

# EDA\_CollegeInsight\_Polina

September 28, 2021

## LOAD THE DATA AND PACKAGES

```
[1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load CSV
df = pd.read_csv('collegeinsight_data_nolabel_ICs_by_year.csv',
↳encoding="ISO-8859-1") # default encoding couldn't ready all characters
```

/opt/conda/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3165:  
DtypeWarning: Columns (7,16) have mixed types.Specify dtype option on import or  
set low\_memory=False.

has\_raised = await self.run\_ast\_nodes(code\_ast.body, cell\_name,

## EDA AND DATA CLEANING

```
[2]: df.head()
```

```
[2]: data_yr_string  entity_id  entity_type  name \
0      2000-01    111100654           1    ALABAMA A & M UNIVERSITY
1      2000-01    111100663           1  UNIVERSITY OF ALABAMA AT BIRMINGHAM
2      2000-01    111100706           1  UNIVERSITY OF ALABAMA IN HUNTSVILLE
3      2000-01    111100724           1    ALABAMA STATE UNIVERSITY
4      2000-01    111100751           1    UNIVERSITY OF ALABAMA

      city state  state_fips  cong_dist  webaddr  sector  ...  \
0    NORMAL   AL         1.0        NaN  WWW.AAMU.EDU    1.0  ...
1  BIRMINGHAM   AL         1.0        NaN  www.uab.edu    1.0  ...
2  HUNTSVILLE   AL         1.0        NaN  www.uah.edu    1.0  ...
3  MONTGOMERY   AL         1.0        NaN  www.alasu.edu    1.0  ...
4  TUSCALOOSA   AL         1.0        NaN  www.ua.edu     1.0  ...

      compl_rpy_7yr_n  noncom_rpy_7yr_n  dep_rpy_7yr_n  ind_rpy_7yr_n  \
0                NaN                NaN            NaN            NaN
1                NaN                NaN            NaN            NaN
2                NaN                NaN            NaN            NaN
3                NaN                NaN            NaN            NaN
```

	pell_rpy_7yr_n	nopell_rpy_7yr_n	female_rpy_7yr_n	male_rpy_7yr_n	\
0	NaN	NaN	NaN	NaN	
1	NaN	NaN	NaN	NaN	
2	NaN	NaN	NaN	NaN	
3	NaN	NaN	NaN	NaN	
4	NaN	NaN	NaN	NaN	

	firstgen_rpy_7yr_n	notfirstgen_rpy_7yr_n
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN

[5 rows x 408 columns]

```
[3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 139149 entries, 0 to 139148
Columns: 408 entries, data_yr_string to notfirstgen_rpy_7yr_n
dtypes: float64(396), int64(3), object(9)
memory usage: 433.1+ MB
```

```
[4]: # all columns
my_list = list(df)

print (my_list)
```

```
['data_yr_string', 'entity_id', 'entity_type', 'name', 'city', 'state',
'state_fips', 'cong_dist', 'webaddr', 'sector', 'control', 'level', 'title_iv',
'hbcu', 'unitid', 'opeid6', 'opeid8', 'reporting_type', 'deggrant', 'ugoffer',
'carnegie2000', 'carnegie2010', 'carnegie2015', 'carnegie2018',
'fa_loans_debt_p', 'fa_loans_debt_avg_d', 'fa_loans_debt_pc_d',
'fa_loans_fed_p', 'fa_loans_fed_avg_d', 'fa_loans_fed_pc_d',
'fa_loans_nfed_pc_d', 'fa_loans_fed_vol_p', 'fa_loans_nfed_vol_p',
'fa_loans_debt_cohort_n', 'fa_loans_debt_n', 'fa_ftft_loans_any_n',
'fa_ftft_loans_any_p', 'fa_ftft_loans_any_avg_d', 'fa_ug_nd_n', 'fa_ug_nd_p',
'fa_ug_nd_met_p', 'fa_ug_award_n', 'fa_ug_nd_full_n', 'fa_ug_nd_full_p',
'fa_ftft_aid_n', 'fa_ftft_aid_p', 'fa_grants_tot_d', 'fa_grants_fed_tot_d',
'fa_grants_fed_nb_d', 'fa_grants_fed_nnb_d', 'fa_grants_state_tot_d',
'fa_grants_state_nb_d', 'fa_grants_state_nnb_d', 'fa_grants_inst_tot_d',
'fa_grants_inst_nb_d', 'fa_grants_inst_nnb_d', 'fa_grants_inst_nb_p',
'fa_grants_ext_tot_d', 'fa_grants_ext_nb_d', 'fa_grants_ext_nnb_d',
'fa_ug_avg_grant_d', 'grad_debt_mdn', 'wdraw_debt_mdn', 'fa_fws_n',
'fa_fws_tot_d', 'fa_fws_avg_d', 'fa_ows_n', 'fa_ows_tot_d', 'fa_ows_avg_d',
```

