

*MAP 553 - Statistical learning and nonparametric  
estimation*

Numerical project

**Report**

EZZAKI Mahmoud

ROHRMANN Till

December 28, 2012

# 1 Problem description

# 2 Data exploration

# 3 Interpretation

## A Annex

```
# q.R
#
# This file creates the graphs for our problem. It calculates the regression
# function with the GDP per capita as its dependent variable for different
# independent variables:
# - Export ratio to GDP per capita
# - Import ratio to GDP per capita
# - Worker ratio to total population
# - Life expectation
# - Average school years
# - Political instability
# - Political rights
# - Civil liberties
#
# Moreover, the following 2 dimensional regression functions are calculated with
# the GDP per capita as its dependent variable:
# - Export ratio to GDP per capita and import ratio to GDP per capita
# - Export ratio to GDP per capita and life expectation
# - Export ratio to GDP per capita and worker ratio to total population
# - Political rights and civil liberties

source("normalize.R")
source("gaussianNoyau.R")
source("rectangularNoyau.R")
source("localPolynomsEstimation.R")
source("leastSquaresEstimation.R")
source("plotHelper1D.R")
source("plotHelper2D.R")
source("plotHelperD1D.R")
source("trigonometricDictionary.R")
source("derive1D.R")
source("transformEstimator1D.R")

data <- read.csv("../donnees/BARLEE/base.csv")

cols <- c(3:7, 9:13, 15:29, 31:40, 55:64, 41:48)
workdata <- data[, cols]
GDP <- c(workdata$GDP.65, workdata$GDP.70, workdata$GDP.75, workdata$GDP.80, workdata$GDP.85)
Worker <- c(workdata$Worker.65, workdata$Worker.70, workdata$Worker.75, workdata$Worker.80,
  workdata$Worker.85)
Life.expectation <- c(workdata$Life.expectation.65, workdata$Life.expectation.70,
  workdata$Life.expectation.75, workdata$Life.expectation.80, workdata$Life.expectation.85)
Average.school.years <- c(workdata$Average.school.years.65, workdata$Average.school.years.70,
  workdata$Average.school.years.75, workdata$Average.school.years.80, workdata$Average.school.years.85)
Export <- c(workdata$Export.60.65, workdata$Export.65.70, workdata$Export.70.75,
  workdata$Export.75.80, workdata$Export.80.85)
Import <- c(workdata$Import.60.65, workdata$Import.65.70, workdata$Import.70.75,
  workdata$Import.75.80, workdata$Import.80.85)
PoInst <- c(workdata$Political.Instability.60.65, workdata$Political.Instability.65.70,
  workdata$Political.Instability.70.75, workdata$Political.Instability.75.80,
  workdata$Political.Instability.80.85)

PoRights <- c(workdata$Political.rights.70.75, workdata$Political.rights.75.80,
  workdata$Political.rights.80.85, workdata$Political.rights.85.90)
CivLibs <- c(workdata$Civil.Liberties.70.75, workdata$Civil.Liberties.75.80, workdata$Civil.Liberties.80.85,
  workdata$Civil.Liberties.85.90)
GDPRights <- c(workdata$GDP.70, workdata$GDP.75, workdata$GDP.80, workdata$GDP.85)

