

Manqiong

2017 State of Data Science - Kaggle survey

▲
0

voters

↳ forked from 2017 State of Data Science - Kaggle survey by Anisotropic (+94/-173)

last run a minute ago · Python notebook · 10 views

using data from [multiple data sources](#) · 👁 Public

Tags

multiple data sources

Add Tag

Notebook

In [1]:

```
import numpy as np
import pandas as pd
import scipy as sp
import csv
%matplotlib inline
import matplotlib.pyplot as plt

import pandas as pd
import matplotlib.pyplot as plt; plt.rcParamsdefaults()
import numpy as np
import matplotlib.pyplot as plt

import seaborn as sns
import matplotlib.ticker as ticker
```

In [2]:

```
mcr=pd.read_csv('../input/kaggle-survey-2017/multipleChoiceResponses.csv',
encoding="ISO-8859-1", low_memory=False)
mcr.head()
```

Out[2]:

	GenderSelect	Country	Age	EmploymentStatus	StudentStatus	LearningDataScience	CodeWri
0	Non-binary, genderqueer, or gender non- conforming	NaN	NaN	Employed full-time	NaN	NaN	Yes
1	Female	United States	30.0	Not employed, but looking for work	NaN	NaN	NaN
2	Male	Canada	28.0	Not employed, but looking for work	NaN	NaN	NaN
3	Male	United States	56.0	Independent contractor, freelancer, or self- em...	NaN	NaN	Yes
4	Male	Taiwan	38.0	Employed full-time	NaN	NaN	Yes

5 rows × 228 columns

In [3]:

```
cr=pd.read_csv('../input/kaggle-survey-2017/conversionRates.csv', encoding
="ISO-8859-1", low_memory=False)
cr.head()
```

Out[3]:

	Unnamed: 0	originCountry	exchangeRate
0	1	USD	1.000000
1	2	EUR	1.195826
2	3	INR	0.015620
3	4	GBP	1.324188
4	5	BRL	0.321350

In [4]:

```
fr=pd.read_csv('../input/kaggle-survey-2017/freeformResponses.csv', encodi
ng="ISO-8859-1", low_memory=False)
fr.head()
```

Out[4]:

	GenderFreeForm	KaggleMotivationFreeForm	CurrentJobTitleFreeForm	MLToolNextYearFreeForm	M
0	NaN	NaN	NaN	NaN	N
1	NaN	NaN	NaN	NaN	N
2	NaN	NaN	teacher	NaN	N
3	NaN	NaN	NaN	NaN	N
4	NaN	NaN	NaN	NaN	N

5 rows × 62 columns

In [5]:

```
sch=pd.read_csv('../input/kaggle-survey-2017/schema.csv', encoding="ISO-88
59-1", low_memory=False)
sch.head()
```

Out[5]:

	Column	Question	Asked
0	GenderSelect	Select your gender identity. - Selected Choice	All

1	GenderFreeForm	Select your gender identity. - A different ide...	All
2	Country	Select the country you currentiy live in.	All
3	Age	What's your age?	All
4	EmploymentStatus	What's your current employment status?	All

In [6]:

```
sch.tail()
```

Out[6]:

	Column	Question	Asked
285	JobFactorRemote	How are you assessing potential job opportunit...	Learners
286	JobFactorIndustry	How are you assessing potential job opportunit...	Learners
287	JobFactorLeaderReputation	How are you assessing potential job opportunit...	Learners
288	JobFactorDiversity	How are you assessing potential job opportunit...	Learners
289	JobFactorPublishingOpportunity	How are you assessing potential job opportunit...	Learners

In [7]:

```
sch.Asked.value_counts()
```

Out[7]:

```
CodingWorker      161
All                70
Learners          41
Worker1           6
CodingWorker-NC   5
OnlineLearners    2
Worker            2
Non-worker        2
Non-switcher      1
Name: Asked, dtype: int64
```

In [8]:

```
sch.groupby('Asked').size()
```

Out[8]:

```
Asked
All      70
CodingWorker  161
```

```

CodingWorker      161
CodingWorker-NC   5
Learners          41
Non-switcher      1
Non-worker        2
OnlineLearners    2
Worker            2
Worker1           6
dtype: int64

```

In [9]:

```

#determine the number of categories in Asked
Piedata=sch.Asked.value_counts()

labels=[]
for i in Piedata.index:
    labels.append('{0}'.format(i))
labels

x = labels
y = Piedata[:]

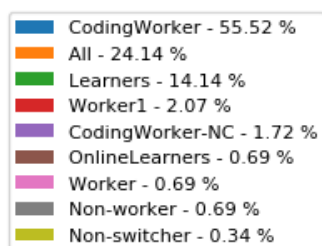
percent = 100.*y/y.sum()

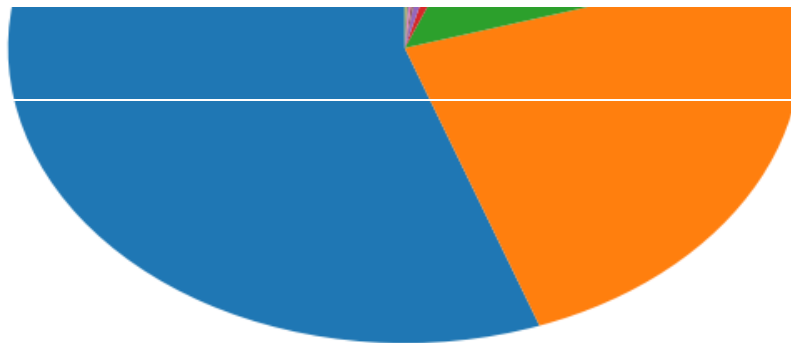
patches, texts = plt.pie(y, startangle=90, radius=1.2)
labels = ['{0} - {1:1.2f} %'.format(i,j) for i,j in zip(x, percent)]

sort_legend = True
if sort_legend:
    patches, labels, dummy = zip(*sorted(zip(patches, labels, y),
                                           key=lambda x: x[2],
                                           reverse=True))

plt.legend(patches, labels, loc='best', bbox_to_anchor=(-0.1, 1.),
           fontsize=8)
plt.show()

```





In [10]:

```
# Group by Asked on MultipleChoiceResponse and FreeFormResponse
Asked_N=range(len(sch.Asked.value_counts().index))
Askgp_mcr=[] for z in Asked_N]
Askgp_fr=[] for z in Asked_N]
for i in Asked_N:
    Askgp=list(sch['Column'][sch['Asked']==sch.Asked.value_counts().index[
i]])
    for j in Askgp:
        if j not in fr:
            Askgp_mcr[i].append(j)
        else:
            Askgp_fr[i].append(j)
```

In [11]:

```
mcr_perc_resp=[]

for i in range(len(sch.Asked.value_counts().index)):

    # Percentage of Response of each questions on MultipleChoiceResponse
    perc_resp=mcr[Askgp_mcr[i]].count()/(mcr[Askgp_mcr[i]].count()+mcr[Ask
gp_mcr[i]].isnull().sum())
    perc_resp_sort=perc_resp.sort_values(ascending=False)
    df_perc_resp=perc_resp_sort.to_frame() # You need to create a datafram
e first, then you can append to a list.
    mcr_perc_resp.append(df_perc_resp)

    y = mcr_perc_resp[i].iloc[:10][0]
    N = len(y)
    x = range(N)
    width = 1/1.5
```

```

plt.title('Percentage of Response on MultipleChoiceResponse of {0} (top
p 10 or less)'.format(sch.Asked.value_counts().index[i]))
plt.xticks(x, list(mcr_perc_resp[i].index)[:10],rotation='vertical')
ax=plt.bar(x, y, width)

plt.show()

for j in range(len(y)):
    plt.figure(figsize=(12,8))
    ax2 = sns.countplot(x=y.index[j], data=mcr, order=mcr[y.index[j]].
value_counts().index[:10])
    plt.title('{0} (top 10 or less)'.format(list(sch['Question'][sch[
'Column']==y.index[j]])[0]))
    plt.xlabel(y.index[j])
    plt.xticks(rotation=90)

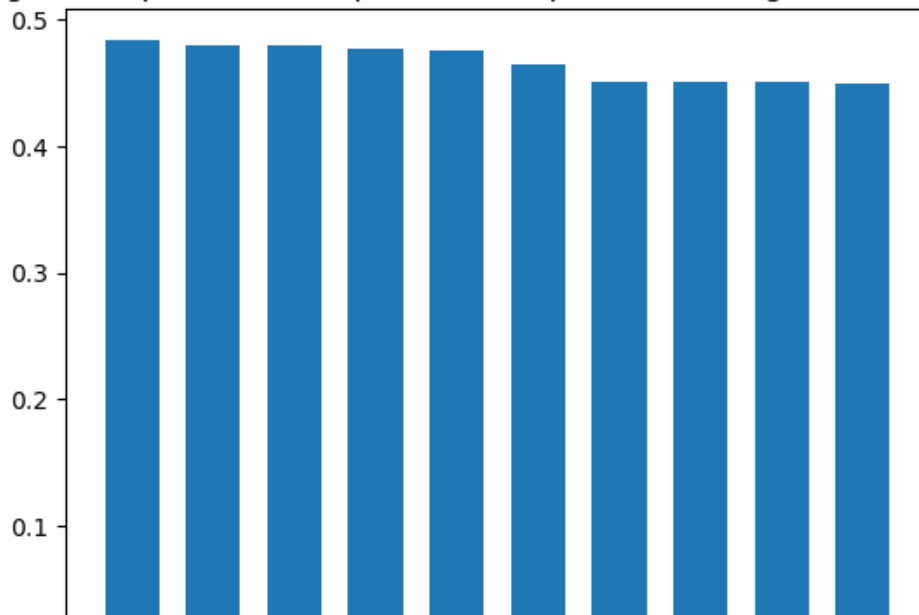
    ncount = len(mcr[y.index[j]].dropna())

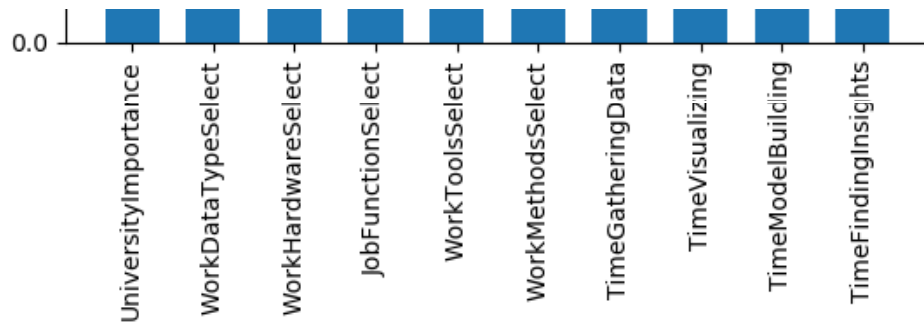
    for p in ax2.patches:
        x2=p.get_bbox().get_points()[:,0]
        y2=p.get_bbox().get_points()[1,1]
        ax2.annotate('{:.1f}%'.format(100.*y2/ncount), (x2.mean(), y2
),
                    ha='center', va='bottom') # set the alignment of the t
ext

plt.show()

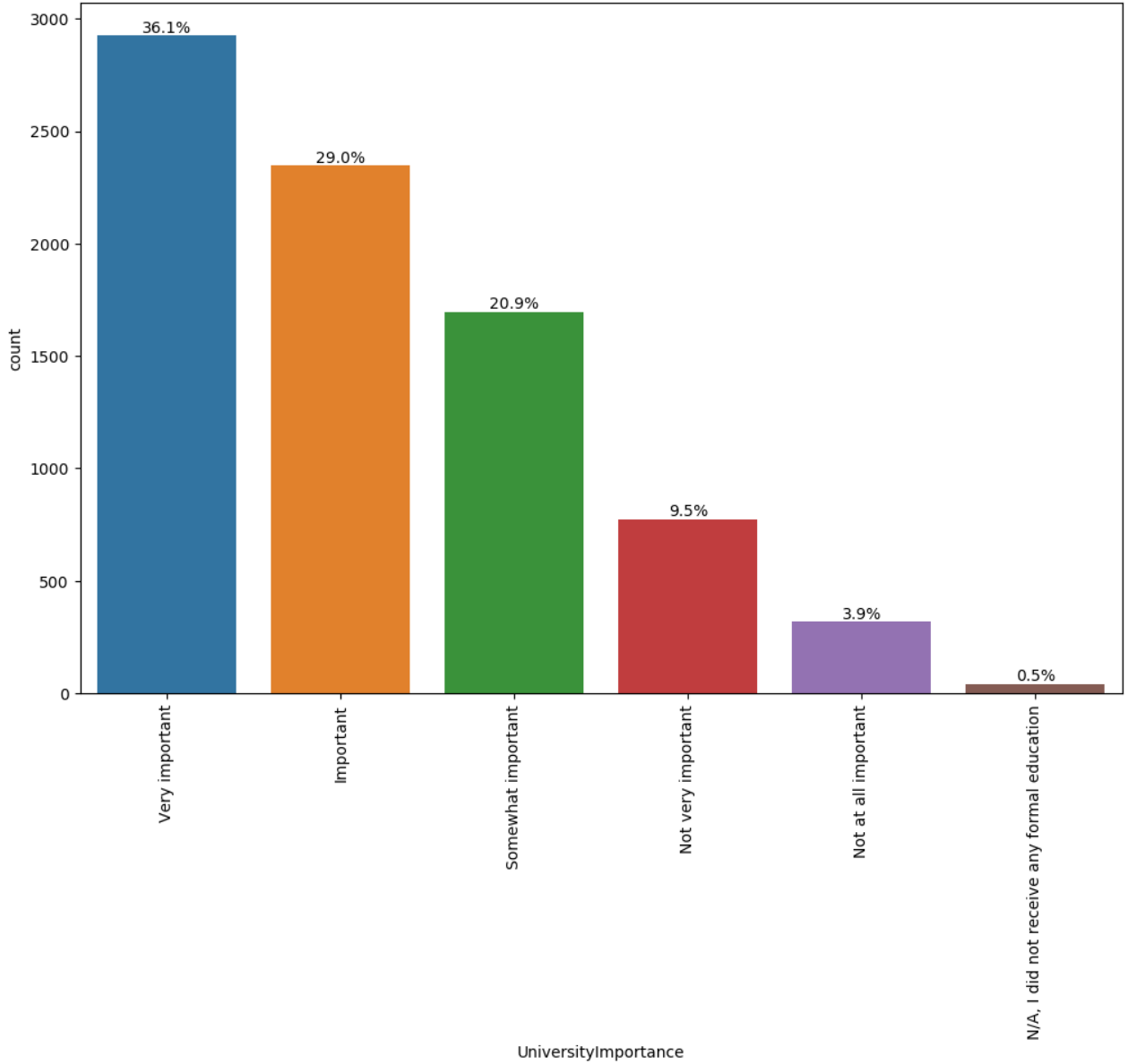
```

Percentage of Response on MultipleChoiceResponse of CodingWorker (top 10 or less)

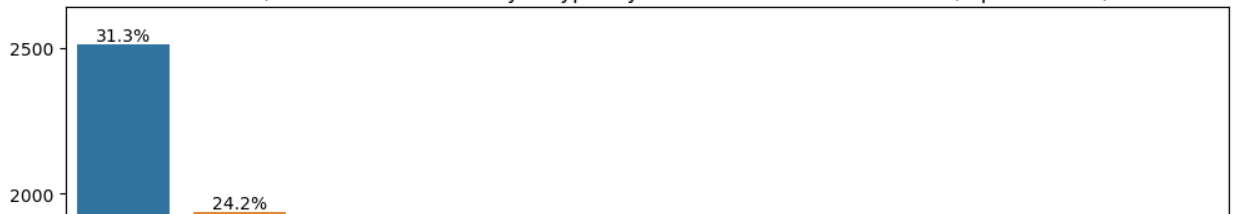


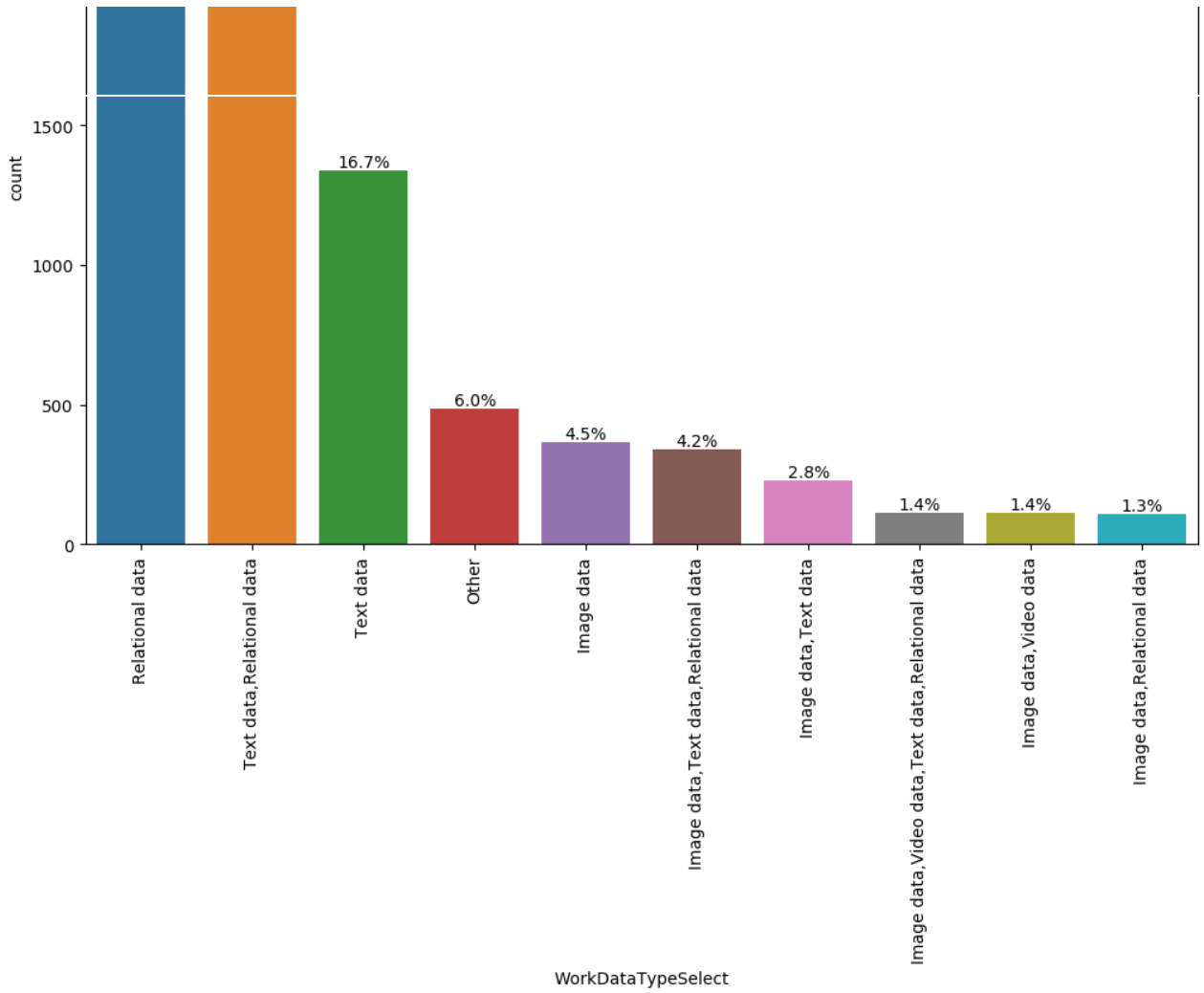


How important was your formal education or degree to your career success analyzing data? (top 10 or less)

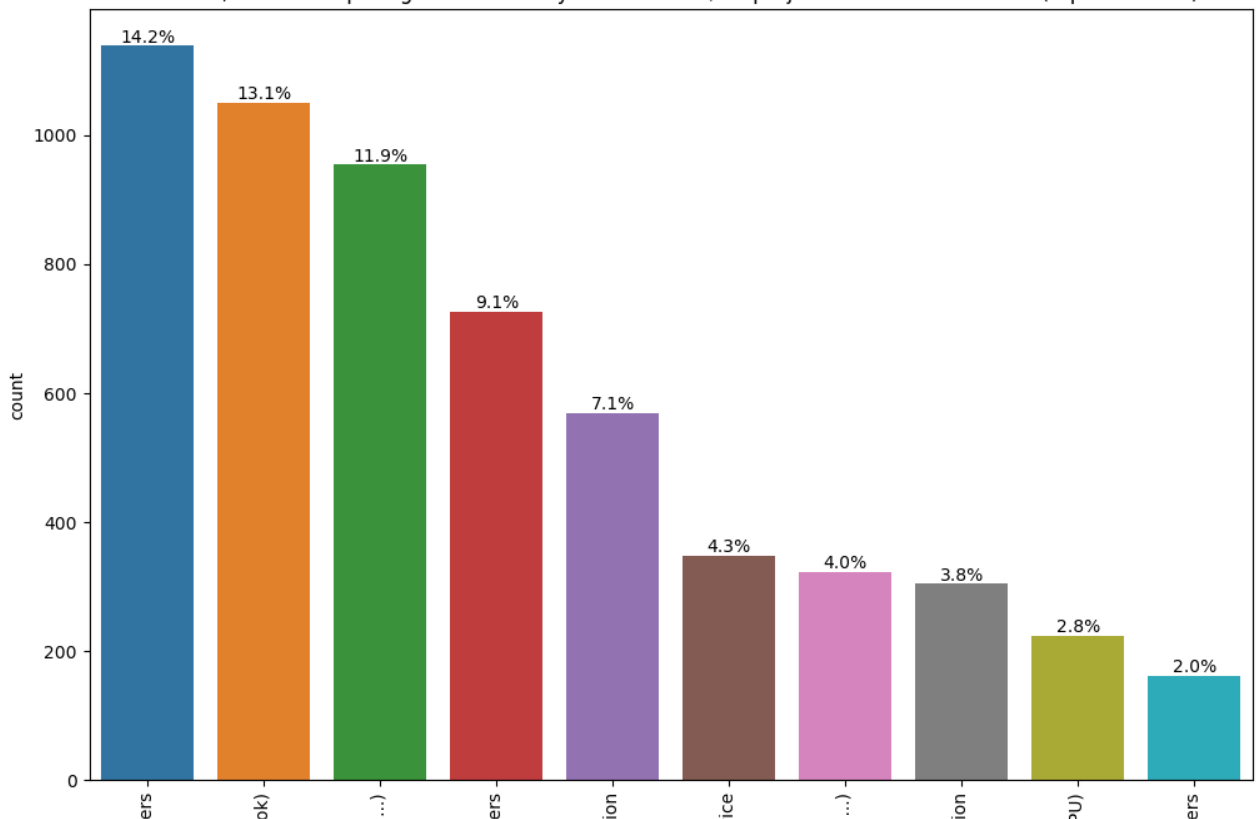


At work, which kind of data do you typically work with? - Selected Choice (top 10 or less)





At work, which computing hardware do you use for ML/DS projects? - Selected Choice (top 10 or less)



Notebook

Code

Data (2)

Comments (0)

Log

Versions (4)