'fa\_ftft\_grants\_any\_n', 'fa\_ftft\_grants\_any\_p', 'fa\_ftft\_grants\_any\_avg\_d',  
 'fa\_ftft\_grants\_any\_pc\_d', 'fa\_ftft\_grants\_fed\_n', 'fa\_ftft\_grants\_fed\_p',  
 'fa\_ftft\_grants\_fed\_avg\_d', 'fa\_ftft\_grants\_fed\_pc\_d', 'fa\_ftft\_grants\_sta\_n',  
 'fa\_ftft\_grants\_sta\_p', 'fa\_ftft\_grants\_sta\_avg\_d', 'fa\_ftft\_grants\_sta\_pc\_d',  
 'fa\_ftft\_grants\_inst\_n', 'fa\_ftft\_grants\_inst\_p', 'fa\_ftft\_grants\_inst\_avg\_d',  
 'fa\_ftft\_grants\_inst\_pc\_d', 'fa\_ftft\_loans\_any\_pc\_d', 'fa\_ftft\_loans\_fed\_n',  
 'fa\_ftft\_loans\_fed\_p', 'fa\_ftft\_loans\_fed\_avg\_d', 'fa\_ftft\_loans\_fed\_pc\_d',  
 'fa\_ftft\_loans\_nfed\_n', 'fa\_ftft\_loans\_nfed\_p', 'fa\_ftft\_loans\_nfed\_avg\_d',  
 'fa\_ftft\_loans\_nfed\_pc\_d', 'coa\_tuit\_fees\_d', 'coa\_books\_supp\_d',  
 'coa\_on\_room\_board\_d', 'coa\_on\_other\_d', 'coa\_on\_tcoa\_d', 'coa\_off\_tcoa\_d',  
 'coa\_cipcode\_prog1', 'coa\_length\_prog1', 'coa\_off\_room\_board\_d',  
 'coa\_off\_other\_d', 'coa\_room\_board\_d', 'coa\_other\_d', 'coa\_tcoa\_d',  
 'fa\_sfa\_pell\_p', 'fa\_sfa\_pell\_avg\_d', 'fa\_sfa\_pell\_n', 'fa\_sfa\_pell\_tot\_d',  
 'en\_amerind\_n', 'en\_amerind\_p', 'en\_asn\_pac\_n', 'en\_asn\_pac\_p', 'en\_black\_n',  
 'en\_black\_p', 'en\_hispanic\_n', 'en\_hispanic\_p', 'en\_white\_n', 'en\_white\_p',  
 'en\_intl\_n', 'en\_intl\_p', 'en\_multiracial\_n', 'en\_multiracial\_p',  
 'en\_race\_unknown\_n', 'en\_race\_unknown\_p', 'en\_12mo\_amerind\_n',  
 'en\_12mo\_amerind\_p', 'en\_12mo\_asn\_pac\_n', 'en\_12mo\_asn\_pac\_p',  
 'en\_12mo\_black\_n', 'en\_12mo\_black\_p', 'en\_12mo\_hispanic\_n',  
 'en\_12mo\_hispanic\_p', 'en\_12mo\_white\_n', 'en\_12mo\_white\_p', 'en\_12mo\_intl\_n',  
 'en\_12mo\_intl\_p', 'en\_12mo\_multiracial\_n', 'en\_12mo\_multiracial\_p',  
 'en\_12mo\_race\_unknown\_n', 'en\_12mo\_race\_unknown\_p', 'en\_12mo\_ipeds\_tot\_n',  
 'en\_12mo\_ipeds\_ug\_n', 'en\_12mo\_ipeds\_grad\_n', 'en\_fall\_tot\_n',  
