

# KID

## Shade Extraction

Fabian Blasch

08.08.2021

### Packages

```
# Packages
get.package <- function(package){

  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)
  })
}

# exec
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, ...))
invis.lapply <- function(x, f, ...) invisible(lapply(x, f, ...))
```

### Import KIDs

```
# set
setwd("C:/Users/blasc/OneDrive/Documents/GitHub/KID/KIDs")

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

# extract text split by linebreak
lapply(file_names, function(x){

  # split by line break and extract text
  strsplit(pdftools::pdf_text(x), "\n")

}) -> pdf.text.list
```

```

# import PDF and convert to Png
lapply(file_names, function(x){

  # convert first page of pdf to bitmap
  pdftools::pdf_render_page(x, page = 1, dpi = 50)

}) -> bitmap.list

# to JPG
jpeg::writeJPEG(bitmap.list[[1]], "test1.jpeg")

# JPEG
imt <- jpeg::readJPEG("test.jpeg")

```

## Extract SSRI

### Bitmap

```

# bitmap of first fund
bitmp1 <- bitmap.list[[1]]

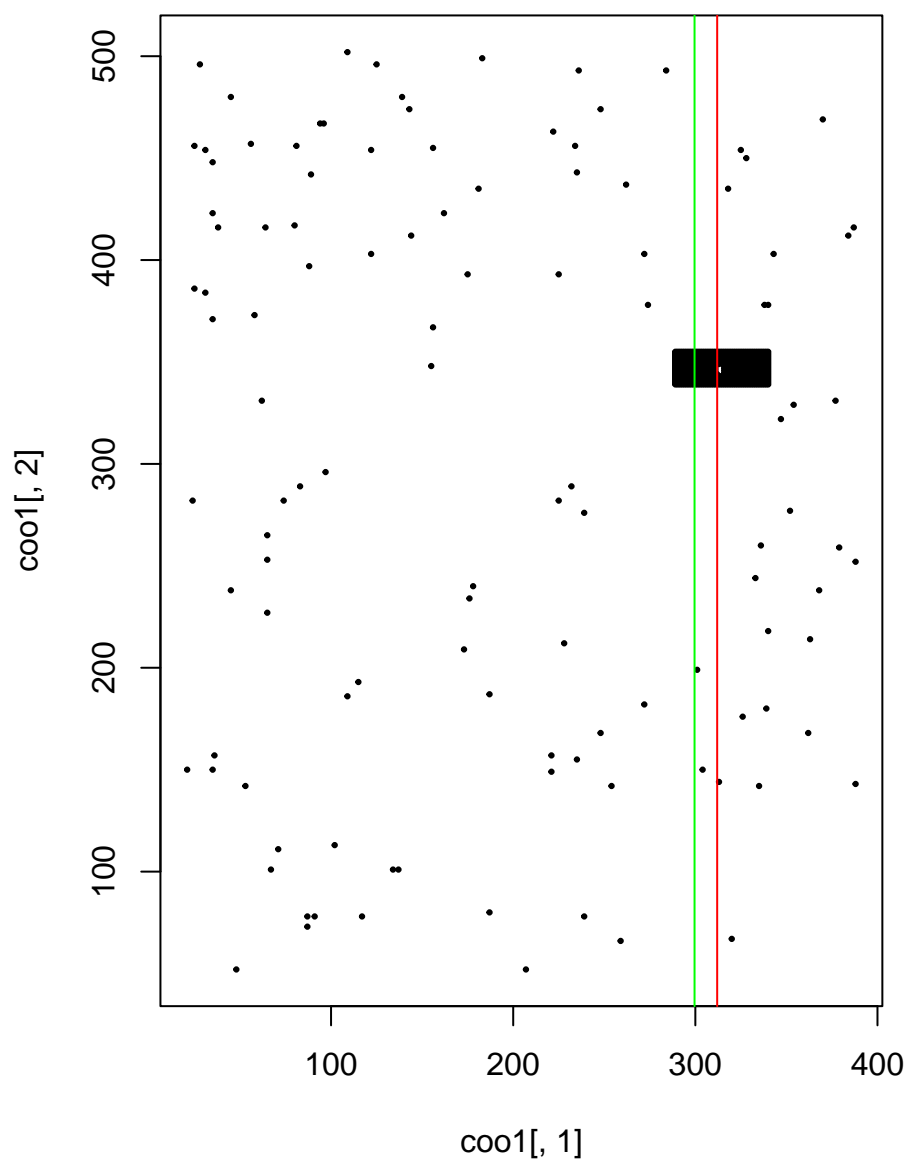
# col
# bitmp1[,587, 690]

# coordinates
coo1 <- which(bitmp1[1,,] == "a6" & bitmp1[2,,] == "a6" & bitmp1[3,,] == "a6" &
  bitmp1[4,,] == "ff", arr.ind = T)

# plot
{
plot(coo1[,1], coo1[, 2], main = "Pixels: a6a6a6ff", pch = 19, cex = 0.3)
abline(v = median(coo1[, 1]), col = "red")
abline(v = mean(coo1[, 1]), col = "green")
}