Options

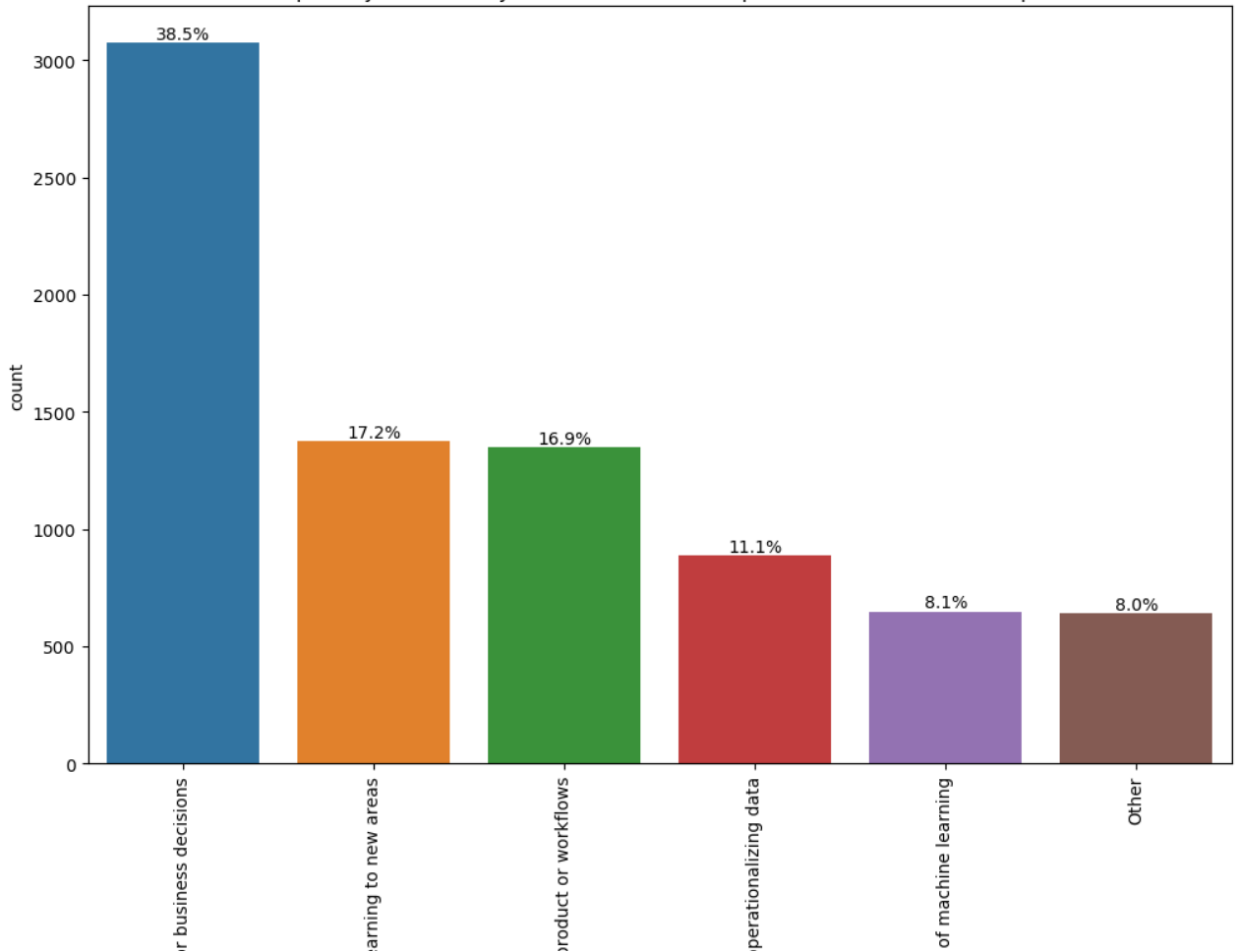
Fork Notebook

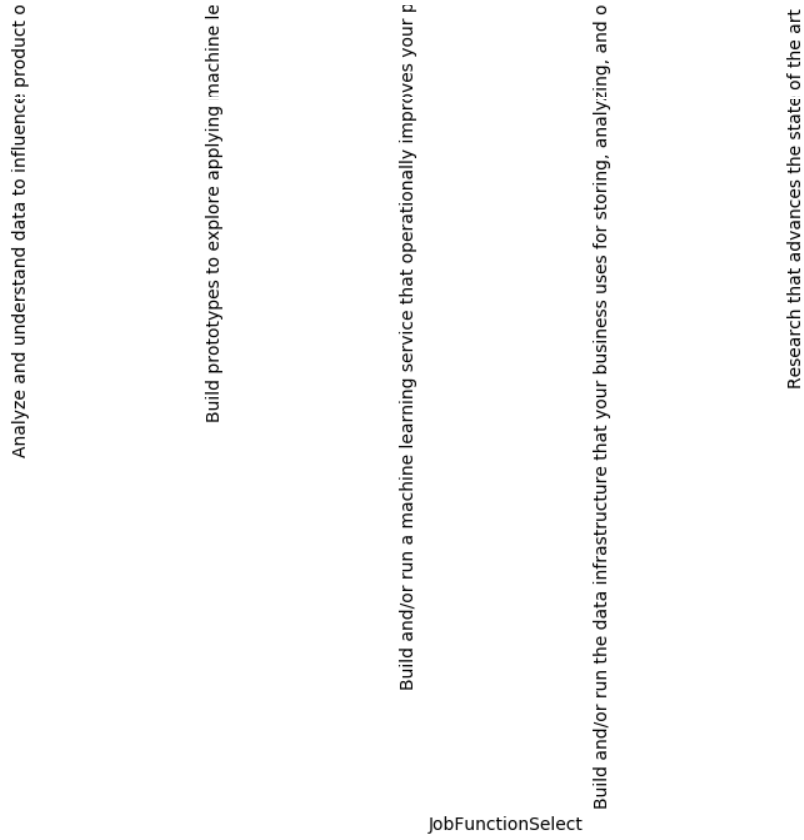
Edit Note

Laptop or Workst	private datacent
Laptop + Clouc	ic laptop (Macboi
Laptop or Workstation	AWS, Azure, GCE
	IT supported serv
	ditional Workstat
	ion + Cloud serv
	AWS, Azure, GCE
	lerated Workstat
	CUDA capable Gf
	private datacent

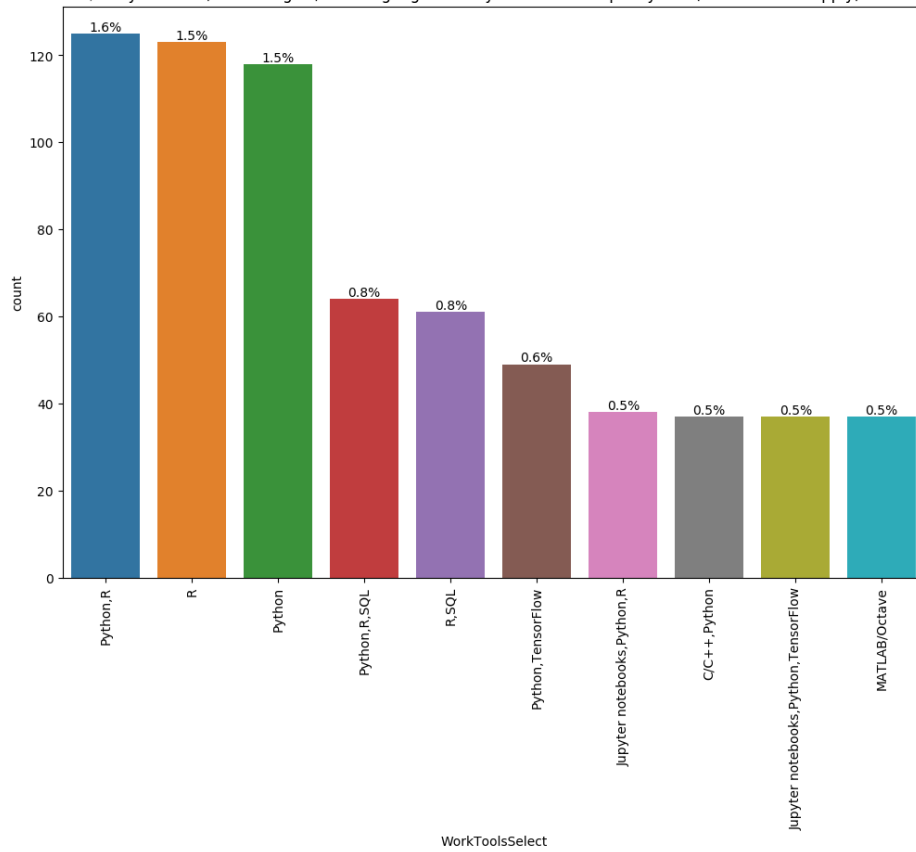
WorkHardwareSelect

What is the primary function of your role? (Select one option) - Selected Choice (top 10 or less)



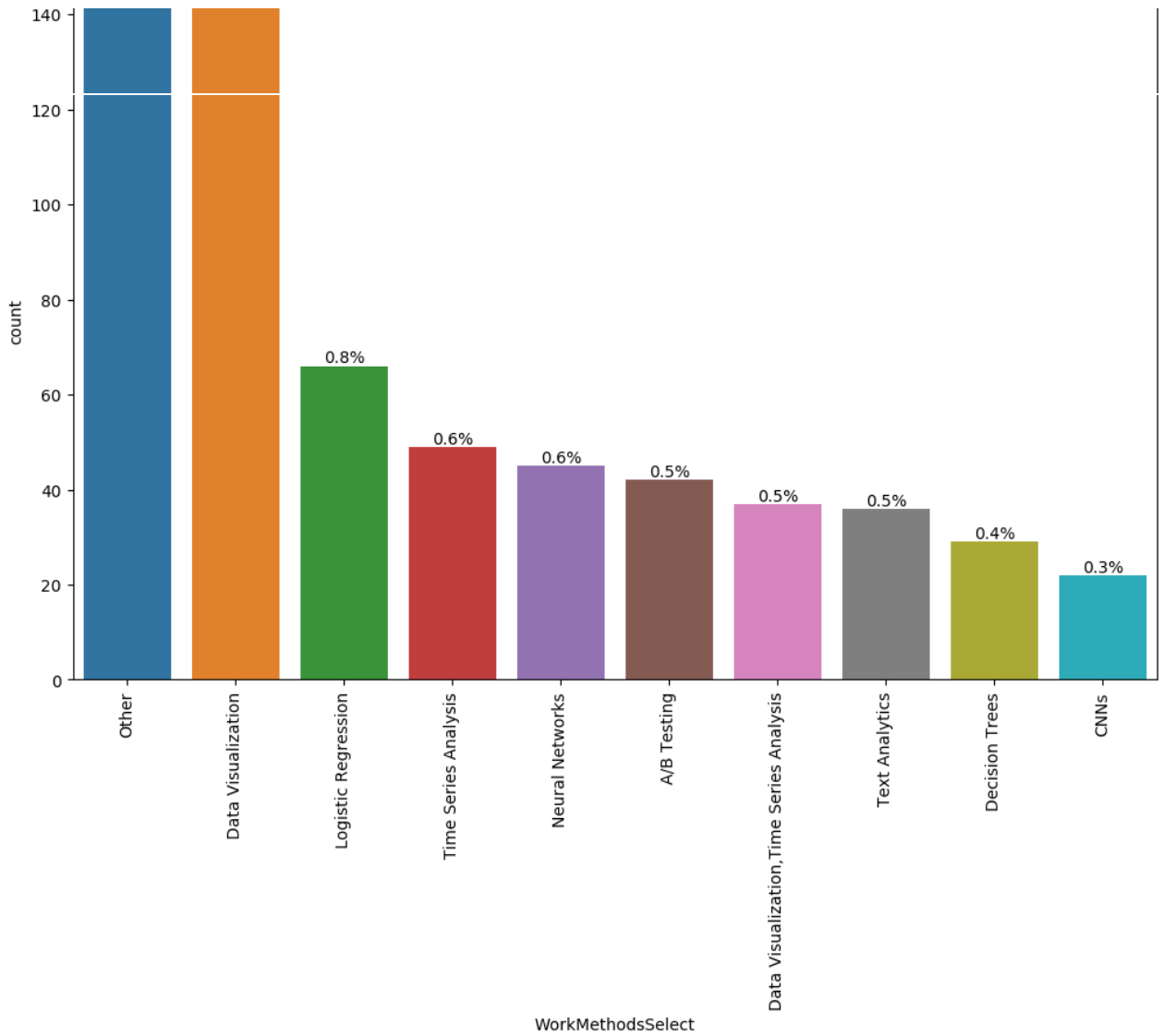


For work, which data science/analytics tools, technologies, and languages have you used in the past year? (Select all that apply) - Selected Choice (top 10 or less)

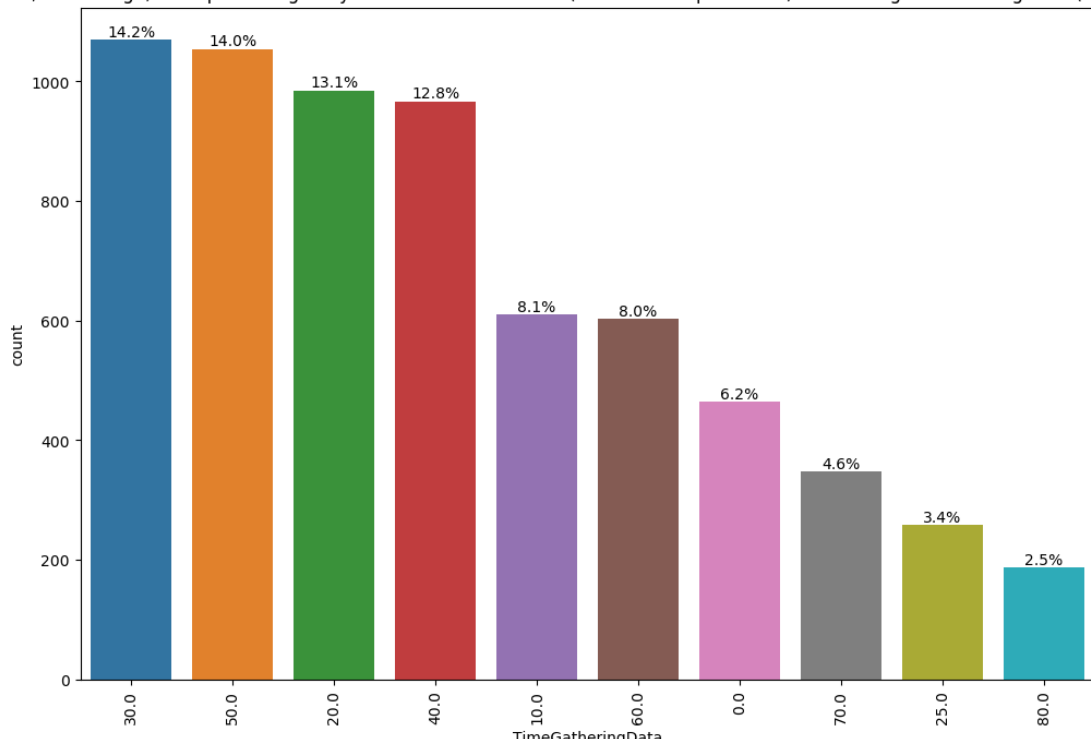


At work, which data science methods do you use? (Select all that apply) - Selected Choice (top 10 or less)

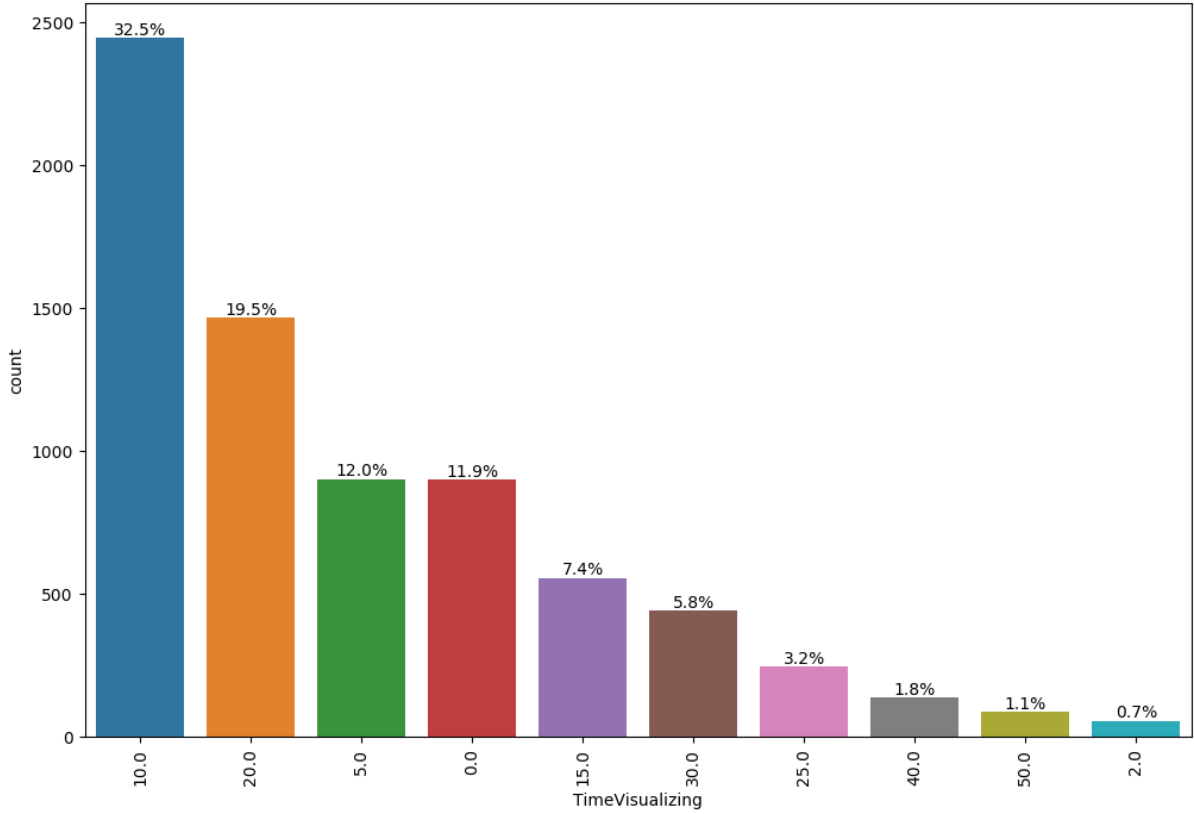




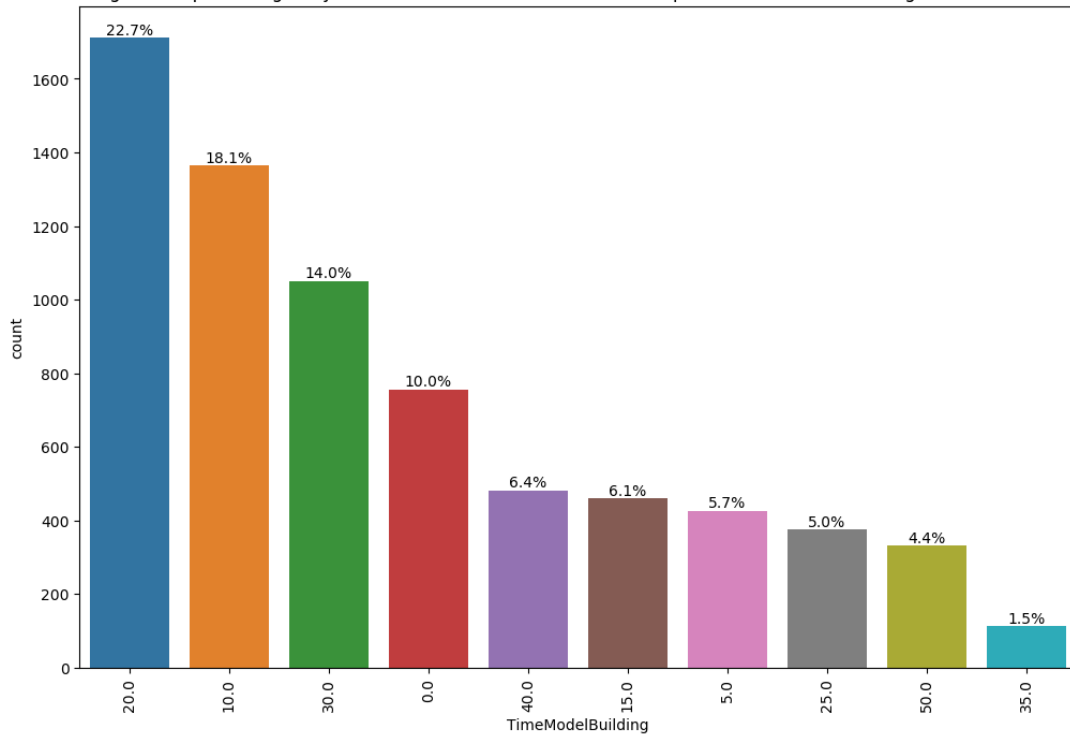
At work, on average, what percentage of your time is devoted to: (Total must equal 100%) - Gathering and cleaning data (top 10 or less)



At work, on average, what percentage of your time is devoted to: (Total must equal 100%) - Visualizing data (top 10 or less)

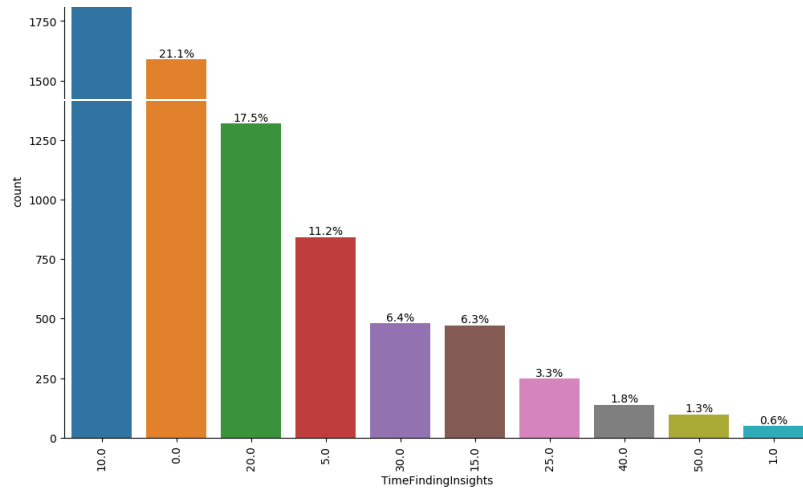


At work, on average, what percentage of your time is devoted to: (Total must equal 100%) - Model building/model selection (top 10 or less)

