

Descriptive Statistics: the Education edition

M P Gururajan, Hina A Gokhale and Dayadeep Monder

Indian Institute of Technology Bombay, Mumbai

In this session, we are going to use a data set downloaded from data.gov.in site on education for our analysis. The dataset is called `Education.csv` and is stored in `csv` format.

Of course, we know how to load the data set. Before we do that, we will also remove all data from the previous sessions. And, after reading the data, we will check it once before proceeding further.

```
unlink("~/RData")
X <- read.csv("../Data/Education.csv",header=TRUE)
str(X)
```

```
## 'data.frame': 36 obs. of 16 variables:
## $ State : chr "Andhra Pradesh" "Arunachal Pradesh" "Assam" "Bihar" "Chhattisgarh" "Goa" "Gujarat" "Haryana" "Karnataka" "Kerala" "Madhya Pradesh" "Maharashtra" "Manipur" "Meghalaya" "Mizoram" "Nagaland" "Odisha" "Punjab" "Rajasthan" "Sikkim" "Tamil Nadu" "Telangana" "Tripura" "Uttar Pradesh" "West Bengal"
## $ Population.in..Cr..of.Total.Census.2001 : num 7.62 0.11 2.67 8.3 2.08 0.09 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01
## $ Population.in..Cr..of.6.14.age.2004 : num 1.3 0.02 0.53 1.88 0.41 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.01
## $ Literacy.Rate : num 60.5 54.3 63.3 47 64.7 82.1 82.1 82.1 82.1 82.1 82.1 82.1 82.1 82.1 82.1 82.1
## $ Gross.Enrollment.Ratio..Classes...I.VIII. : num 87 106.7 91.9 65.2 112.6 112.6 112.6 112.6 112.6 112.6 112.6 112.6 112.6 112.6 112.6 112.6
## $ Drop.out.Classes..I.X. : num 63.7 70.8 75 83.1 0 ...
## $ Pupil.Teacher.Ratio : int 33 34 42 104 48 21 35 44 44 44 44 44 44 44 44 44
## $ Pupil.Teacher.Ratio.Upper.Primary : int 31 30 16 75 46 17 39 30 30 30 30 30 30 30 30 30
## $ Elementary.School.per.lakh.population : int 99 163 137 57 203 76 73 61 61 61 61 61 61 61 61 61
## $ Sec.Hr.Sec.Schools.per.lakh.population : int 22 19 19 4 12 31 14 23 37 37 37 37 37 37 37 37
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Rs.cr. : num 1807 187 1165 2479 1440 1440 1440 1440 1440 1440 1440 1440 1440 1440 1440 1440
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure.as.percentage.of.Total : num 5.09 0.53 3.28 6.99 4.06 4.06 4.06 4.06 4.06 4.06 4.06 4.06 4.06 4.06 4.06 4.06
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age : int 1387 8179 2186 1319 3513 3513 3513 3513 3513 3513 3513 3513 3513 3513 3513 3513
## $ Lakh.of.Population.per.Institution..University. : num 31.6 11.4 46.5 46.2 43.8 43.8 43.8 43.8 43.8 43.8 43.8 43.8 43.8 43.8 43.8 43.8
## $ Lakh.of.Population.per.Institution..College. : num 0.59 1.14 0.88 1.18 1.03 1.03 1.03 1.03 1.03 1.03 1.03 1.03 1.03 1.03 1.03 1.03
## $ Lakh.of.Population.per.Institution..Technical. : num 1.97 5.7 25.34 43.87 18.2 18.2 18.2 18.2 18.2 18.2 18.2 18.2 18.2 18.2 18.2 18.2
```