 'en\_fall\_ug\_tot\_n', 'en\_fall\_grad\_n', 'en\_fall\_ug\_ft\_n', 'en\_fall\_ug\_pt\_n',  
 'en\_fall\_ug\_ft\_p', 'en\_fresh\_n', 'en\_fresh\_ft\_n', 'en\_fall\_fte\_n',  
 'en\_under25\_n', 'en\_under25\_p', 'en\_25plus\_n', 'en\_25plus\_p',  
 'comp\_fresh\_ret\_p', 'fa\_ftft\_ug\_n', 'fa\_ftft\_ug\_p', 'fa\_dep\_n',  
 'ef\_undg\_allstudentsenrolled', 'ef\_undg\_students\_distanceonly',  
 'ef\_undg\_students\_distanceonly\_p', 'ef\_undg\_distance\_samestate',  
 'ef\_undg\_instate\_dist\_p', 'ef\_undg\_students\_distancesome',  
 'ef\_undg\_students\_distancesome\_p', 'ef\_undg\_students\_distancenone',  
 'ef\_undg\_students\_distancenone\_p', 'ef\_undg\_distance\_US\_notsamestate',  
 'ef\_undg\_distance\_US\_stateunknown', 'ef\_undg\_distance\_outside\_US',  
 'ef\_undg\_distance\_unknown', 'ef\_all\_allstudentsenrolled',  
 'ef\_all\_students\_distanceonly', 'ef\_all\_students\_distanceonly\_p',  
 'ef\_all\_distance\_samestate', 'ef\_all\_instate\_dist\_p',  
 'ef\_all\_students\_distancesome', 'ef\_all\_students\_distancesome\_p',  
 'ef\_all\_students\_distancenone', 'ef\_all\_students\_distancenone\_p',  
 'ef\_all\_distance\_US\_notsamestate', 'ef\_all\_distance\_US\_stateunknown',  
 'ef\_all\_distance\_outside\_US', 'ef\_all\_distance\_unknown',  
 'fa\_sfa\_living\_onscamp\_p', 'fa\_sfa\_living\_offfam\_p', 'fa\_sfa\_living\_offnotfam\_p',  
 'fa\_sfa\_living\_unkn\_p', 'comp\_grad\_p', 'pell\_comp\_gr\_p', 'nonpell\_comp\_gr\_p',  
 'ba\_gr\_p', 'pell\_ba\_gr\_p', 'nonpell\_ba\_gr\_p', 'comp\_certif\_n', 'comp\_assoc\_n',  
 'comp\_bach\_n', 'comp\_grad\_certif\_n', 'comp\_ma\_n', 'comp\_phd\_n', 'om4\_total\_p',  
 'om4\_pell\_p', 'om4\_nonpell\_p', 'ftft\_om4\_total\_p', 'ftpt\_om4\_total\_p',  
 'nftft\_om4\_total\_p', 'nftpt\_om4\_total\_p', 'om6\_total\_p', 'om6\_pell\_p',  
 'om6\_nonpell\_p', 'ftft\_om6\_total\_p', 'ftpt\_om6\_total\_p', 'nftft\_om6\_total\_p',  
 'nftpt\_om6\_total\_p', 'om8\_total\_p', 'om8\_pell\_p', 'om8\_nonpell\_p',

