

A decorative graphic on the left side of the slide, consisting of white lines and circles on a blue gradient background, resembling a circuit board or a stylized tree structure.

# READING PDF FILES

REBECCA DUONG

- Goals:
  - Extract the tables and the data from Pdf
  - Find the country with the highest % of Child Marriage by age 15 and 18 (using the data from the textbook by Jacqueline Kazil, Katherine Jarmul)
- Python Libraries that are available for reading PDF files:
  - pdfminer
  - pdftotext
  - tabula-py
  - PyPDF2
  - pdftables
- The library I will use for the tutorial:
  - PyPDF2
  - pdftables
  - matplotlib
  - pandas

**TABLE 9 | CHILD PROTECTION**

Countries and areas	Child labour (%) <sup>a</sup> 2005–2012*			Child marriage (%) 2005–2012*		Birth registration (%) <sup>a</sup> 2005–2012*	Female genital mutilation/cutting (%) <sup>a</sup> 2002–2012*			Justification of wife beating (%) 2005–2012*		Violent discipline (%) <sup>a</sup> 2005–2012*		
							prevalence		attitudes					
	total	male	female	married by 15	married by 18		total	women <sup>a</sup>	girls <sup>a</sup>	support for the practice <sup>a</sup>	male	female	total	male
Afghanistan	10	11	10	15	40	37	–	–	–	–	90	74	75	74
Albania	12	14	9	0	10	99	–	–	–	36	30	75	78	71
Algeria	5 y	6 y	4 y	0	2	99	–	–	–	–	68	88	89	87
Andorra	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Angola	24 x	22 x	25 x	–	–	36 x	–	–	–	–	–	–	–	–
Antigua and Barbuda	–	–	–	–	–	–	–	–	–	–	–	–	–	–
Argentina	7 y	8 y	5 y	–	–	99 y	–	–	–	–	–	–	–	–
Armenia	4	5	3	0	7	100	–	–	–	20	9	70	72	67
Australia	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Austria	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Azerbaijan	7 y	8 y	5 y	1	12	94	–	–	–	58	49	75	79	71
Bahamas	–	–	–	–	–	–	–	–	–	–	–	–	–	–
Bahrain	5 x	6 x	3 x	–	–	–	–	–	–	–	–	–	–	–
Bangladesh	13	18	8	29	65	31	–	–	–	–	33 y	–	–	–
Barbados	–	–	–	–	–	–	–	–	–	–	–	–	–	–
Belarus	1	1	2	0	3	100 y	–	–	–	4	4	65 y	67 y	62 y
Belgium	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Belize	6	7	5	3	26	95	–	–	–	–	9	71	71	70
Benin	46	47	45	8	34	80	13	2 y	1	14	47	–	–	–
Bhutan	3	3	3	6	26	100	–	–	–	–	68	–	–	–
Bolivia (Plurinational State of)	26 y	28 y	24 y	3	22	76 y	–	–	–	–	16	–	–	–
Bosnia and Herzegovina	5	7	4	0	4	100	–	–	–	6	5	55	60	50
Botswana	9 y	11 y	7 y	–	–	72	–	–	–	–	–	–	–	–
Brazil	9 y	11 y	6 y	11	36	93 y	–	–	–	–	–	–	–	–
Brunei Darussalam	–	–	–	–	–	–	–	–	–	–	–	–	–	–
Bulgaria	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Burkina Faso	39	42	36	10	52	77	76	13	9	34	44	83	84	82
Burundi	26	26	27	3	20	75	–	–	–	44	73	–	–	–
Cabo Verde	3 x,y	4 x,y	3 x,y	3	18	91	–	–	–	16 y	17	–	–	–
Cambodia	36 y	36 y	36 y	2	18	62	–	–	–	22 y	46 y	–	–	–
Cameroon	42	43	40	13	38	61	1	1 y	7	39	47	93	93	93
Canada	–	–	–	–	–	100 v	–	–	–	–	–	–	–	–
Central African Republic	29	27	30	29	68	61	24	1	11	80 y	80	92	92	92
Chad	26	25	28	29	68	16	44	18 y	38	–	62	84	85	84
Chile	3 x	3 x	2 x	–	–	100 y	–	–	–	–	–	–	–	–
China	–	–	–	–	–	–	–	–	–	–	–	–	–	–
Colombia	13 y	17 y	9 y	6	23	97	–	–	–	–	–	–	–	–
Comoros	27 x	26 x	28 x	–	–	88 x	–	–	–	–	–	–	–	–
Congo	25	24	25	7	33	91	–	–	–	–	76	–	–	–