We know that this data, in addition to the statewise information, also contains all India information. In some cases, we want to remove the India information; we can generate a new data frame in which the last line is removed as follows. Note that we are storing one row of data corresponding to India in a different variable.

```
Y <- X[-36,]
Z <- X[36,]
str(Y)

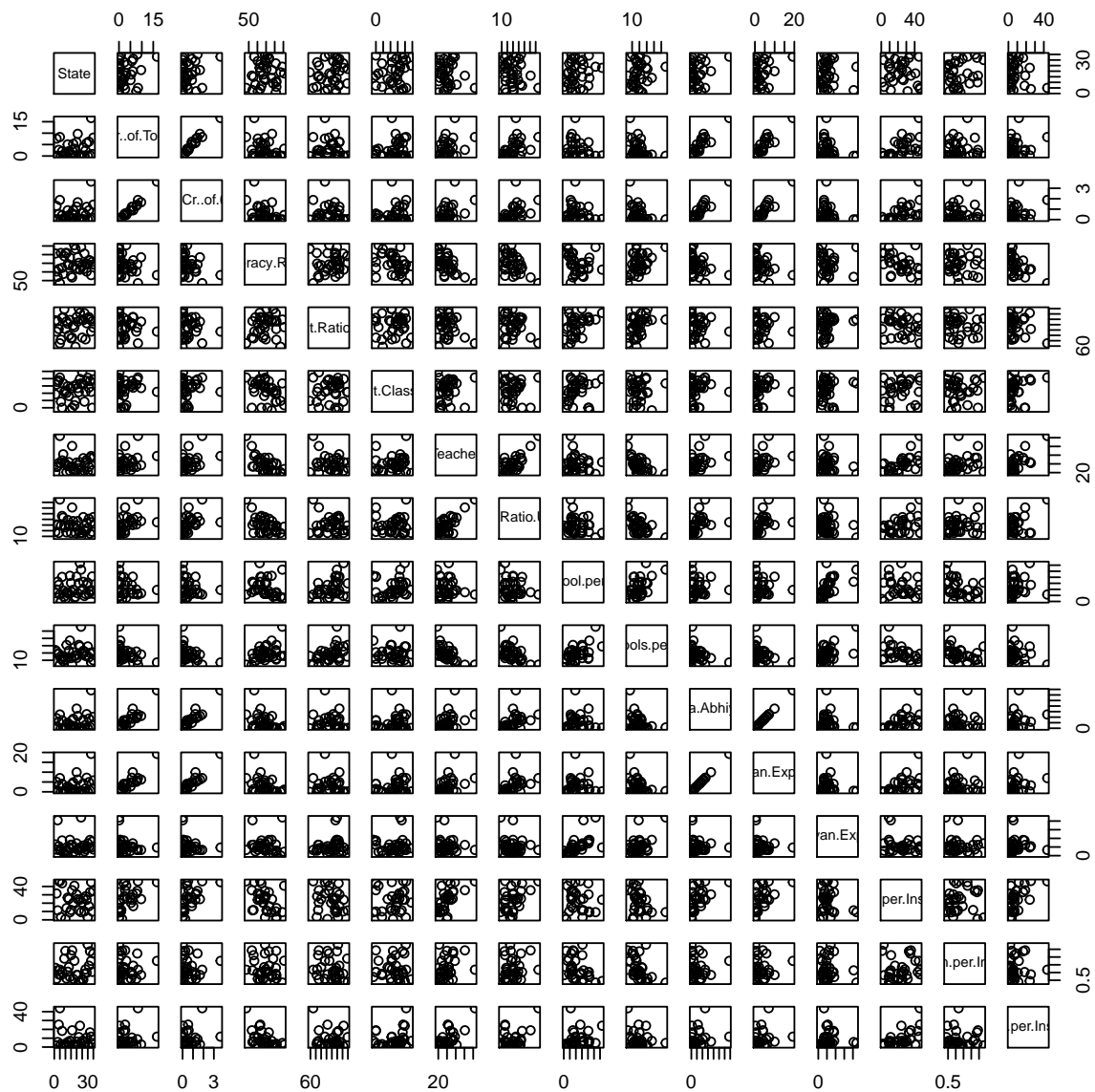
## 'data.frame': 35 obs. of 16 variables:
## $ State : chr "Andhra Pradesh" "Arunachal Pradesh" "Assam" "Bihar" "Chhattisgarh" "Goa" "Gujarat" "Haryana" "Himachal Pradesh" "Jammu and Kashmir" "Jharkhand" "Karnataka" "Kerala" "Madhya Pradesh" "Maharashtra" "Manipur" "Meghalaya" "Mizoram" "Nagaland" "Odisha" "Punjab" "Rajasthan" "Sikkim" "Tamil Nadu" "Telangana" "Tripura" "Uttar Pradesh" "Uttarakhand" "West Bengal"
## $ Population.in..Cr..of.Total.Census.2001 : num 7.62 0.11 2.67 8.3 2.08 0.003 0.004 0.005 0.006 0.007 0.008 0.009 0.01 0.011 0.012 0.013 0.014 0.015 0.016 0.017 0.018 0.019 0.02 0.021 0.022 0.023 0.024 0.025 0.026 0.027 0.028 0.029 0.03 0.031 0.032 0.033 0.034 0.035 0.036 0.037 0.038 0.039 0.04 0.041 0.042 0.043 0.044 0.045 0.046 0.047 0.048 0.049 0.05 0.051 0.052 0.053 0.054 0.055 0.056 0.057 0.058 0.059 0.06 0.061 0.062 0.063 0.064 0.065 0.066 0.067 0.068 0.069 0.07 0.071 0.072 0.073 0.074 0.075 0.076 0.077 0.078 0.079 0.08 0.081 0.082 0.083 0.084 0.085 0.086 0.087 0.088 0.089 0.09 0.091 0.092 0.093 0.094 0.095 0.096 0.097 0.098 0.099 0.1 0.101 0.102 0.103 0.104 0.105 0.106 0.107 0.108 0.109 0.11 0.111 0.112 0.113 0.114 0.115 0.116 0.117 0.118 0.119 0.12 0.121 0.122 0.123 0.124 0.125 0.126 0.127 0.128 0.129 0.13 0.131 0.132 0.133 0.134 0.135 0.136 0.137 0.138 0.139 0.14 0.141 0.142 0.143 0.144 0.145 0.146 0.147 0.148 0.149 0.15 0.151 0.152 0.153 0.154 0.155 0.156 0.157 0.158 0.159 0.16 0.161 0.162 0.163 0.164 0.165 0.166 0.167 0.168 0.169 0.17 0.171 0.172 0.173 0.174 0.175 0.176 