'ftft\_om8\_total\_p', 'ftpt\_om8\_total\_p', 'nftft\_om8\_total\_p',  
 'nftpt\_om8\_total\_p', 'fa\_sfa\_np\_avg\_d', 'fa\_sfa\_np\_0\_30\_avg\_d',  
 'fa\_sfa\_np\_30\_48\_avg\_d', 'fa\_sfa\_np\_48\_75\_avg\_d', 'fa\_sfa\_np\_75\_110\_avg\_d',  
 'fa\_sfa\_np\_110plus\_avg\_d', 'fsa\_dl\_sub\_recip\_n', 'fsa\_dl\_sub\_orig\_n',  
 'fsa\_dl\_sub\_orig\_amt', 'fsa\_dl\_sub\_dis\_n', 'fsa\_dl\_sub\_dis\_amt',  
 'fsa\_dl\_sub\_ug\_recip\_n', 'fsa\_dl\_sub\_ug\_orig\_n', 'fsa\_dl\_sub\_ug\_orig\_amt',  
 'fsa\_dl\_sub\_ug\_dis\_n', 'fsa\_dl\_sub\_ug\_dis\_amt', 'fsa\_dl\_sub\_grad\_recip\_n',  
 'fsa\_dl\_sub\_grad\_orig\_n', 'fsa\_dl\_sub\_grad\_orig\_amt', 'fsa\_dl\_sub\_grad\_dis\_n',  
 'fsa\_dl\_sub\_grad\_dis\_amt', 'fsa\_dl\_unsub\_recip\_n', 'fsa\_dl\_unsub\_orig\_n',  
 'fsa\_dl\_unsub\_orig\_amt', 'fsa\_dl\_unsub\_dis\_n', 'fsa\_dl\_unsub\_dis\_amt',  
 'fsa\_dl\_unsub\_ug\_recip\_n', 'fsa\_dl\_unsub\_ug\_orig\_n', 'fsa\_dl\_unsub\_ug\_orig\_amt',  
 'fsa\_dl\_unsub\_ug\_dis\_n', 'fsa\_dl\_unsub\_ug\_dis\_amt', 'fsa\_dl\_unsub\_grad\_recip\_n',  
 'fsa\_dl\_unsub\_grad\_orig\_n', 'fsa\_dl\_unsub\_grad\_orig\_amt',  
 'fsa\_dl\_unsub\_grad\_dis\_n', 'fsa\_dl\_unsub\_grad\_dis\_amt',  
 'fsa\_dl\_parent\_plus\_recip\_n', 'fsa\_dl\_parent\_plus\_orig\_n',  
 'fsa\_dl\_parent\_plus\_orig\_amt', 'fsa\_dl\_parent\_plus\_dis\_n',  
 'fsa\_dl\_parent\_plus\_dis\_amt', 'fsa\_dl\_grad\_plus\_recip\_n',  
 'fsa\_dl\_grad\_plus\_orig\_n', 'fsa\_dl\_grad\_plus\_orig\_amt',  
 'fsa\_dl\_grad\_plus\_dis\_n', 'fsa\_dl\_grad\_plus\_dis\_amt', 'fsa\_gr\_pell\_recip\_n',  
 'fsa\_gr\_pell\_dis\_amt', 'fsa\_gr\_acg\_recip\_n', 'fsa\_gr\_acg\_dis\_amt',  
 'fsa\_gr\_smart\_recip\_n', 'fsa\_gr\_smart\_dis\_amt', 'fsa\_gr\_teach\_recip\_n',  
 'fsa\_gr\_teach\_dis\_amt', 'fsa\_gr\_iasg\_recip\_n', 'fsa\_gr\_iasg\_dis\_amt',  
 'fsa\_cba\_fws\_recip\_n', 'fsa\_cba\_fws\_fedaw\_amt', 'fsa\_cba\_fws\_dis\_amt',  
 'fsa\_cba\_pl\_recip\_n', 'fsa\_cba\_pl\_fedaw\_amt', 'fsa\_cba\_pl\_dis\_amt',  
 'fsa\_cba\_fseog\_recip\_n', 'fsa\_cba\_fseog\_fedaw\_amt', 'fsa\_cba\_fseog\_dis\_amt',  
 'fsa\_ffel\_sub\_recip\_n', 'fsa\_ffel\_sub\_orig\_n', 'fsa\_ffel\_sub\_orig\_amt',  
 'fsa\_ffel\_sub\_dis\_n', 'fsa\_ffel\_sub\_dis\_amt', 'fsa\_ffel\_unsub\_recip\_n',  
 'fsa\_ffel\_unsub\_orig\_n', 'fsa\_ffel\_unsub\_orig\_amt', 'fsa\_ffel\_unsub\_dis\_n',  
 'fsa\_ffel\_unsub\_dis\_amt', 'fsa\_ffel\_parent\_plus\_recip\_n',  
 'fsa\_ffel\_parent\_plus\_orig\_n', 'fsa\_ffel\_parent\_plus\_orig\_amt',  
 'fsa\_ffel\_parent\_plus\_dis\_n', 'fsa\_ffel\_parent\_plus\_dis\_amt',  
