

Class Activity 21

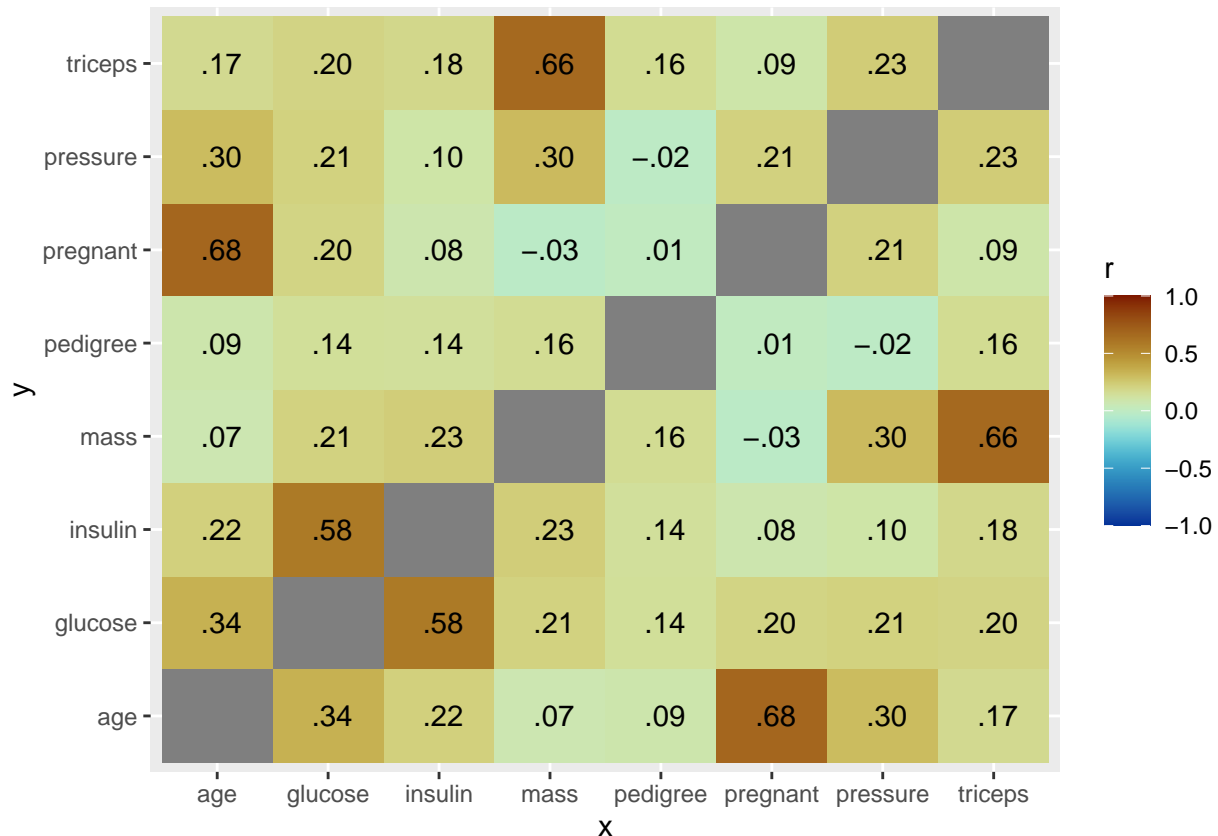
Arthur Viegas Eguia

May 13 2024

Group Activity 1

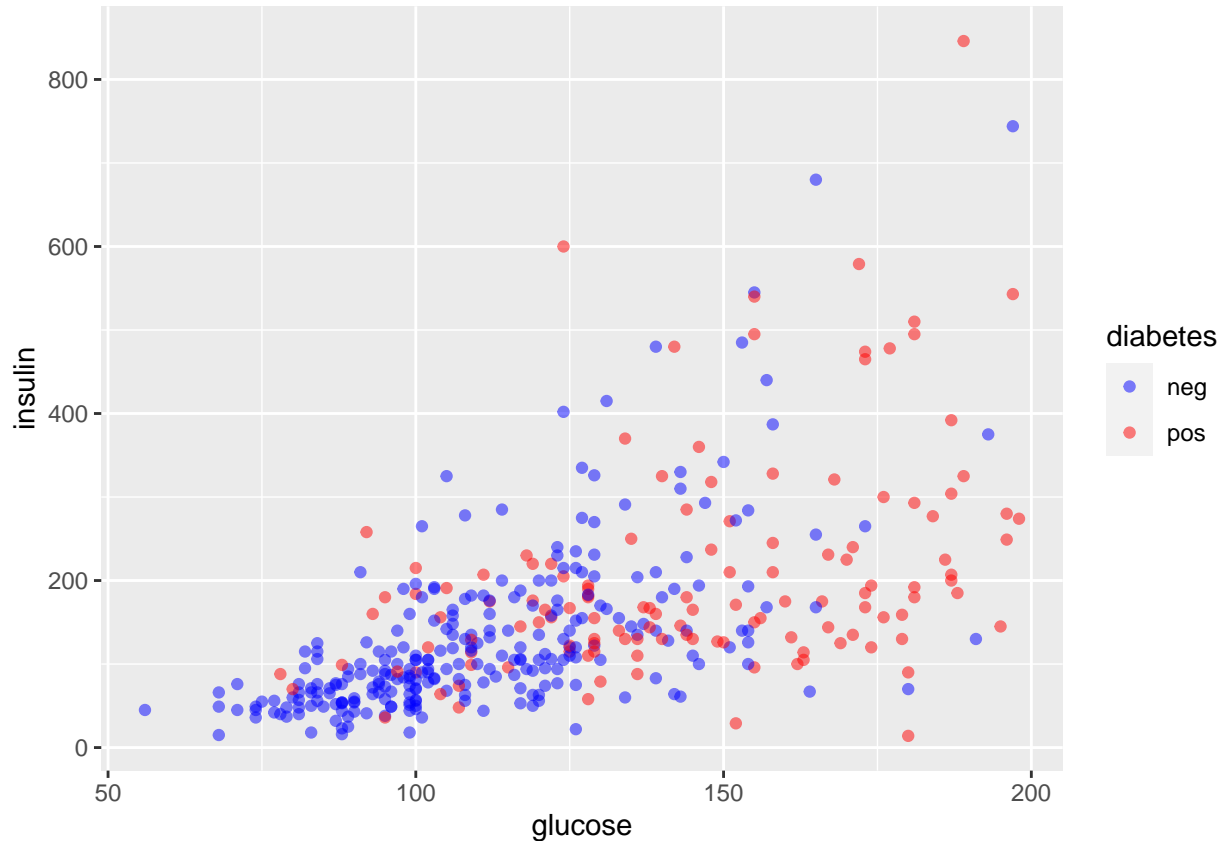
```
# Load the data
data(PimaIndiansDiabetes2)
db <- PimaIndiansDiabetes2 %>% tidyr::drop_na()

# correlation plot of the variables
db %>%
  select(-diabetes) %>% # only numerical variables
  correlate() %>%
  stretch() %>%
  ggplot(aes(x, y, fill = r)) +
  geom_tile() +
  geom_text(aes(label = as.character(fashion(r)))) +
  scale_fill_paletteer_c("scico::roma", limits = c(-1, 1), direction = -1)
```



- a. Create a scatter plot using ggplot2 to visualize the classification of diabetes status based on glucose and insulin levels, color-coding negative cases in blue and positive cases in red.

```
ggplot(db, aes(x = glucose, y = insulin, color = diabetes)) +
  geom_point(alpha = 0.5) +
  scale_color_manual(values = c("neg" = "blue", "pos" = "red"))
```