0.177 0.178 0.179 0.18 0.181 0.182 0.183 0.184 0.185 0.186 0.187 0.188 0.189 0.19 0.191 0.192 0.193 0.194 0.195 0.196 0.197 0.198 0.199 0.2 0.201 0.202 0.203 0.204 0.205 0.206 0.207 0.208 0.209 0.21 0.211 0.212 0.213 0.214 0.215 0.216 0.217 0.218 0.219 0.22 0.221 0.222 0.223 0.224 0.225 0.226 0.227 0.228 0.229 0.23 0.231 0.232 0.233 0.234 0.235 0.236 0.237 0.238 0.239 0.24 0.241 0.242 0.243 0.244 0.245 0.246 0.247 0.248 0.249 0.25 0.251 0.252 0.253 0.254 0.255 0.256 0.257 0.258 0.259 0.26 0.261 0.262 0.263 0.264 0.265 0.266 0.267 0.268 0.269 0.27 0.271 0.272 0.273 0.274 0.275 0.276 0.277 0.278 0.279 0.28 0.281 0.282 0.283 0.284 0.285 0.286 0.287 0.288 0.289 0.29 0.291 0.292 0.293 0.294 0.295 0.296 0.297 0.298 0.299 0.3 0.301 0.302 0.303 0.304 0.305 0.306 0.307 0.308 0.309 0.31 0.311 0.312 0.313 0.314 0.315 0.316 0.317 0.318 0.319 0.32 0.321 0.322 0.323 0.324 0.325 0.326 0.327 0.328 0.329 0.33 0.331 0.332 0.333 0.334 0.335 0.336 0.337 0.338 0.339 0.34 0.341 0.342 0.343 0.344 0.345 0.346 0.347 0.348 0.349 0.35 0.351 0.352 0.353 0.354 0.355 0.356 0.357 0.358 0.359 0.36 0.361 0.362 0.363 0.364 0.365 0.366 0.367 0.368 0.369 0.37 0.371 0.372 0.373 0.374 0.375 0.376 0.377 0.378 0.379 0.38 0.381 0.382 0.383 0.384 0.385 0.386 0.387 0.388 0.389 0.39 0.391 0.392 0.393 0.394 0.395 0.396 0.397 0.398 0.399 0.4 0.401 0.402 0.403 0.404 0.405 0.406 0.407 0.408 0.409 0.41 0.411 0.412 0.413 0.414 0.415 0.416 0.417 0.418 0.419 0.42 0.421 0.422 0.423 0.424 0.425 0.426 0.427 0.428 0.429 0.43 0.431 0.432 0.433 0.434 0.435 0.436 0.437 0.438 0.439 0.44 0.441 0.442 0.443 0.444 0.445 0.446 0.447 0.448 0.449 0.45 0.451 0.452 0.453 0.454 0.455 0.456 0.457 0.458 0.459 0.46 0.461 0.462 0.463 0.464 0.465 0.466 0.467 0.468 0.469 0.47 0.471 0.472 0.473 0.474 0.475 0.476 0.477 0.478 0.479 0.48 0.481 0.482 0.483 0.484 0.485 0.486 0.487 0.488 0.489 0.49 0.491 0.492 0.493 0.494 0.495 0.496 0.497 0.498 0.499 0.5 0.501 0.502 0.503 