domainGDP = c(min(GDP, na.rm = TRUE), max(GDP, na.rm = TRUE))
GDP <- t(normalize(GDP))

domainWorker = c(min(Worker, na.rm = TRUE), max(Worker, na.rm = TRUE))
Worker <- t(normalize(Worker))

domainLifeExpectation = c(min(Life.expectation, na.rm = TRUE), max(Life.expectation,
  na.rm = TRUE))
Life.expectation <- t(normalize(Life.expectation))

domainAvgSchoolYears = c(min(Average.school.years, na.rm = TRUE), max(Average.school.years,
  na.rm = TRUE))
Average.school.years <- t(normalize(Average.school.years))

domainExport <- c(min(Export, na.rm = TRUE), max(Export, na.rm = TRUE))
Export <- t(normalize(Export))

domainImport <- c(min(Import, na.rm = TRUE), max(Import, na.rm = TRUE))
Import <- t(normalize(Import))

domainPoInst <- c(min(PoInst, na.rm = TRUE), max(PoInst, na.rm = TRUE))
PoInst <- t(normalize(PoInst))

domainCivLibs <- c(min(CivLibs, na.rm = TRUE), max(CivLibs, na.rm = TRUE))
CivLibs <- t(normalize(CivLibs))

domainPoRights <- c(min(PoRights, na.rm = TRUE), max(PoRights, na.rm = TRUE))
```

```

PoRights <- t(normalize(PoRights))

domainGDPRights <- c(min(GDPRights, na.rm = TRUE), max(GDPRights, na.rm = TRUE))
GDPRights <- t(normalize(GDPRights))

final <- data.frame(GDP = GDP, Worker = Worker, Life.expectation = Life.expectation,
  Average.school.years = Average.school.years, Export = Export, Import = Import)

dict <- trigonometricDictionary(dimension, 0, 1)
noyau <- gaussianNoyau(dimension)

ordre <- 3
dimension <- 1
blocksize <- 10
h <- 0.01

final <- data.frame(GDP = GDP, Export = Export)
final <- final[complete.cases(final), ]

estimatorGDPEXP <- localPolynomsEstimation(final$Export, final$GDP, ordre, dimension,
  noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$Export, domainExport, "Export_ratio", final$GDP, domainGDP, "GDP",
  estimatorGDPEXP, "Export_ratio_to_GDP")

estimatorGDPEXP_LS <- leastSquaresEstimation(final$Export, final$GDP, dimension, dict,
  blocksize, 1, 20)
plotHelper1D(final$Export, domainExport, "Export_ratio", final$GDP, domainGDP, "GDP",
  estimatorGDPEXP_LS, "Export_ratio_to_GDP_LS")

estimatorGDPEXP_D <- derive1D(estimatorGDPEXP, h)
plotHelperD1D(domainExport, "Export_ratio", domainGDP, "Derived_GDP", estimatorGDPEXP_D,
  "Derivation_export_ratio_to_GDP")

final <- data.frame(Import = Import, GDP = GDP)
final <- final[complete.cases(final), ]

estimatorGDPIMP <- localPolynomsEstimation(final$Import, final$GDP, ordre, dimension,
  noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$Import, domainImport, "Import_ratio", final$GDP, domainGDP, "GDP",
  estimatorGDPIMP, "Import_ratio_to_GDP")

estimatorGDPIMP_LS <- leastSquaresEstimation(final$Import, final$GDP, dimension, dict,
  blocksize, 1, 20)
plotHelper1D(final$Import, domainImport, "Import_ratio", final$GDP, domainGDP, "GDP",
