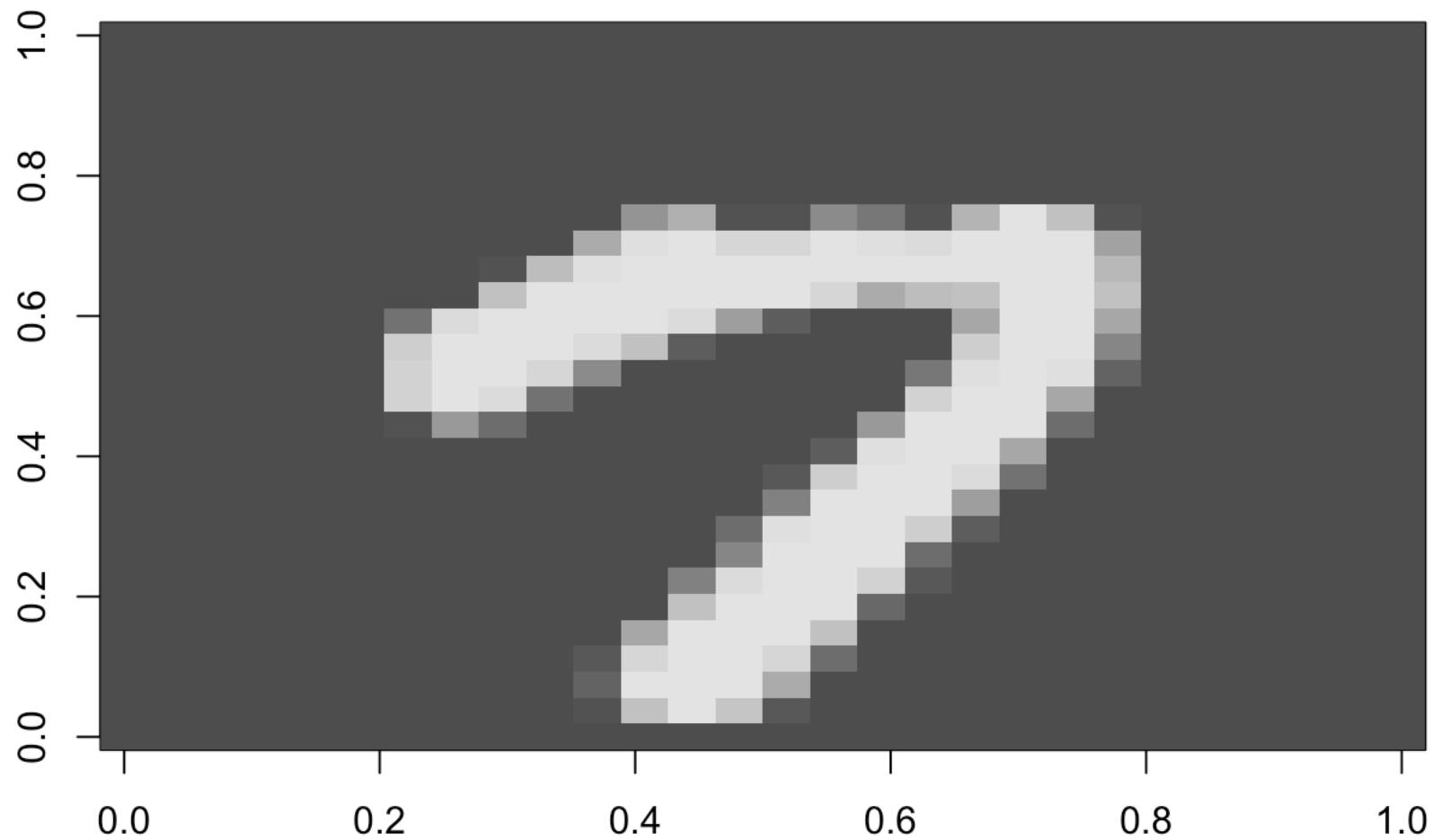


3

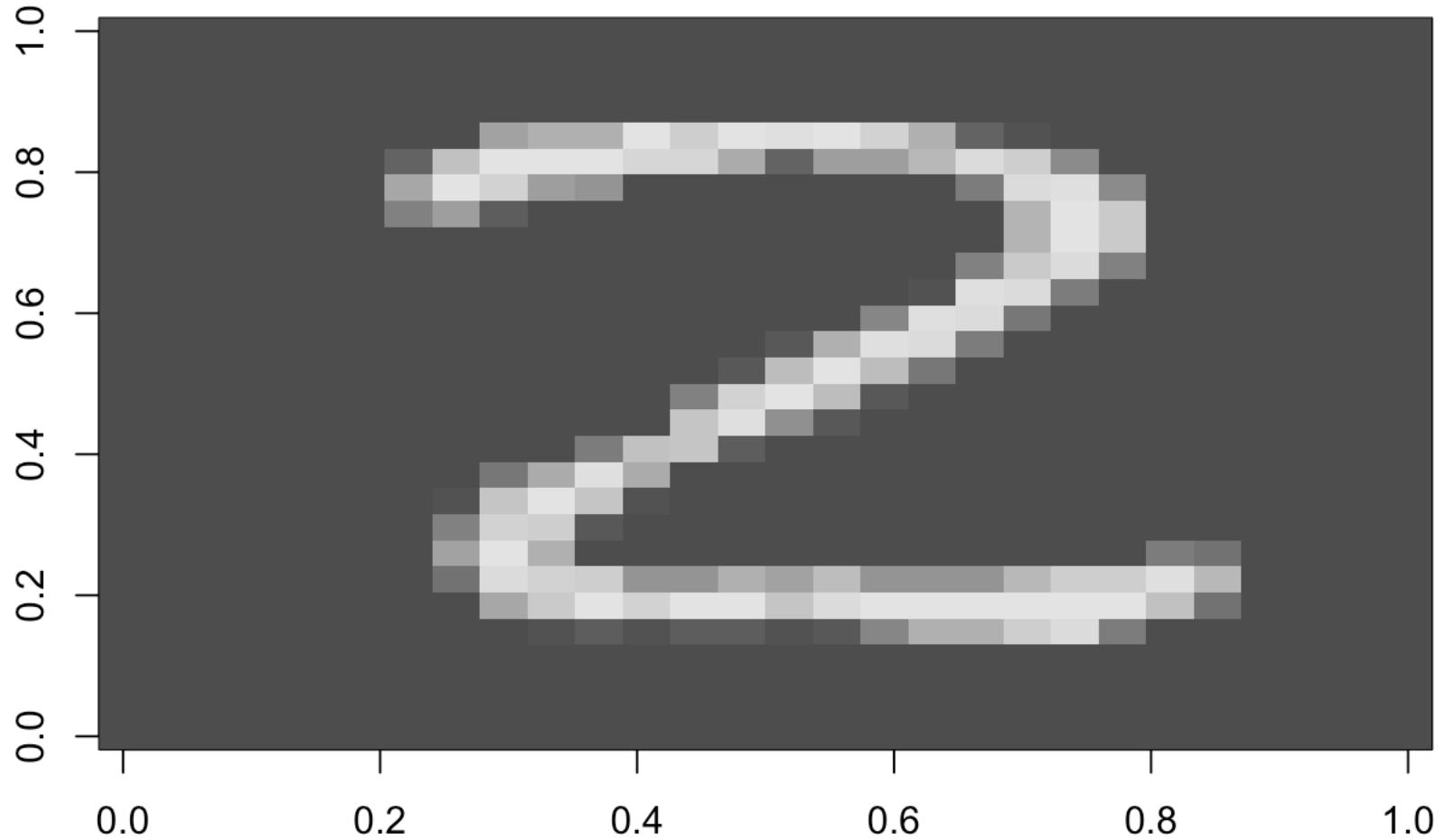