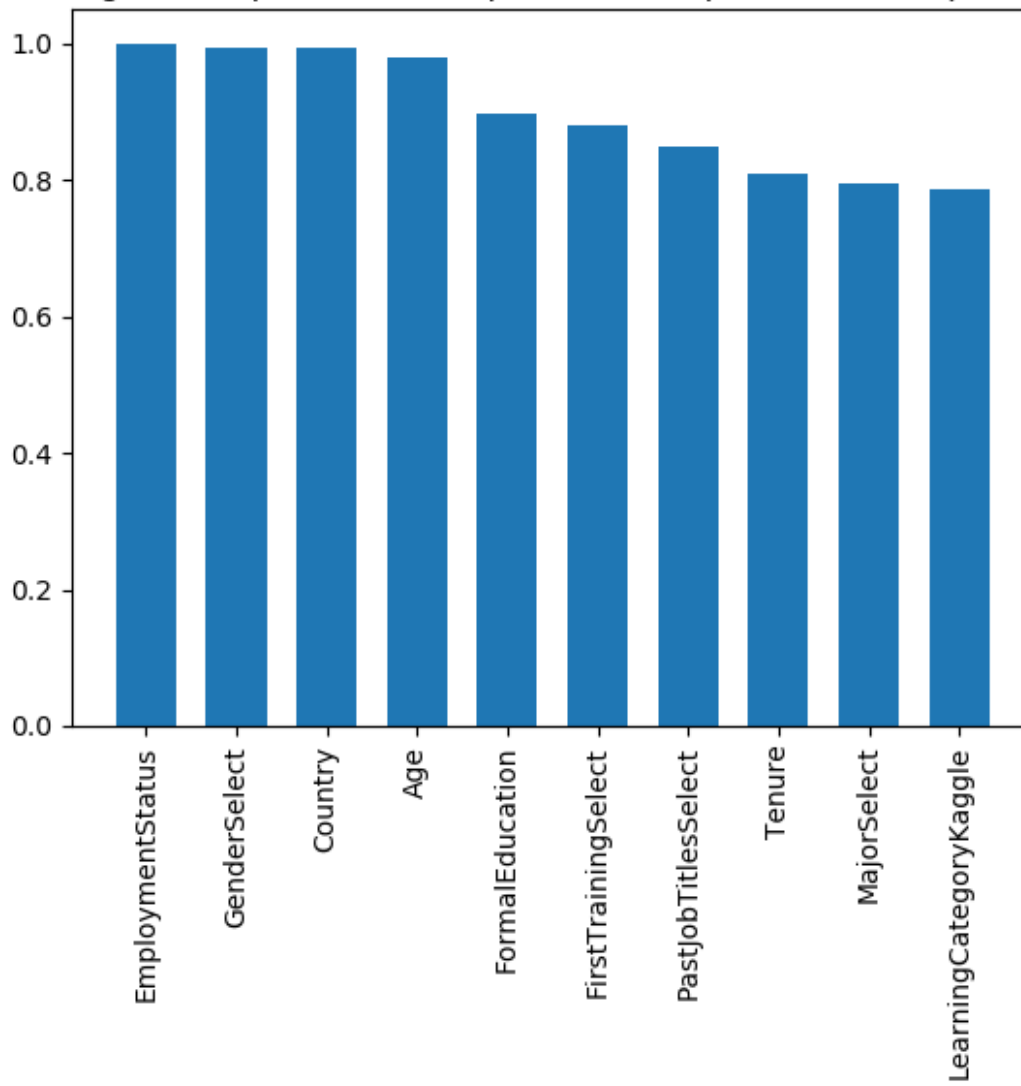


At work, on average, what percentage of your time is devoted to: (Total must equal 100%) - Finding insights in the data and communicating these to relevant stakeholders (top 10 or less)

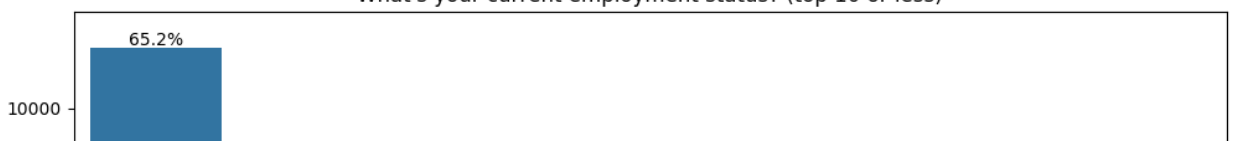


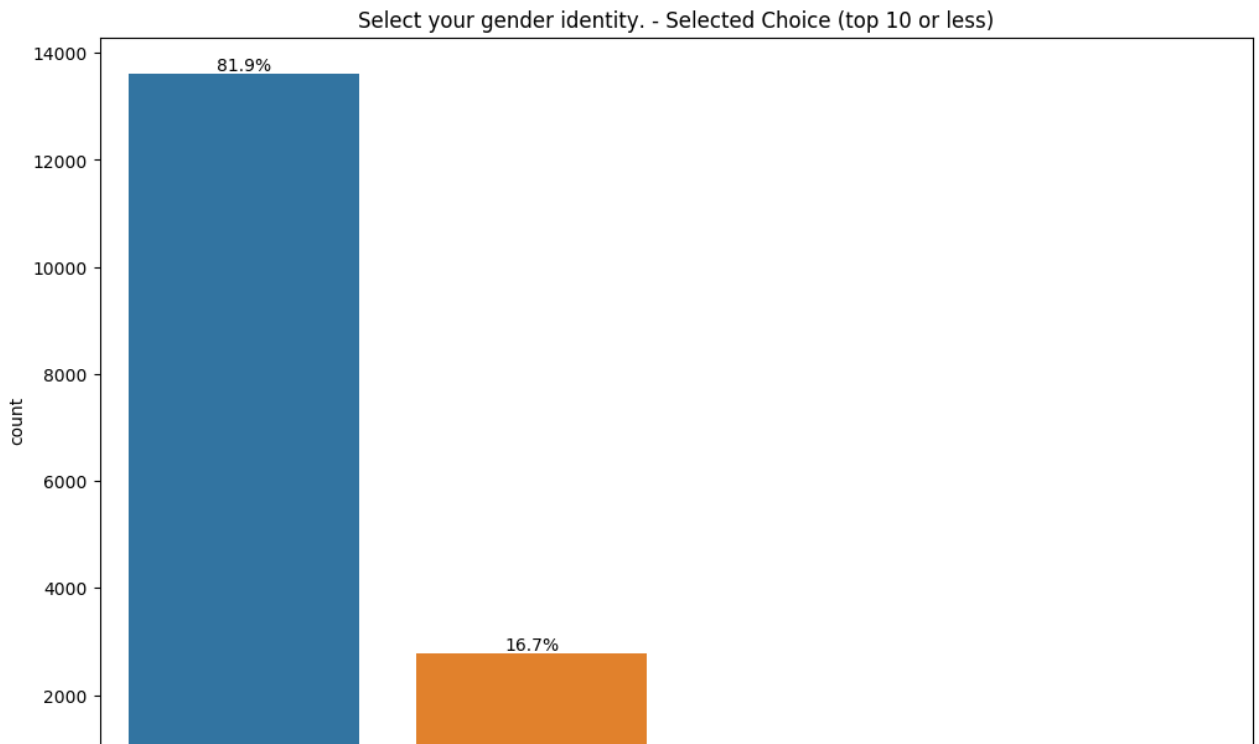
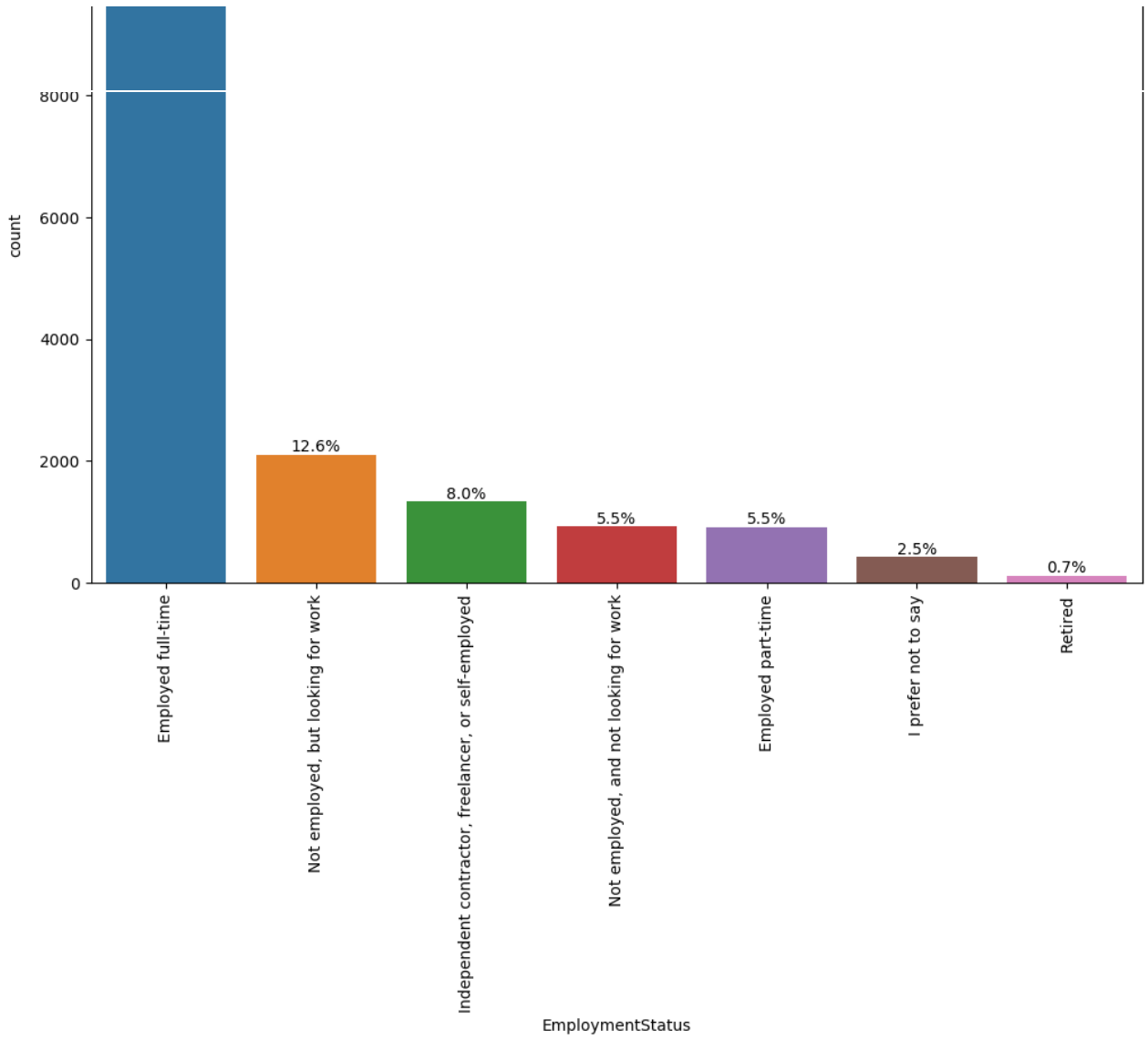


Percentage of Response on MultipleChoiceResponse of All (top 10 or less)

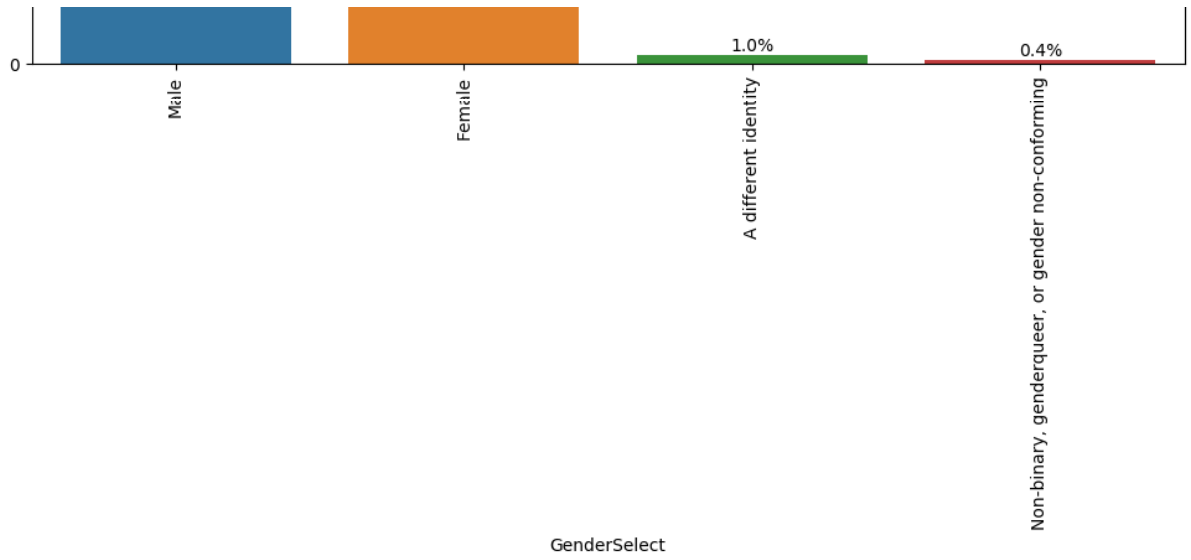


What's your current employment status? (top 10 or less)

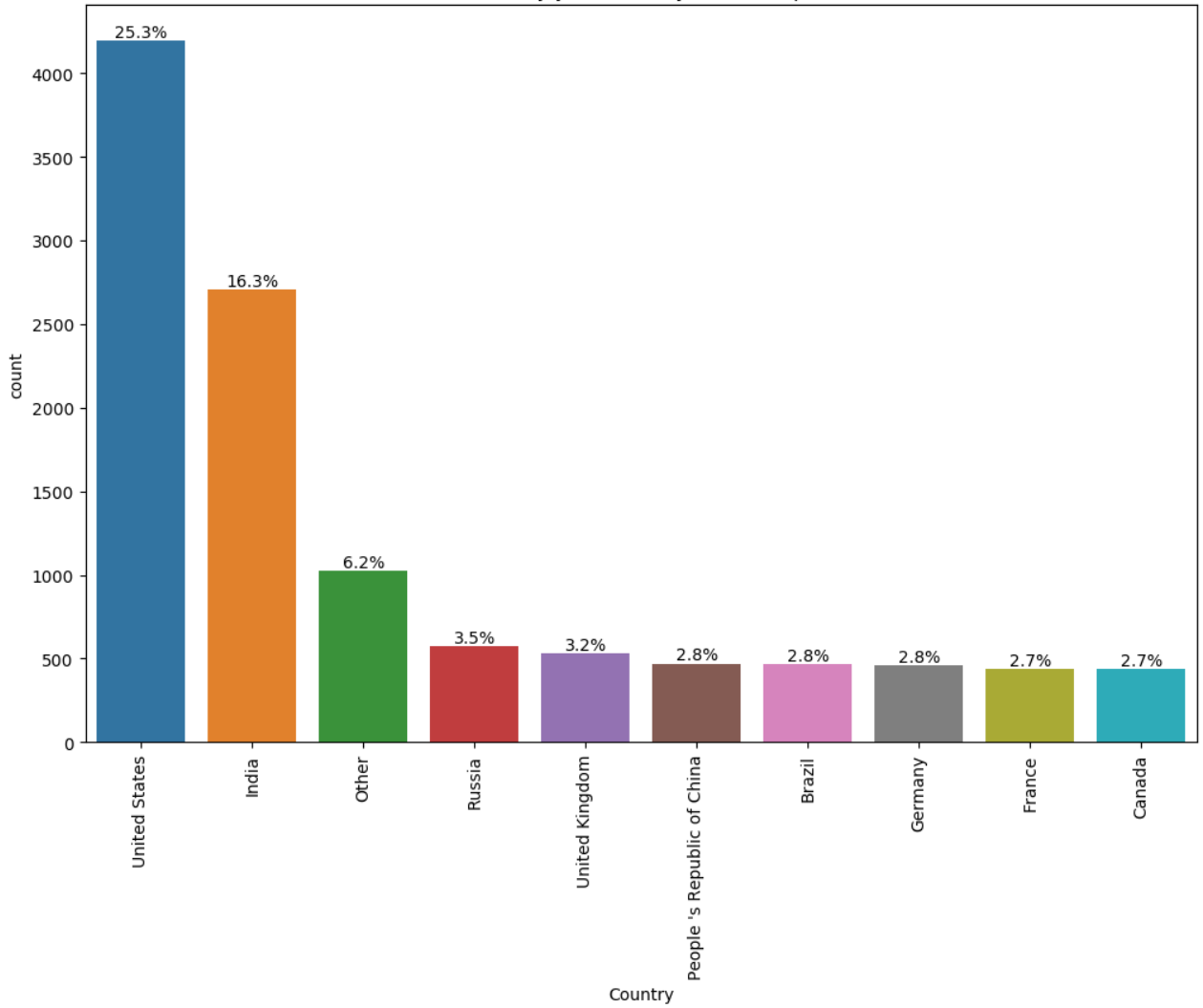




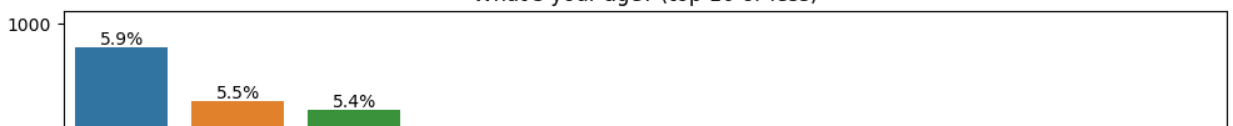
2017 State of Data Science - Kaggle survey



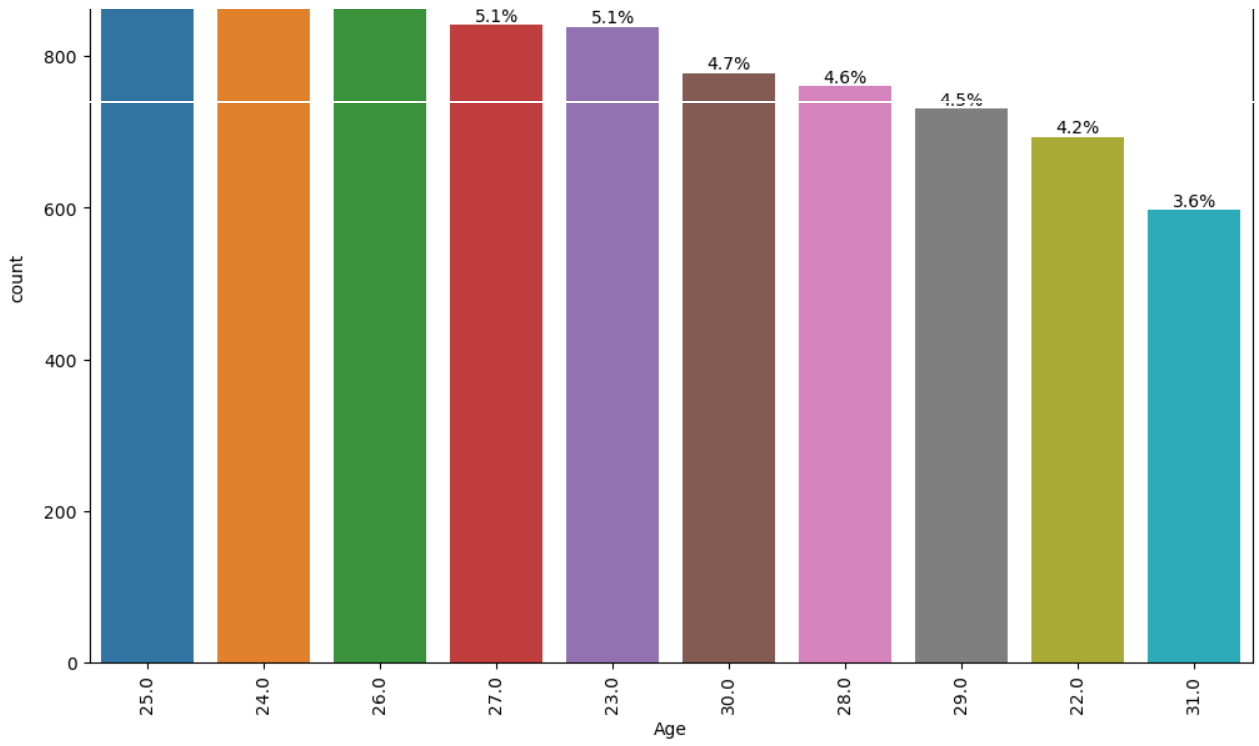
Select the country you currently live in. (top 10 or less)



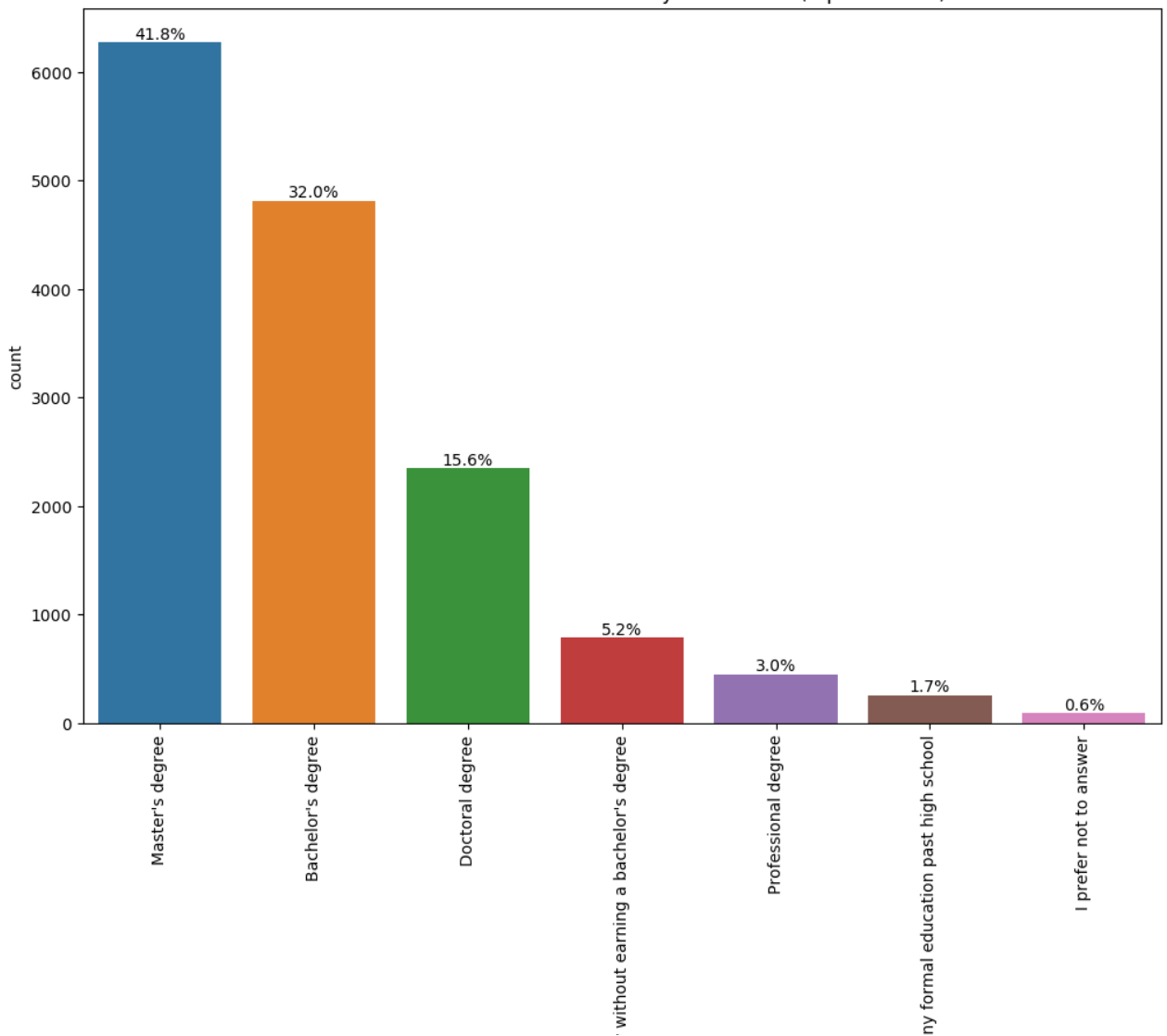
What's your age? (top 10 or less)



