# Understanding theTrends in Indian Education

## Imports

```
In [1]:   1  import numpy as np
          2  import pandas as pd
          3  import seaborn as sns
          4  import matplotlib.pyplot as plt
          5  import random
```

```
In [2]:   1  random.seed(a=42)
```

## Import Data

```
In [3]:   1  xls_06_07 = pd.ExcelFile('Data/processed/SRC_Rawdata_2006-07_mod.xls')
          2  basic_06_07 = pd.read_excel(xls_06_07, 'Basic Data')
          3  facilities_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
          4  condition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
          5  enrolment_06_07 = pd.read_excel(xls_06_07, 'Enrolment')
          6  teacher_06_07 = pd.read_excel(xls_06_07, 'Teacher')
```

```
In [4]:   1  xls_07_08 = pd.ExcelFile('Data/processed/SRC_Rawdata_2007-08_mod.xls')
          2  basic_07_08 = pd.read_excel(xls_07_08, 'Basic Data')
          3  facilities_07_08 = pd.read_excel(xls_07_08, 'School Facilities')
          4  condition_07_08 = pd.read_excel(xls_07_08, 'School Condition')
          5  enrolment_07_08 = pd.read_excel(xls_07_08, 'Enrolment')
          6  teacher_07_08 = pd.read_excel(xls_07_08, 'Teacher')
```

```
In [5]:   1  xls_08_09 = pd.ExcelFile('Data/processed/SRC_Rawdata_2008-09_mod.xls')
          2  basic_08_09 = pd.read_excel(xls_08_09, 'Basic Data')
          3  facilities_08_09 = pd.read_excel(xls_08_09, 'School Facilities')
          4  condition_08_09 = pd.read_excel(xls_08_09, 'School Condition')
          5  enrolment_08_09 = pd.read_excel(xls_08_09, 'Enrolment')
          6  teacher_08_09 = pd.read_excel(xls_08_09, 'Teacher')
```

```
In [6]:   1  xls_09_10 = pd.ExcelFile('Data/processed/SRC_Rawdata_2009-10_mod.xls')
          2  basic_09_10 = pd.read_excel(xls_09_10, 'Basic Data')
          3  facilities_09_10 = pd.read_excel(xls_09_10, 'School Facilities')
          4  condition_09_10 = pd.read_excel(xls_09_10, 'School Condition')
          5  enrolment_09_10 = pd.read_excel(xls_09_10, 'Enrolment')
          6  teacher_09_10 = pd.read_excel(xls_09_10, 'Teacher')
```

```
In [7]:   1  xls_10_11 = pd.ExcelFile('Data/processed/SRC_Rawdata_2010-11_mod.xls')
          2  basic_10_11 = pd.read_excel(xls_10_11, 'Basic Data')
          3  facilities_10_11 = pd.read_excel(xls_10_11, 'School Facilities')
          4  condition_10_11 = pd.read_excel(xls_10_11, 'School Condition')
          5  enrolment_10_11 = pd.read_excel(xls_10_11, 'Enrolment')
          6  teacher_10_11 = pd.read_excel(xls_10_11, 'Teacher')
```

```
In [8]:   1  xls_11_12 = pd.ExcelFile('Data/processed/SRC_Rawdata_2011-12_mod.xls')
          2  basic_11_12 = pd.read_excel(xls_11_12, 'Basic Data')
          3  facilities_11_12 = pd.read_excel(xls_11_12, 'School Facilities')
          4  condition_11_12 = pd.read_excel(xls_11_12, 'School Condition')
          5  enrolment_11_12 = pd.read_excel(xls_11_12, 'Enrolment')
          6  teacher_11_12 = pd.read_excel(xls_11_12, 'Teacher')
```

```
In [9]:   1  xls_12_13 = pd.ExcelFile('Data/processed/SRC_Rawdata_2012-13_mod.xls')
          2  basic_12_13 = pd.read_excel(xls_12_13, 'Basic Data')
          3  facilities_12_13 = pd.read_excel(xls_12_13, 'School Facilities')
          4  condition_12_13 = pd.read_excel(xls_12_13, 'School Condition')
          5  enrolment_12_13 = pd.read_excel(xls_12_13, 'Enrolment')
          6  teacher_12_13 = pd.read_excel(xls_12_13, 'Teacher')
```
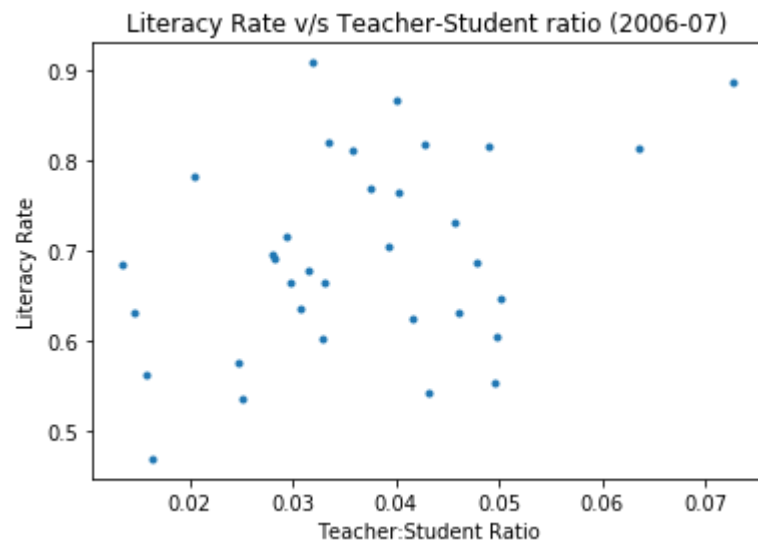
```
In [10]:   1  xls_13_14 = pd.ExcelFile('Data/processed/SRC_Rawdata_2013-14_mod.xls')
           2  basic_13_14 = pd.read_excel(xls_13_14, 'Basic Data')
           3  facilities_13_14 = pd.read_excel(xls_13_14, 'School Facilities')
           4  condition_13_14 = pd.read_excel(xls_13_14, 'School Condition')
           5  enrolment_13_14 = pd.read_excel(xls_13_14, 'Enrolment')
           6  teacher_13_14 = pd.read_excel(xls_13_14, 'Teacher')
```