  estimatorGDPIMP_LS, "Import_ratio_to_GDP_LS")

estimatorGDPIMP_D <- derive1D(estimatorGDPIMP, h)
plotHelperD1D(domainImport, "Import_ratio", domainGDP, "Derived_GDP", estimatorGDPIMP_D,
  "Derivation_import_ratio_to_GDP")

final <- data.frame(GDP = GDP, Life.expectation = Life.expectation)
final <- final[complete.cases(final), ]

estimatorGDPLIFE <- localPolynomsEstimation(final$Life.expectation, final$GDP, ordre,
  dimension, noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$Life.expectation, domainLifeExpectation, "Life_expectation", final$GDP,
  domainGDP, "GDP", estimatorGDPLIFE, "Life_expectation_to_GDP")

estimatorGDPLIFE_LS <- leastSquaresEstimation(final$Life.expectation, final$GDP, dimension,
  dict, blocksize, 1, 20)
plotHelper1D(final$Life.expectation, domainLifeExpectation, "Life_expectation", final$GDP,
  domainGDP, "GDP", estimatorGDPLIFE_LS, "Life_expectation_to_GDP_LS")

estimatorGDPLIFE_D <- derive1D(estimatorGDPLIFE, h)
plotHelperD1D(domainLifeExpectation, "Life_expectation", domainGDP, "Derived_GDP",
  estimatorGDPLIFE_D, "Derivation_life_expectation_to_GDP")

final <- data.frame(GDP = GDP, Worker = Worker)
final <- final[complete.cases(final), ]

estimatorGDPWORKER <- localPolynomsEstimation(final$Worker, final$GDP, ordre, dimension,
  noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$Worker, domainWorker, "Worker_ratio_to_population", final$GDP,
  domainGDP, "GDP", estimatorGDPWORKER, "Worker_ratio_to_GDP")

estimatorGDPWORKER_LS <- leastSquaresEstimation(final$Worker, final$GDP, dimension,
  dict, blocksize, 1, 20)
plotHelper1D(final$Worker, domainWorker, "Worker_ratio_to_population", final$GDP,
  domainGDP, "GDP", estimatorGDPWORKER_LS, "Worker_ratio_to_GDP_LS")

estimatorGDPWORKER_D <- derive1D(estimatorGDPWORKER, h)
plotHelperD1D(domainWorker, "Worker_ratio", domainGDP, "Derived_GDP", estimatorGDPWORKER_D,
  "Derivation_worker_ratio_to_GDP")

final <- data.frame(GDP = GDP, PoInst = PoInst)
final <- final[complete.cases(final), ]

estimatorGDPPoInst <- localPolynomsEstimation(final$PoInst, final$GDP, ordre, dimension,
  noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$PoInst, domainPoInst, "Political_instability", final$GDP, domainGDP,
  "GDP", estimatorGDPPoInst, "Political_instability_to_GDP")

estimatorGDPPoInst_D <- derive1D(estimatorGDPPoInst, h)
plotHelperD1D(domainPoInst, "Political_instability", domainGDP, "Derived_GDP", estimatorGDPPoInst_D,
  "Derivation_political_instability_to_GDP")

final <- data.frame(GDP = GDPRights, CivLibs = CivLibs)
final <- final[complete.cases(final), ]

estimatorGDPCivLibs <- localPolynomsEstimation(final$CivLibs, final$GDP, ordre, dimension,
  noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$CivLibs, domainCivLibs, "Civil_liberties", final$GDP, domainGDPRights,
  "GDP", estimatorGDPCivLibs, "Civil_liberties_to_GDP")

