# KID

## Shade Extraction

### Anonymized for Review

27.04.2022

### **Packages**

```
# qit install opencv
# devtools::install_github("ropenscilabs/opencu", force = T)
# Packages
get.package <- function(package){</pre>
  lapply(package, \(x){
    # check if packages are installed and if not install them
    if(!require(x, character.only = T)){
       install.packages(x)
    }
    # call package
    library(x, character.only = T)
  })
}
get.package(c("png", "jpeg", "tabulizer", "pdftools", "raster", "rgdal", "sp",
              "cluster"))
\# since I will use Map() / lapply() alot for plotting I will wrap them in invisible()
invis.Map <- function(f, ...) invisible(Map(f, ...))</pre>
invis.lapply <- function(x, f, ...) invisible(lapply(x, f, ...))</pre>
```

## Import KIDs

```
# set
setwd("./../KIDs/Erste")

# all PDF files in the current directory
file_names <- list.files(pattern = ".pdf")

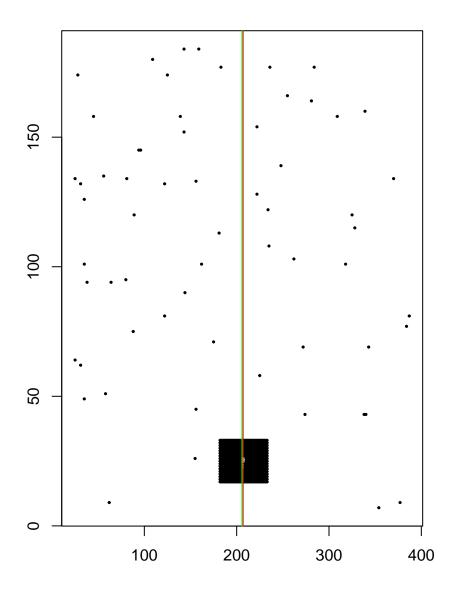
# extract text split by linebreak
lapply(file_names, function(x){</pre>
```

```
# split by line break and extract text
  strsplit(pdftools::pdf_text(x), "\n")
}) -> pdf.text.list
# identofy line that contains "Risiko- und Ertragsprofil"
lapply(pdf.text.list, function(x){
  # obtain relative line on page
  sapply(x, function(y){
    \# id
    tmp1 <- grep("Risiko- und Ertragsprofil", y) / length(y)</pre>
    # return
    if(length(tmp1) == 0){
      return(NA)
    } else {
      return(tmp1)
    }
  })
}) -> pos.list
# find page that contains desired text
sapply(pos.list, function(x) which(!is.na(x))) -> page.oi
# import PDF and convert to Png
Map(function(x, y, z){
  # convert first page of pdf to bitmap
  tmp1 <- pdftools::pdf_render_page(x, page = y, dpi = 50)</pre>
  # subset array
  tmp2 <- round(dim(tmp1)[3] * z)</pre>
  # return
  tmp1[ , , -c(tmp2:1)]
}, file_names, page.oi, sapply(pos.list, "[", 1)) -> bitmap.list
# import PDF and convert to Png
# to JPG
jpeg::writeJPEG(bitmap.list[[1]], "test1.jpeg")
# JPEG
imt <- jpeg::readJPEG("test.jpeg")</pre>
```

### **Extract SSRI**

### Bitmap

## Pixels: a6a6a6ff

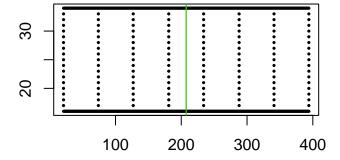


```
# find page margin for now use same color later should be switched to black
# left side margin
lsm <- min(coo1[, 1])
rsm <- max(coo1[, 1])

# scale
int_leng <- (rsm - lsm) / 7

# midpoints
scale <- setNames(cumsum(c(lsm + int_leng / 2, rep(int_leng, 6))), 1:7)</pre>
```

## Pixels: #005290



# We realize the page is upside down in the bitmap but not mirrored!

#### **JPEG**

Pixels: a6a6a6ff Noise: Low

Noise: Low

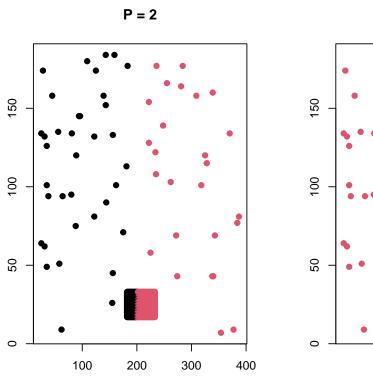
Noise: High

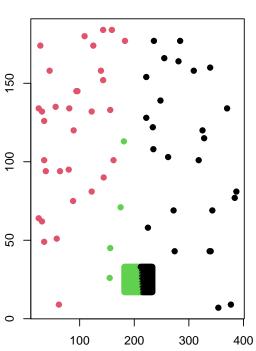
Pixels: a6a6a6ff

```
#off
#dev.off()
```

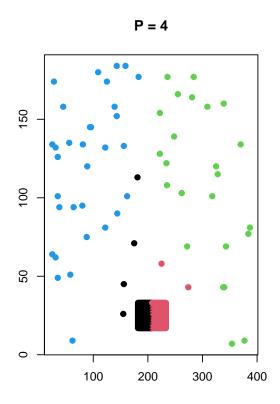
#### Classify utilizing k-means.