## Literacy Rate v/s Teacher:Student ratio

*Literacy rate against Teacher-Student ratio was plotted for different years. As we can see, the Teacher-Student ratio is increasing through the years as we look from 2006 to 2013. This implies that there has been an increase in the number of teachers. Also, the literacy rates have also increased over these years implying that the increase in the number of teachers has a positive influence on the literacy rate.*
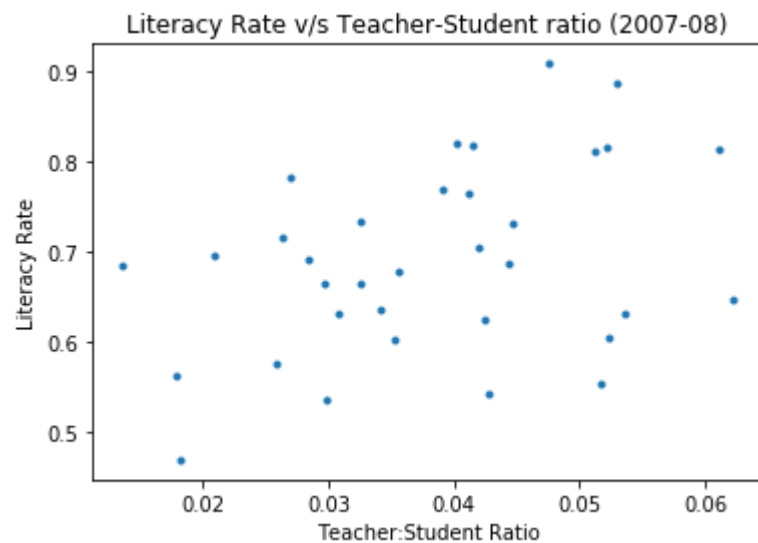
```
In [11]:   1  lit_06_07 = basic_06_07['literacy_rate']/100
           2  tch_stud_06_07 = (teacher_06_07['govt_tch_3'] + teacher_06_07['pvt_tch_3']) / (enrolment_06_07['govt_enr_3'
```

In [12]:
```python
1  plt.plot(tch_stud_06_07, lit_06_07, '.')
2  plt.xlabel('Teacher:Student Ratio')
3  plt.ylabel('Literacy Rate')
4  plt.title('Literacy Rate v/s Teacher-Student ratio (2006-07)')
5  plt.show()
```