0.504 0.505 0.506 0.507 0.508 0.509 0.51 0.511 0.512 0.513 0.514 0.515 0.516 0.517 0.518 0.519 0.52 0.521 0.522 0.523 0.524 0.525 0.526 0.527 0.528 0.529 0.53 0.531 0.532 0.533 0.534 0.535 0.536 0.537 0.538 0.539 0.54 0.541 0.542 0.543 0.544 0.545 0.546 0.547 0.548 0.549 0.55 0.551 0.552 0.553 0.554 0.555 0.556 0.557 0.558 0.559 0.56 0.561 0.562 0.563 0.564 0.565 0.566 0.567 0.568 0.569 0.57 0.571 0.572 0.573 0.574 0.575 0.576 0.577 0.578 0.579 0.58 0.581 0.582 0.583 0.584 0.585 0.586 0.587 0.588 0.589 0.59 0.591 0.592 0.593 0.594 0.595 0.596 0.597 0.598 0.599 0.6 0.601 0.602 0.603 0.604 0.605 0.606 0.607 0.608 0.609 0.61 0.611 0.612 0.613 0.614 0.615 0.616 0.617 0.618 0.619 0.62 0.621 0.622 0.623 0.624 0.625 0.626 0.627 0.628 0.629 0.63 0.631 0.632 0.633 0.634 0.635 0.636 0.637 0.638 0.639 0.64 0.641 0.642 0.643 0.644 0.645 0.646 0.647 0.648 0.649 0.65 0.651 0.652 0.653 0.654 0.655 0.656 0.6
```

```
## $ Literacy.Rate : num 60.5 54.3 63.3 47 64.7 82
## $ Gross.Enrollment.Ratio..Classes...I.VIII. : num 87 106.7 91.9 65.2 112.6
## $ Drop.out.Classes..I.X. : num 63.7 70.8 75 83.1 0 ...
## $ Pupil.Teacher.Ratio : int 33 34 42 104 48 21 35 44
## $ Pupil.Teacher.Ratio.Upper.Primary : int 31 30 16 75 46 17 39 30 3
## $ Elementary.School.per.lakh.population : int 99 163 137 57 203 76 73 6
## $ Sec.Hr.Sec.Schools.per.lakh.population : int 22 19 19 4 12 31 14 23 37
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Rs.cr. : num 1807 187 1165 2479 1440
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure.as.percentage.of.Total : num 5.09 0.53 3.28 6.99 4.06
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age : int 1387 8179 2186 1319 3513
## $ Lakh.of.Population.per.Institution..University. : num 31.6 11.4 46.5 46.2 43.8
## $ Lakh.of.Population.per.Institution..College. : num 0.59 1.14 0.88 1.18 1.03
## $ Lakh.of.Population.per.Institution..Technical. : num 1.97 5.7 25.34 43.87 18.2
```