# USING PDFTABLES

```
from pdftables import get_tables
import pprint
```

```
# The get_tables function return each pages as its own table
# Each of those tables has a list of rows with a contained list of columns
all_tables = get_tables(open('EN-FINAL Table 9.pdf', 'rb'))
```

```
# Trying to find the titles we can use for our columns
# Here we are looking at the first page's first 6 rows
print all_tables[0][:6]
```

# NOTE: we can see that the titles are included in the first 3 list and they are messy

```
# Manually setting up our titles
# So we add all of our headers, including the country names, to one list
headers = ['Country', 'Child Labor 2005-2012 (%) total',
           'Child Labor 2005-2012 (%) male',
           'Child Labor 2005-2012 (%) female',
           'Child Marriage 2005-2012 (%) married by 15',
           'Child Marriage 2005-2012 (%) married by 18',
           'Birth registration 2005-2012 (%)',
           'Female Genital mutilation 2002-2012 (prevalence), women',
           'Female Genital mutilation 2002-2012 (prevalence), girls',
           'Female Genital mutilation 2002-2012 (support)',
           'Justification of wife beating 2005-2012 (%) male',
           'Justification of wife beating 2005-2012 (%) female',
           'Violent discipline 2005-2012 (%) total',
           'Violent discipline 2005-2012 (%) male',
           'Violent discipline 2005-2012 (%) female']
```



```
final_data = []
```

```
# isolates only the rows for each page we want; slice from the 5th index onward
```

```
#if data row is missing index 0, it has no country name and is a blank row
```

```
if roww[0] == " " or roww[0][0].isdigit():
```

```
# if the data row is missing index 2, we know this is probably the first part of a country name
```

```
elif row[2] == ":
```

```
continue
```

```
# Manipulate the country name entry in the row if it has a first_name
```

# Put the second part of a country name back together with the first name

```
if first_name:
```

```
row[0] = u'{} {}'.format(first_name, row[0])
```

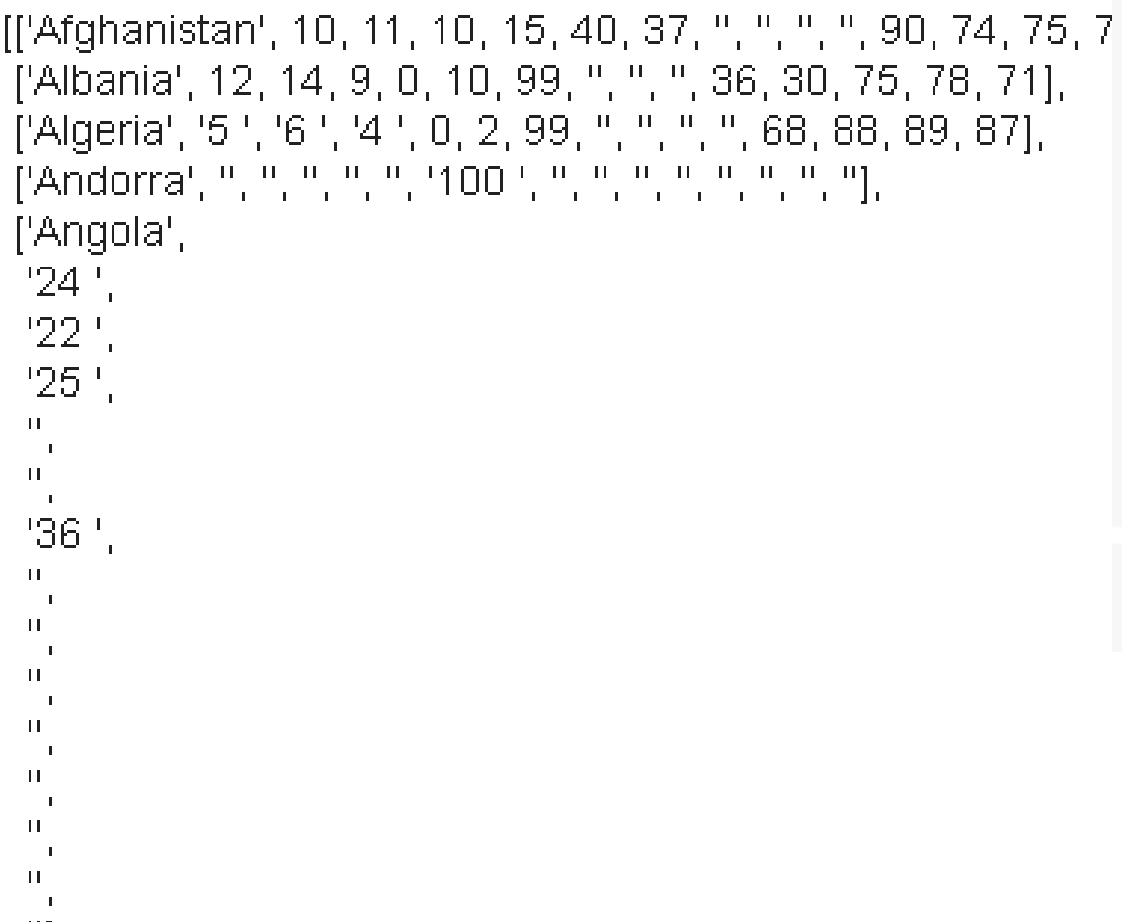
```
# set first_name back to False, so our next iteration operates properly
```

```
first_name = False
```