In [13]:
```python
1  lit_07_08 = basic_07_08['literacy_rate']/100
2  tch_stud_07_08 = (teacher_07_08['govt_tch_3'] + teacher_07_08['pvt_tch_3']) / (enrolment_07_08['govt_enr_3'
```

```
In [14]:    1  plt.plot(tch_stud_07_08, lit_07_08, '.')
            2  plt.xlabel('Teacher:Student Ratio')
            3  plt.ylabel('Literacy Rate')
            4  plt.title('Literacy Rate v/s Teacher-Student ratio (2007-08)')
            5  plt.show()
```
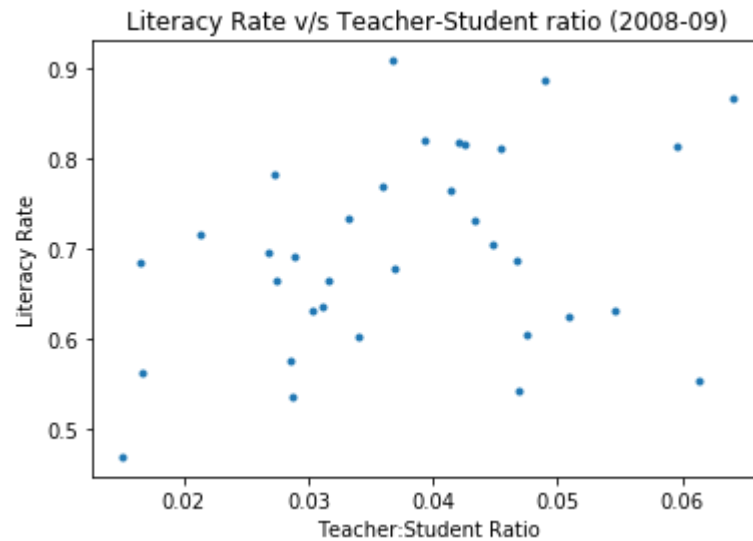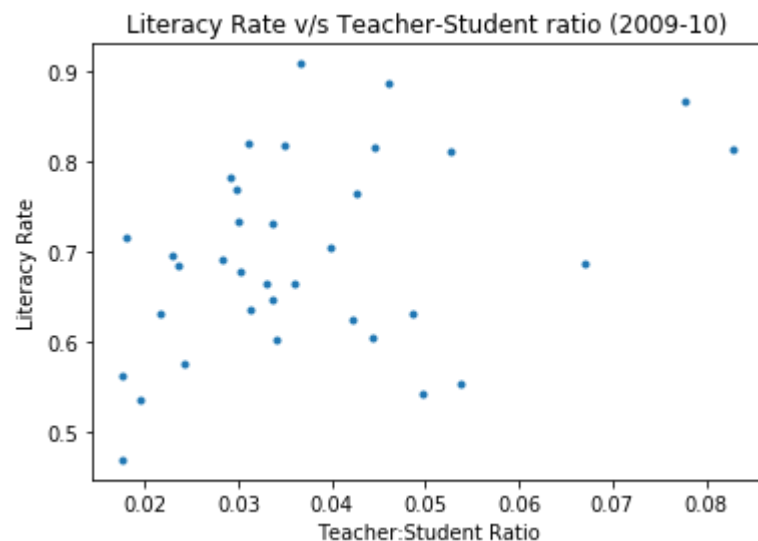


```
In [15]:    1  lit_08_09 = basic_08_09['literacy_rate']/100
            2  tch_stud_08_09 = (teacher_08_09['govt_tch_3'] + teacher_08_09['pvt_tch_3']) / (enrolment_08_09['govt_enr_3'
```

In [16]:
```python
1  plt.plot(tch_stud_08_09, lit_08_09, '.')
2  plt.xlabel('Teacher:Student Ratio')
3  plt.ylabel('Literacy Rate')
4  plt.title('Literacy Rate v/s Teacher-Student ratio (2008-09)')
5  plt.show()
```



In [17]:
```python
1  lit_09_10 = basic_09_10['literacy_rate']/100
2  tch_stud_09_10 = (teacher_09_10['govt_tch_3'] + teacher_09_10['pvt_tch_3']) / (enrolment_09_10['govt_enr_3'
```