1 Correlation

It is possible to plot the complete data – each column against all the other columns. This can be done by simply calling the plot function:

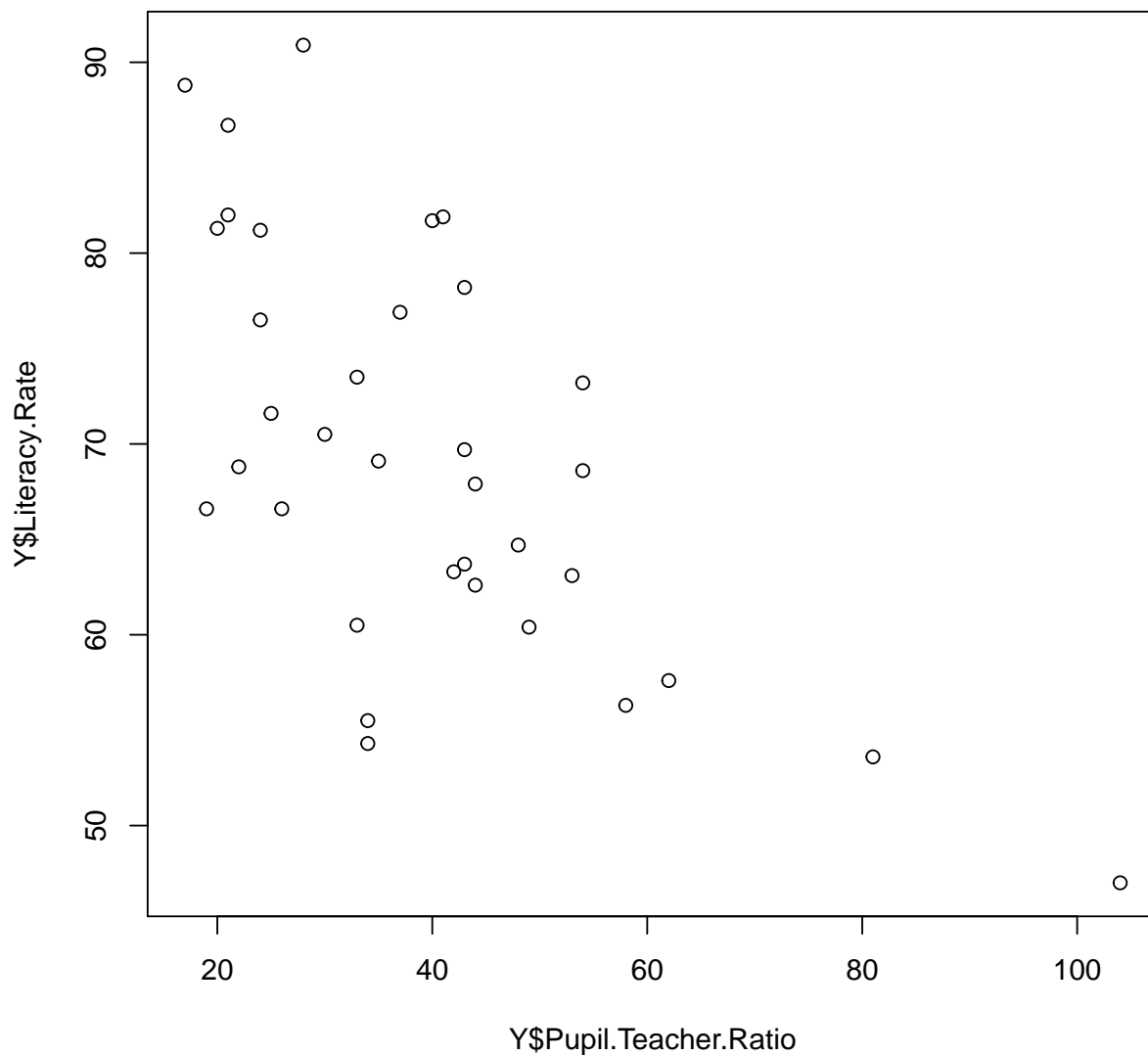
```
plot(Y)
```



From these plots, once you zoom in and see, it is easy to notice the trends and correlations, if any!

Of course, one can also plot specific variables. For example, is literacy rate related to the pupil teacher ratio? Let us check!

```
plot(Y$Pupil.Teacher.Ratio,Y$Literacy.Rate)
```



The overall trend seems to be that the literacy rate decreases as the pupil teacher ratio increases.

How about the correlation coefficient? As you have learnt from the lecture, the correlation coefficient is defined as follows:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

So, let us write an R script that does this. We are going to calculate the coefficient of correlation in two different ways. In one case, we are going to use the `sum` function. In the other, we are going to explicitly run a loop.

```
x <- Y$Pupil.Teacher.Ratio
y <- Y$Literacy.Rate
xbar = mean(x) # Calculate the mean x
ybar = mean(y) # Calculate the mean y
```

```

xp <- x-xbar # Define xprime = x - xbar
yp <- y-ybar # Define yprime = y - ybar
r <- sum(xp*yp)/sqrt(sum(xp*xp)*sum(yp*yp))
r

## [1] -0.6509299

# Method 2
n <- length(x)
sum1 = 0
sum2 = 0
sum3 = 0
for(i in 1:n){
  sum1 = sum1 + (x[i]-xbar)*(y[i]-ybar)
  sum2 = sum2 + (x[i]-xbar)*(x[i]-xbar)
  sum3 = sum3 + (y[i]-ybar)*(y[i]-ybar)
}
r = sum1/(sqrt(sum2*sum3))
r

## [1] -0.6509299

```

As one can see, there is correlation and it is negative.