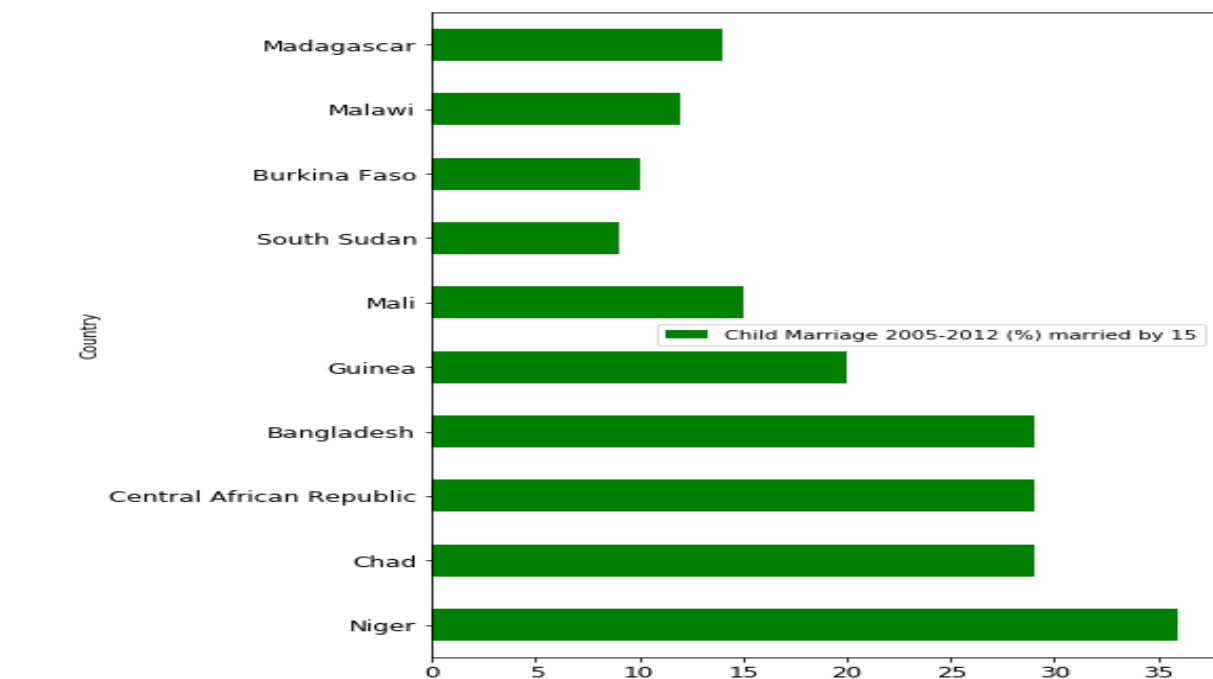
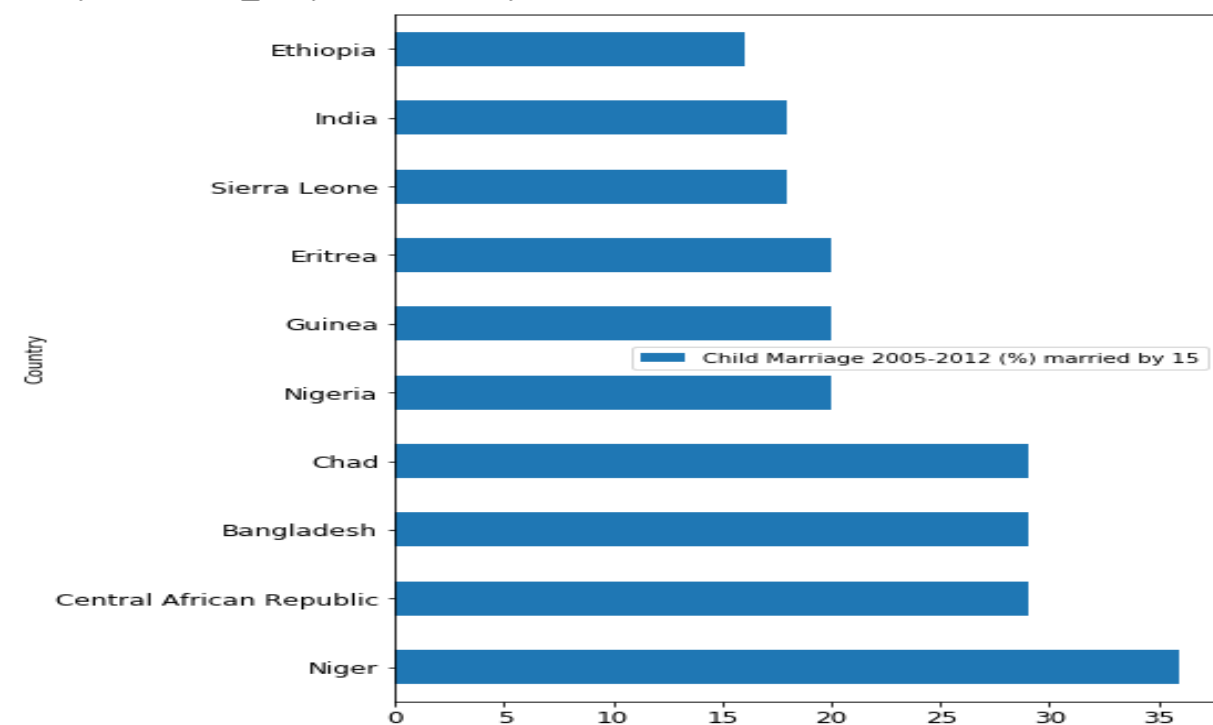
```
# Before cleaning up the messy data...
```

```
# Note there are in unicode
```

```
# print(row)
```



```
# prints out a list of each row in the table
pprint.pprint(final_data)
```



```
import pandas as pd
```

```
# convert the list of dict into a dataframe and set the 'Country' column as the index
df = pd.DataFrame(final_data, columns = headers)
df = df.set_index("Country")
df = df.apply(pd.to_numeric)
# df.columns
df
```

```
import matplotlib.pyplot as plt
```

```
# Find the country with the highest % of child marriage in 2005-2012 of age 15 and 18
```

```
fig, axes = plt.subplots(nrows = 1, ncols = 2)
```

```
# adjust the space between the 2 plots; set it a little more farther apart
plt.subplots_adjust(wspace = 1)
```

```
# Sort the dataframe according to its % of child marriage by age 15 and 18; the highest % will be at the top
# store these into a new variable
df_1 = df.sort_values('Child Marriage 2005-2012 (%) married by 15', ascending = False)
# df_1
df_2 = df.sort_values('Child Marriage 2005-2012 (%) married by 18', ascending = False)
# df_2
```

```
# Looking at the data, it seems like Niger the highest % of Child Marriage in 2005-2012 for both age 15 and 18
df_1.head(10).plot(ax = axes[0], kind = 'barh', y = ['Child Marriage 2005-2012 (%) married by 15'], fontdict = {'color': 'red'})
df_2.head(10).plot(color = 'g', ax = axes[1], kind = 'barh', y = ['Child Marriage 2005-2012 (%) married by 18'], fontdict = {'color': 'red'})
```

```
# Convert the dataframe to csv
df.to_csv('table_9.csv')
```

