

Ques 1 Hierarchical clustering on the Iris dataset

Algorithm

- ① load data
- ② standardize the features to zero mean and unit variance
- ③ compute Euclidean distance matrix
- ④ apply hclust with method = "ward.D2"
- ⑤ plot the dendrogram
- ⑥ cut off tree at  $k \leftarrow$  decide
- ⑧ plot clusters

code

```
library(stats)
library(cluster)
library(ggplot2)

X ← scale(iris[, 1:4])
d ← dist(X, method = "euclidean")
hc ← hclust(d, method = "ward.D2")

plot(as.dendrogram(hc), main = "practical", ylab = "Height",
      xlab = "Samples")

k ← 3

cl ← cutree(hc, k = k)
table(cl)
```

```
sil_k <- sapply(2:6, function(kk) {
  ss <- silhouette(cutree(hc, k=kk), d)
  mean(ss[, "sil_width"])
})
```

```
pca <- prcomp(X, center = FALSE, scale = FALSE)
pc <- as.data.frame(pca$x[, 1:2])
pc$cluster <- factor(cl)
```

```
ggplot(pc, aes(PC1, PC2, cluster)) +
  geom_point(size = 2) +
  labs(title = "Iris") +
  theme_minimal()
```

## Output

```
cl
 1  2  3
49 30 71
```

## Dendrogram