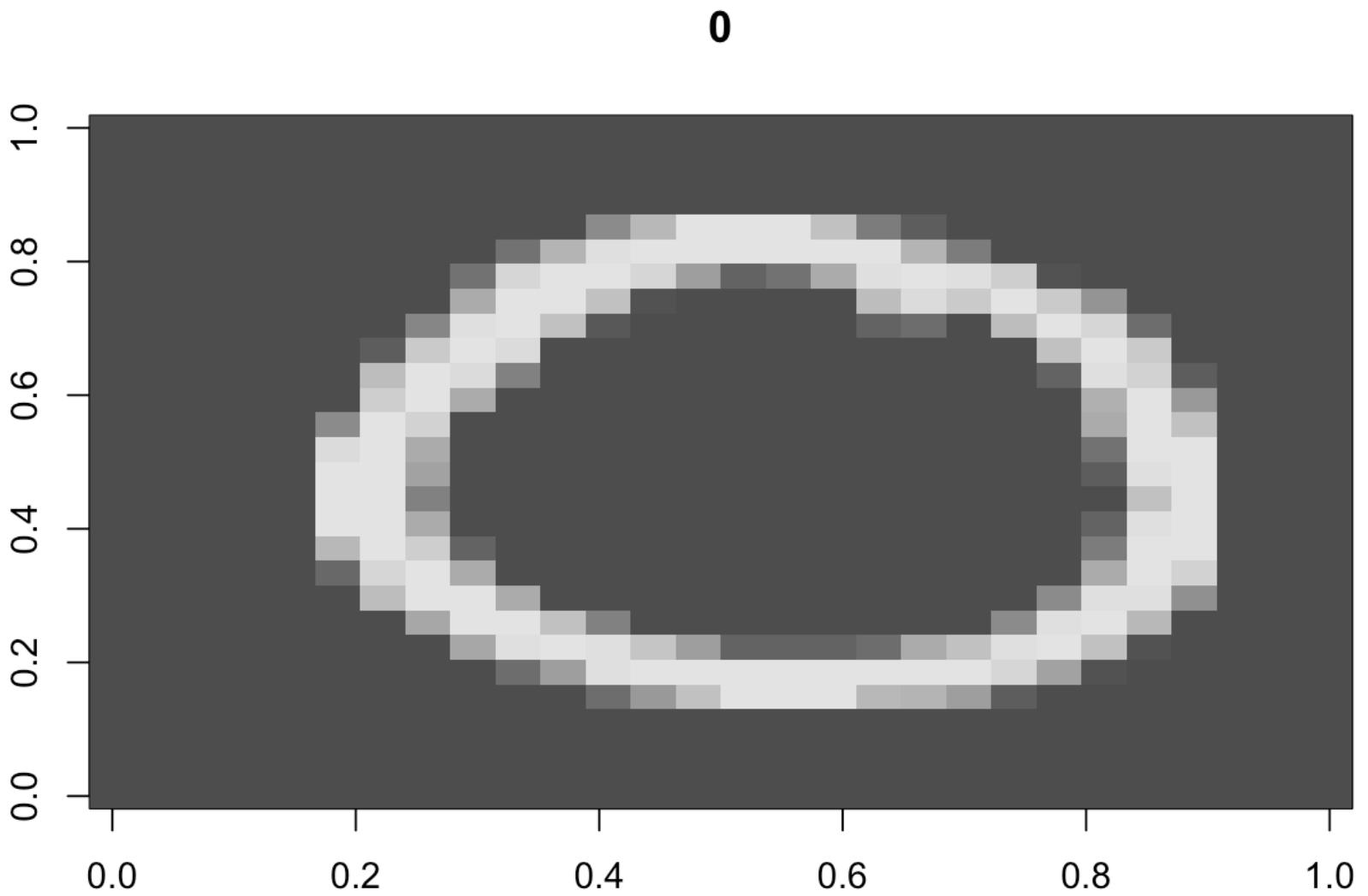
6



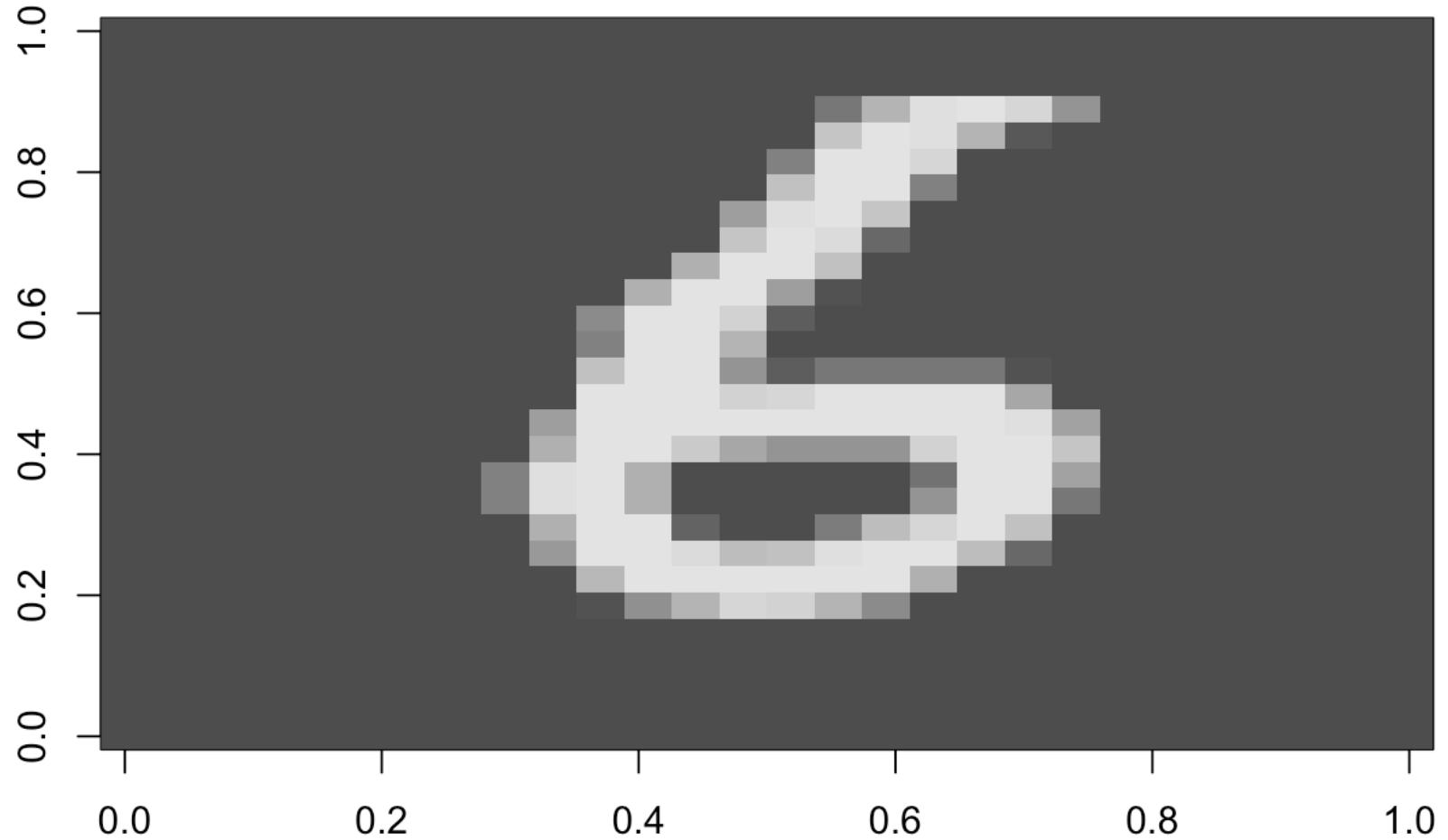


2