# USING PYPDF2 LIBRARY

Countries and areas	Child labour (%) 2005-2012*			Child marriage (%) 2005-2012*			Birth registration (%) 2005-2012*	Female genital mutilation/cutting (%) 2005-2012*			Justification of wife beating (%) 2005-2012*			Violent discipline (%) 2005-2012*		
	total	male	female	married by 15	married by 18	total		women*	girls*	support for the practice*	male	female	total	male	female	
Afghanistan	10	11	10	15	40	37	—	—	—	90	36	74	75	74		
Albania	12	14	9	0	10	99	—	—	—	30	78	71	89	87		
Algeria	5	6	4	0	2	99	—	—	—	68	88	87	88	89		
Andorra	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Angola	24	22	25	—	—	36	—	—	—	—	—	—	—	—		
Antigua and Barbuda	—	—	—	—	—	—	—	—	—	—	—	—	—	—		
Argentina	7	8	5	—	—	99	—	—	—	—	—	—	—	—		
Armenia	4	5	3	0	7	100	—	—	—	20	70	67	72	72		
Australia	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Austria	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Azerbaijan	7	8	5	1	12	94	—	—	—	58	49	71	75	79		
Bahamas	—	—	—	—	—	—	—	—	—	—	—	—	—	—		
Bahrain	5	6	3	—	—	—	—	—	—	—	—	—	—	—		
Bangladesh	13	18	8	29	65	31	—	—	—	33	—	—	—	—		
Barbados	—	—	—	—	—	—	—	—	—	—	—	—	—	—		
Belarus	1	1	2	0	3	100	—	—	—	4	4	62	65	67		
Belgium	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Belize	6	7	5	3	26	95	—	—	—	9	71	70	71	71		
Benin	46	47	45	8	34	80	—	—	—	14	47	—	—	—		
Bhutan	3	3	3	6	26	100	—	—	—	68	—	—	—	—		
Bolivia (Plurinational State of)	26	28	24	3	22	76	—	—	—	—	16	—	—	—		
Bosnia and Herzegovina	5	7	4	0	4	100	—	—	—	6	5	50	55	60		
Botswana	9	11	7	—	—	72	—	—	—	—	—	—	—	—		
Brazil	9	11	6	—	—	93	—	—	—	—	—	—	—	—		
Brunei Darussalam	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Bulgaria	—	—	—	—	—	—	—	—	—	—	—	—	—	—		
Burkina Faso	39	42	36	10	52	77	—	—	—	44	44	82	83	84		
Burundi	26	26	27	3	20	75	—	—	—	44	73	—	—	—		
Cabo Verde	3	4	3	3	18	91	—	—	—	16	17	—	—	—		
Cambodia	36	36	36	2	18	62	—	—	—	22	46	—	—	—		
Cameroon	42	43	40	13	38	61	—	—	—	39	47	93	93	93		
Canada	—	—	—	—	—	100	—	—	—	—	—	—	—	—		
Central African Republic	29	27	30	29	68	61	—	—	—	80	80	92	92	92		
Chad	26	25	28	29	68	16	—	—	—	—	—	84	85	84		
Chile	3	3	2	—	—	100	—	—	—	—	—	—	—	—		
China	—	—	—	—	—	—	—	—	—	—	—	—	—	—		
Colombia	13	17	9	6	23	97	—	—	—	—	—	—	—	—		
Comoros	27	26	28	—	—	88	—	—	—	—	—	—	—	—		
Congo	25	24	25	7	33	91	—	—	—	—	76	—	—	—		