```

## 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)
```

```
# only using the median to predict SSRI  
which.min(abs(median(coo[, 1]) - scale))
```

```
## 6
```

```
## 6
```

```
# in this case we correctly predict the SSRI, without further classification!
```

## JPEG

```
# convert by alpha value in this case 255
```

```
imt <- imt * 255
```

```
# coordinates
```

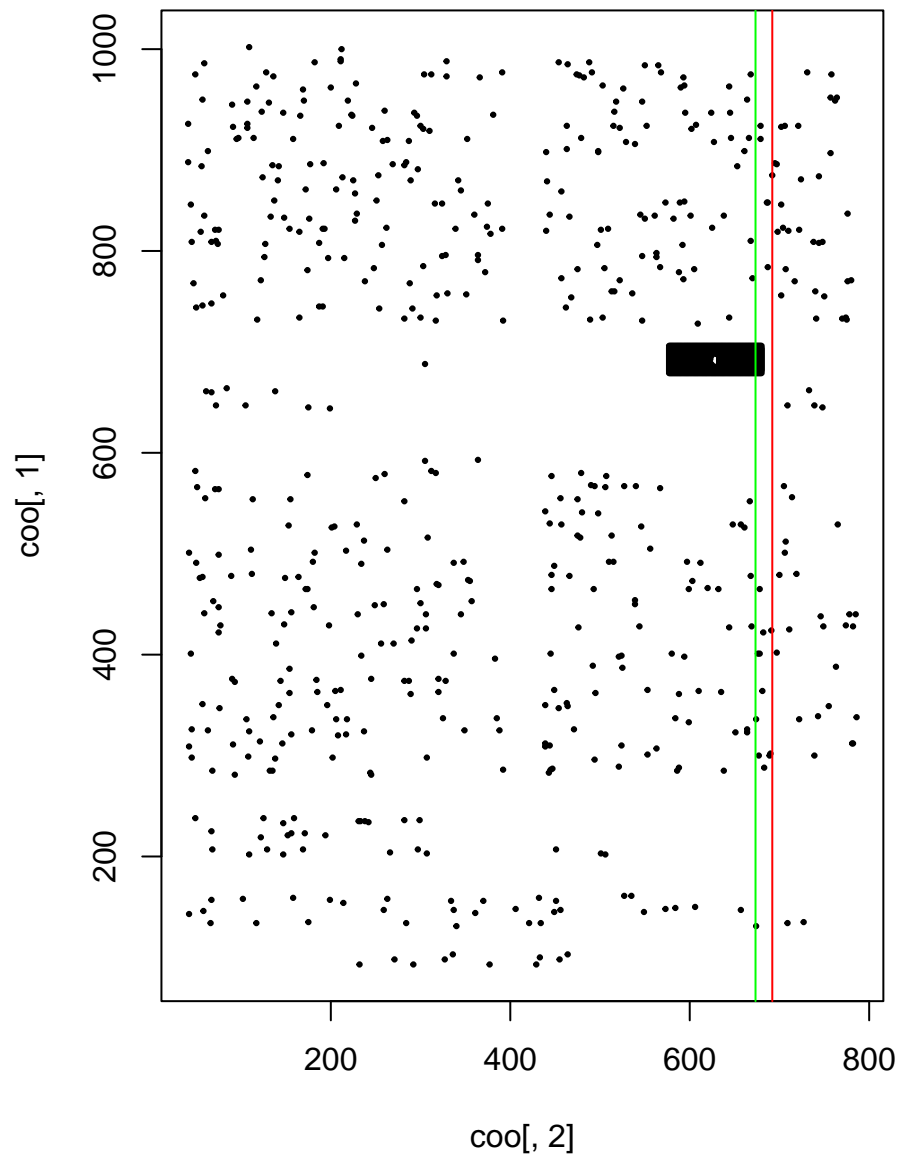
```
  # coo <- which((0 < imt[, 1] & imt[, 1] < 10) &  
  #               (75 < imt[, 2] & imt[, 2] < 85) &  
  #               (130 < imt[, 3] & imt[, 3] < 150), arr.ind = T)
```