 'fsa\_ffel\_grad\_plus\_recip\_n', 'fsa\_ffel\_grad\_plus\_orig\_n',  
 'fsa\_ffel\_grad\_plus\_orig\_amt', 'fsa\_ffel\_grad\_plus\_dis\_n',  
 'fsa\_ffel\_grad\_plus\_dis\_amt', 'md\_earn\_wne\_p10', 'gt\_25k\_p10', 'gt\_28k\_p10',  
 'md\_earn\_wne\_p8', 'gt\_25k\_p8', 'gt\_28k\_p8', 'md\_earn\_wne\_p6', 'gt\_25k\_p6',  
 'gt\_28k\_p6', 'cdr2', 'cdr3', 'rpy\_3yr\_rt\_supp', 'lo\_inc\_rpy\_3yr\_rt\_supp',  
 'md\_inc\_rpy\_3yr\_rt\_supp', 'hi\_inc\_rpy\_3yr\_rt\_supp', 'compl\_rpy\_3yr\_rt\_supp',  
 'noncom\_rpy\_3yr\_rt\_supp', 'dep\_rpy\_3yr\_rt\_supp', 'ind\_rpy\_3yr\_rt\_supp',  
 'pell\_rpy\_3yr\_rt\_supp', 'nopell\_rpy\_3yr\_rt\_supp', 'female\_rpy\_3yr\_rt\_supp',  
 'male\_rpy\_3yr\_rt\_supp', 'firstgen\_rpy\_3yr\_rt\_supp',  
 'notfirstgen\_rpy\_3yr\_rt\_supp', 'rpy\_5yr\_rt', 'lo\_inc\_rpy\_5yr\_rt',  
 'md\_inc\_rpy\_5yr\_rt', 'hi\_inc\_rpy\_5yr\_rt', 'compl\_rpy\_5yr\_rt',  
 'noncom\_rpy\_5yr\_rt', 'dep\_rpy\_5yr\_rt', 'ind\_rpy\_5yr\_rt', 'pell\_rpy\_5yr\_rt',  
 'nopell\_rpy\_5yr\_rt', 'female\_rpy\_5yr\_rt', 'male\_rpy\_5yr\_rt',  
 'firstgen\_rpy\_5yr\_rt', 'notfirstgen\_rpy\_5yr\_rt', 'rpy\_7yr\_rt',  
 'lo\_inc\_rpy\_7yr\_rt', 'md\_inc\_rpy\_7yr\_rt', 'hi\_inc\_rpy\_7yr\_rt',  
 'compl\_rpy\_7yr\_rt', 'noncom\_rpy\_7yr\_rt', 'dep\_rpy\_7yr\_rt', 'ind\_rpy\_7yr\_rt',  
 'pell\_rpy\_7yr\_rt', 'nopell\_rpy\_7yr\_rt', 'female\_rpy\_7yr\_rt', 'male\_rpy\_7yr\_rt',

```
'firstgen_rpy_7yr_rt', 'notfirstgen_rpy_7yr_rt', 'cdr2_denom', 'cdr3_denom',
' rpy_3yr_n', 'lo_inc_rpy_3yr_n', 'md_inc_rpy_3yr_n', 'hi_inc_rpy_3yr_n',
' compl_rpy_3yr_n', 'noncom_rpy_3yr_n', 'dep_rpy_3yr_n', 'ind_rpy_3yr_n',
' pell_rpy_3yr_n', 'nopell_rpy_3yr_n', 'female_rpy_3yr_n', 'male_rpy_3yr_n',
' firstgen_rpy_3yr_n', 'notfirstgen_rpy_3yr_n', 'rpy_5yr_n', 'lo_inc_rpy_5yr_n',
' md_inc_rpy_5yr_n', 'hi_inc_rpy_5yr_n', 'compl_rpy_5yr_n', 'noncom_rpy_5yr_n',
' dep_rpy_5yr_n', 'ind_rpy_5yr_n', 'pell_rpy_5yr_n', 'nopell_rpy_5yr_n',
' female_rpy_5yr_n', 'male_rpy_5yr_n', 'firstgen_rpy_5yr_n',
' notfirstgen_rpy_5yr_n', 'rpy_7yr_n', 'lo_inc_rpy_7yr_n', 'md_inc_rpy_7yr_n',
' hi_inc_rpy_7yr_n', 'compl_rpy_7yr_n', 'noncom_rpy_7yr_n', 'dep_rpy_7yr_n',
' ind_rpy_7yr_n', 'pell_rpy_7yr_n', 'nopell_rpy_7yr_n', 'female_rpy_7yr_n',
' male_rpy_7yr_n', 'firstgen_rpy_7yr_n', 'notfirstgen_rpy_7yr_n']
```

```
[5]: print(df['data_yr_string'])
```

```
0      2000-01
1      2000-01
2      2000-01
3      2000-01
4      2000-01
```

```
...
139144    2018-19
139145    2018-19
139146    2018-19
139147    2018-19
139148    2018-19
```