TABLE 9 | CHILD PROTECTION

```
# A function for rotating the page
def PDFrotate(origFileName, newFileName, rotation):
```

```
# creating a pdf File object of original pdf
pdfFileObj = open(origFileName, 'rb')
```

```
# creating a pdf Reader object
pdfReader = PyPDF2.PdfFileReader(pdfFileObj)
```

```
# creating a pdf writer object for new pdf
# rotated pages will be written to a new pdf
pdfWriter = PyPDF2.PdfFileWriter()
```

```
# iterate through each page of original pdf
for page in range(pdfReader.numPages):
```

```
# creating rotated page object
pageObj = pdfReader.getPage(page)
```

```
# now we rotate the page by rotateClockwise() method
pageObj.rotateClockwise(rotation)
```

```
# adding rotated page object to pdf writer
pdfWriter.addPage(pageObj)
```

```
# open new pdf file object
newFile = open(newFileName, 'wb')
```

```
# writing rotated pages to new file using write() method
pdfWriter.write(newFile)
```

```
# closing the original pdf file object
pdfFileObj.close()
```

```
# closing the new pdf file object
newFile.close()
```







Algeria	5 y	8 y	4 y	8	2	99	--	--	--	--	88	88	88	87
Andorra	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Angola	24 x	22 x	26 x	--	--	26 x	--	--	--	--	--	--	--	--
Antigua and Barbuda	--	--	--	--	--	--	--	--	--	--	--	--	--	--
Argentina	7 y	8 y	5 y	--	--	99 y	--	--	--	--	--	--	--	--
Armenia	4	5	3	8	7	100	--	--	--	20	9	78	72	87
Australia	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Austria	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Azerbaijan	7 y	8 y	5 y	1	12	94	--	--	--	58	49	75	79	71
Bahamas	--	--	--	--	--	--	--	--	--	--	--	--	--	--
Bahrain	5 x	8 x	3 x	--	--	--	--	--	--	--	--	--	--	--
Bangladesh	12	18	8	28	65	21	--	--	--	22 y	--	--	--	--
Barbados	--	--	--	--	--	--	--	--	--	--	--	--	--	--
Belarus	1	1	2	8	3	100 y	--	--	--	4	4	66 y	67 y	60 y
Belgium	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Belize	8	7	5	3	26	95	--	--	--	9	--	71	71	76
Benin	48	47	45	8	34	80	12	2 y	1	14	47	--	--	--
Bhutan	3	3	3	8	26	100	--	--	--	68	--	--	--	--
Bolivia (Plurinational State of)	28 y	28 y	28 y	3	22	76 y	--	--	--	16	--	--	--	--
Bosnia and Herzegovina	5	7	4	8	4	100	--	--	--	9	5	54	60	58
Botswana	9 y	11 y	7 y	--	--	72	--	--	--	--	--	--	--	--
Brazil	9 y	11 y	8 y	11	26	93 y	--	--	--	--	--	--	--	--
Brunei Darussalam	--	--	--	--	--	--	--	--	--	--	--	--	--	--
Bulgaria	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Burkina Faso	28	42	26	10	52	77	76	12	9	26	44	82	84	82
Burundi	28	28	27	3	20	75	--	--	--	44	72	--	--	--
Cabo Verde	2 x,y	4 x,y	3 x,y	3	18	91	--	--	--	16 y	17	--	--	--
Cambodia	28 y	28 y	28 y	2	18	82	--	--	--	22 y	58 y	--	--	--
Cameroon	42	42	40	12	28	61	1	1 y	7	29	47	92	92	92
Canada	--	--	--	--	--	100 v	--	--	--	--	--	--	--	--
Central African Republic	28	27	28	28	68	61	24	1	11	80 y	80	92	92	92
Chad	28	25	28	28	68	16	44	18 y	28	--	82	84	85	84
Chile	2 x	2 x	2 x	--	--	100 y	--	--	--	--	--	--	--	--
China	--	--	--	--	--	--	--	--	--	--	--	--	--	--
Colombia	12 y	17 y	9 y	8	22	97	--	--	--	--	--	--	--	--
Comoros	27 x	28 x	28 x	--	--	38 x	--	--	--	--	--	--	--	--
Congo	25	24	25	7	22	91	--	--	--	--	--	--	--	--

78 THE STATE OF THE WORLD'S CHILDREN 2014 IN NUMBERS

TABLE 9 CHILD PROTECTION &gt;&gt;