2017 State of Data Science - Kaggle survey

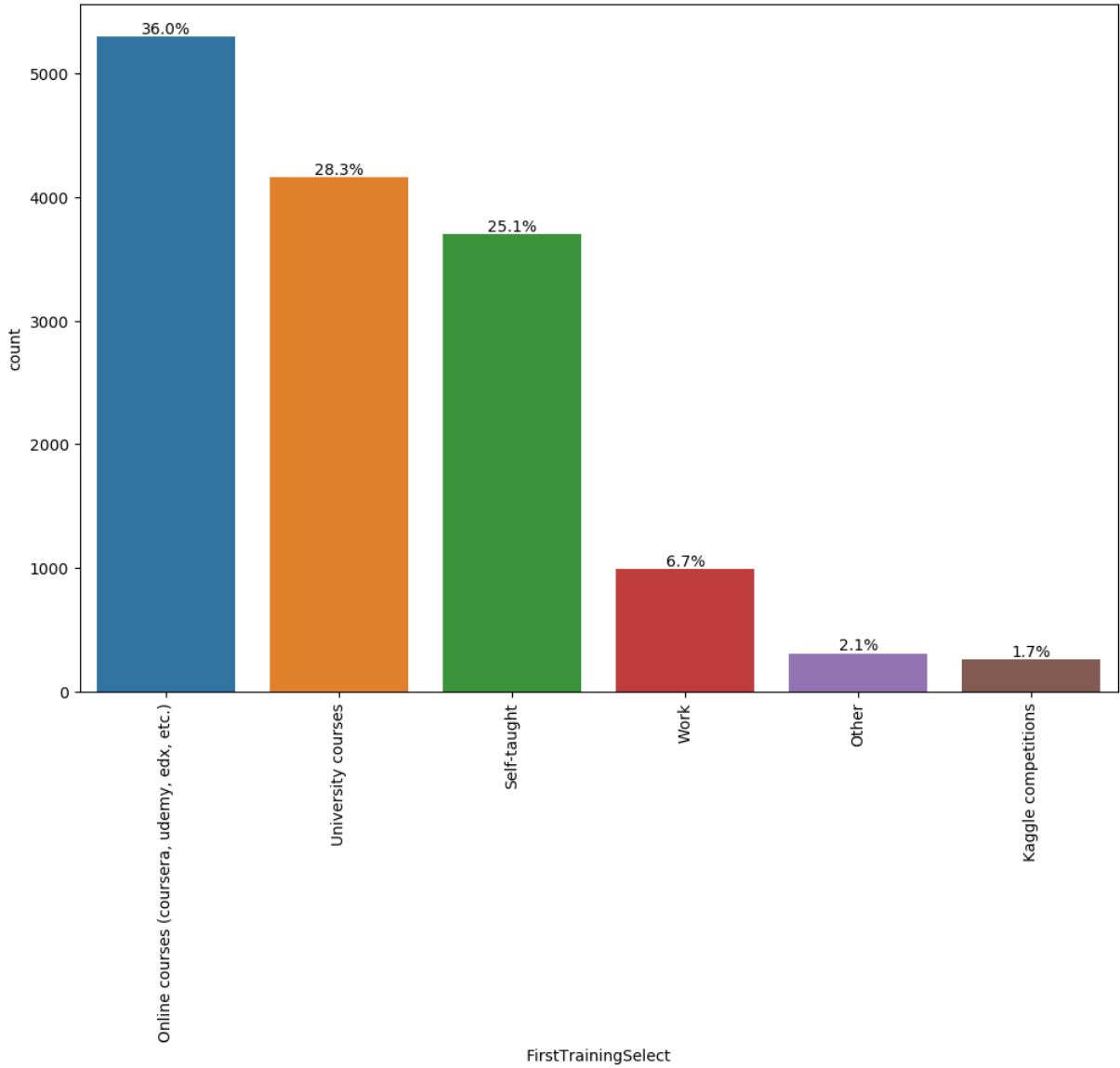


Which level of formal education have you attained? (top 10 or less)

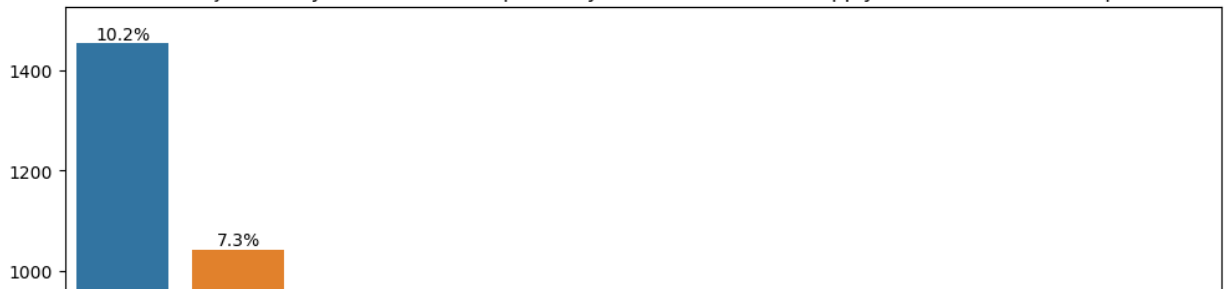


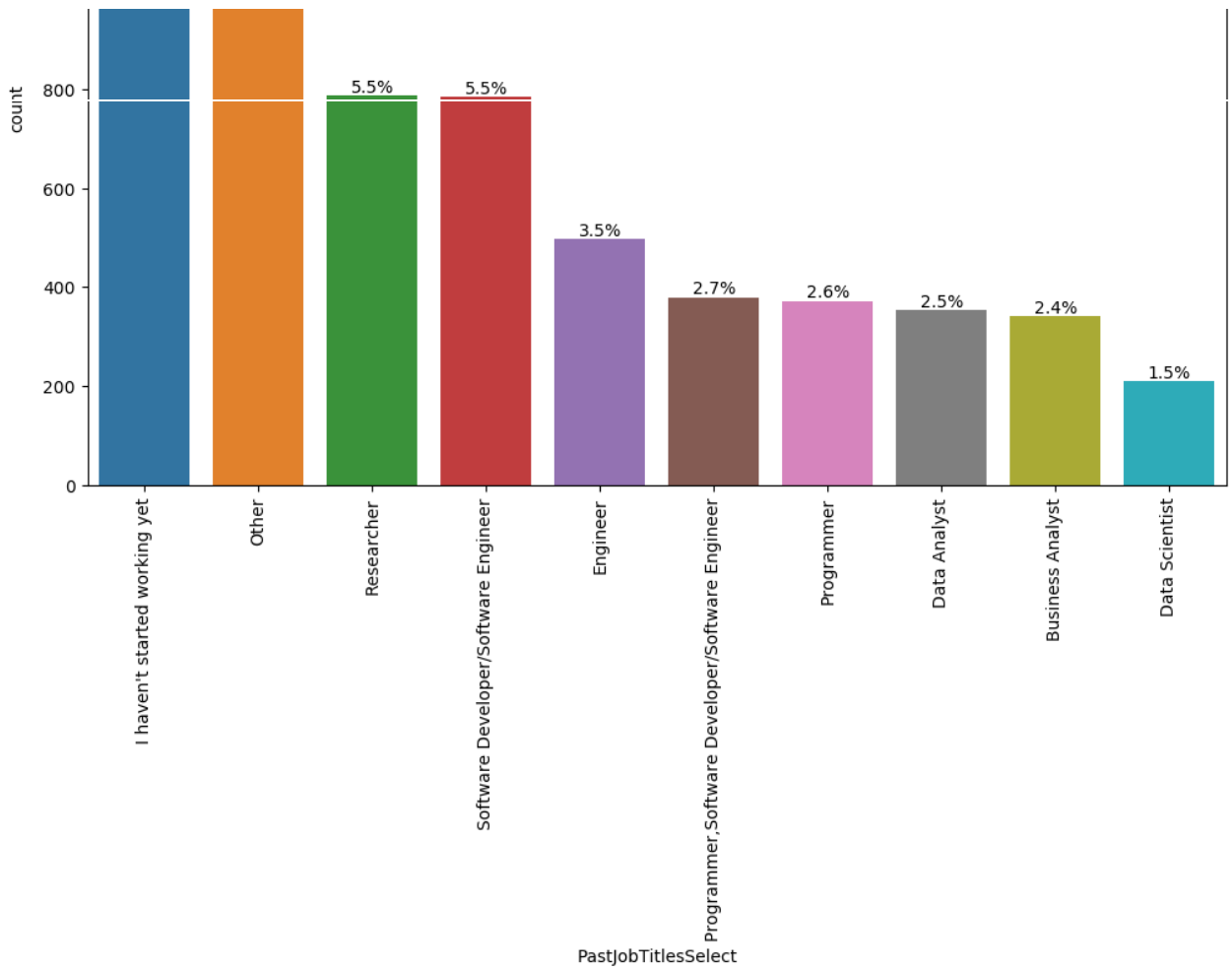
Some college/university study
FormalEducation
I did not complete a

How did you first start your machine learning / data science training? (Select one option) - Selected Choice (top 10 or less)

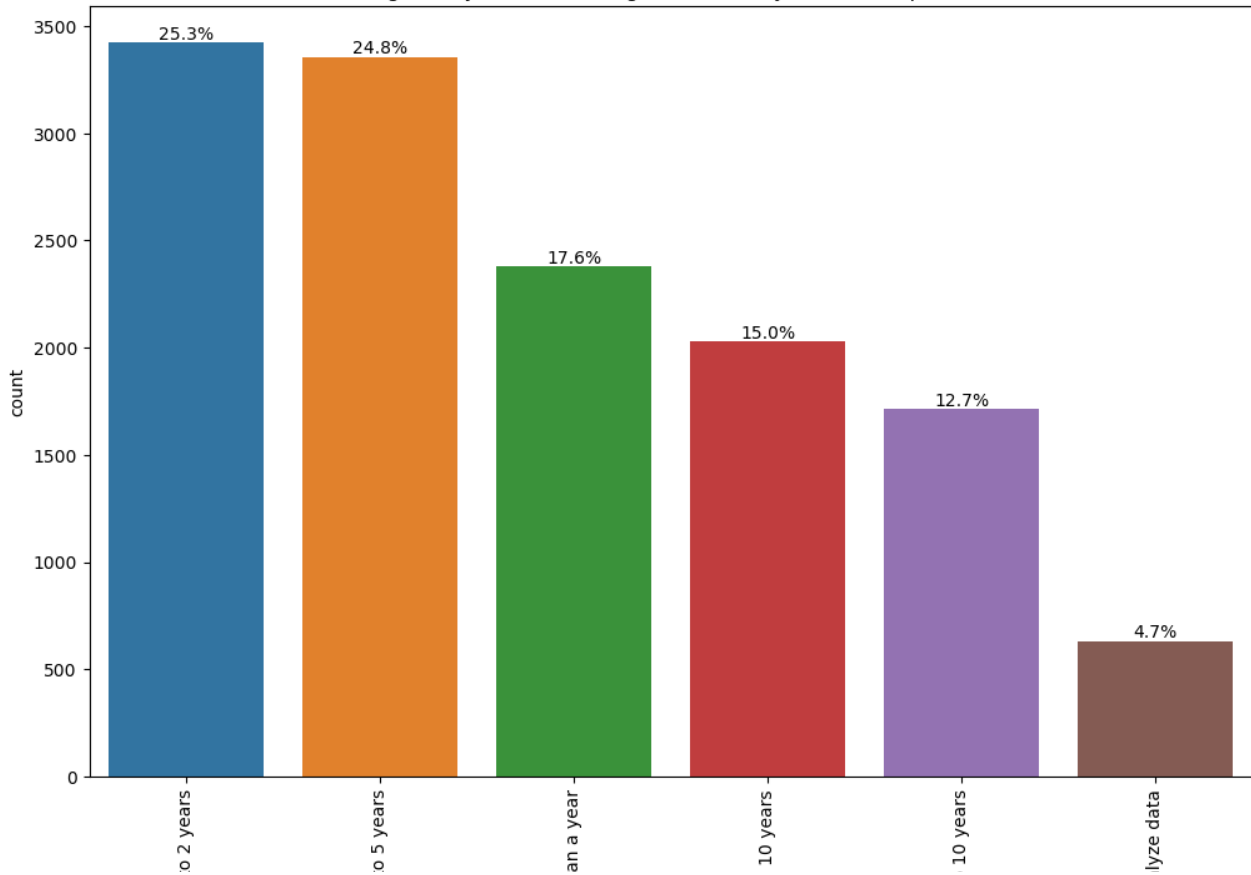


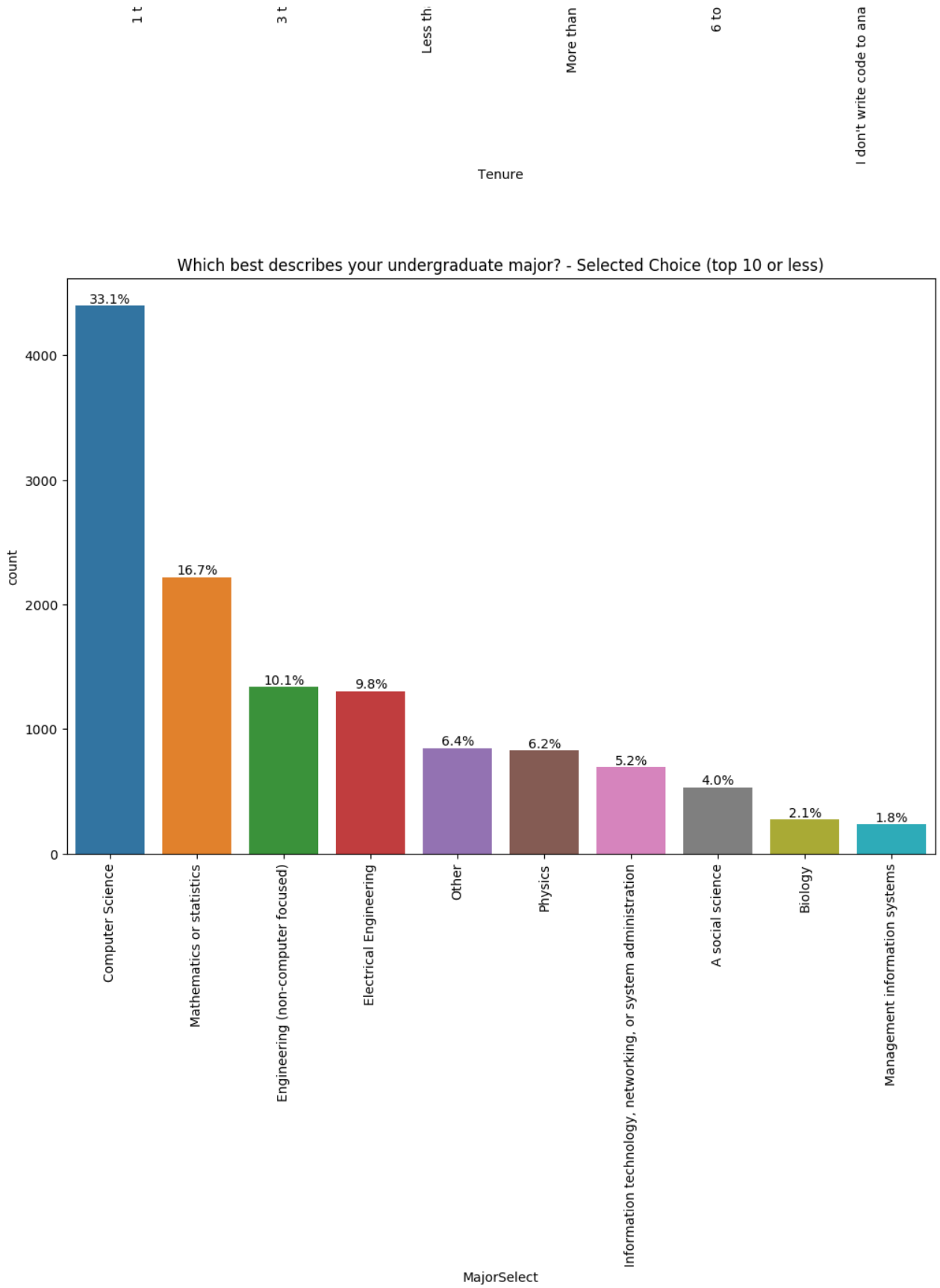
Select all other job titles you've held in the past 10 years? (Select all that apply) - Selected Choice (top 10 or less)



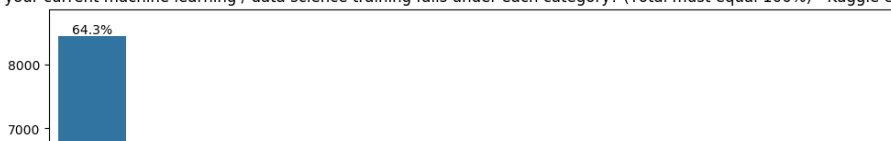


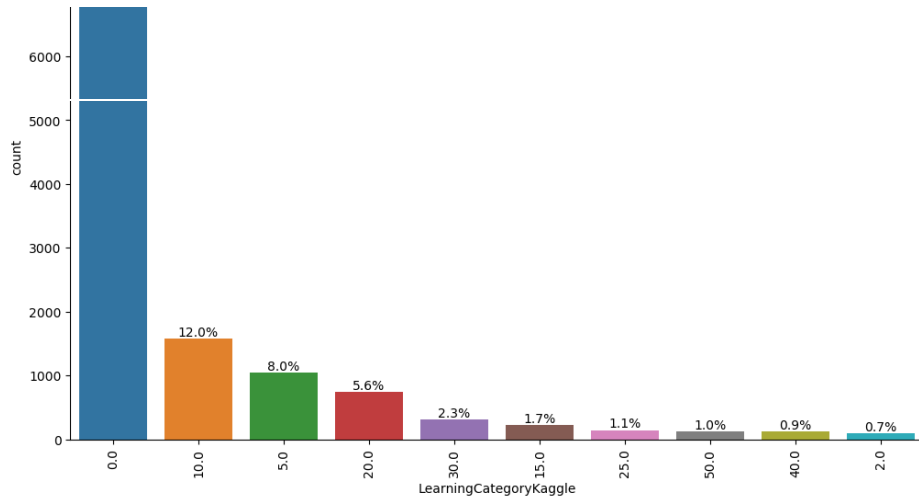
How long have you been writing code to analyze data? (top 10 or less)



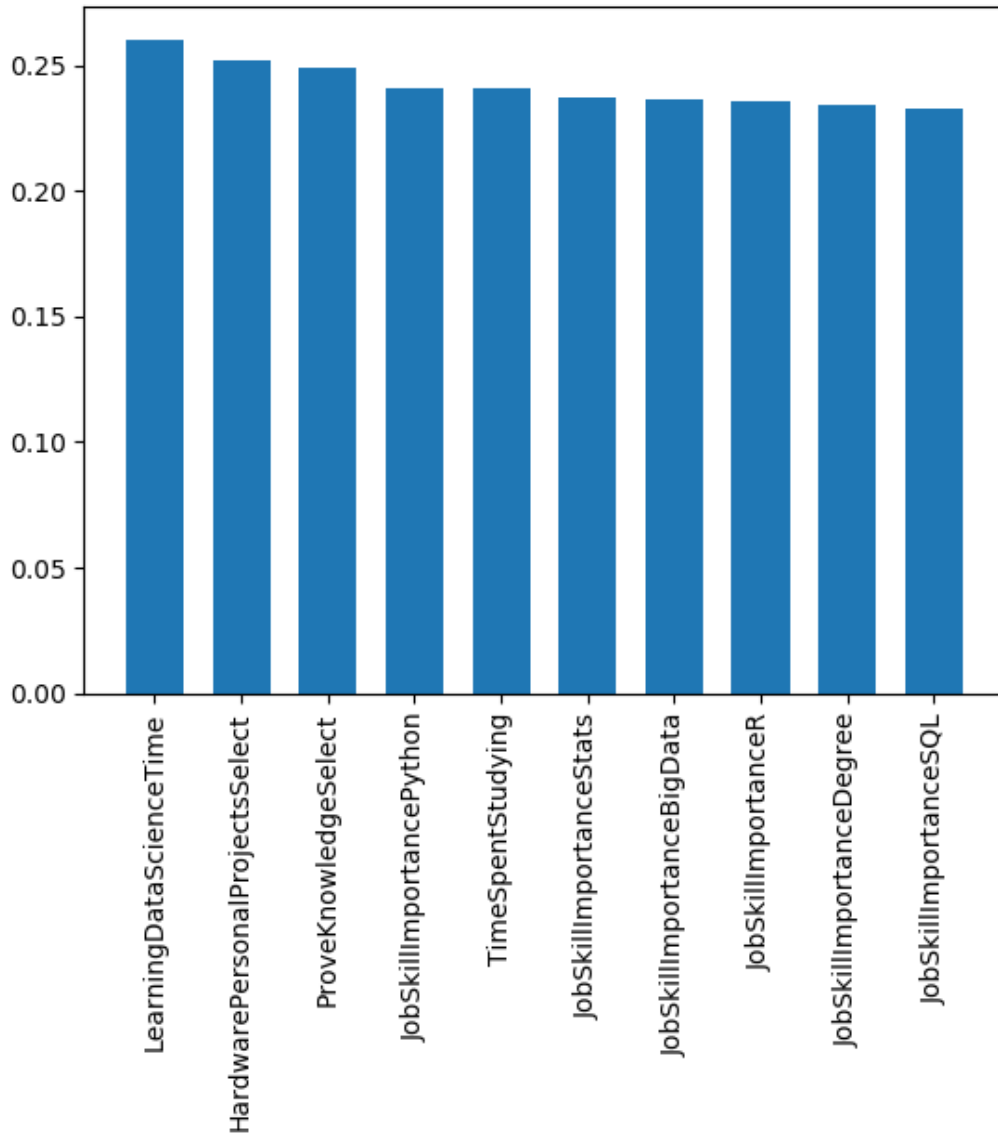


What percentage of your current machine learning / data science training falls under each category? (Total must equal 100%) - Kaggle competitions (top 10 or less)