In [18]:
```python
1  plt.plot(tch_stud_09_10, lit_09_10, '.')
2  plt.xlabel('Teacher:Student Ratio')
3  plt.ylabel('Literacy Rate')
4  plt.title('Literacy Rate v/s Teacher-Student ratio (2009-10)')
5  plt.show()
```
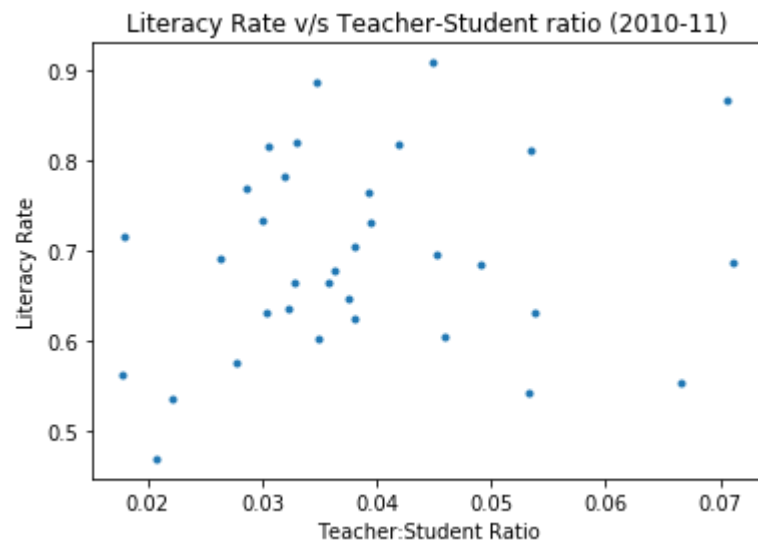


In [19]:
```python
1  lit_10_11 = basic_10_11['Literacy Rate']/100
2  tch_stud_10_11 = (teacher_10_11['Govt Tch 3'] + teacher_10_11['Pvt Tch 3']) / (enrolment_10_11['govt_enr_3'
```

In [20]:
```python
1  plt.plot(tch_stud_10_11, lit_10_11, '.')
2  plt.xlabel('Teacher:Student Ratio')
3  plt.ylabel('Literacy Rate')
4  plt.title('Literacy Rate v/s Teacher-Student ratio (2010-11)')
5  plt.show()
```
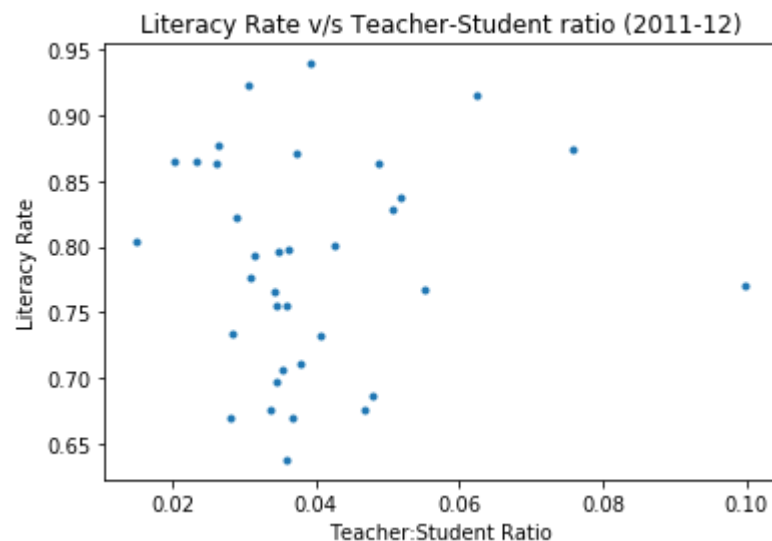


In [21]:
```python
1  lit_11_12 = basic_11_12['literacy_rate']/100
2  tch_stud_11_12 = (teacher_11_12['govt_tch_3'] + teacher_11_12['pvt_tch_3']) / (enrolment_11_12['govt_enr_3'
```

In [22]:
```python
1  tch_stud_11_12 = tch_stud_11_12[1:]
```

In [23]:
```python
plt.plot(tch_stud_11_12, lit_11_12, '.')
plt.xlabel('Teacher:Student Ratio')
plt.ylabel('Literacy Rate')
plt.title('Literacy Rate v/s Teacher-Student ratio (2011-12)')
plt.show()
```



In [24]:
```python
lit_12_13 = basic_12_13['literacy_rate']/100
tch_stud_12_13 = (teacher_12_13['govt_tch_3'] + teacher_12_13['pvt_tch_3']) / (enrolment_12_13['govt_enr_3'
```