```

```

final <- data.frame(GDP = GDPRights, PoRights = PoRights)
final <- final[complete.cases(final), ]

estimatorGDPPoRights <- localPolynomsEstimation(final$PoRights, final$GDP, ordre,
  dimension, noyau, blocksize, 0.005, 1, 0.005)
plotHelper1D(final$PoRights, domainPoRights, "Political_rights", final$GDP, domainGDPRights,
  "GDP", estimatorGDPPoRights, "Political_rights_to_GDP")

dimension <- 2
ordre <- 1
blocksize <- 25

final <- data.frame(GDP = GDP, Import = Import, Export = Export)
final <- final[complete.cases(final), ]
plot(final$Export * (domainExport[2] - domainExport[1]) + domainExport[1], final$Import *
  (domainImport[2] - domainImport[1]) + domainImport[1], xlab = "Export_ratio_to_GDP",
  ylab = "Import_ratio_to_GDP")
title(main = "Export_ratio_and_import_ratio_points")

exportImport <- rbind(final$Export, final$Import)
estimatorGDPEXPimp <- localPolynomsEstimation(exportImport, final$GDP, ordre, dimension,
  gaussianNoyau(dimension), blocksize, 0.005, 1, 0.005)
plotHelper2D(domainExport, "\nExport_ratio", domainImport, "\nImport_ratio", estimatorGDPEXPimp,
  domainGDP, "\n\nGDP", "Export_ratio_and_Import_ratio_to_GDP", xstart = 0,
  xend = 1, ystart = 0, yend = 1)

final <- data.frame(GDP = GDP, Export = Export, Life.expectation = Life.expectation)
final <- final[complete.cases(final), ]
plot(final$Export * (domainExport[2] - domainExport[1]) + domainExport[1], final$Life.expectation *
  (domainLifeExpectation[2] - domainLifeExpectation[1]) + domainLifeExpectation[1],
  xlab = "Export_ratio", ylab = "Life_expectation", main = "Export_ratio_and_life_expectation_points")
exportLife <- rbind(final$Export, final$Life.expectation)
estimatorGDPEXPLife <- localPolynomsEstimation(exportLife, final$GDP, ordre, dimension,
  gaussianNoyau(dimension), blocksize, 0.005, 1, 0.005)
plotHelper2D(domainExport, "\nExport_ratio", domainLifeExpectation, "\nLife_expectation",
  estimatorGDPEXPLife, domainGDP, "\n\nGDP", "Export_ratio_and_life_expectation_to_GDP")

final <- data.frame(GDP = GDP, Export = Export, Worker = Worker)
final <- final[complete.cases(final), ]
plot(final$Export * (domainExport[2] - domainExport[1]) + domainExport[1], final$Worker *
  (domainWorker[2] - domainWorker[1]) + domainWorker[1], xlab = "Export_ratio",
  ylab = "Worker_ratio", main = "Export_ratio_and_Worker_ratio_points")
exportWorker <- rbind(final$Export, final$Worker)
estimatorGDPEXPWorker <- localPolynomsEstimation(exportWorker, final$GDP, ordre,
  dimension, gaussianNoyau(dimension), blocksize, 0.005, 1, 0.005)
plotHelper2D(domainExport, "\nExport_ratio", domainWorker, "\nWorker_ratio", estimatorGDPEXPWorker,
  domainGDP, "\n\nGDP", "Export_ratio_and_worker_ratio_to_GDP")

final <- data.frame(GDP = GDPRights, PoRights = PoRights, CivLibs = CivLibs)
final <- final[complete.cases(final), ]
plot(final$PoRights * (domainPoRights[2] - domainPoRights[1]) + domainPoRights[1],
  final$CivLibs * (domainCivLibs[2] - domainCivLibs[1]) + domainCivLibs[1], xlab = "Political_rights",
  ylab = "Civil_liberties", main = "Political_rights_and_civil_liberties_points")
poRightsCivLibs <- rbind(final$PoRights, final$CivLibs)
estimatorGDPPoRivL <- localPolynomsEstimation(poRightsCivLibs, final$GDP, ordre,
  dimension, gaussianNoyau(dimension), blocksize, 0.005, 1, 0.005)
plotHelper2D(domainPoRights, "\nPolitical_rights", domainCivLibs, "\nCivLibLiberties",
  estimatorGDPPoRivL, domainGDPRights, "\n\nGDP", "Political_rights_and_civil_liberties_to_GDP")

# localPolynomsEstimation.R
#
# This file contains the function localPolynomsEstimation which realizes
# a local polynoms estimation.

source("crossValidation.R")

# function localPolynomsEstimation
#
# This function calculates via the method of cross validation a local polynoms
# estimator for the observations (xdatapoints, ydatapoints). For this purpose,
# it calculates for every h, starting at hstart, stepsize hstep and ending at hend,
# the estimation for the risk and selects the estimator with the minimal risk.
#
# @param xdatapoints multidimensional x-values --> "dimensions" x "number observations"
# @param ydatapoints 1 dimensional y-values --> 1 x "number observations"
# @param ordre Ordre of the local polynoms
# @param dimension Dimension of the data
# @param noyau Function pointer to a kernel
# @param blocksize Block size which is used for the cross validation
# @param hstart Start value for h for the incremental search procedure
# @param hend End value for h for the incremental search procedure
# @param hstep Step width for the incremental search procedure of h
#
# @return Function pointer to the estimator using local polynoms. It takes as
# argument a vector of size 1 x "dimension" for which it calculates the estimated
# value
localPolynomsEstimation <- function(xdatapoints, ydatapoints, ordre, dimension, noyau,
  blocksize, hstart, hend, hstep) {

  minValue <- 2^109
  minParameter <- 0

  for (h in seq(hstart, hend, hstep)) {
    value <- crossValidation(xdatapoints, ydatapoints, localPolynomsEstimatorGenerator(ordre,
      h, noyau, dimension), dimension, blocksize)
    if (value < minValue) {
      minValue <- value
      minParameter <- h
    }
  }
  cat("h:", h, "_value:", value, "\n")
}