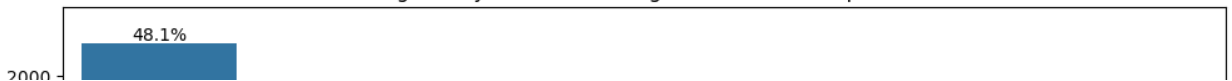


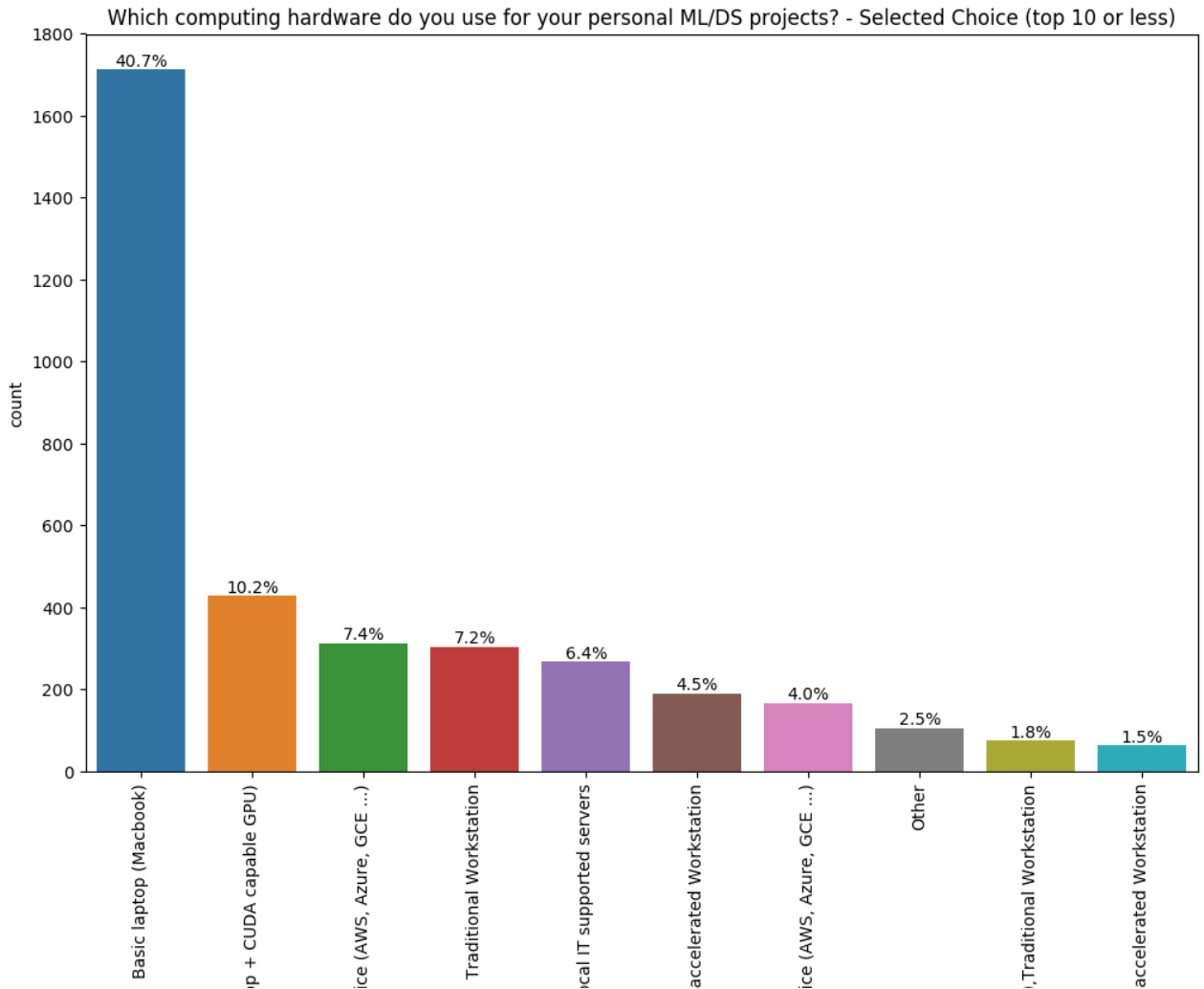
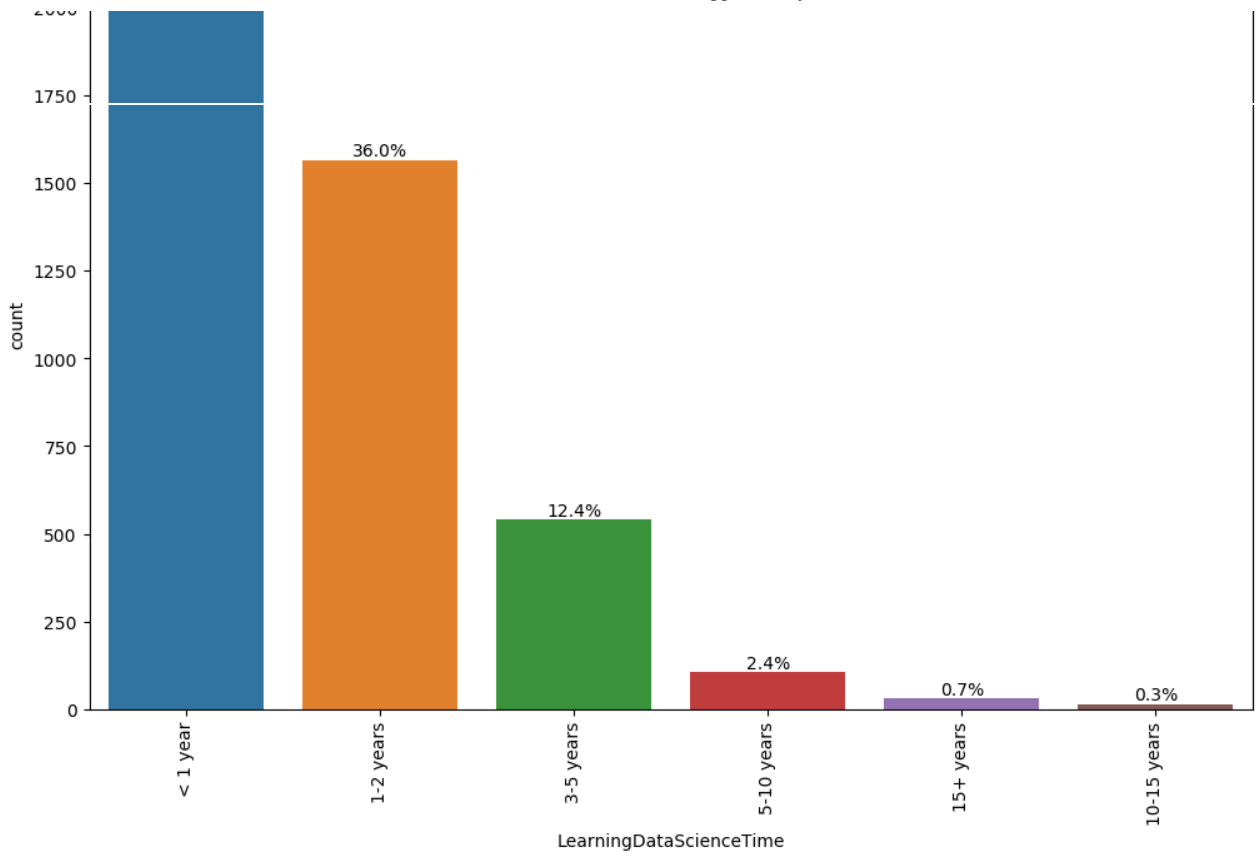


Percentage of Response on MultipleChoiceResponse of Learners (top 10 or less)



How long have you been learning data science? (top 10 or less)

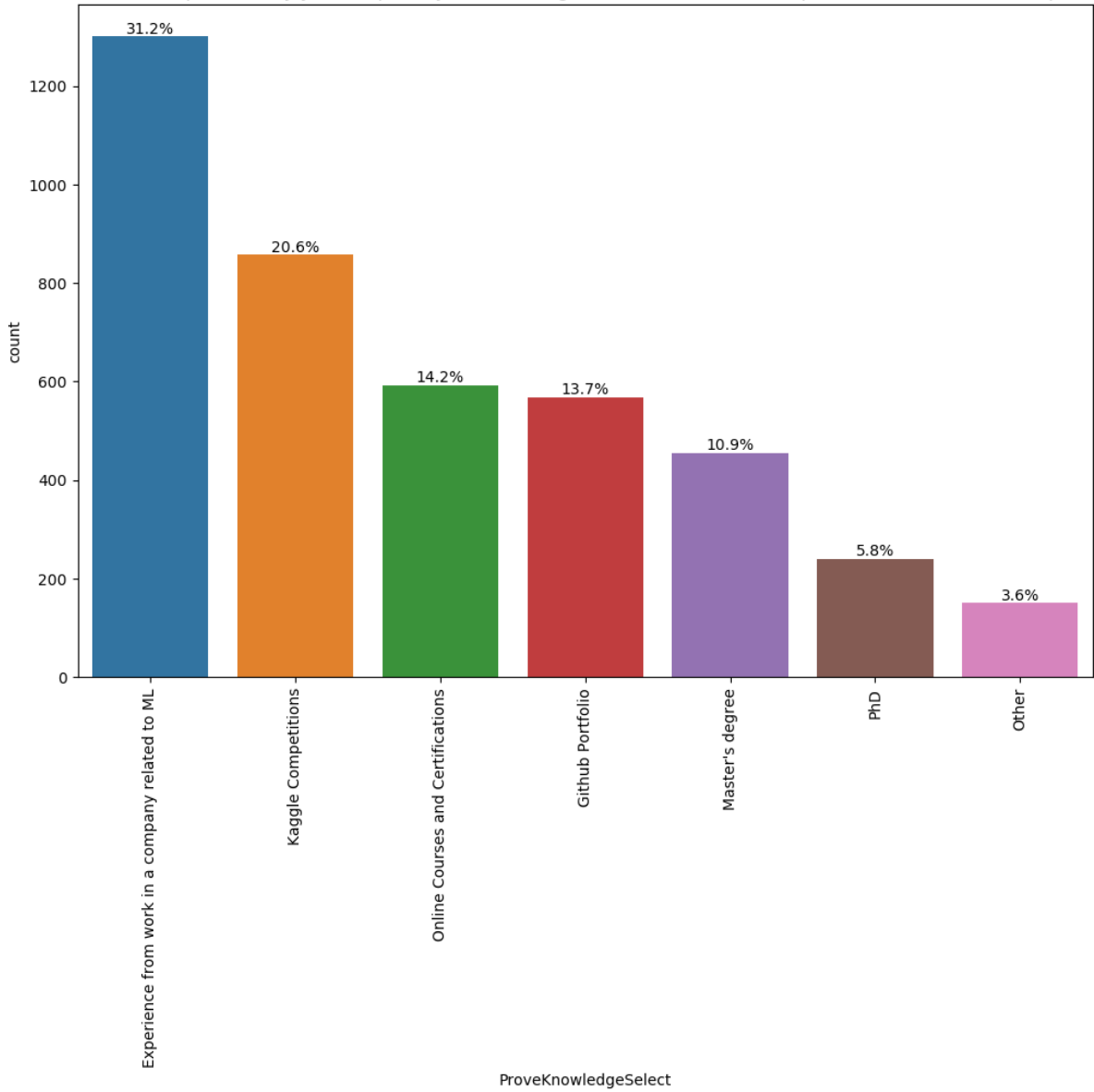




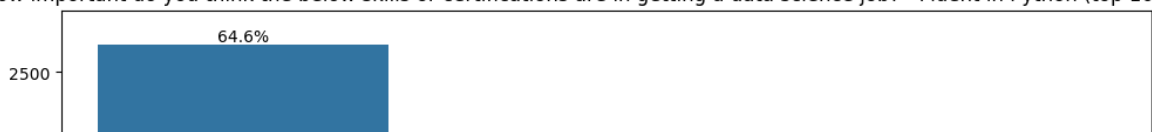
Gaming Laptop (Lapto
Laptop + Cloud servi
Laptop or Workstation and lo
GPU
Basic laptop (Macbook),Laptop + Cloud servi
Basic laptop (Macbook)
Basic laptop (Macbook),GPU

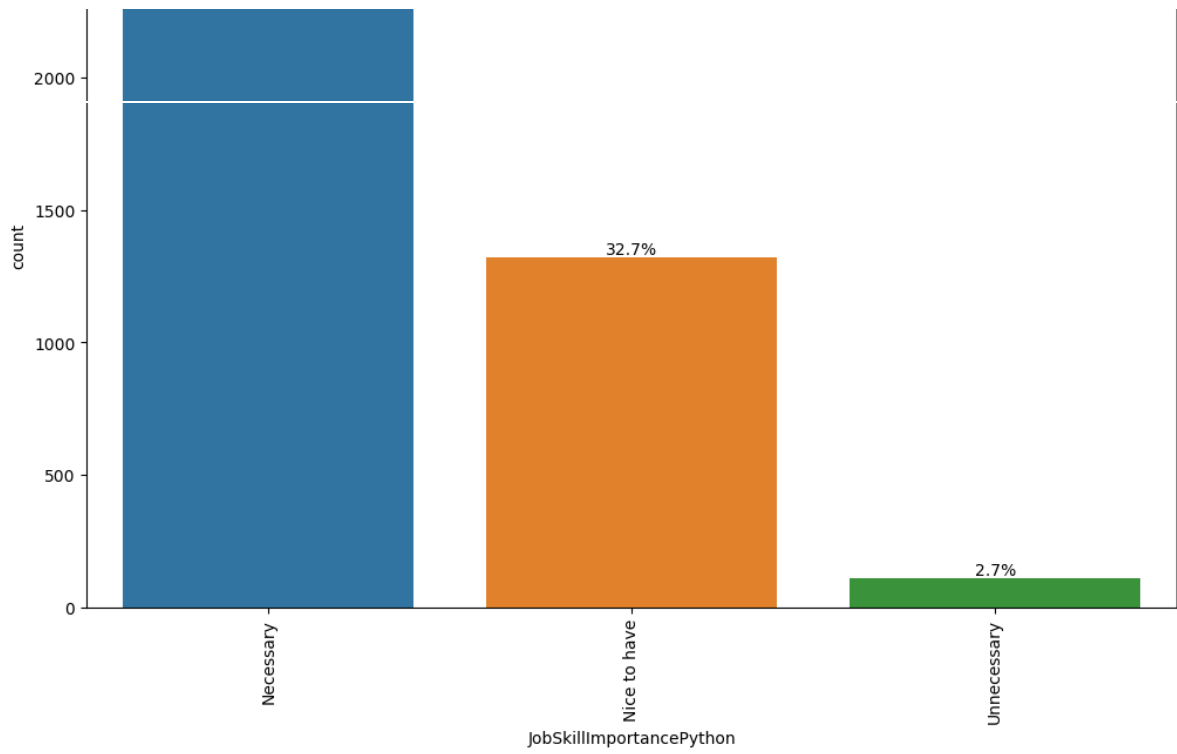
HardwarePersonalProjectsSelect

What's the most important way you can prove your knowledge of ML/DS? (Select one option) - Selected Choice (top 10 or less)

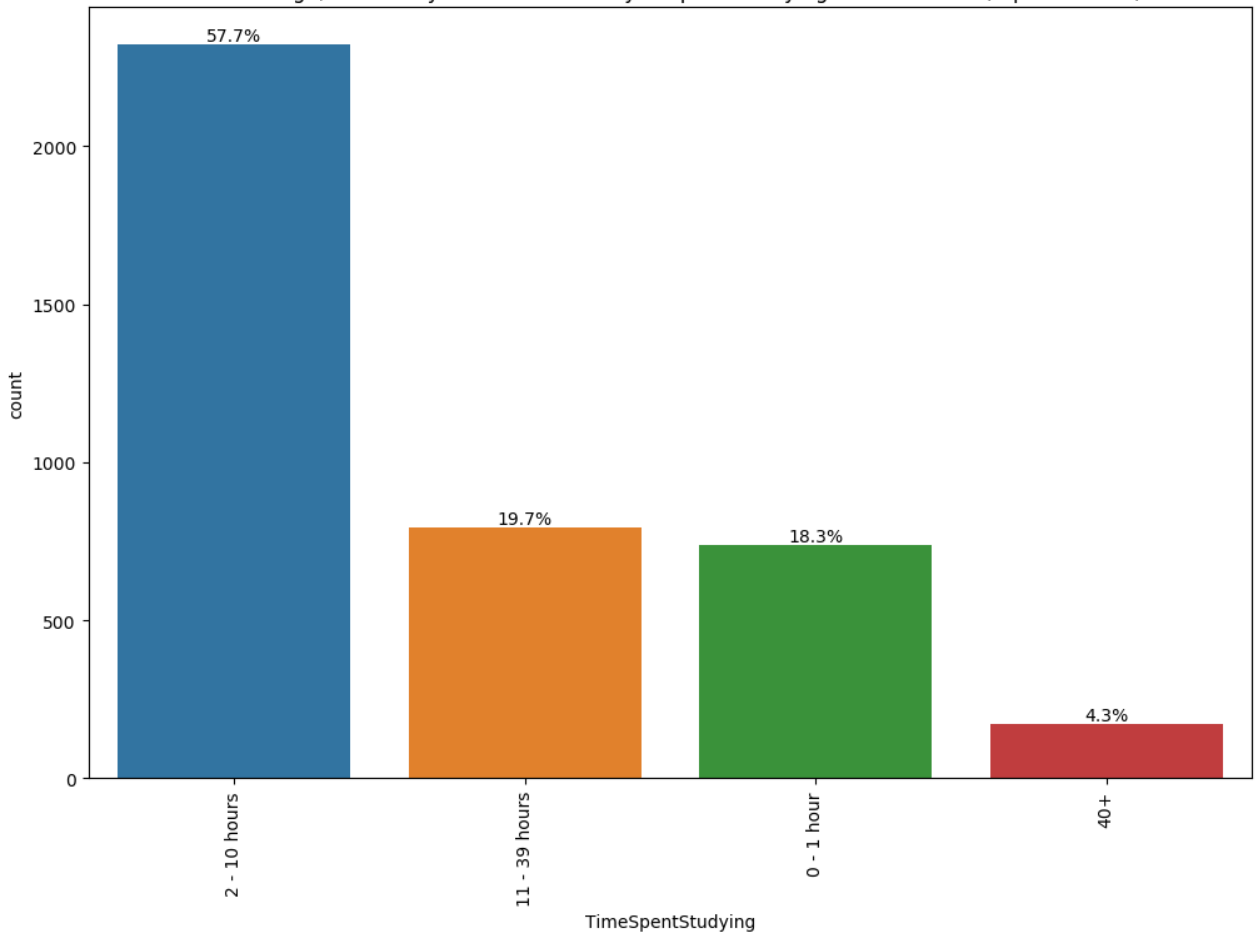


How important do you think the below skills or certifications are in getting a data science job? - Fluent in Python (top 10 or less)



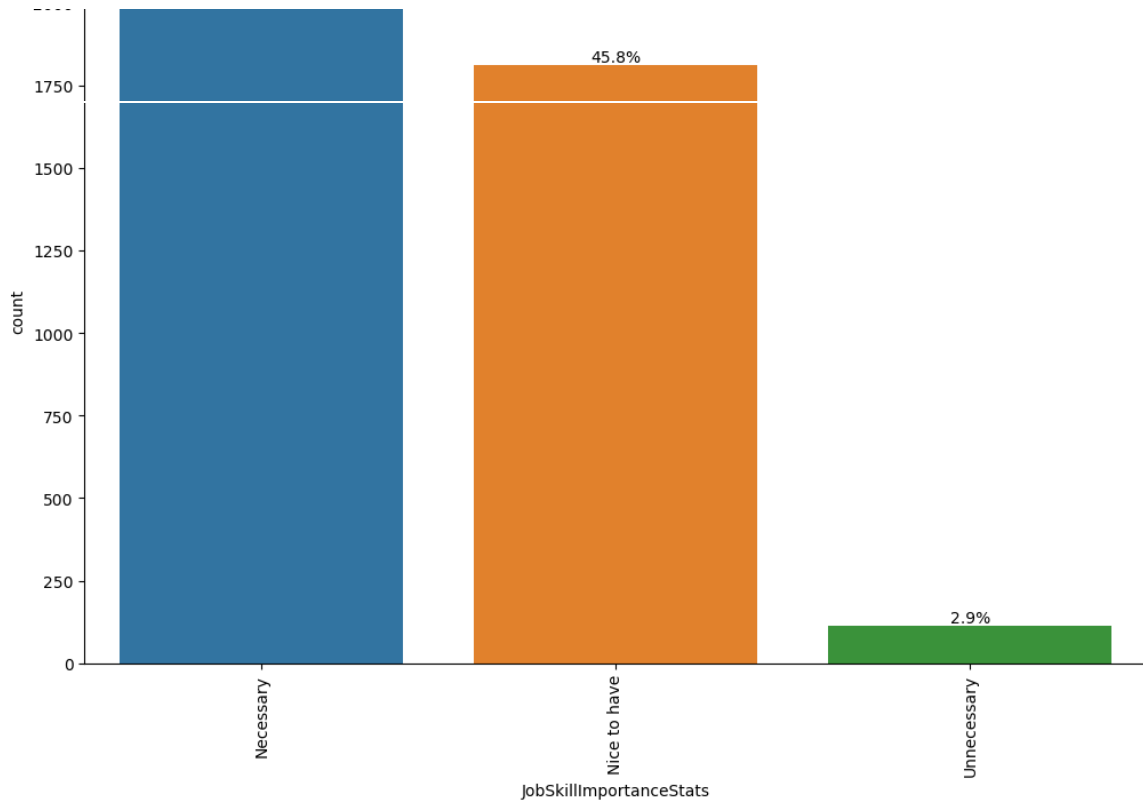


On average, how many hours a week do you spend studying data science? (top 10 or less)

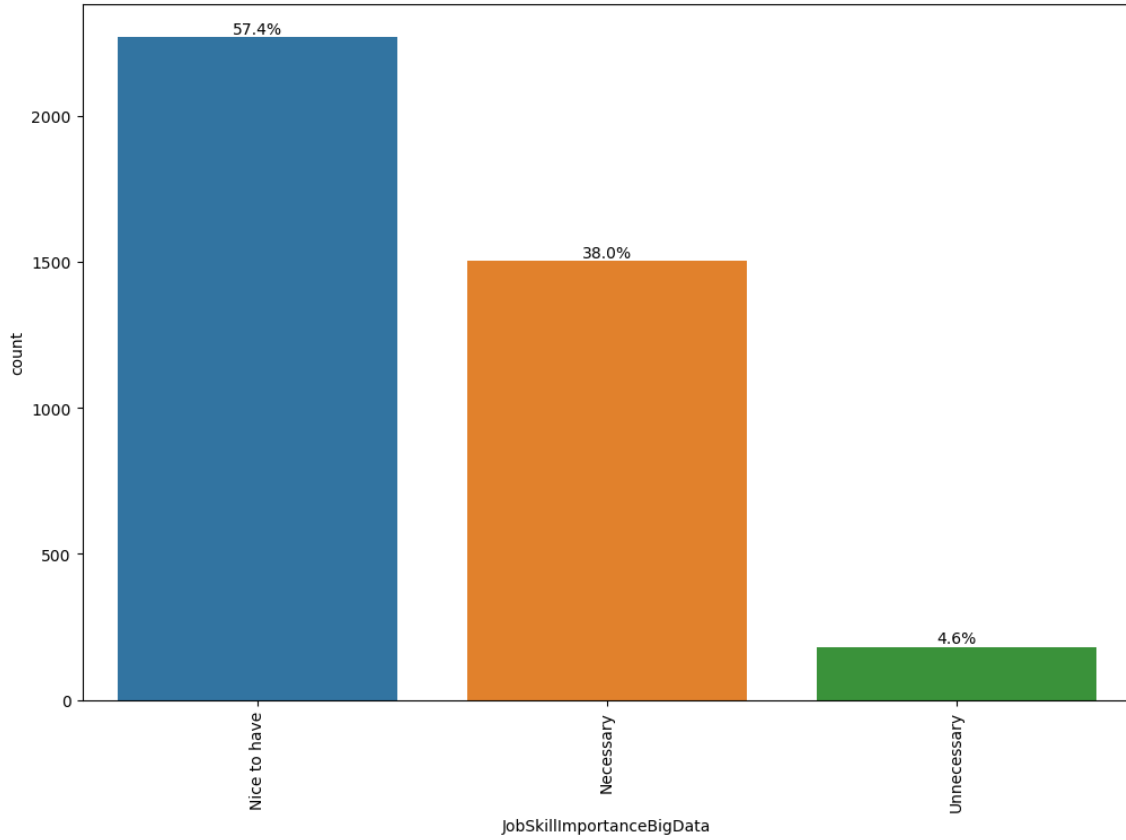


How important do you think the below skills or certifications are in getting a data science job? - Advanced Statistics (top 10 or less)