```
coo <- which(imt[,1] == 166 & imt[,2] == 166 & imt[,3] == 166, arr.ind = T)
```

```
# Pixel / Coordinateplot
```

```
{  
plot(coo[,2], coo[, 1], main = "Pixels: r = 166, g = 166, b = 166", pch = 19, cex = 0.3)  
abline(v = median(coo[, 1]), col = "red")  
abline(v = mean(coo[, 1]), col = "green")  
}
```

**Pixels: r = 166, g = 166, b = 166**



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))

# return SC and Data
list(
  "SC" = max(tapply(tmp1[, "sil_width"], tmp1[, "cluster"], mean)),
  "dat" = dat)
}) -> dat.kmeans

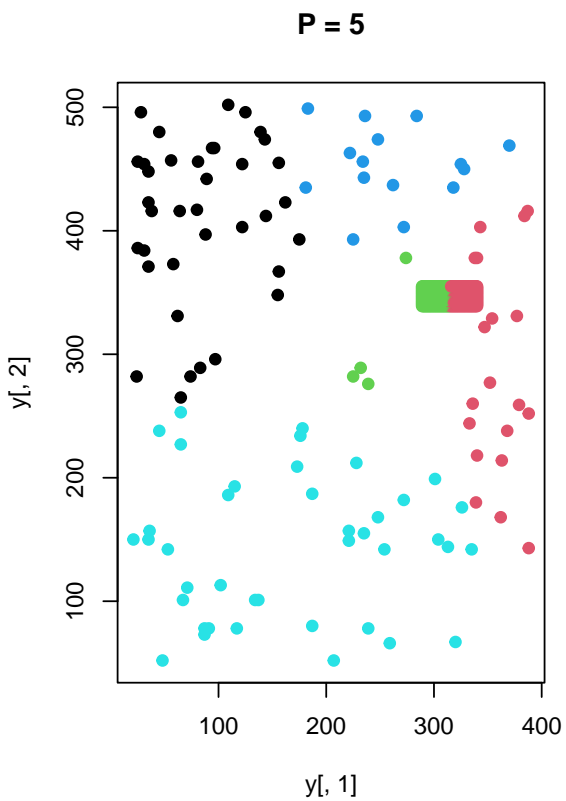
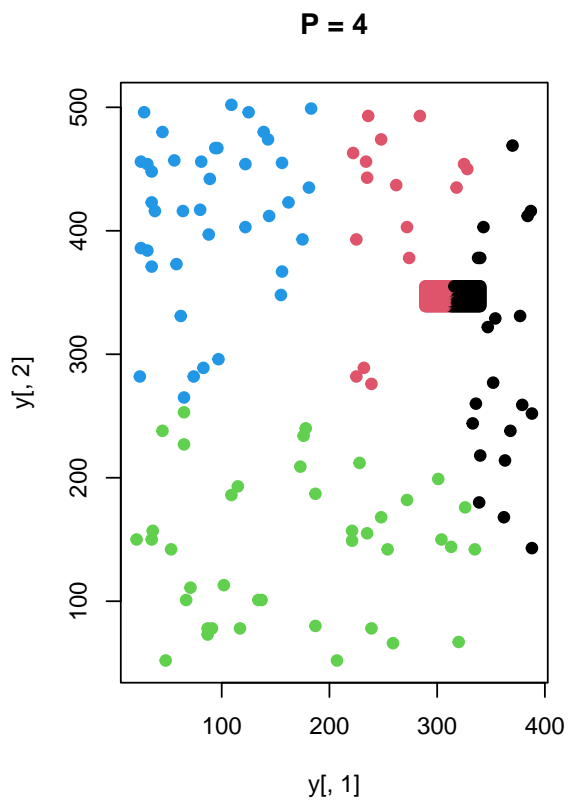
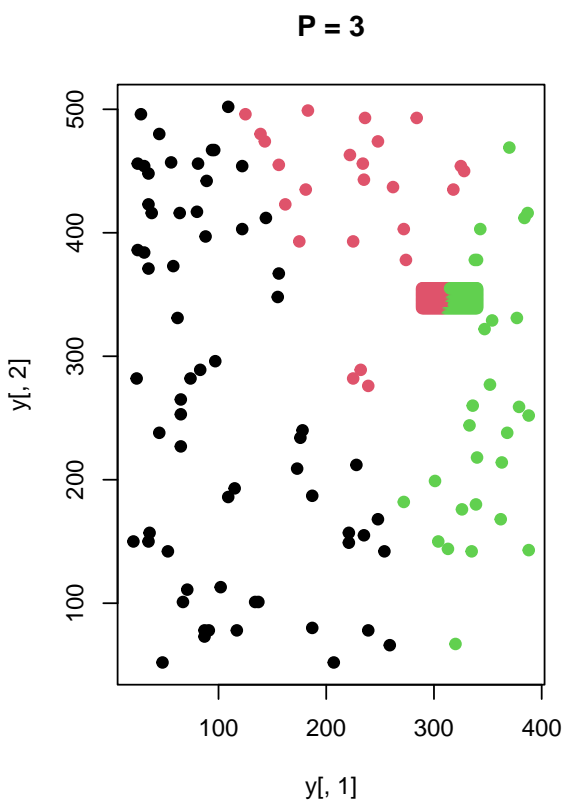
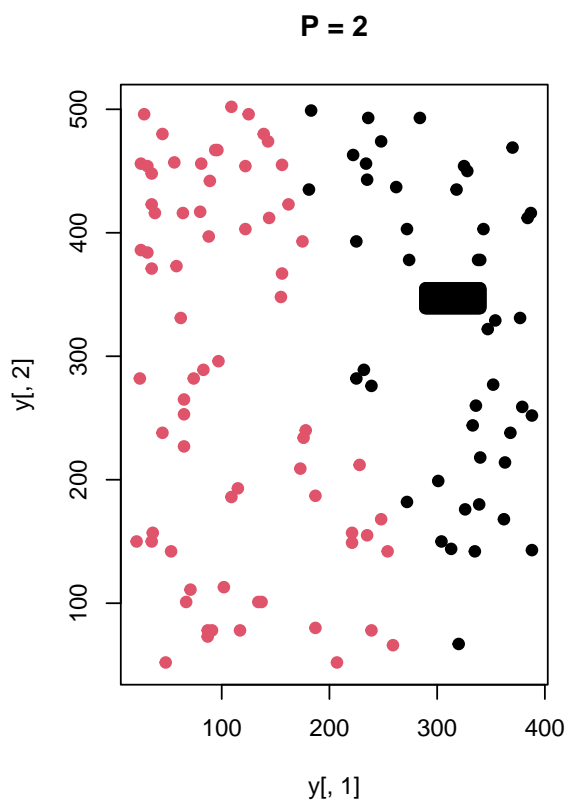
# sil coef
sapply(dat.kmeans, "[[", 1)

## [1] 0.8870700 0.4330818 0.4989778 0.5861058

# data.frames with the classification of different amt. of groups
# arrange
par(mfrow = c(2, 2))

# plot
invisible(Map(function(x, y){
  plot(y[, 1], y[, 2], col = y[, ncol(y)], pch = 19, main = paste("P =", x))
}, p, lapply(dat.kmeans, "[[", 2))

```



## Classify utilizing hierarchical clustering

```
# methods
meth <- c("single", "average", "complete")

# 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)

    # restrict amnt of groups
    tmp2 <- cutree(tmp1, k = y)

    # bind
    tmp3 <- cbind(coo1, tmp2)

    # calculate coefficients
    tmp4 <- silhouette(tmp3[, ncol(tmp3)], dist(tmp3[, 2:3]))

    # SC
    SC <- max(tapply(tmp4[, "sil_width"], tmp4[, "cluster"], mean))

    # 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
## 0.8863065 0.8745861 0.8745861 0.9211543
##
## $average
##      P = 2      P = 3      P = 4      P = 5
## 0.8736758 0.8292062 0.8724868 0.7837196
##
## $complete
##      P = 2      P = 3      P = 4      P = 5
```



```
## 0.9135547 0.8963935 0.8955598 0.8861353
# plot
# arrange
par(mfrow = c(4, 3))

invisible.lapply(paste("P =", p), \ (x){

  # over p
  invisible.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))

  })
})
```