In [25]:
```python
plt.plot(tch_stud_12_13, lit_12_13, '.')
plt.xlabel('Teacher:Student Ratio')
plt.ylabel('Literacy Rate')
plt.title('Literacy Rate v/s Teacher-Student ratio (2012-13)')
plt.show()
```



In [26]:
```python
lit_13_14 = basic_13_14['literacy_rate']/100
tch_stud_13_14 = (teacher_13_14['govt_tch_3'] + teacher_13_14['pvt_tch_3']) / (enrolment_13_14['govt_enr_3'
```
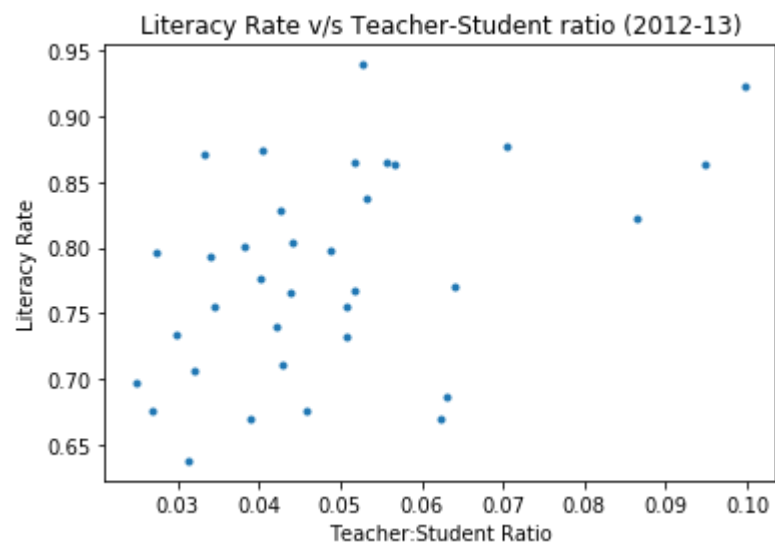
In [27]:
```python
1  plt.plot(tch_stud_13_14, lit_13_14, '.')
2  plt.xlabel('Teacher:Student Ratio')
3  plt.ylabel('Literacy Rate')
4  plt.title('Literacy Rate v/s Teacher-Student ratio (2013-14)')
5  plt.show()
```
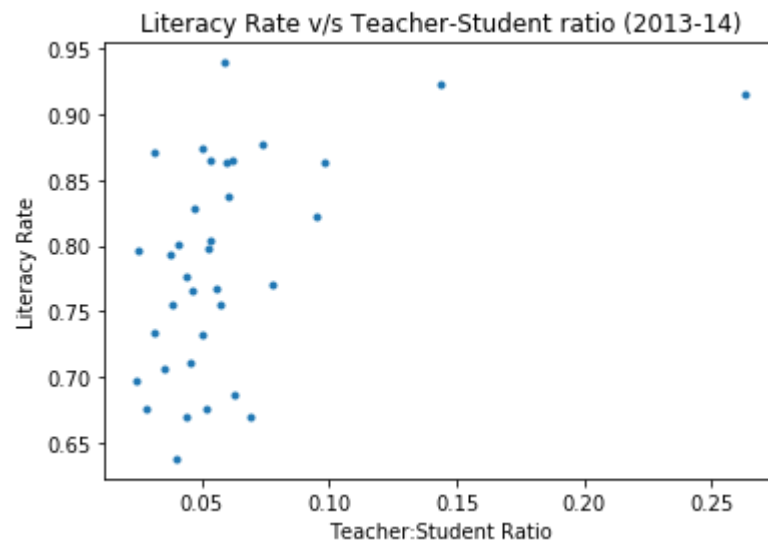


# Data Cleaning

```
In [28]:    1  df1 = pd.DataFrame()
            2  states = list(basic_13_14['statname'])
            3  literacy = list(basic_13_14['literacy_rate'])
```

## Basic

```
In [29]:    1  cols = ['area_sqkm', 'schools', 'tot_population', 'literacy_rate']
```

```
In [30]:    1  df1.insert(0, "area_sqkm", basic_13_14['area_sqkm'])
            2  df1.insert(1, "schools", basic_13_14['schools'])
            3  df1.insert(2, "tot_population", basic_13_14['tot_population'])
            4  df1.insert(3, "literacy_rate", basic_13_14['literacy_rate'])
```