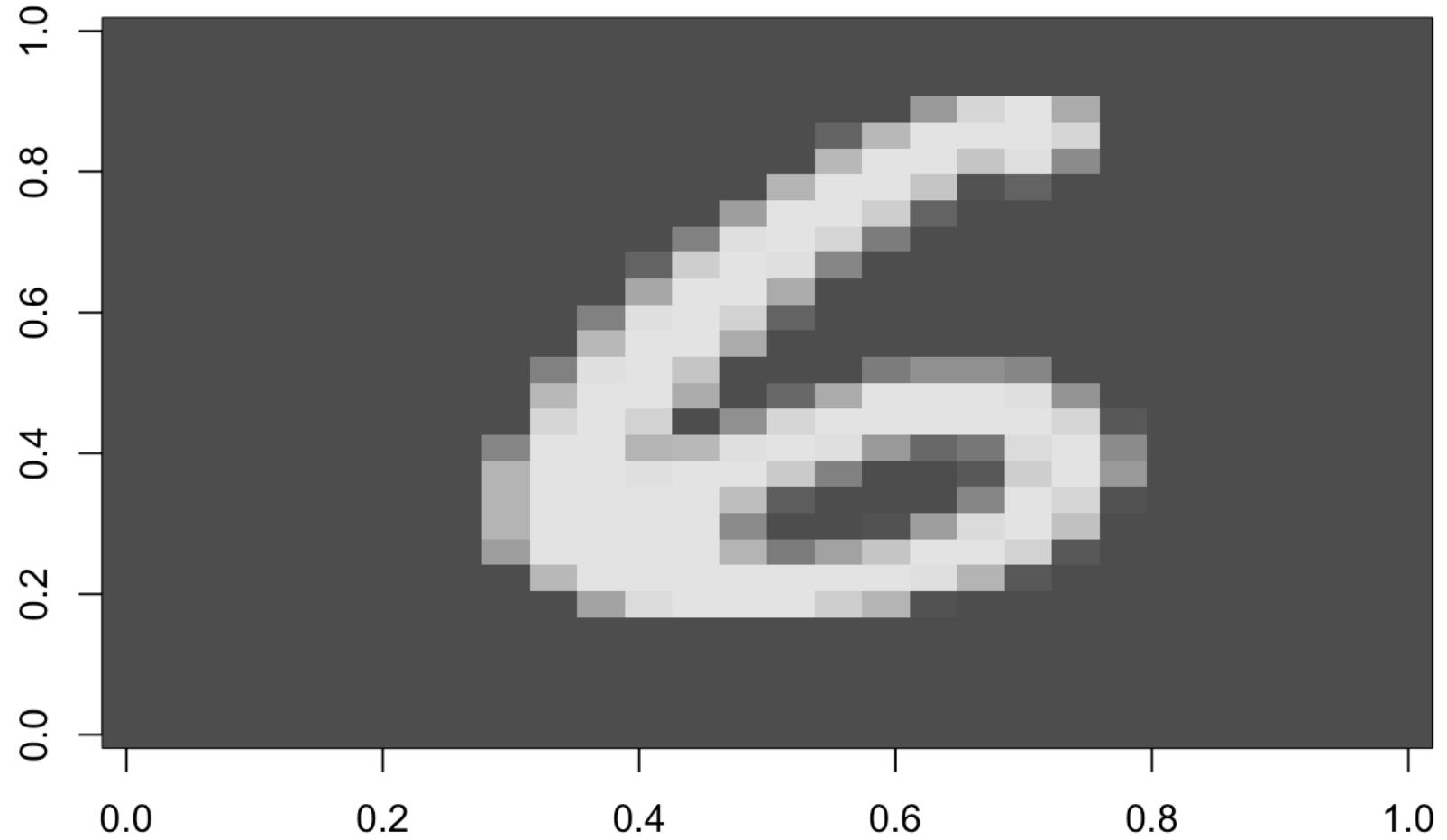




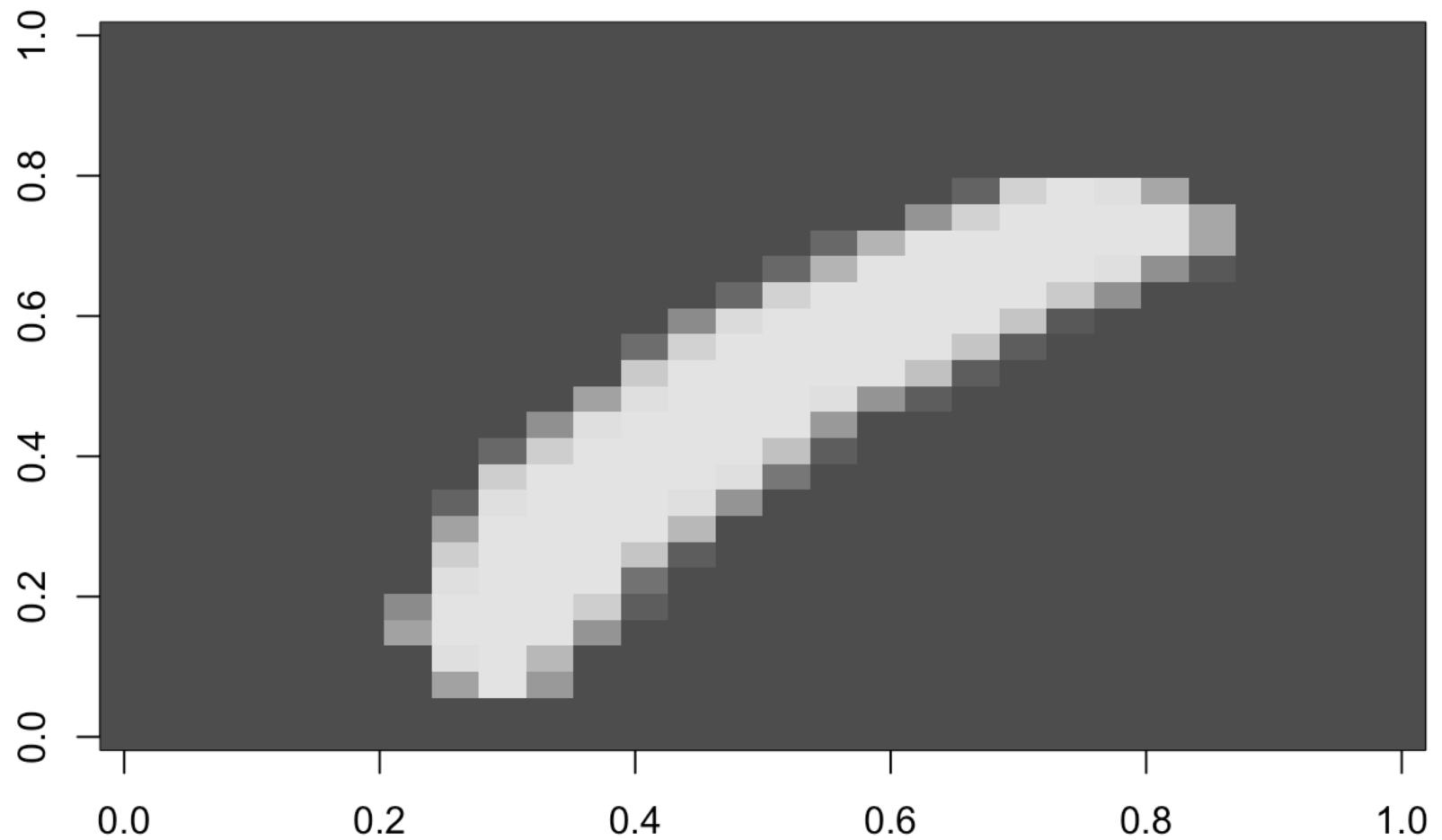
6



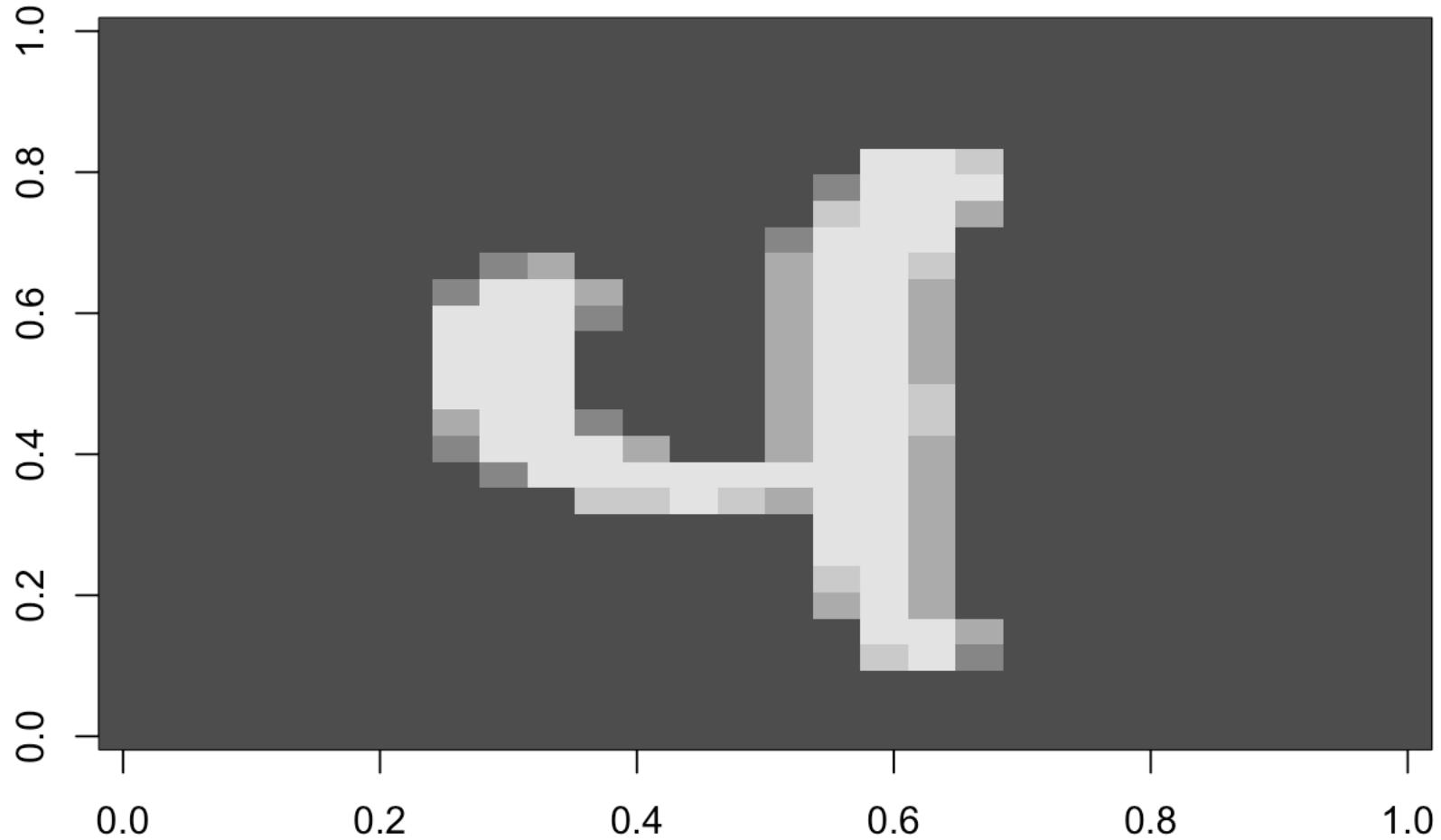