Homework

- Plot the literacy rate against the elementary schools per lakh population. From the plot, can you see if they are correlated? Calculate the correlation coefficient and confirm your answer.
- Plot the literacy rate against the tenth plan Sarva Siksha Abhiyan expenditure (per capita in the 6 to 14 age group). From the plot, can you see if they are correlated? Calculate the correlation coefficient and confirm your answer.
- In terms of spending per capita, number of elementary schools and pupil to teachers ratio, which is the most effective in terms of improving literacy rate? Order them in light of the given data and your analysis.

2 Ordering and plotting data

Let us pick the top ten states / union territories in terms of the literacy rate and plot their literacy rate against the spending (per capita in the relevant age group). In this plot, we also want to label the data according to their state. The following script does that.

```

TopTen <- head(Y[order(Y$Literacy.Rate,decreasing=TRUE),],10) # Pick the top ten states
str(TopTen)

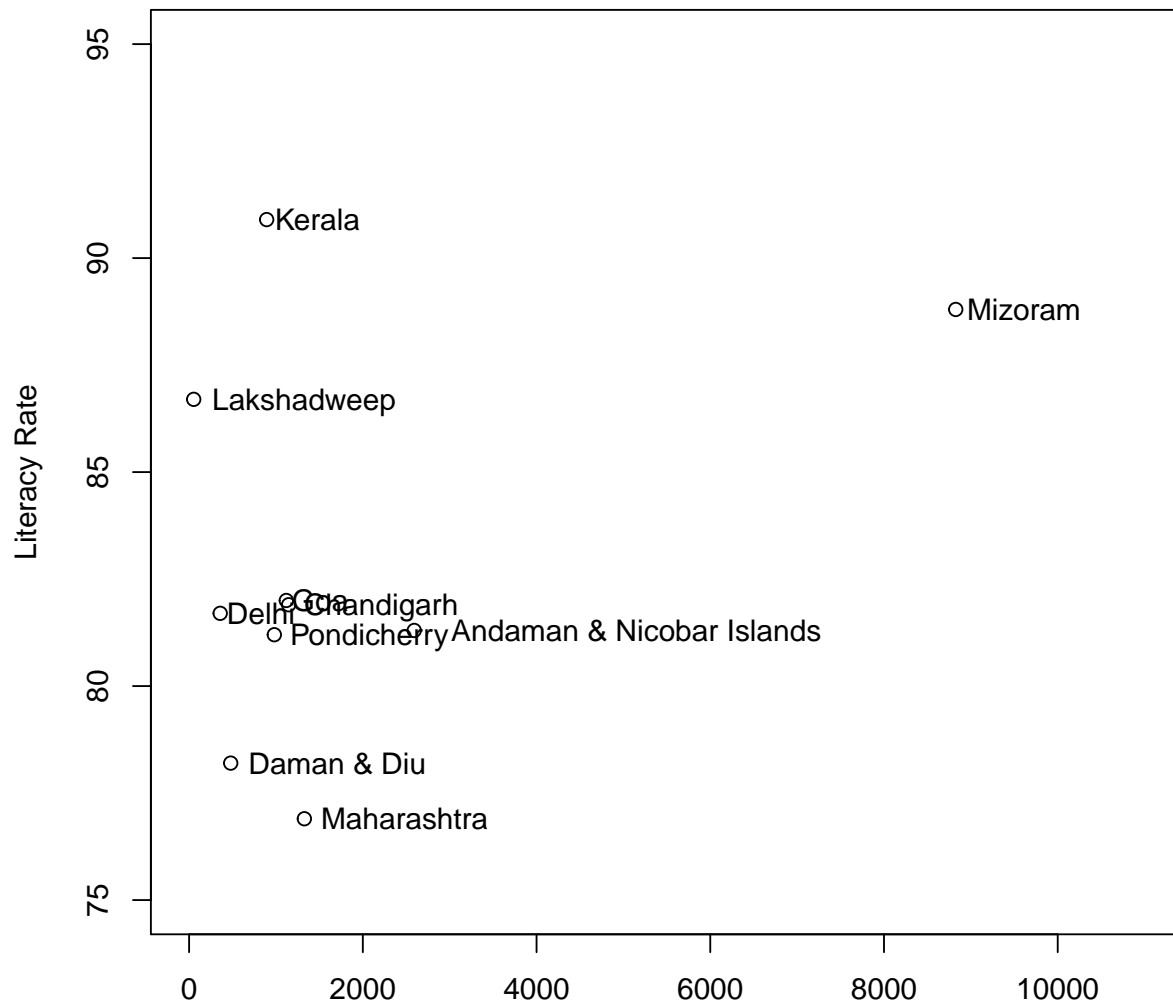
## 'data.frame': 10 obs. of 16 variables:
## $ State                                     : chr  "Kerala" "Mizoram" "Lakshadweep"
## $ Population.in..Cr..of.Total.Census.2001 : num  3.18 0.09 0.01 0.13 0.09
## $ Population.in..Cr..of.6.14.age.2004      : num  0.44 0.02 0 0.02 0.01 0.2
## $ Literacy.Rate                             : num  90.9 88.8 86.7 82 81.9 81