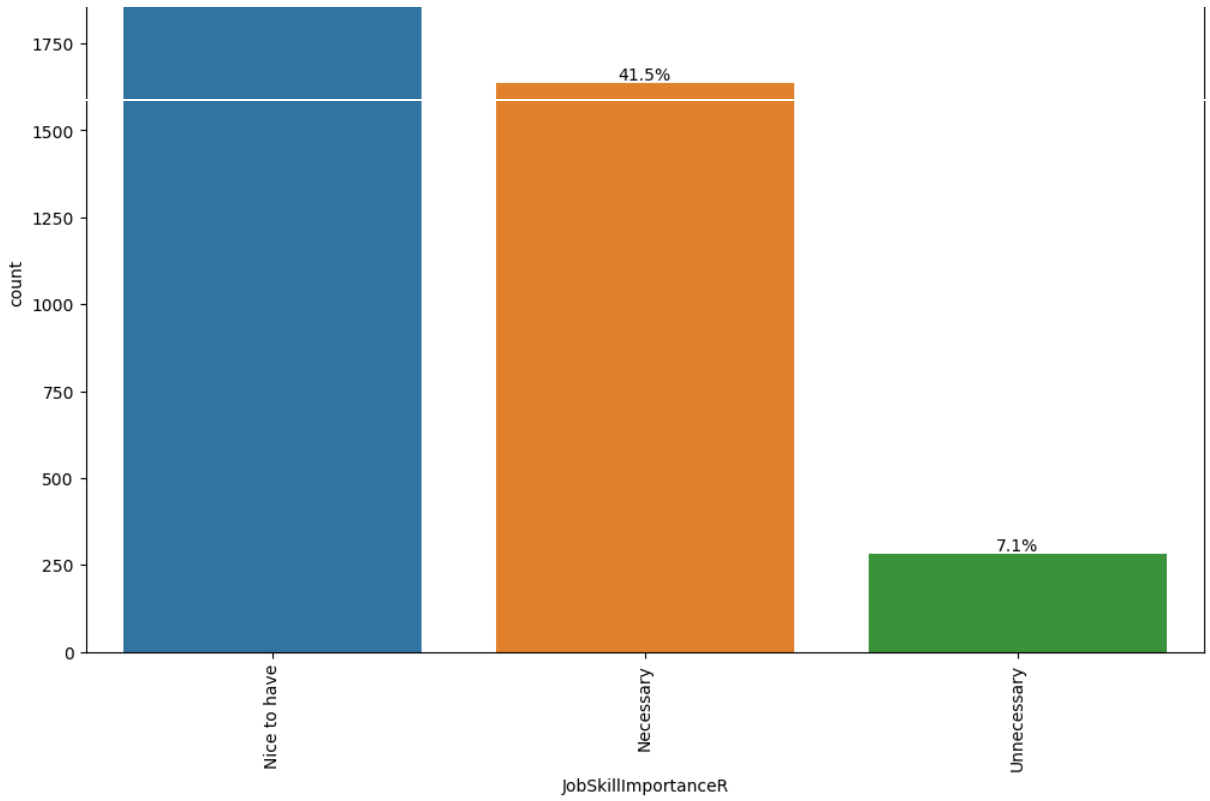


How important do you think the below skills or certifications are in getting a data science job? - 'Big Data' technology (top 10 or less)

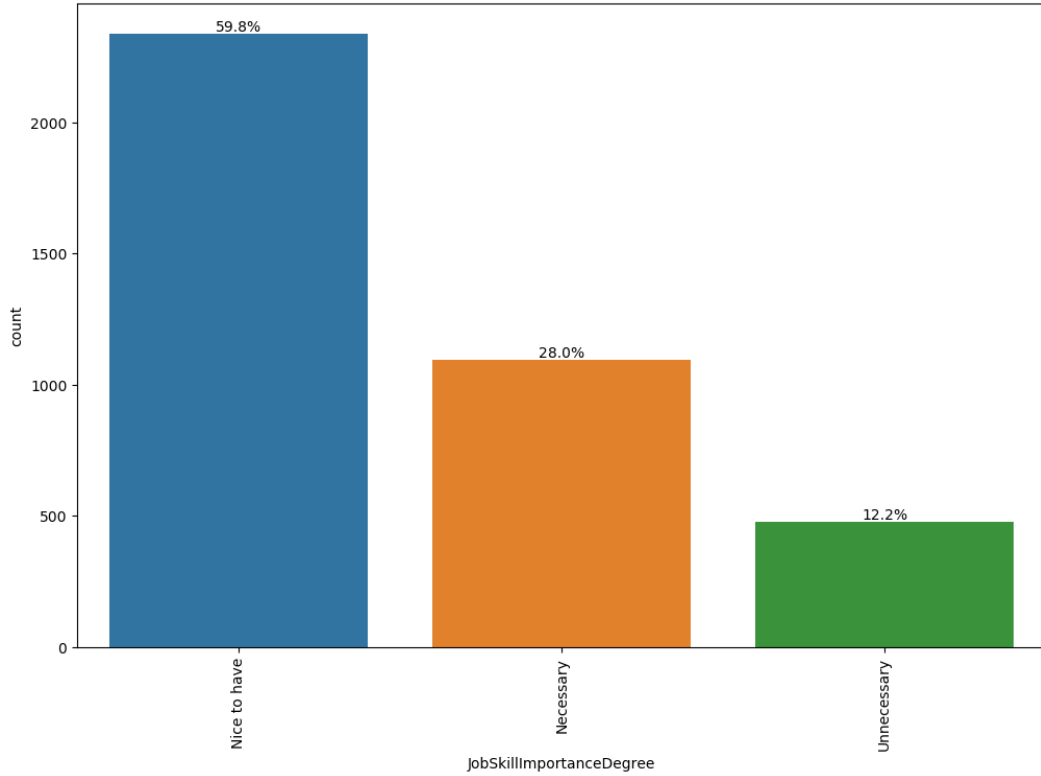


How important do you think the below skills or certifications are in getting a data science job? - Fluent in R (top 10 or less)

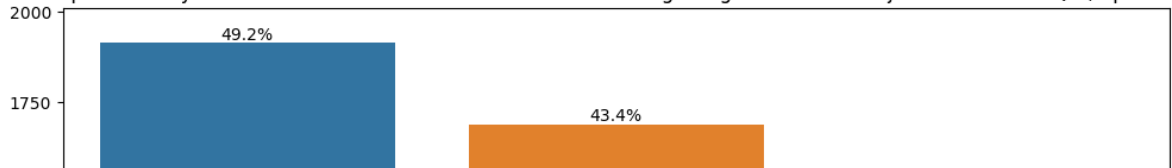


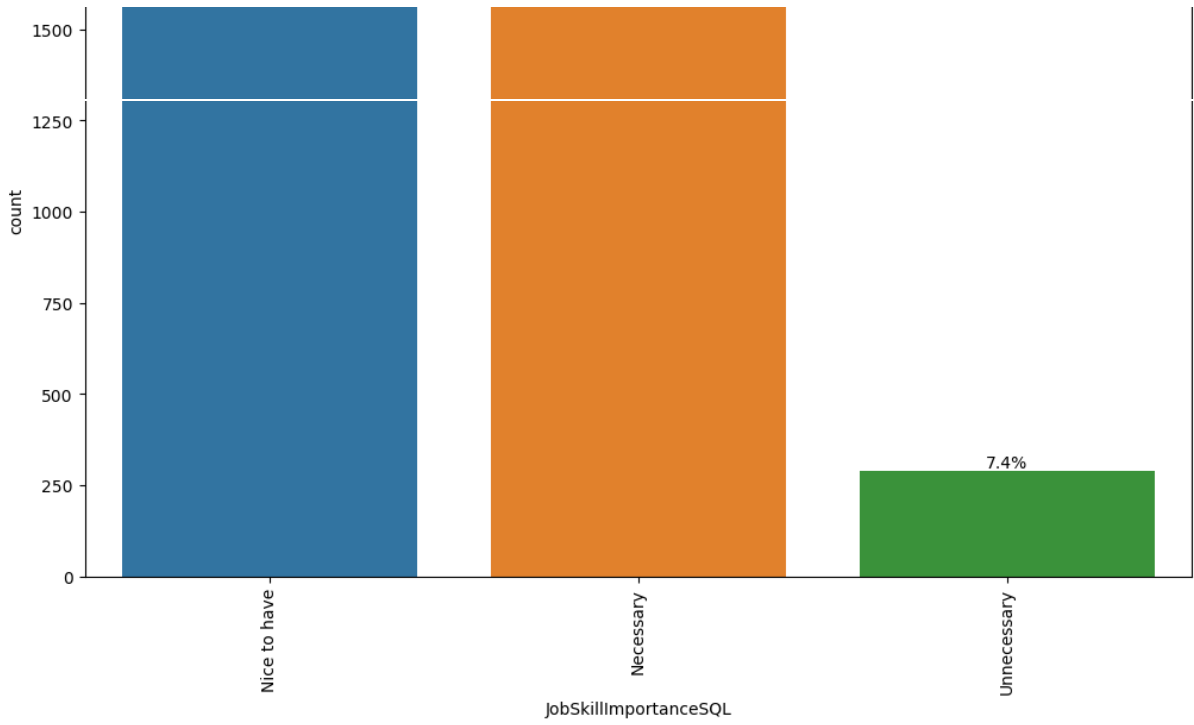


How important do you think the below skills or certifications are in getting a data science job? - Academic degree in related field (top 10 or less)

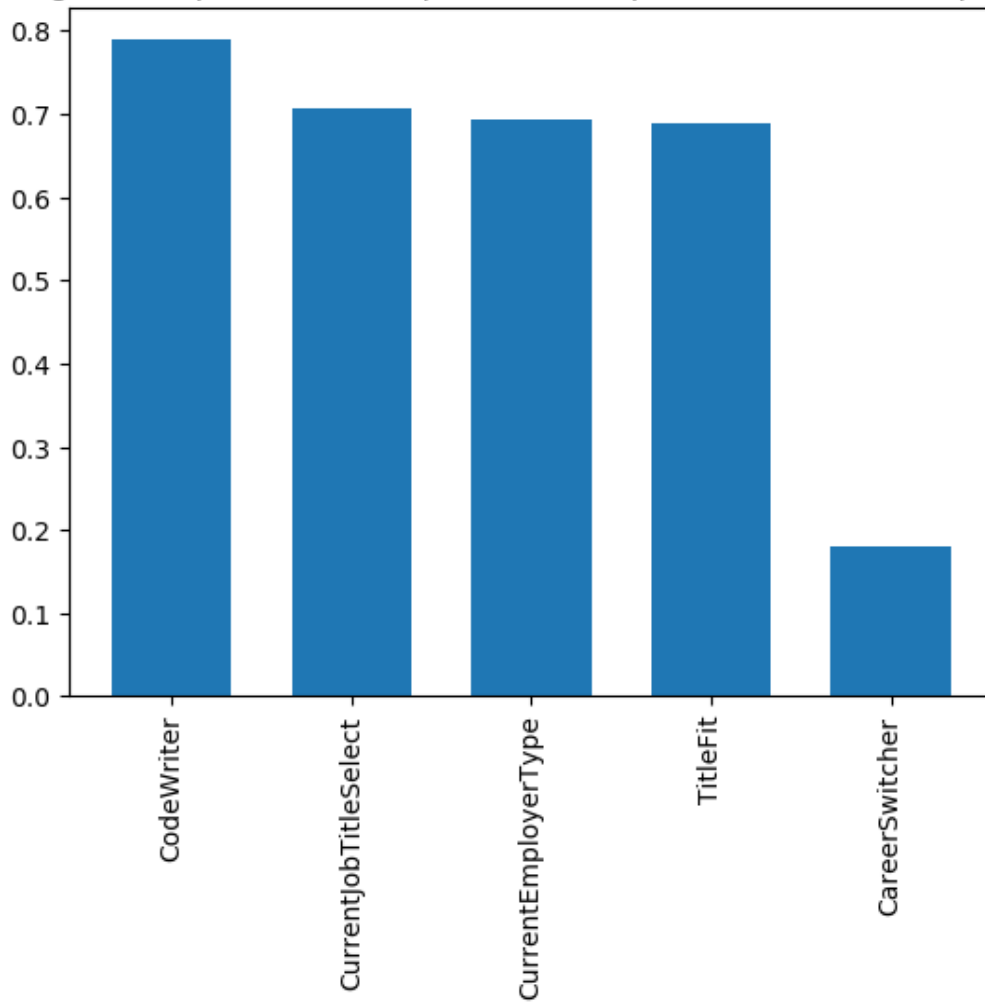


How important do you think the below skills or certifications are in getting a data science job? - Fluent in SQL (top 10 or less)



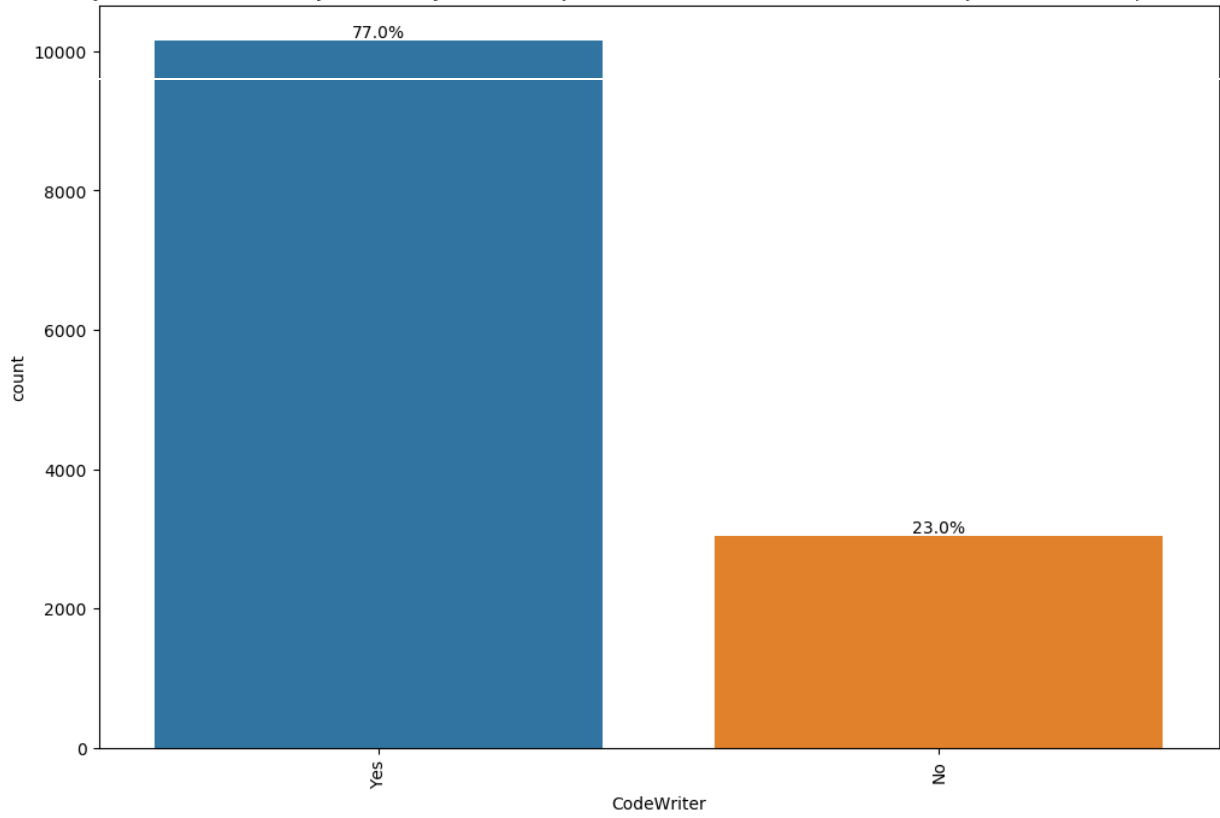


Percentage of Response on MultipleChoiceResponse of Worker1 (top 10 or less)

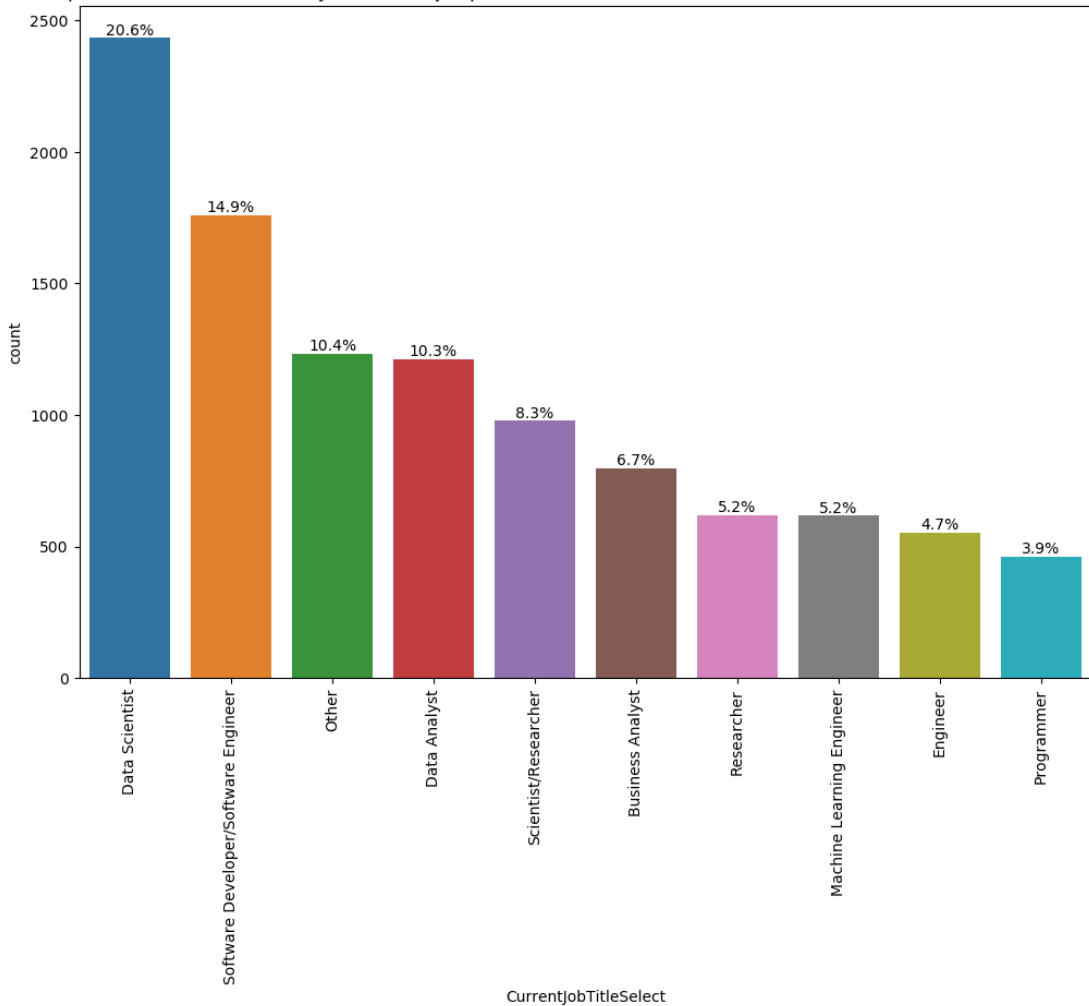


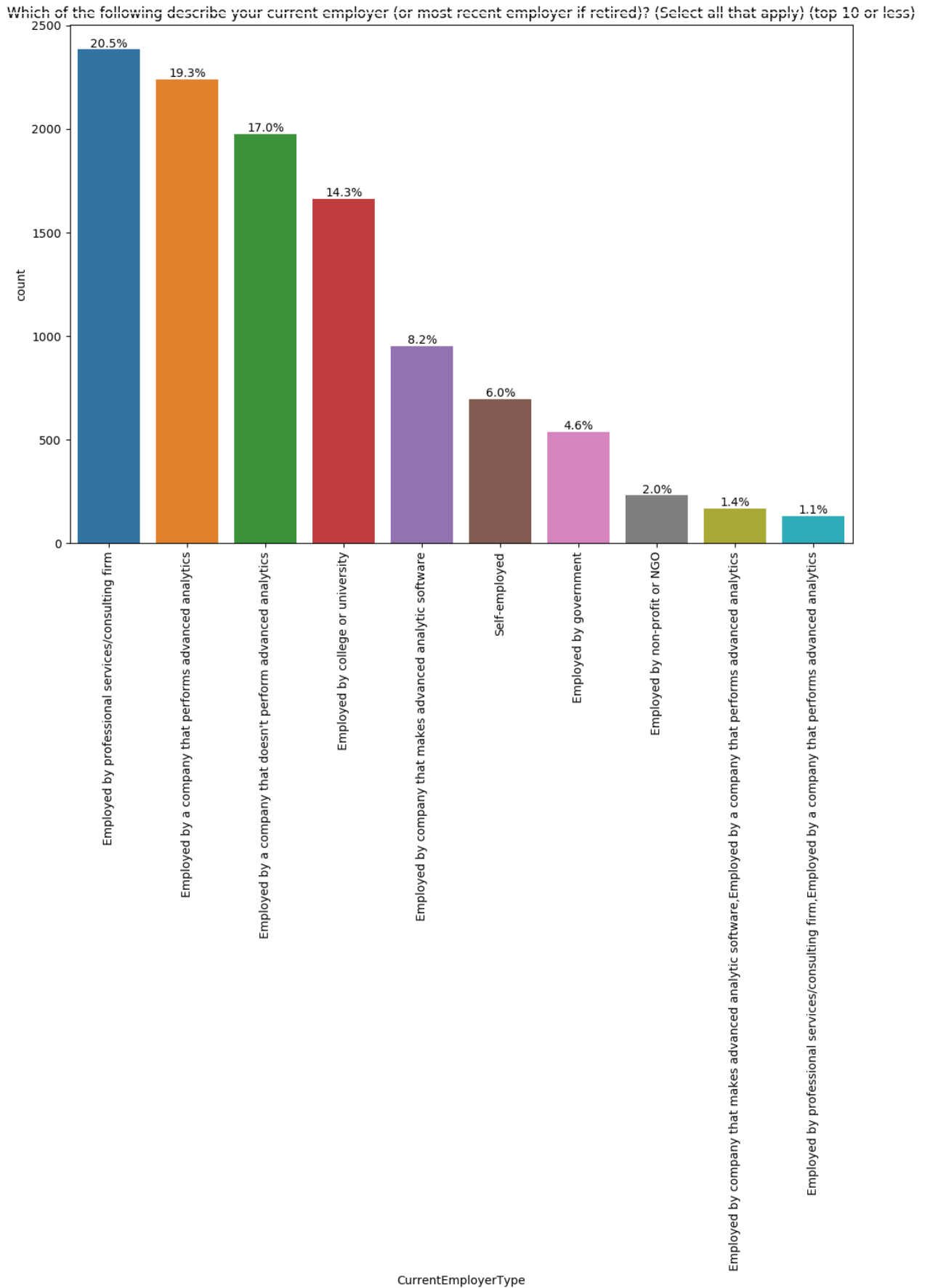
Do you write code to analyze data in your current job, freelance contracts, or most recent job if retired? (top 10 or less)

Do you write code to analyze data in your current job, freelance contracts, or most recent job if retired? (top 10 or less)