6



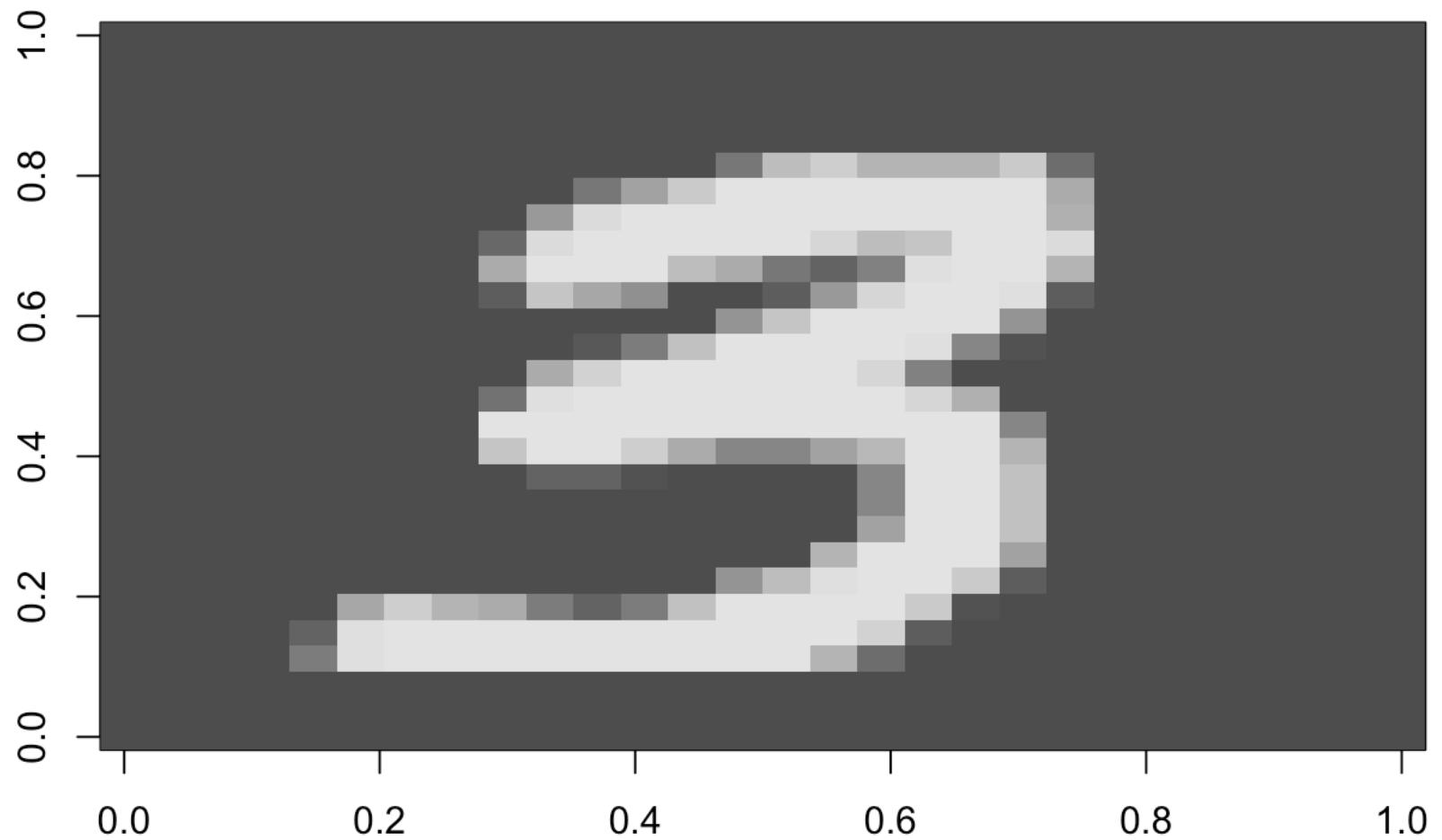
1



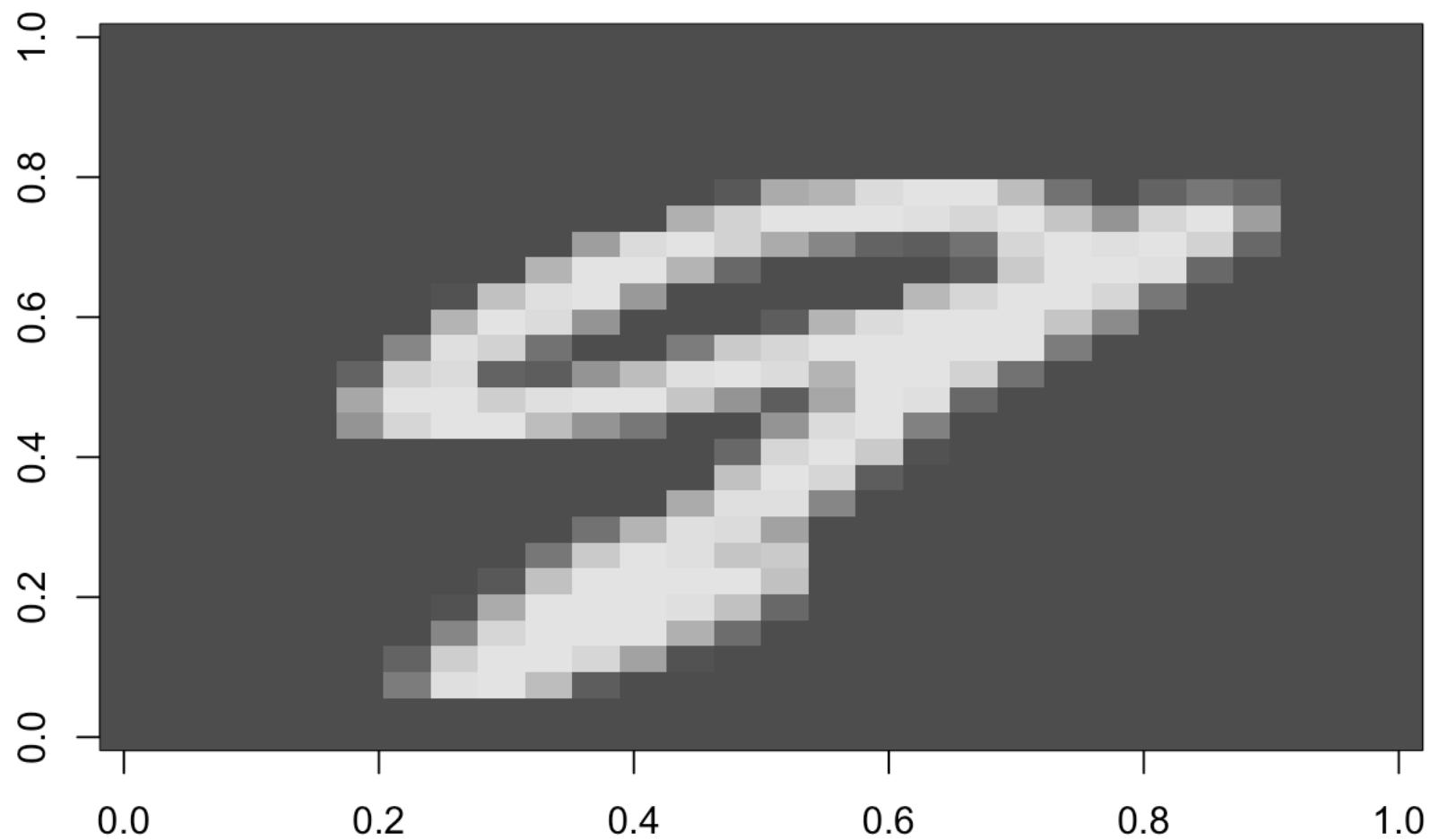