Name: data\_yr\_string, Length: 139149, dtype: object

```
[6]: df[['data_yr_string1','data_yr_string2']] = df['data_yr_string'].str.
      ↪split("-",expand=True,)
print(df['data_yr_string2'])
# I use df['data_yr_string2'] cus it is a School year coded as for example ↪
↪2018-19,
# so I am using the last year
```

```
0      01
1      01
2      01
3      01
4      01
```

```
..
139144    19
139145    19
139146    19
139147    19
139148    19
```

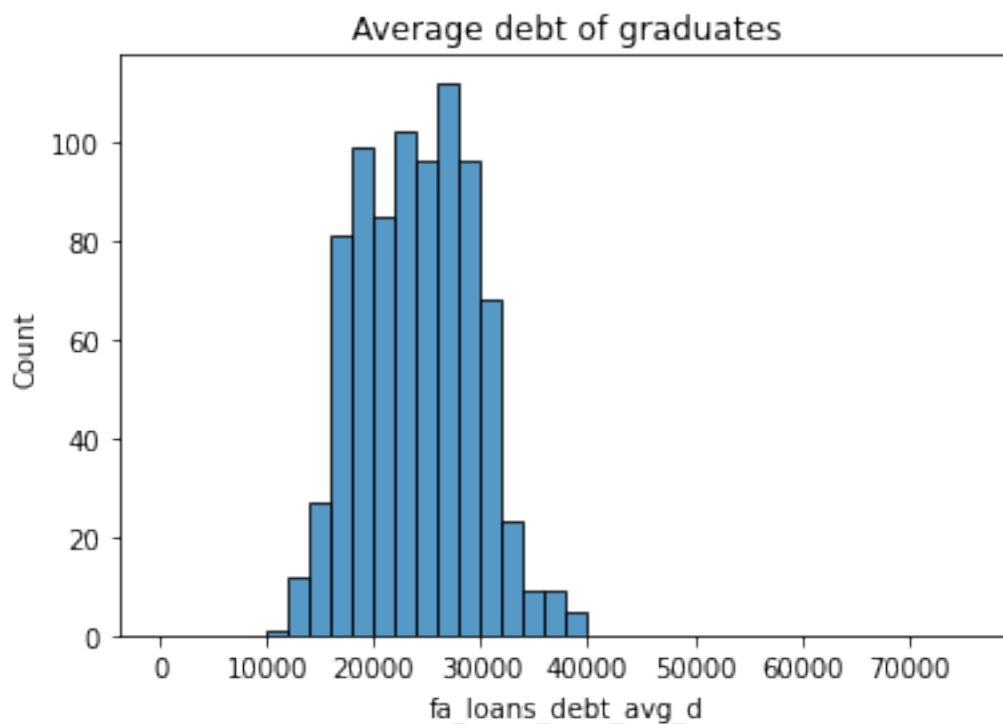
Name: data\_yr\_string2, Length: 139149, dtype: object

```
[7]: # entity type 6 is state total and entity type 10 is nation total
x = df[df['entity_type'] == 6]
y = df[df['entity_type'] == 10]
new_df = x + y
new_df.head()
new_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 867 entries, 145 to 139132
Columns: 410 entries, data_yr_string to data_yr_string2
dtypes: float64(396), int64(3), object(11)
memory usage: 2.7+ MB
```

```
[8]: sns.histplot(data=x, x='fa_loans_debt_avg_d', binrange=(0,75000),
    ↪binwidth=2000).set_title('Average debt of graduates')
```

```
[8]: Text(0.5, 1.0, 'Average debt of graduates')
```



```
[9]: x['fa_loans_debt_p'] = x['fa_loans_debt_p'] * 100
```