```

```

    }
    cat("min_h:", minParameter, "\n")
    localPolynomsEstimatorGenerator(ordre, minParameter, noyau, dimension)(xdatapoints,
    ydatapoints)
}

# leastSquaresEstimation.R
#
# This file contains a helper function for the least squares estimation
source("crossValidation.R")

leastSquaresEstimation <- function(xdatapoints, ydatapoints, dimension, dict, blocksize,
mstart, mend) {
    minValue <- 2^109
    minParameter <- 0

    for (i in seq(mstart, mend)) {
        value <- crossValidation(xdatapoints, ydatapoints, leastSquaresEstimatorGenerator(dict,
        i), dimension, blocksize)
        if (value < minValue) {
            minValue <- value
            minParameter <- i
        }
        cat("m:", i, "_value:", value, "\n")
    }
    cat("min_m:", minParameter, "\n")
    leastSquaresEstimatorGenerator(dict, minParameter)(xdatapoints, ydatapoints)
}

# crossValidation.R
#
# function crossValidation
#
# This function performs a cross validation to estimate the risk for a given
# estimator generator.
#
# @param xdatapoints x-values of the observations
# @param ydatapoints y-values of the observations
# @param estimatorGenerator Generator function for the estimator
# @param dimension dimension of the x-values
# @param blocksize size of the blocks for the cross-validation
#
# @return estimated risk value
crossValidation <- function(xdatapoints, ydatapoints, estimatorGenerator, dimension,
blocksize = 1) {
    risk <- 0
    numdatapoints <- length(xdatapoints)/dimension
    dim(xdatapoints) = c(dimension, numdatapoints)

    # mix the data
    xdatapoints <- xdatapoints[, sample.int(numdatapoints)]

    dim(xdatapoints) = c(dimension, numdatapoints)

    # for every block of test data do
    for (index in seq(1, numdatapoints, by = blocksize)) {
        xtestpoints <- xdatapoints[, index:min(numdatapoints, index + blocksize -
        1)]
        ytestpoints <- ydatapoints[index:min(numdatapoints, index + blocksize - 1)]

        dim(xtestpoints) = c(dimension, min(numdatapoints - index + 1, blocksize))
        xlearningpoints <- xdatapoints[, -(index:min(numdatapoints, index + blocksize -
        1))]
        ylearningpoints <- ydatapoints[-(index:min(numdatapoints, index + blocksize -
        1))]

        dim(xlearningpoints) = c(dimension, numdatapoints - min(numdatapoints - index +
        1, blocksize))
        # estimator for current learning data set
        estimator <- estimatorGenerator(xlearningpoints, ylearningpoints)

        numblockdatapoints <- min(blocksize, numdatapoints - index + 1)
        # the risk for the corresponding test data set
        temprisk <- crossprod((ytestpoints - estimator(xtestpoints)), (ytestpoints -
        estimator(xtestpoints)))
        risk <- risk + temprisk/numblockdatapoints
    }

    risk
}

```