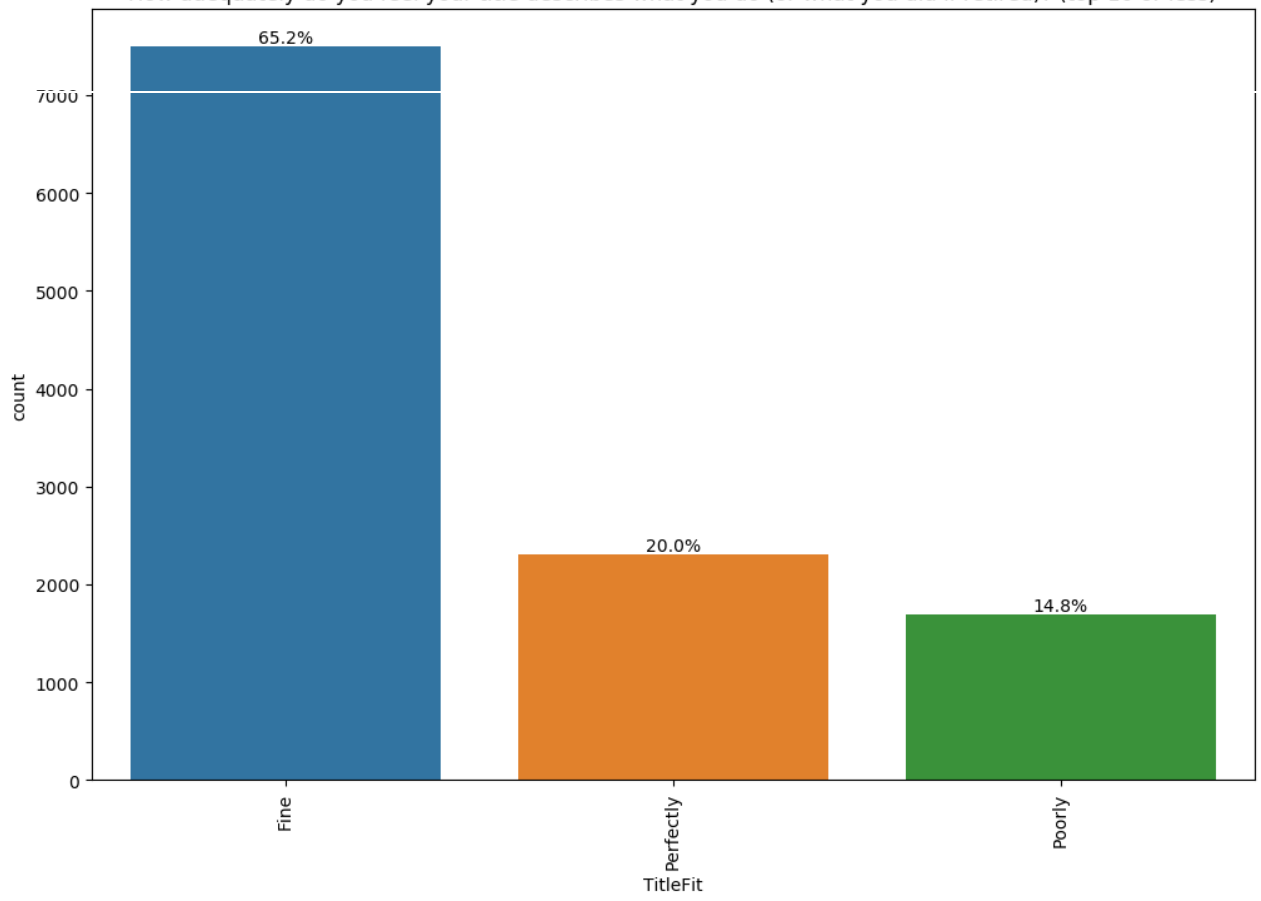


Select the option that's most similar to your current job/professional title (or most recent title if retired). - Selected Choice (top 10 or less)

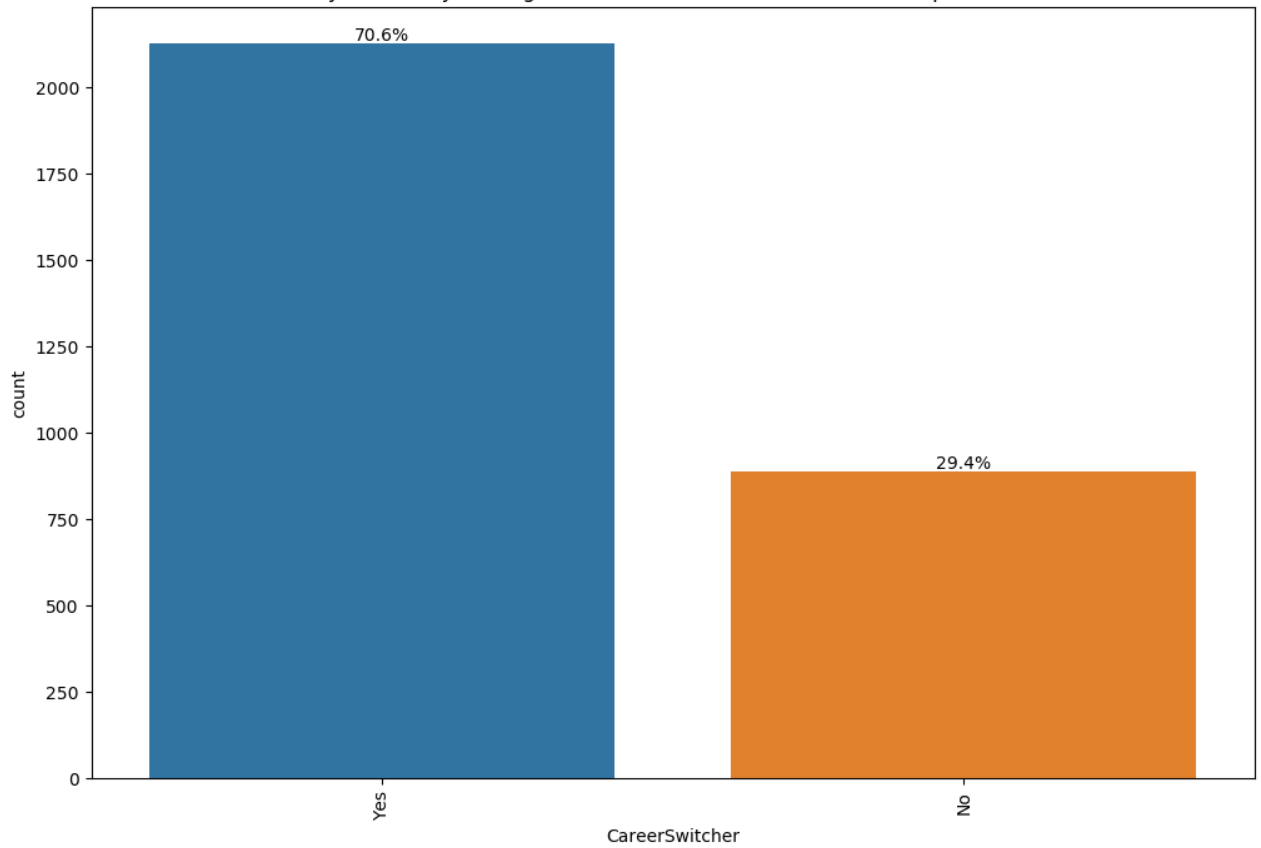




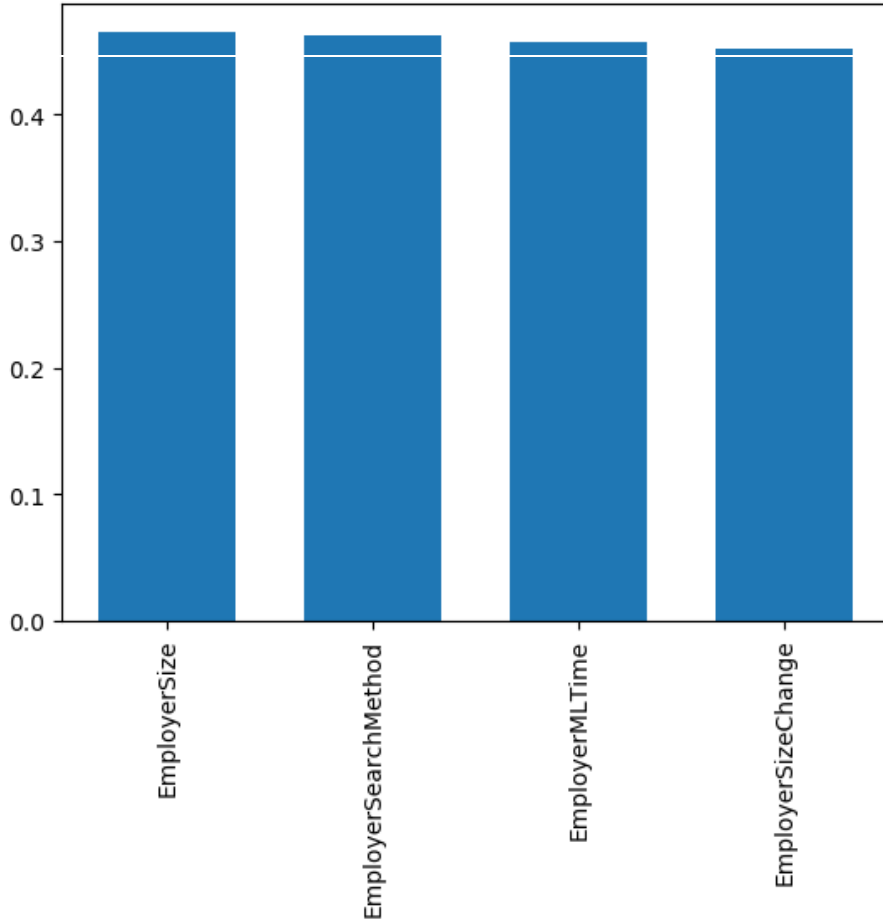
How adequately do you feel your title describes what you do (or what you did if retired)? (top 10 or less)



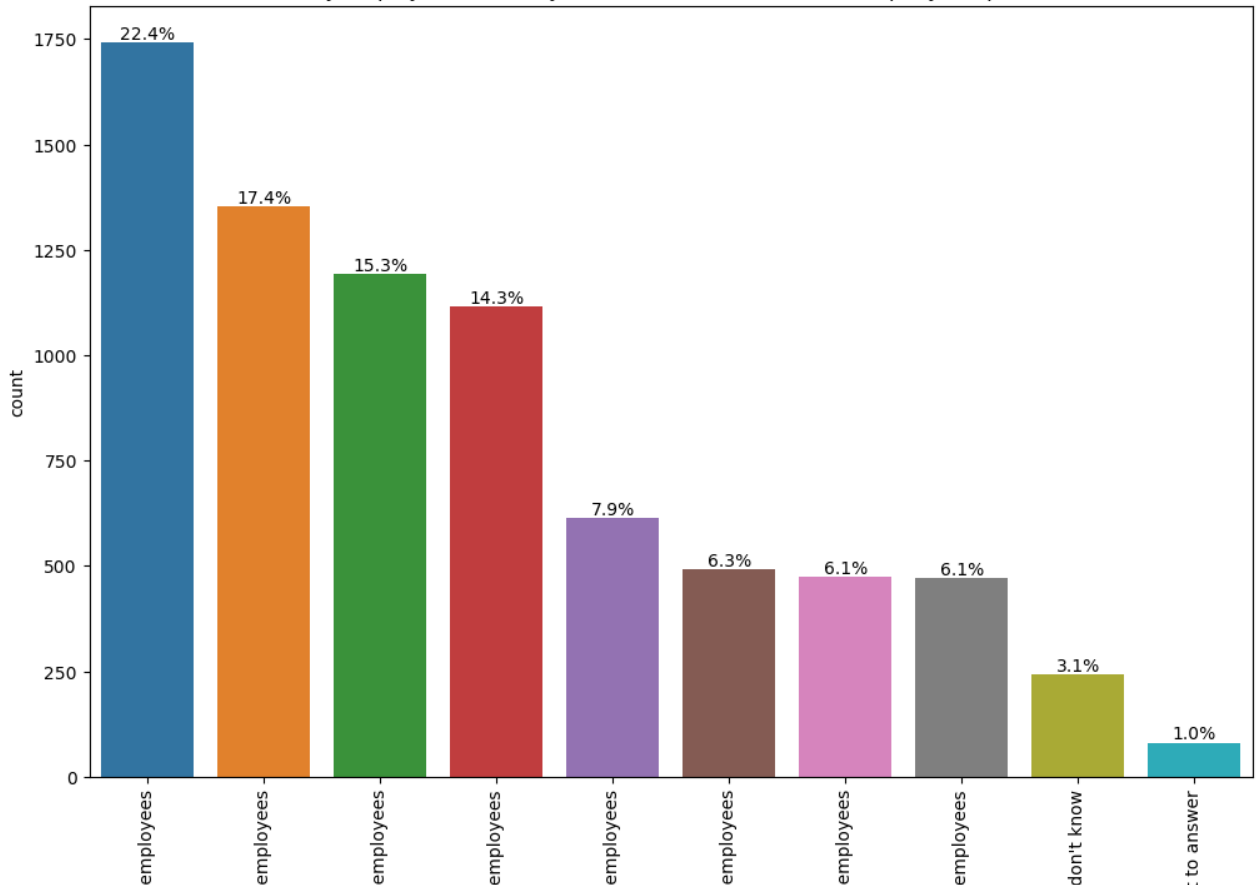
Are you actively looking to switch careers to data science? (top 10 or less)

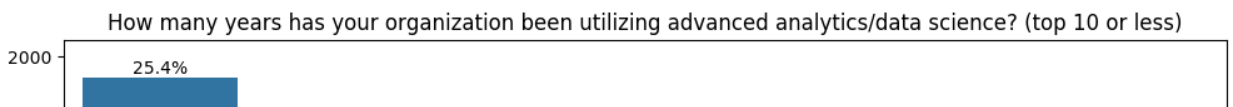
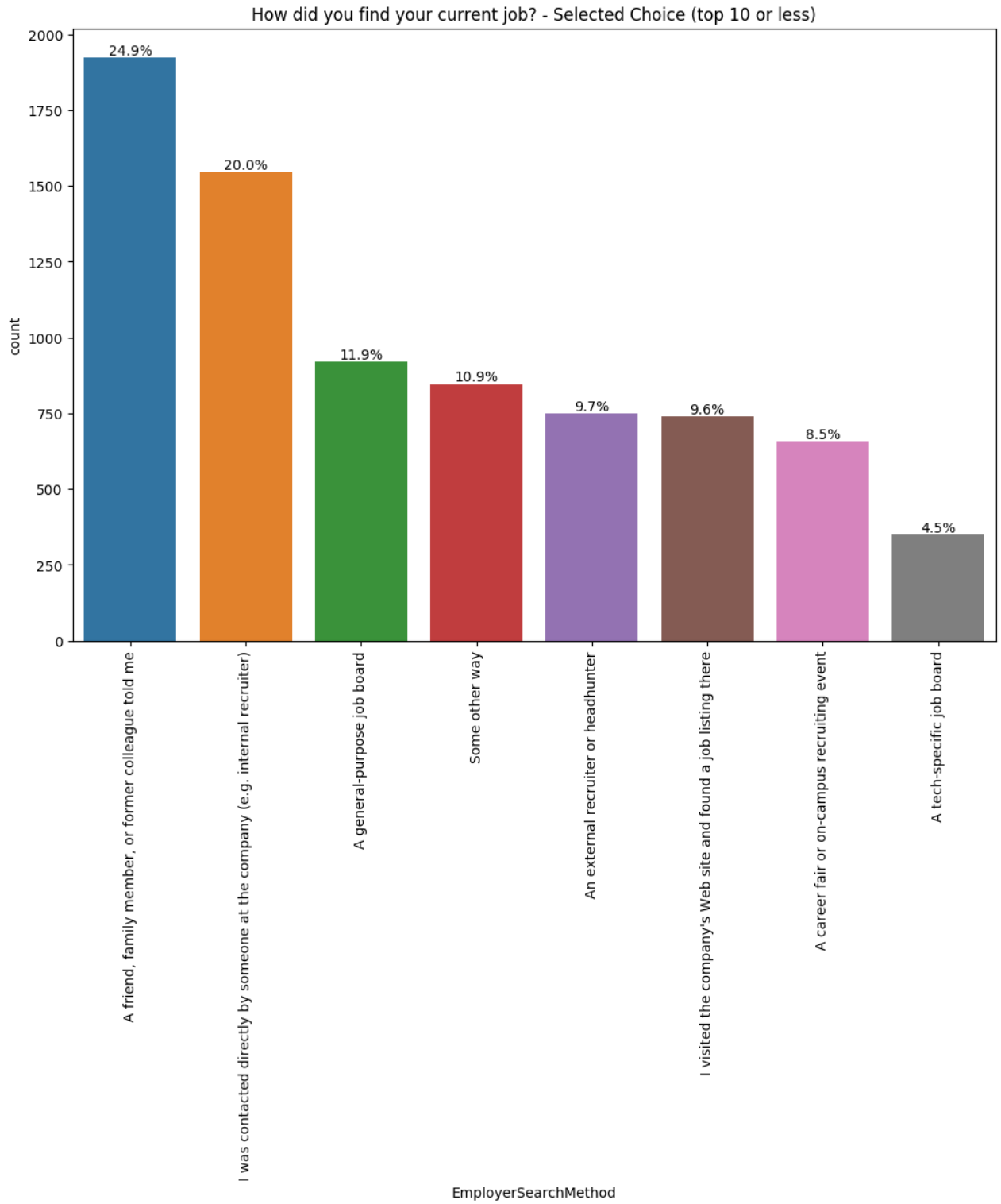


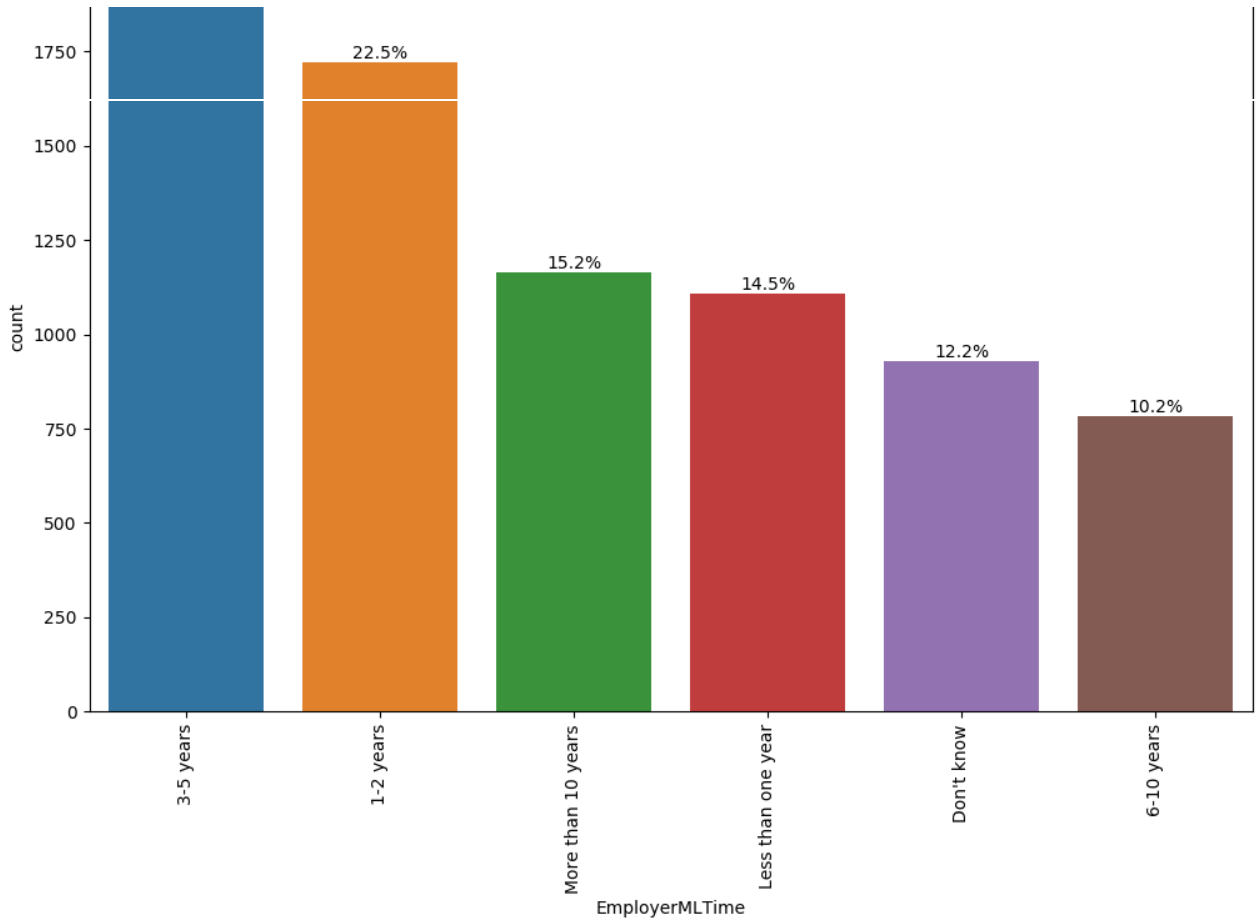
Percentage of Response on MultipleChoiceResponse of CodingWorker-NC (top 10 or less)



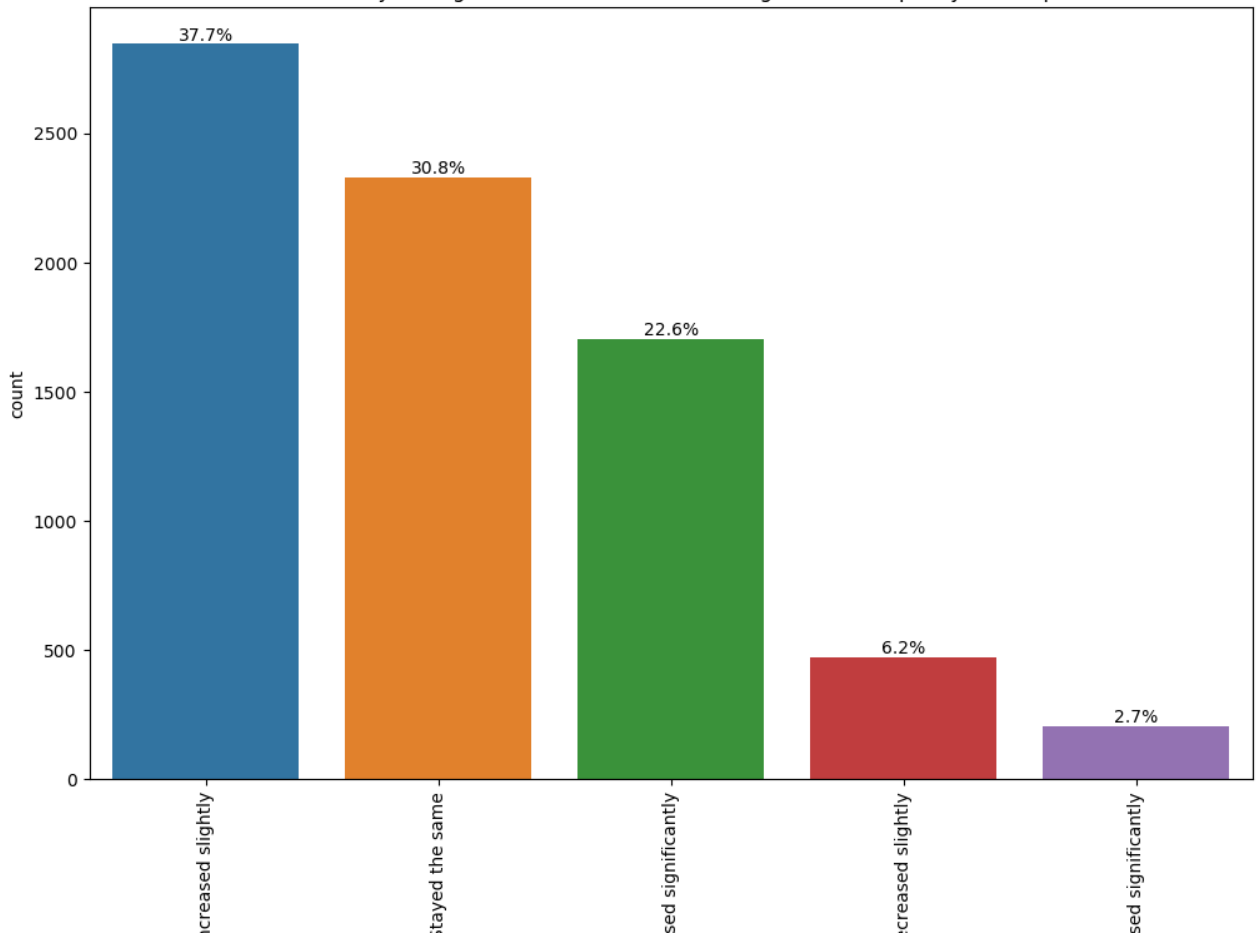
How many employees work at your current or most recent company? (top 10 or less)