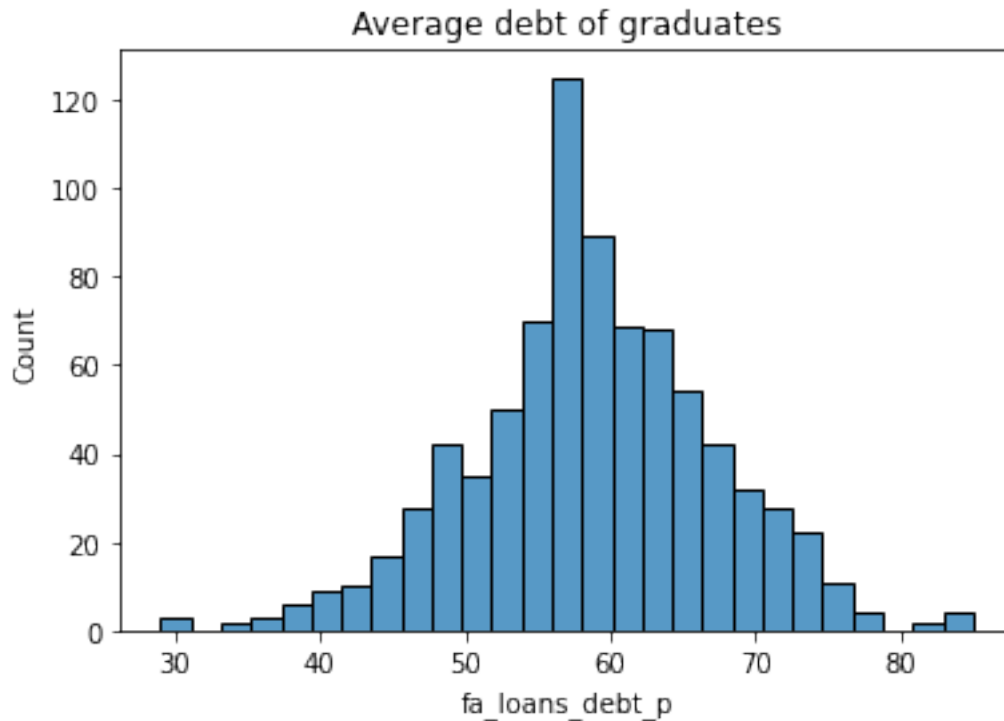
```
<ipython-input-9-891b78c0a8cd>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
x['fa_loans_debt_p'] = x['fa_loans_debt_p'] * 100
```