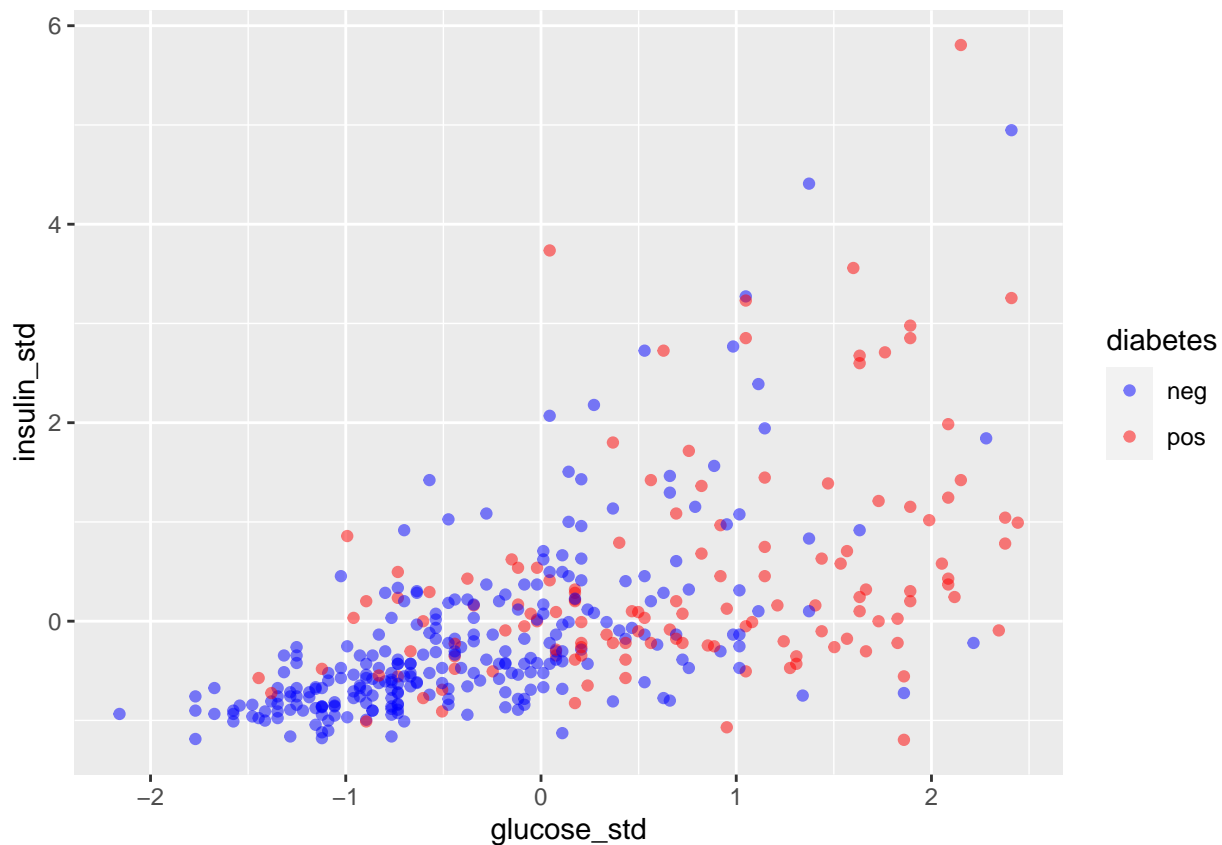


- b. Using the provided standardization function, apply it to both the glucose and insulin columns of your dataset to create new standardized columns, then plot these standardized values to analyze diabetes status.

```
standardize <- function(x, na.rm = FALSE) {
  (x - mean(x, na.rm = na.rm)) /
  sd(x, na.rm = na.rm)
}
```

```
db_std <- db %>%
  mutate(glucose_std = standardize(glucose, na.rm = TRUE),
         insulin_std = standardize(insulin, na.rm = TRUE))
```

```
ggplot(db_std, aes(x = glucose_std, y = insulin_std, color = diabetes)) +
  geom_point(alpha = 0.5) +
  scale_color_manual(values = c("neg" = "blue", "pos" = "red"))
```



c. Let's perform all the steps involved in classifying whether a patient with certain glucose and insulin would have diabetes or not.

```
# 1 Prepare raw data
db_raw <- db %>% drop_na() %>% select(glucose, insulin, diabetes)

# 2 Create a recipe for data pre-processing
db_recipe <- recipe(diabetes ~ ., data = db_raw) %>%
  step_scale(all_predictors()) %>%
  step_center(all_predictors()) %>%
  prep()

# 3 Apply the recipe to the data set
db_scaled <- bake(db_recipe, db_raw)

# 4 Create a model specification
knn_spec <- nearest_neighbor(mode = "classification",
                             engine = "kkn",
                             neighbors = 5)

# 5 Fit the model on the pre-processed data
knn_fit <- knn_spec %>% fit(diabetes ~ ., data = db_scaled)
knn_fit
parsnip model object

Call:
kkn::train.kkn(formula = diabetes ~ ., data = data, ks = min_rows(5, data, 5))
```

```
Type of response variable: nominal
Minimal misclassification: 0.2780612
Best kernel: optimal
Best k: 5
```

```
# 6 Classify
# These are standardized value!!
new_observations <- tibble(glucose = c(1, 2), insulin = c(-1, 1))

predict(knn_fit, new_data = new_observations)
# A tibble: 2 x 1
  .pred_class
  <fct>
1 neg
2 pos
```

- d. We already know the labels of the patients in the dataset. How well does the model predict their diabetes status? We will see more of this in the coming lectures.

```
#scaled_observations <-
```

What is the accuracy percentage?

Answer:

```
db_scaled %>%
  predict(knn_fit, new_data = .) %>%
  bind_cols(db_scaled, .) %>%
  mutate(correct_prediction = diabetes == .pred_class) %>%
  summarize(accuracy = mean(correct_prediction))
# A tibble: 1 x 1
  accuracy
  <dbl>
1      0.860
```

- e. Repeat part d. with a different model that consists of all the available features fitted with different number of neighbors. See if the accuracy percentage change in this new setting.

```
db_new <- db %>% tidyr::drop_na()

db_recipe <- recipe(diabetes ~ ., data = db_new) %>%
  step_scale(all_predictors()) %>%
  step_center(all_predictors()) %>%
  prep()

db_scaled <- bake(db_recipe, db_new)

knn_spec <- nearest_neighbor(mode = "classification",
                             engine = "kkn",
                             neighbors = 25)

knn_fit <- knn_spec %>% fit(diabetes ~ ., data = db_scaled)
db_scaled %>%
  predict(knn_fit, new_data = .) %>%
  bind_cols(db_scaled, .) %>%
  mutate(correct_prediction = diabetes == .pred_class) %>%
  summarize(accuracy = mean(correct_prediction))
```

```
# A tibble: 1 x 1
  accuracy
  <dbl>
1    0.832
```