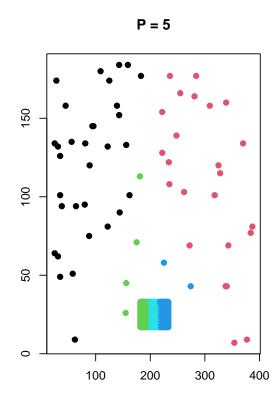
```
# Classify with different amount of groups then check for sil coef
# amt of groups
p <- 2:5
# estimate
lapply(p, function(x){
  # merge cluster into df
  dat <- cbind(coo1, kmeans(coo1, x)$cluster) # try specifying centers as closest to median most top le
  # silhouette
  tmp1 <- cluster::silhouette(dat[, ncol(dat)], dist(coo1))</pre>
  # return SC and Data
  list(
  "SC" = max(tapply(tmp1[, "sil_width"], tmp1[, "cluster"], mean)),
 "dat" = dat)
}) -> dat.kmeans
# sil coef
sapply(dat.kmeans, "[[", 1)
## [1] 0.4937373 0.5905349 0.5577094 0.4888173
# data.frames with the classification of different amt. of groups
# arrange
par(mfrow = c(2, 2))
# plot
invis.Map(function(x, y){
  plot(y[, 1], y[, 2], col = y[, ncol(y)], pch = 19, main = paste("P =", x), xlab = "",
       ylab = "")
}, p, lapply(dat.kmeans, "[[", 2))
```





P = 3





#### Classify utilizing hierarchical clustering

```
meth <- c("single", "average", "complete")</pre>
# ramp up p
p <- 2:5
# estimate
# over methods
lapply(meth, function(x){
  # over p
  lapply(p, function(y){
    # get grouping
    tmp1 <- agnes(coo1, method = x, diss = F)</pre>
    # restrict amnt of groups
    tmp2 \leftarrow cutree(tmp1, k = y)
    # bind
    tmp3 <- cbind(coo1, tmp2)</pre>
    # calculate coefficients
    tmp4 <- silhouette(tmp3[, ncol(tmp3)], dist(tmp3[, 2:3]))</pre>
    SC <- max(tapply(tmp4[, "sil_width"], tmp4[, "cluster"], mean))</pre>
    # return
    list("Data" = tmp3,
    "SC" = SC,
    "tmp.plot.silhouette" = tmp4)
  }) |> setNames(nm = paste("P =", p))
}) |> setNames(nm = meth) -> Group.list
# SC
lapply(Group.list, \(x){
  sapply(x, "[[", "SC")
})
## $single
       P = 2
                 P = 3
                            P = 4
                                      P = 5
## 0.8460797 0.3397616 0.9236410 0.9015994
##
## $average
       P = 2
               P = 3
                        P = 4
                                      P = 5
## 0.8831344 0.9328561 0.8125481 0.8125481
## $complete
      P = 2
               P = 3 	 P = 4
                                     P = 5
```

```
## 0.8787131 0.9219419 0.8719915 0.8497295
```

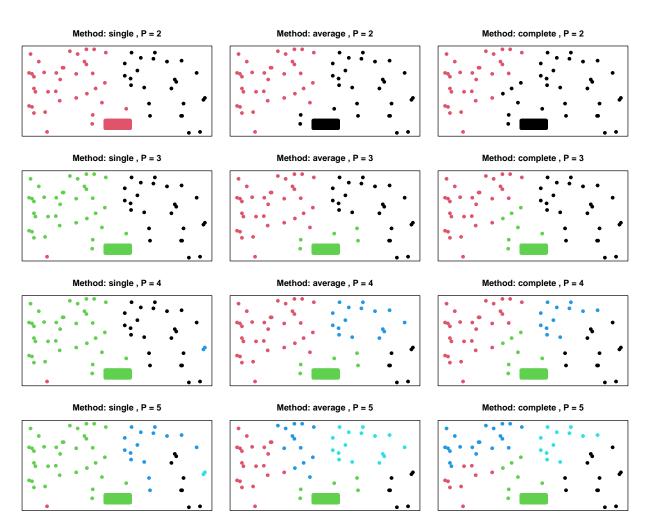
```
# plot
# arrange

#pdf(file = "methodsclust.pdf", height = 8)
par(mfrow = c(4, 3), mar = c(1, 1, 3, 1) + 0.1)
invis.lapply(paste("P =", p), \(x){\}

# over p
invis.lapply(meth, \(y){\}

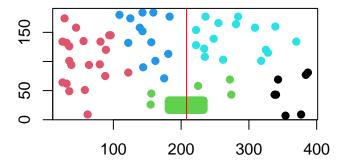
# Data
tmp.plot <- Group.list[[y]][[x]][["Data"]]

# plot
plot(tmp.plot[, 1], tmp.plot[, 2], col = tmp.plot[, ncol(tmp.plot)], pch = 19,
main = paste("Method:", y, ",", x), xlab = "", ylab = "", xaxt = "n", yaxt = "n")
})
})
})</pre>
```



```
#dev.off()
```

## Identifying the cluster and estimating SSRI



```
# classify
dif <- abs(med.rect.grp - scale)

# amt pixels off
dif

## 1 2 3 4 5 6 7
## 157.14286 105.42857 53.71429 2.00000 49.71429 101.42857 153.14286

# which grp
which.min(dif)

## 4
## 4</pre>
```