4

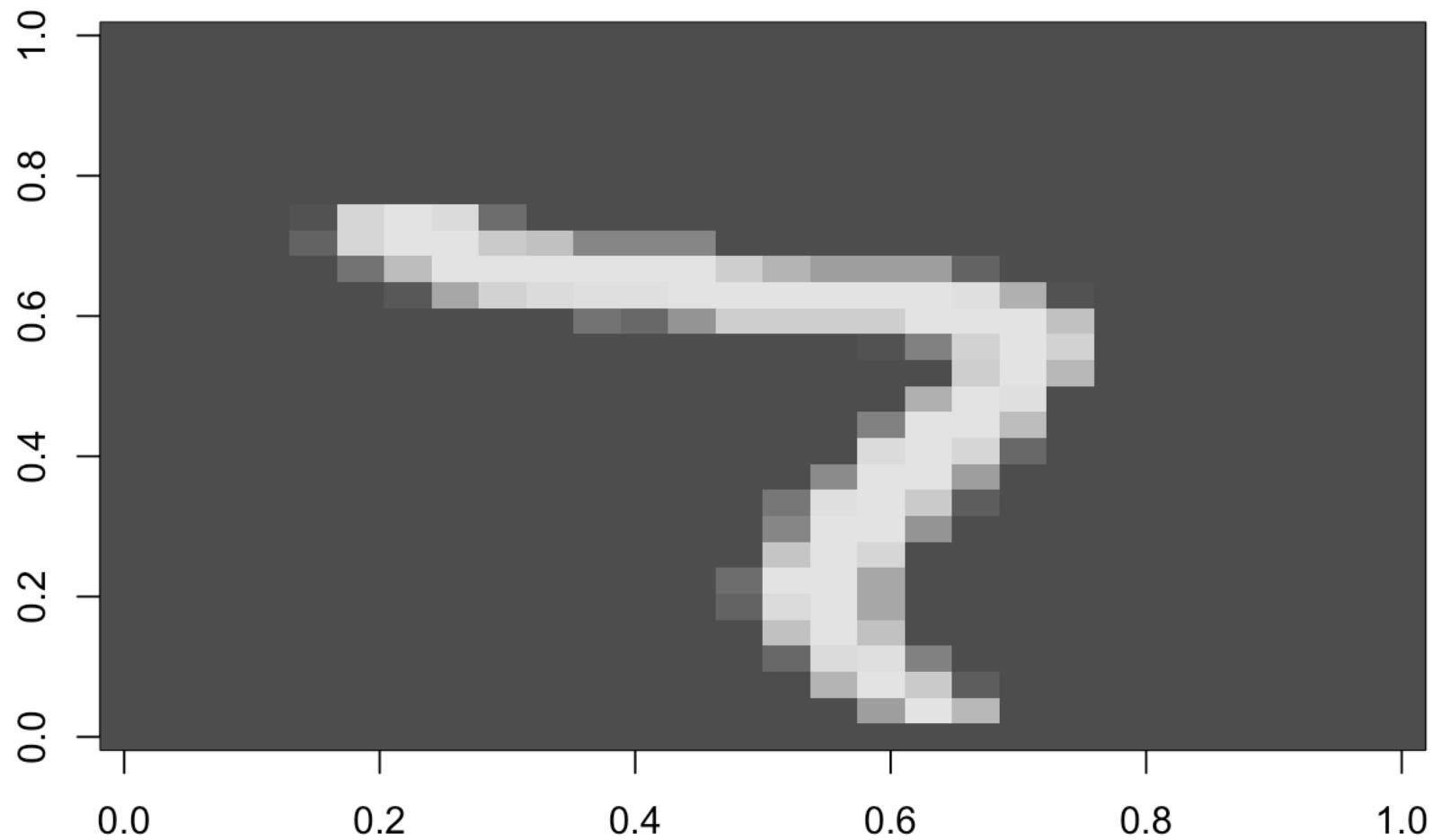


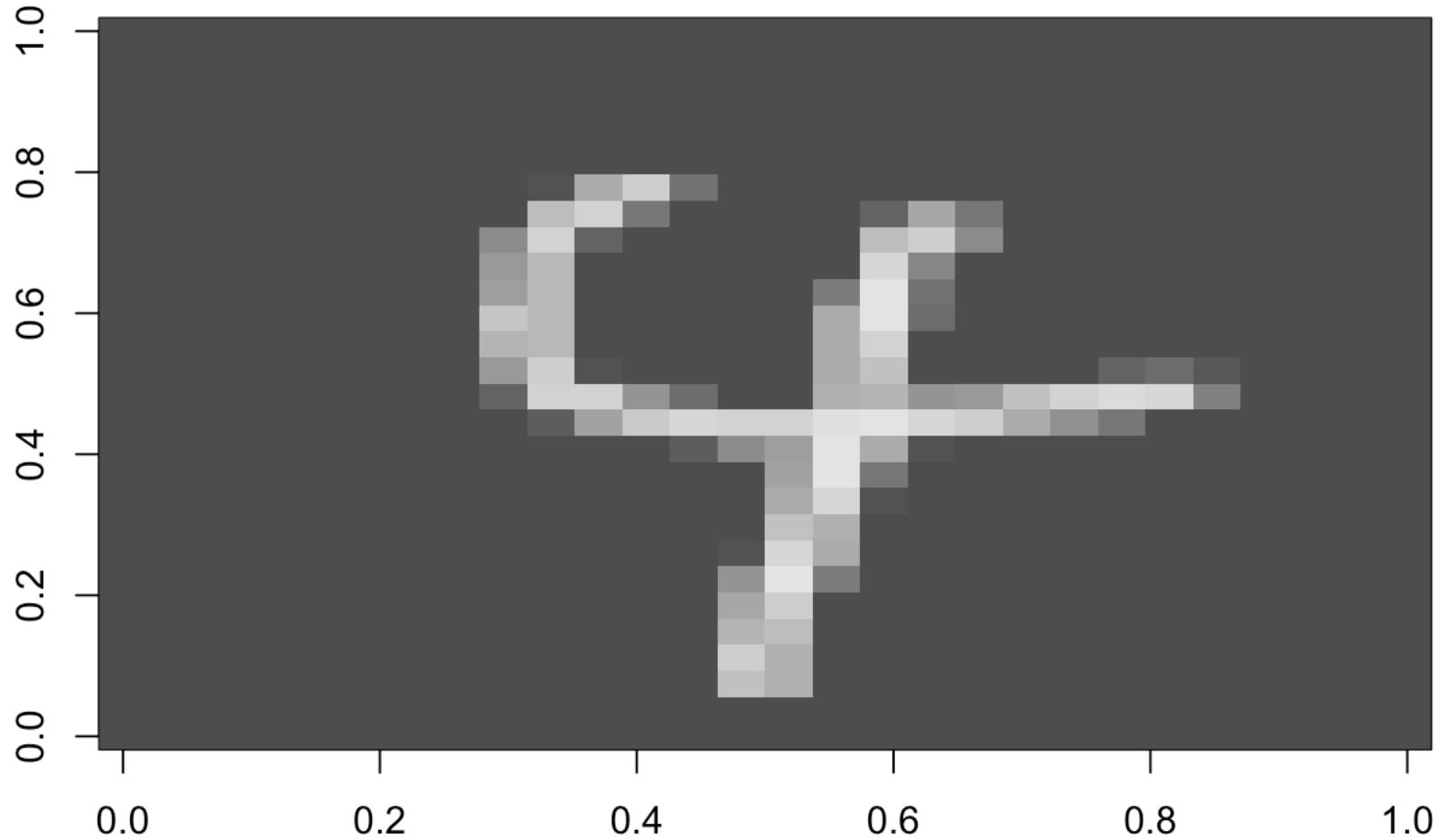
3



9







As we can see here, out of the first 100 predictions, 2 are wrong, which gives the classifier an accuracy rate of 98% on the first 100 observations. This result reaffirms the accuracy of the model at around 90-95% (as we have seen with the training and cross validation dataset).