Assignment 4 Question 2 b)

Sheen Thusoo

30/11/2021

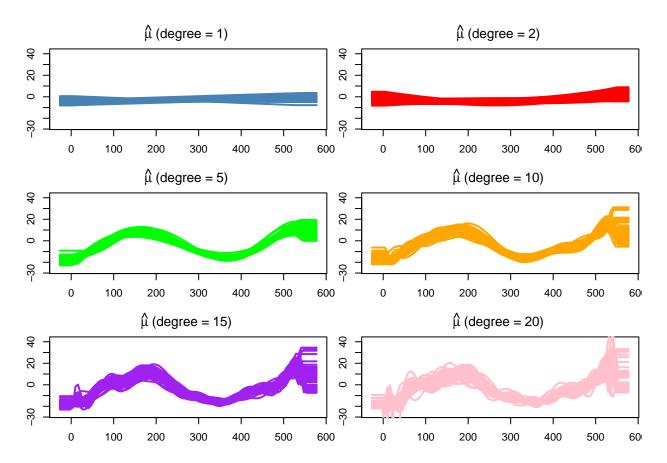
Generate M = 50 samples S_1, S_2, \ldots, S_{50} of size n = 100. You are encouraged (but don't have to) use functions getSampleComp and getXYSample from the lectures. Fit polynomials of degree 1, 2, 5, 10, 15, and 20 to every sample.

```
data <- read.csv("OzoneData.csv")
```

```
ave_y_mu_sq <- function(sample, predfun, na.rm = TRUE){</pre>
  mean((sample$y - predfun(sample$x))^2, na.rm = na.rm)
}
ave_mu_mu_sq <- function(predfun1, predfun2, x, na.rm = TRUE){</pre>
  mean((predfun1(x) - predfun2(x))^2, na.rm = na.rm)
}
getSampleComp <- function(pop, size, replace=FALSE) {</pre>
  N <- popSize(pop)</pre>
  samp <- rep(FALSE, N)</pre>
  samp[sample(1:N, size, replace = replace)] <- TRUE</pre>
  samp
}
getXYSample <- function(xvarname, yvarname, samp, pop) {</pre>
  sampData <- pop[samp, c(xvarname, yvarname)]</pre>
  names(sampData) <- c("x", "y")</pre>
  sampData
}
popSize <- function(pop) {nrow(as.data.frame(pop))}</pre>
sampSize <- function(samp) {popSize(samp)}</pre>
getmuhat <- function(sampleXY, complexity = 1) {</pre>
  formula <- paste0("y ~ ",</pre>
                      if (complexity==0) {
                        111
                        } else
                        pasteO("poly(x, ", complexity, ", raw = FALSE)")
                        #paste0("bs(x, ", complexity, ")")
  )
```

```
fit <- lm(as.formula(formula), data = sampleXY)</pre>
  tx = sampleXY$x
  ty = fit$fitted.values
  range.X = range(tx)
  val.rY = c( mean(ty[tx == range.X[1]]),
               mean(ty[tx == range.X[2]]) )
  ## From this we construct the predictor function
  muhat <- function(x){</pre>
    if ("x" %in% names(x)) {
      ## x is a dataframe containing the variate named
      ## by xvarname
      newdata <- x
    } else
      \#\# x is a vector of values that needs to be a data.frame
    { newdata \leftarrow data.frame(x = x) }
    ## The prediction
    ##
    val = predict(fit, newdata = newdata)
    val[newdata$x < range.X[1]] = val.rY[1]</pre>
    val[newdata$x > range.X[2]] = val.rY[2]
    val
  }
  ## muhat is the function that we need to calculate values
  ## at any x, so we return this function from getmuhat
  muhat
}
getmuFun <- function(pop, xvarname, yvarname){</pre>
  pop = na.omit(pop[, c(xvarname, yvarname)])
  # rule = 2 means return the nearest y-value when extrapolating, same as above.
  # ties = mean means that repeated x-values have their y-values averaged, as above.
  tauFun = approxfun(pop[,xvarname], pop[,yvarname], rule = 2, ties = mean)
  return(tauFun)
}
xnam <- "Dav"
ynam <- "Ozone"</pre>
pop <- data
n <- 100
N_S < -50
set.seed(1) # for reproducibility
samples <- lapply(1:N_S, FUN = function(i) {</pre>
    getSampleComp(pop, n)
})
Ssam <- lapply(samples, FUN = function(Si) {</pre>
    getXYSample(xnam, ynam, Si, pop)
})
Tsam <- lapply(samples, FUN = function(Si) {</pre>
```

```
getXYSample(xnam, ynam, !Si, pop)
})
```



Legend

- Degree 1
 Degree 5
 Degree 10
 Degree 15
 Degree 20