```

```
## $ Gross.Enrollment.Ratio..Classes...I.VIII. : num 95.3 109.5 58.8 106 71.9
## $ Drop.out.Classes..I.X. : num 7.15 66.95 18.88 40.65 16
## $ Pupil.Teacher.Ratio : int 28 17 21 21 41 40 20 24 4
## $ Pupil.Teacher.Ratio.Upper.Primary : int 27 8 16 17 29 26 18 21 29
## $ Elementary.School.per.lakh.population : int 30 263 38 76 3 20 69 45 4
## $ Sec.Hr.Sec.Schools.per.lakh.population : int 16 56 16 31 12 11 24 22
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Rs.cr. : num 390.7 141.2 0.1 19 16 ..
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure.as.percentage.of.Total : num 1.1 0.4 0 0.05 0.05 0.25
## $ Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age : int 893 8826 53 1119 1146 355
## $ Lakh.of.Population.per.Institution..University. : num 40.9 9.2 NA 14.1 3.33 1 M
## $ Lakh.of.Population.per.Institution..College. : num 1.76 0.35 NA 0.61 0.83 2
## $ Lakh.of.Population.per.Institution..Technical. : num 2.68 4.6 NA 1.57 3.33 3.6
```

```
plot(Literacy.Rate~Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age,data=TopTen,main="Literacy Rate vs Tenth Plan Sarva Siksha Abhiyan Expenditure")
text(Literacy.Rate~Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age,data=TopTen,labels=)
```

Expenditure versus Literacy Rate



Tenth plan Sarva Siksha Abhiyan Expenditure (per capita of population in the age group 6–14 y)

From the figure, one can see that Mizoram is an outlier in terms of spending.