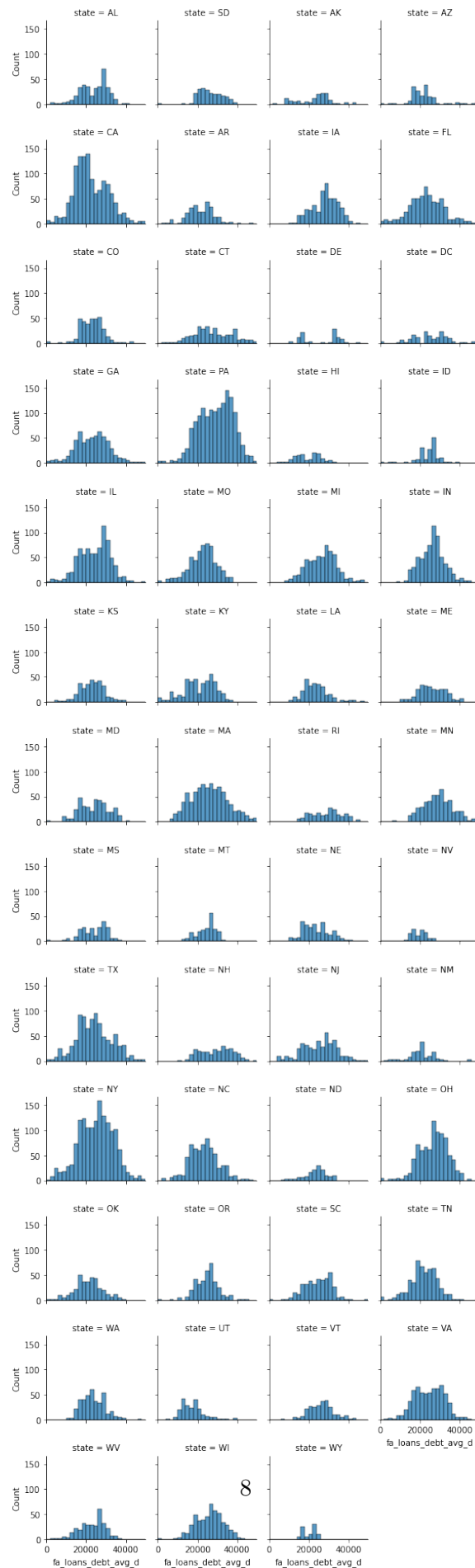
```
[10]: #Percent of graduates from 4-year public and private nonprofit colleges who are
      ↪carrying student debt
      sns.histplot(data=x, x='fa_loans_debt_p').set_title('Average debt of graduates')
```

```
[10]: Text(0.5, 1.0, 'Average debt of graduates')
```



```
[18]: ## Average debt of graduates by state
      g = sns.FacetGrid(df, col="state", col_wrap=4, height=2, xlim=(0,50000))
      g.map(sns.histplot, "fa_loans_debt_avg_d",binrange=(0,50000), binwidth=2000)
```

```
[18]: <seaborn.axisgrid.FacetGrid at 0x7f4096afb700>
```





```
[19]: ## Average debt of graduates by state  
g = sns.FacetGrid(df, col="data_yr_string2", col_wrap=4, height=2,  
→xlim=(0,50000))  
g.map(sns.histplot, "fa_loans_debt_avg_d",binrange=(0,50000), binwidth=2000)
```

```
[19]: <seaborn.axisgrid.FacetGrid at 0x7f408d3032e0>
```