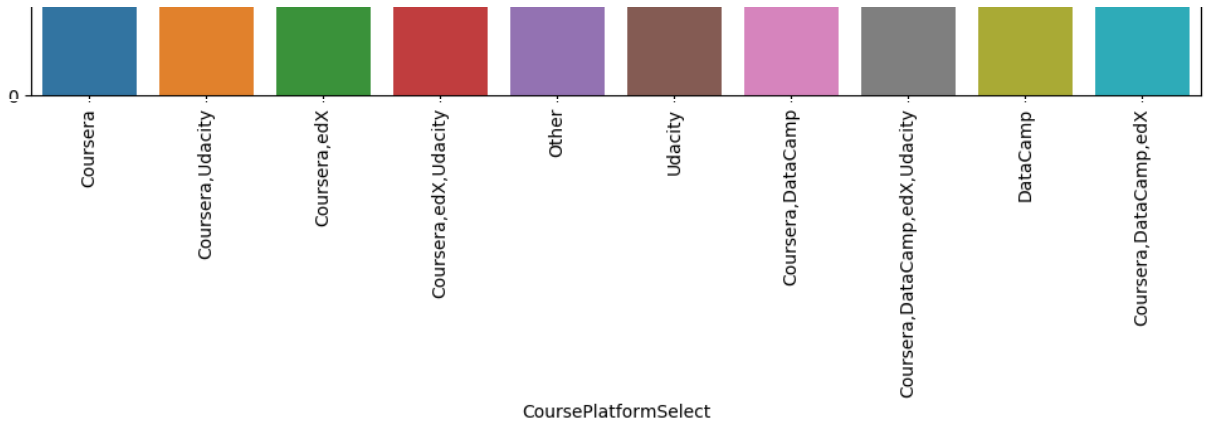




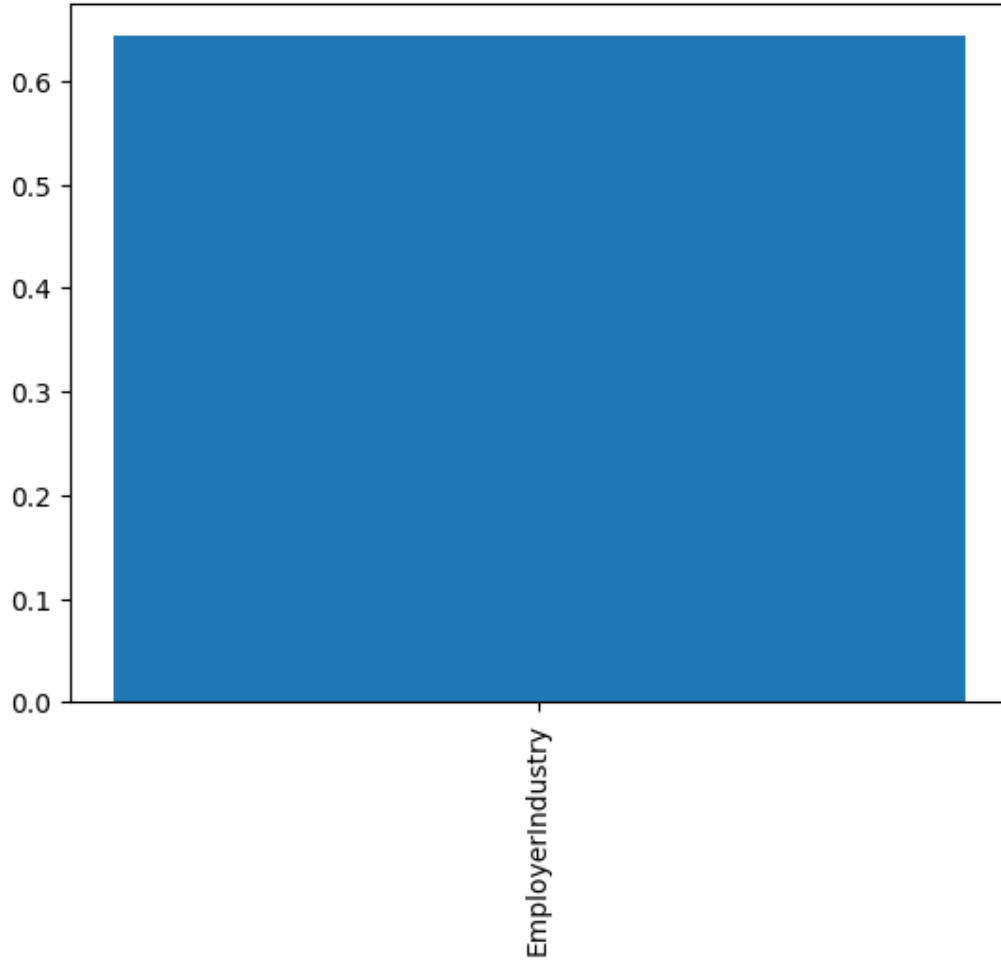


How has the size of your organization's ML/DS staff changed over the past year? (top 10 or less)

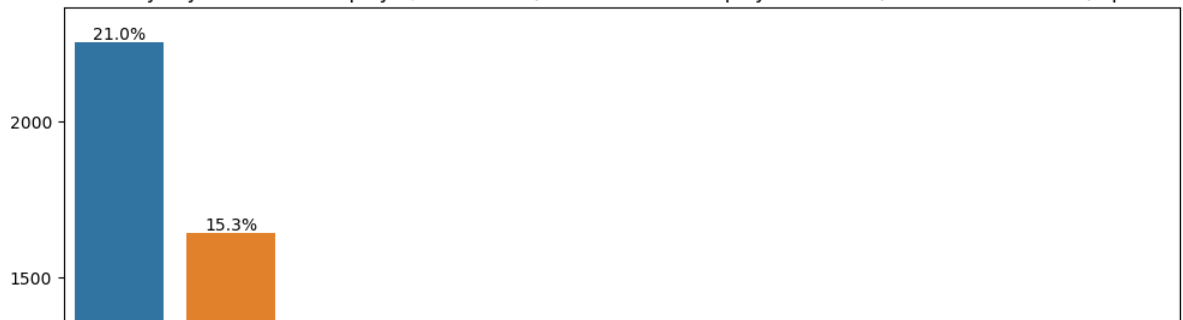


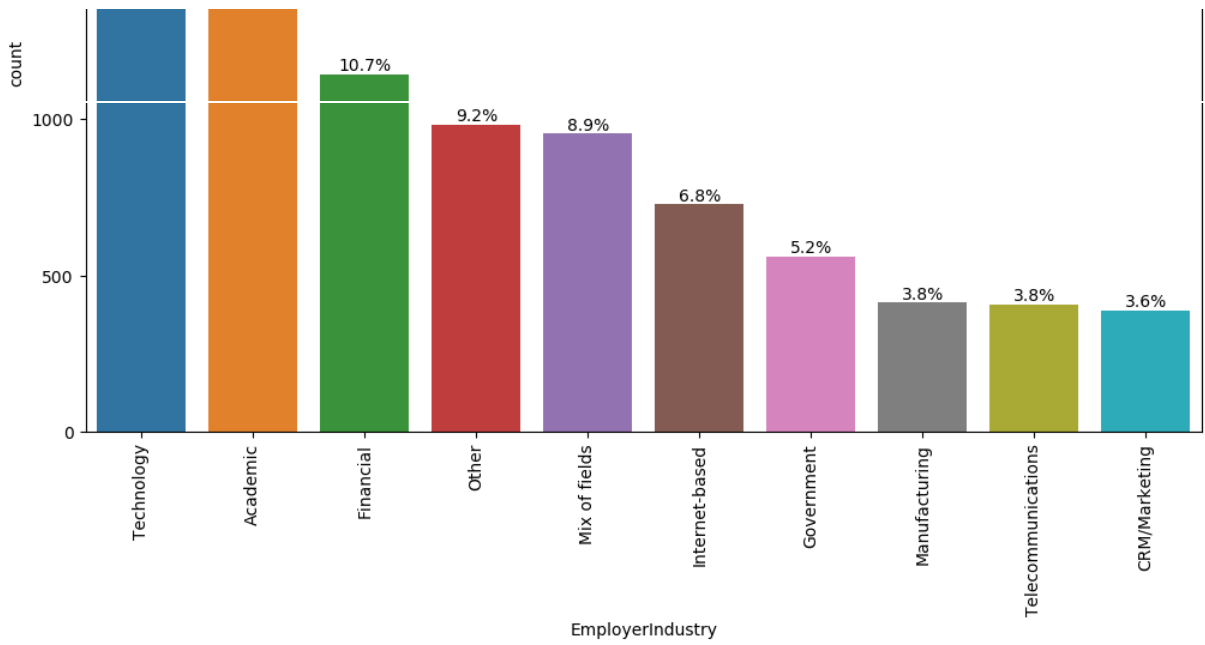


Percentage of Response on MultipleChoiceResponse of Worker (top 10 or less)

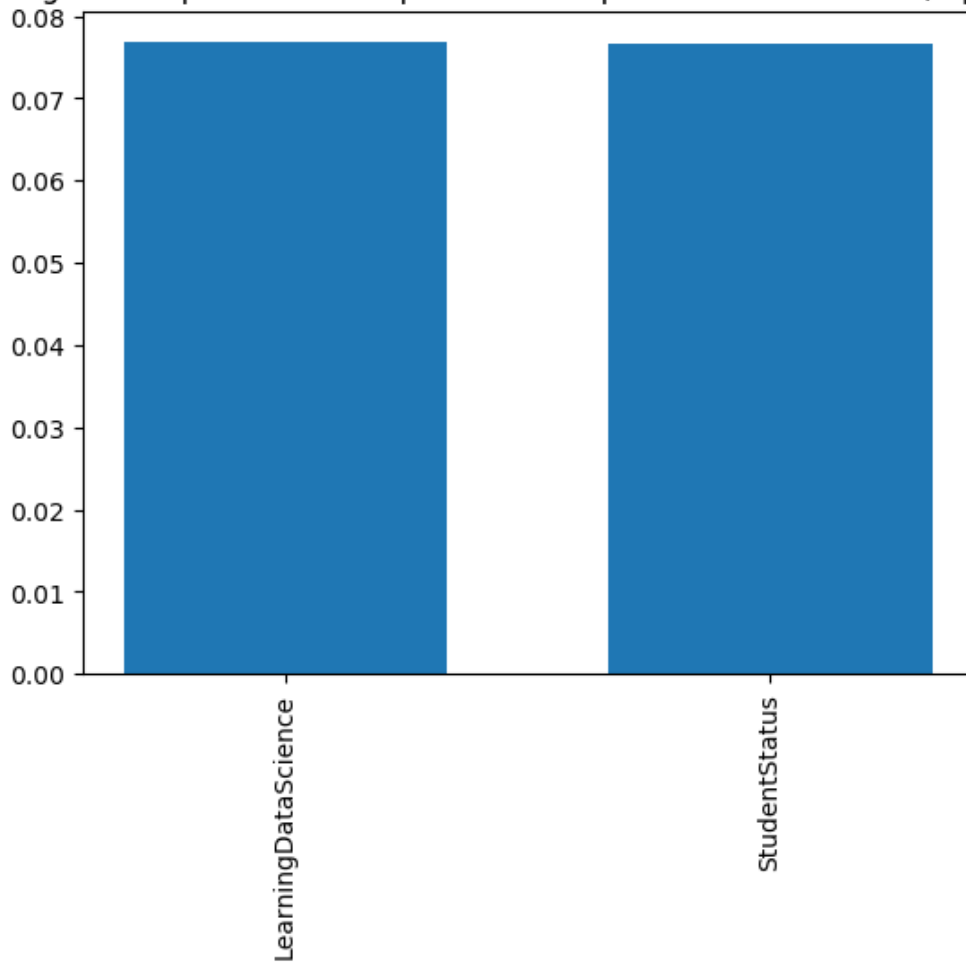


Which industry is your current employer/contract in (or most recent employer if retired)? - Selected Choice (top 10 or less)



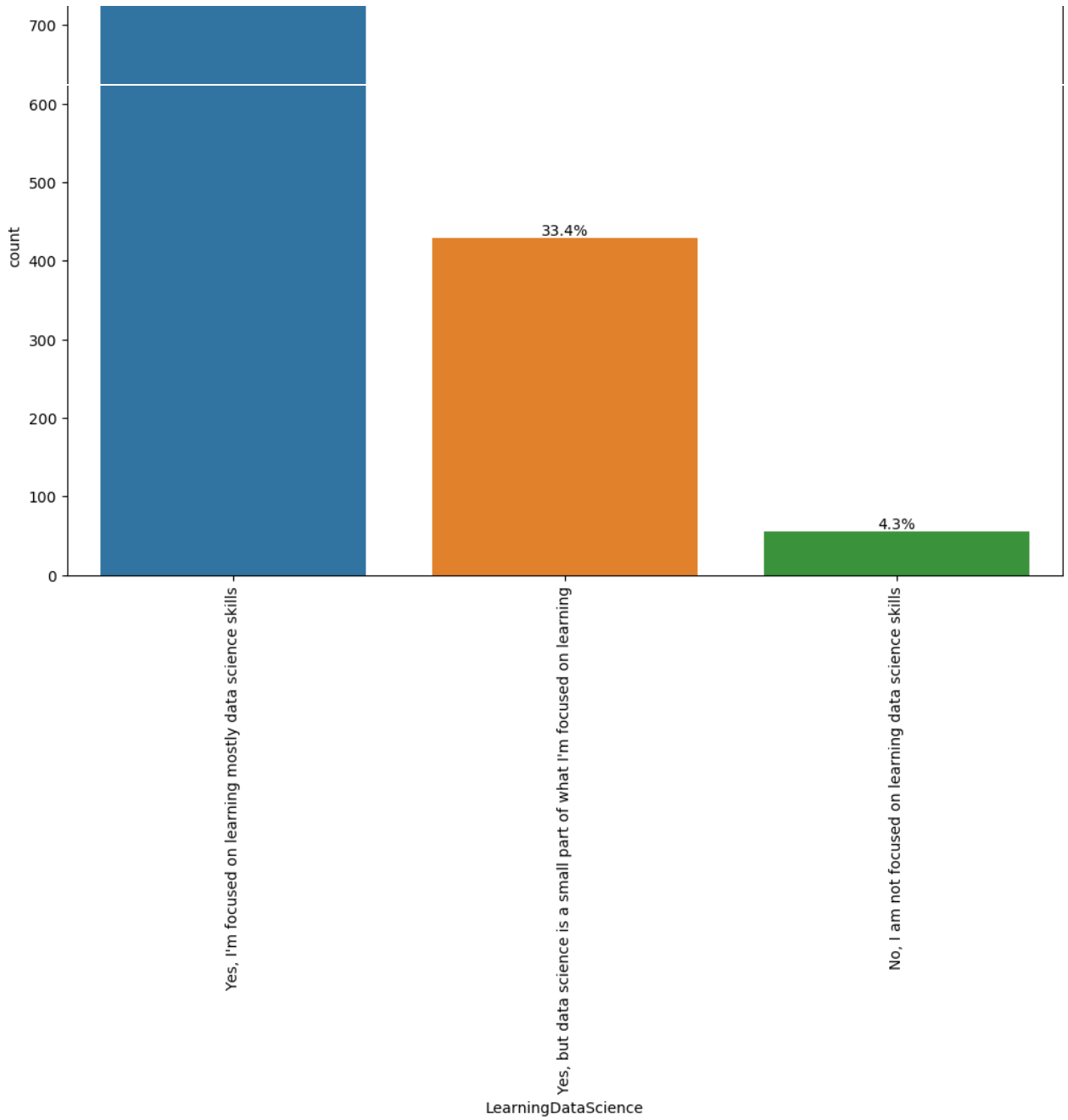


Percentage of Response on MultipleChoiceResponse of Non-worker (top 10 or less)

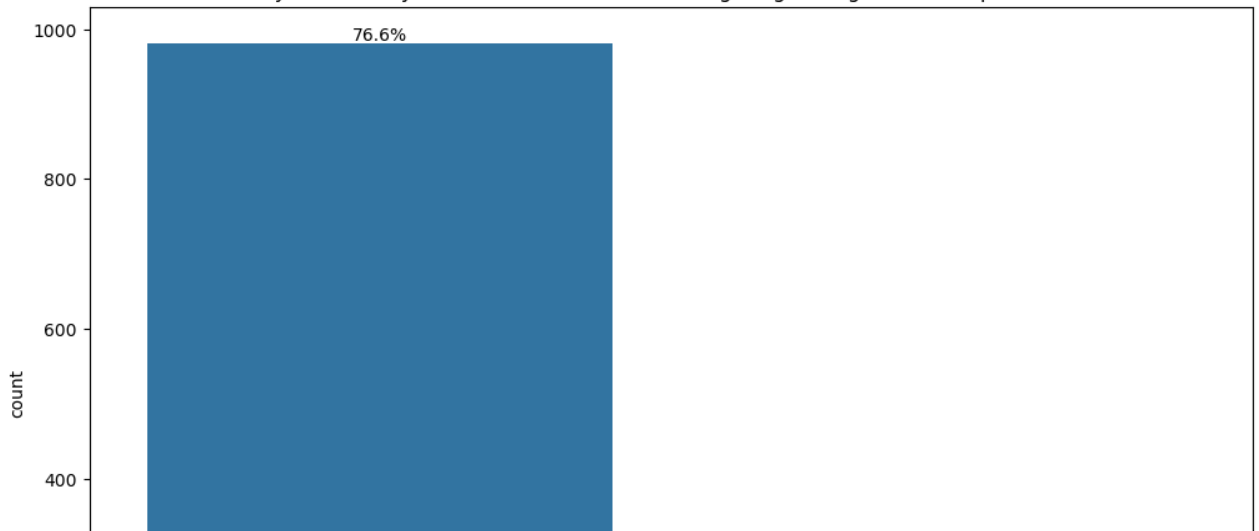


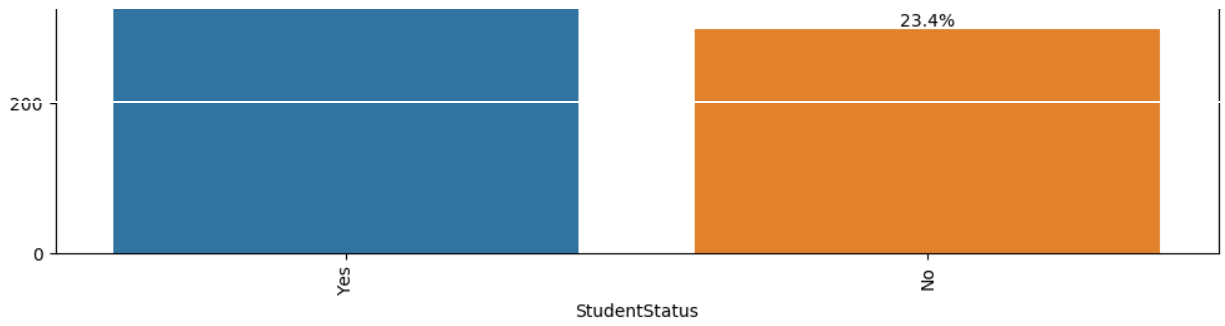
Are you currently focused on learning data science skills either formally or informally? (top 10 or less)





Are you currently enrolled as a student at a degree granting school? (top 10 or less)

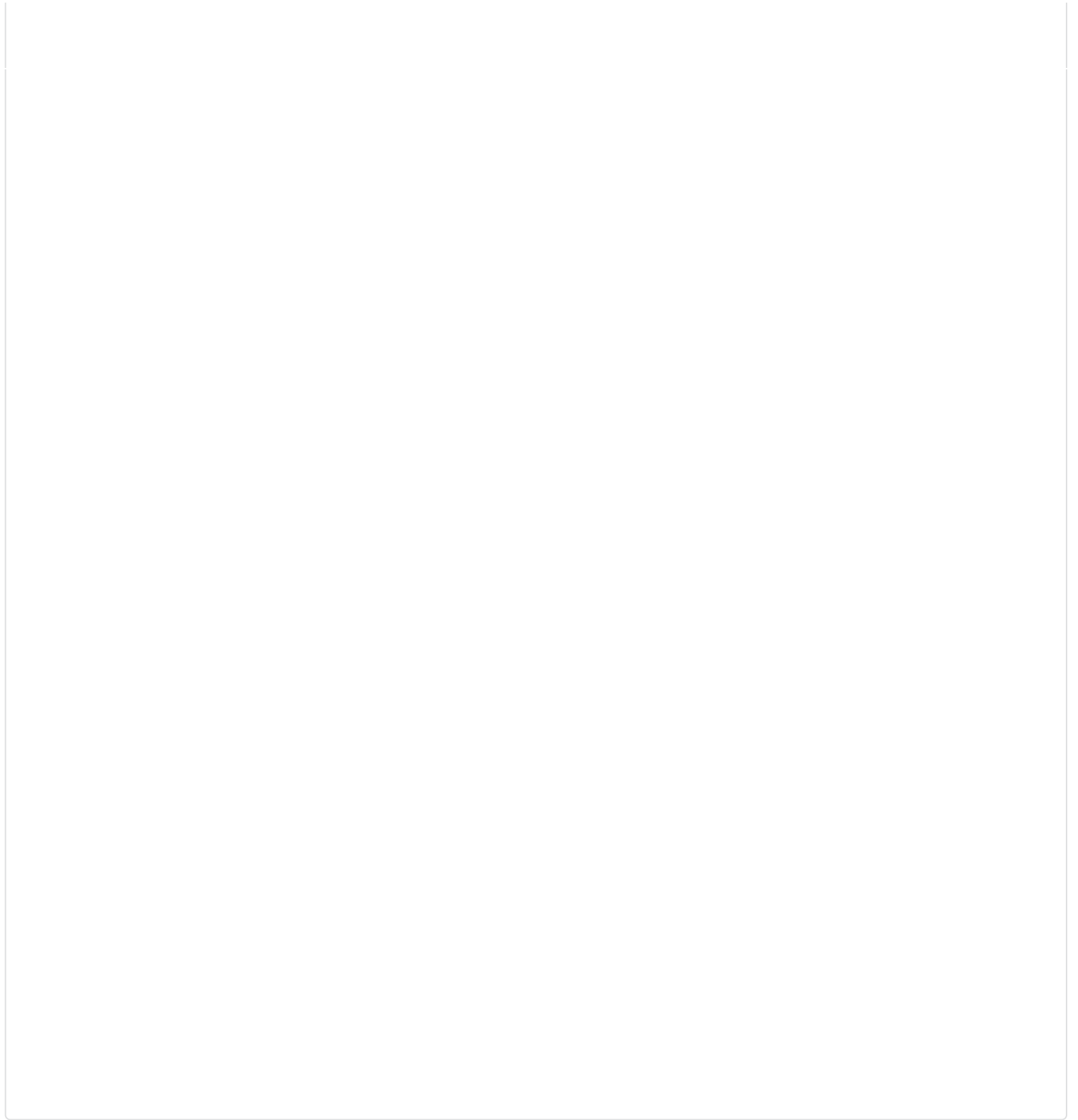




Percentage of Response on MultipleChoiceResponse of Non-switcher (top 10 or less)



In [12]:



Comments (0)

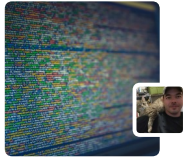
Sort by

Select...

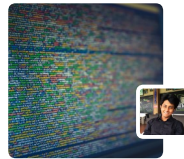


Click here to enter a comment...

Similar Kernels



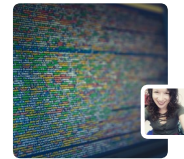
Analyzing the analyzers



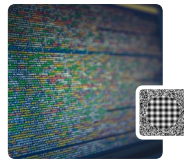
Data Science FAQ



An Interactive Deep Dive into Survey Results



Kaggle 2017 Survey Results



2017 State of Data Science - Kaggle survey