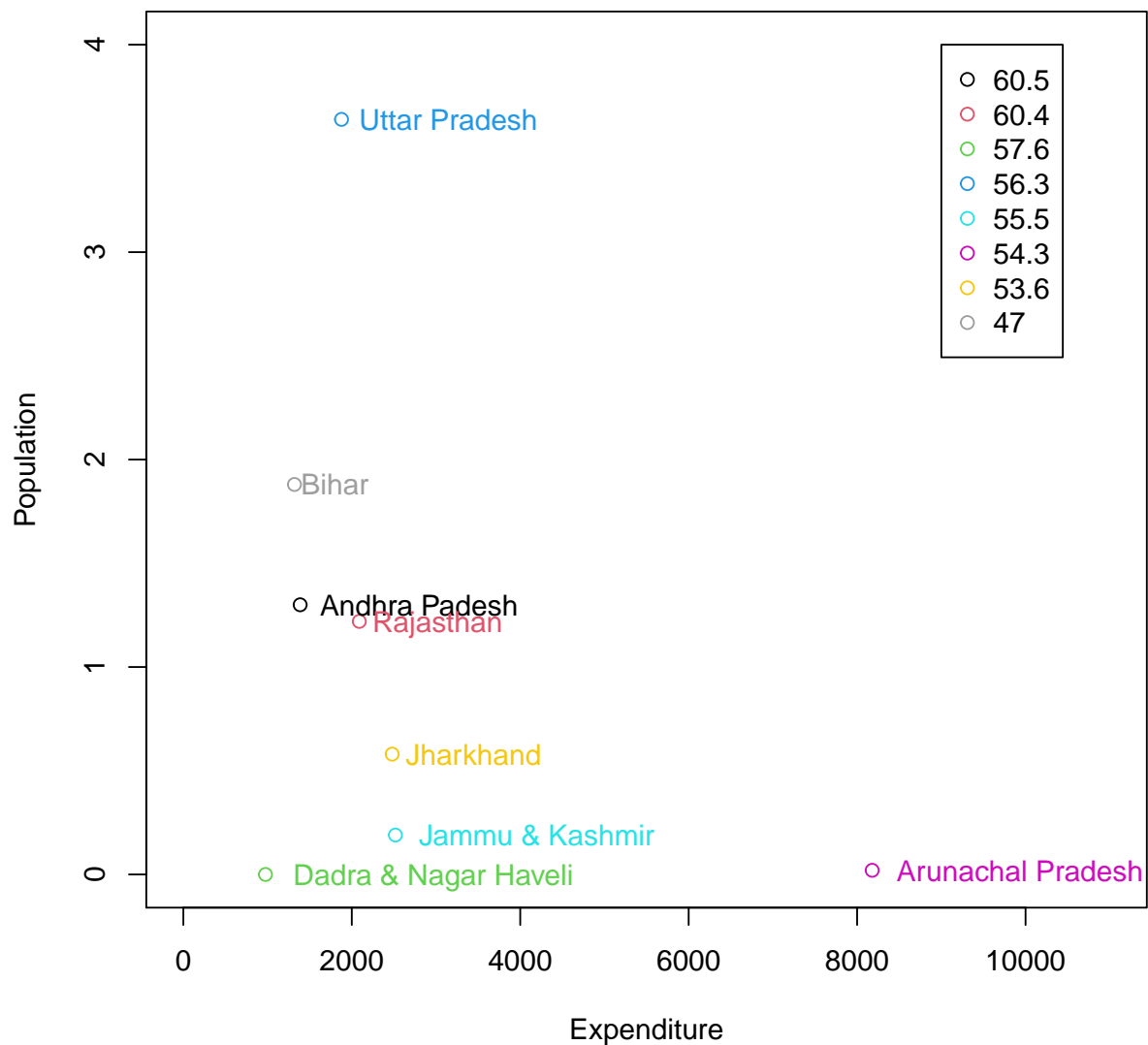
Homework

- In the above script, what happens to the plot if you remove the `xlim` and `ylim` commands?
- Pick the states which are in the bottom ten in terms of literacy and make a plot of spending versus literacy rate. Do you see any outliers?
- Plot two pie-charts next to each other – one of spending and another of literacy rate – for states in the bottom ten. Below them, plot the same two pie-charts for the states in top ten. What do you see?

3 Adding colour!

Adding colour to the plots can help bring in one more dimension to the plots. Here is an example of the expenditure versus population of the bottom eight states plotted with a colour marker which indicated their literacy rate; of course, we have also labeled the data points to indicate the State / UT.

```
BottomEight <- tail(Y[order(Y$Literacy.Rate,decreasing=TRUE),],8)
# Pick the bottom eight states
plot(BottomEight$Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age,
      BottomEight$Population.in..Cr..of.6.14.age.2004,col=1:length(BottomEight$Literacy.Rate),
      xlab="Expenditure",ylab="Population",
      xlim=c(0,11000),ylim=c(0,4))
legend(9000,4,BottomEight$Literacy.Rate,col=1:length(BottomEight$Literacy.Rate),pch=1)
text(BottomEight$Tenth.Plan.Sarva.Siksha.Abhiyan.Expenditure...Per.capita.6.14.age,
      BottomEight$Population.in..Cr..of.6.14.age.2004,
      labels=(BottomEight[1:length(BottomEight$Literacy.Rate),1]),
      adj=c(-0.1,0.5),col=1:length(BottomEight$Literacy.Rate))
```

In the above figure, if I want the names of the states / UTs printed below the point and not next to it, what modification needs to be done?

As a final exercise, plot the histograms of (a) Lakh of Population per University, (b) Tenth plan expenditure, and (c) literacy rate. What can you comment about these distributions? Which data set looks closer to normal, for example? Which data set is skewed? In which direction?

Of course, with the given data set, you can plot many more plots and have much more fun! Happy programming!!