Countries and areas	Child labor (%) 2008-2012*			Child marriage (%) 2008-2012*		Early registration (%) 2008-2012*	Female genital mutilation/cutting (%) 2008-2012*			
	total	male	female	married by 15	married by 18		prevalence	incidence	girls*	attitude support for the practice*
Cook Islands	--	--	--	--	--	--	--	--	--	--
Costa Rica	5 x	8 x	3 x	--	--	--	--	--	--	--
Côte d'Ivoire	26	25	28	10	22	65	28	10	14	4
Croatia	--	--	--	--	--	--	--	--	--	--
Cuba	--	--	--	9	40	100 y	--	--	--	--
Cyprus	27 x	28 x	28 x	--	--	100 v	--	--	--	--
Czech Republic	--	--	--	--	--	100 v	--	--	--	--
Democratic People's Republic of Korea	--	--	--	--	--	100	--	--	--	--
Democratic Republic of the Congo	15	12	17	9	28	28	--	--	--	76
Denmark	--	--	--	--	--	100 v	--	--	--	92
Djibouti	8	8	8	2	5	92	92	48 y	27	72
Dominica	--	--	--	--	--	--	--	--	--	72
Dominican Republic	12	18	8	12	41	82	--	--	--	4
Ecuador	8	7	8	4 x	22 x	90	--	--	--	67
Egypt	9 y	14 y	4 y	2	17	99 y	91	17	54	28 y
El Salvador	18 y	--	--	5	25	99	--	--	--	91
Equatorial Guinea	28 x	28 x	28 x	--	--	27 x	--	--	--	92
Eritrea	--	--	--	25 x	47 x	--	89	82 y	49	71 x
Estonia	--	--	--	--	--	100 v	--	--	--	--
Ethiopia	27	31	25	16	41	7	74	24	31	65
Fiji	--	--	--	--	--	--	--	--	--	72 y
Finland	--	--	--	--	--	100 v	--	--	--	--
France	--	--	--	--	--	100 v	--	--	--	--
Gabon	12	15	12	6	22	90	--	--	--	80
Gambia	18	21	15	7	26	52	76	56	64	76
Georgia	18	20	17	1	14	99	--	--	--	7
Germany	--	--	--	--	--	100 v	--	--	--	87
Ghana	24	24	24	5	21	82	4	1	2	26 y
Greece	--	--	--	--	--	100 v	--	--	--	84
Grenada	--	--	--	--	--	--	--	--	--	94
Guatemala	28 y	25 y	16 y	7	30	97	--	--	--	94
Guinea	42 y	40 y	40 y	20	62	42	95	57 y	69	--

## # Splitting PDF file

## # A function to split the pdf file

```
def PDFsplit(pdf, splits):
    # creating input pdf file object
    pdfFileObj = open(pdf, 'rb')

    # creating pdf reader object
    pdfReader = PyPDF2.PdfFileReader(pdfFileObj)

    # starting index of first slice
    start = 0

    # starting index of last slice
    end = splits[0]

    for i in range(len(splits)+1):
        # creating pdf writer object for (i+1)th split
        pdfWriter = PyPDF2.PdfFileWriter()

        # output pdf file name
        # str(i) is the xth number of split we are currently on
        outputpdf = pdf.split('.')[0] + str(i) + '.pdf'

        # adding pages to pdf writer object
        for page in range(start, end):
            pdfWriter.addPage(pdfReader.getPage(page))

        # writing split pdf pages to pdf file
        with open(outputpdf, "wb") as f:
            pdfWriter.write(f)

    # interchanging page split start position for next split
    start = end

    try:
        # setting split end position for next split
        end = splits[i+1]
    except IndexError:
        # setting split end position for last split
        end = pdfReader.numPages
```

# SOURCES/ SHAREABLE LINK

- Shareable Link of my Tutorial:
  - [https://colab.research.google.com/drive/1Qm7\\_N-UzV9E7Qucywjl2A\\_fpKWYXm3v](https://colab.research.google.com/drive/1Qm7_N-UzV9E7Qucywjl2A_fpKWYXm3v)
- Sources:
  - Jacqueline Kazil, Katherine Jarmul - Data Wrangling with Python (Ch.5)
  - <https://www.geeksforgeeks.org/working-with-pdf-files-in-python/>
  - <https://stackoverflow.com/questions/2365411/convert-unicode-to-ascii-without-errors-in-python/355362>