## School Facilities

```
In [31]:    1  comp_cols = ['statcd', 'ac_year', 'statname', 'govt_sch_9', 'pvt_sch_9', 'govt_sch_r_9', 'pvt_sch_r_9']
```

```
In [32]:    1  for col in comp_cols:
            2      del facilities_13_14[col]
```

```
In [33]:    1  index = 4
            2  for col in list(facilities_13_14.columns.values):
            3      df1.insert(index, col, facilities_13_14[col])
```

## School Condition

```
In [34]:    1  comp_cols = ['statcd', 'ac_year', 'statname', 'sdg', 'tlm']
            2  for col in comp_cols:
            3      del condition_13_14[col]
            4  for col in list(condition_13_14.columns.values):
            5      df1.insert(index, col, condition_13_14[col])
```

## Enrolment

```
In [35]:   1  comp_cols = ['statcd', 'ac_year', 'statname', 'govt_enr_9', 'pvt_enr_9', 'govt_enr_r_9', 'pvt_enr_r_9', 'en
```

```
In [36]:   1  for col in comp_cols:
           2      del enrolment_13_14[col]
```

```
In [37]:   1  for col in list(enrolment_13_14.columns.values):
           2      df1.insert(index, col, enrolment_13_14[col])
```

## Teacher

```
In [38]:   1  comp_cols = ['statcd', 'ac_year', 'statname', 'govt_tch_9', 'pvt_tch_9', 'tch_bs', 'tch_s', 'tch_hs', 'tch_
```

```
In [39]:   1  for col in comp_cols:
           2      del teacher_13_14[col]
```

```
In [40]:   1  for col in list(teacher_13_14.columns.values):
           2      df1.insert(index, col, teacher_13_14[col])
```

# Clustering

### Removing nans & infs

```
In [41]:   1  df1 = df1.fillna(0)
           2  df1 = df1.replace(np.inf, 0)
           3  np.any(np.isnan(df1))
```

Out[41]:  False

```
In [42]:   1  np.all(np.isfinite(df1))
```

Out[42]:  True

## PCA + K-Means

In [43]:
```python
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
from mpl_toolkits.mplot3d import Axes3D
```

In [44]:
```python
pca = PCA(n_components=2)
```

In [45]:
```python
df1_trnsfrm = pca.fit_transform(df1)
```

In [46]:
```python
km_pca = KMeans(n_clusters=3)
```

In [47]:
```python
km_pca.fit(df1_trnsfrm)
```

Out[47]:
```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=None, tol=0.0001, verbose=0)
```

In [48]:
```python
km_pca.labels_
```

Out[48]:
```
array([1, 1, 1, 1, 1, 1, 1, 0, 2, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1])
```

In [49]:
```python
clus0 = [[], []]
clus1 = [[], []]
clus2 = [[], []]
for i, lab in enumerate(km_pca.labels_):
    if lab == 0:
        clus0[0].append(df1_trnsfrm[i, 0])
        clus0[1].append(df1_trnsfrm[i, 1])
    elif lab == 1:
        clus1[0].append(df1_trnsfrm[i, 0])
        clus1[1].append(df1_trnsfrm[i, 1])
    else:
        clus2[0].append(df1_trnsfrm[i, 0])
        clus2[1].append(df1_trnsfrm[i, 1])
```

In [50]:

```python
c1 = []
c2 = []
c3 = []
c4 = []
c5 = []
c6 = []
c7 = []
c8 = []
c9 = []
c10 = []
for index, label in enumerate(km_pca.labels_):
    if label == 0:
        c1.append((states[index], literacy[index]))
    elif label == 1:
        c2.append((states[index], literacy[index]))
    elif label == 2:
        c3.append((states[index], literacy[index]))
#     elif label == 3:
#         c4.append((states[index], literacy[index]))
#     elif label == 4:
#         c5.append((states[index], literacy[index]))
#     elif label == 5:
#         c6.append((states[index], literacy[index]))
#     elif label == 6:
#         c7.append((states[index], literacy[index]))
#     elif label == 7:
#         c8.append((states[index], literacy[index]))
#     elif label == 8:
#         c9.append((states[index], literacy[index]))
#     else:
#         c10.append((states[index], literacy[index]))
```

In [51]:     1  c1

Out[51]:  [('RAJASTHAN                                                                    ',
           67.06),
          ('BIHAR                                                                        ',
           63.82),
          ('WEST BENGAL                                                                  ',
           77.08),
          ('JHARKHAND                                                                    ',
           67.63),
          ('ODISHA                                                                       ',
           73.45),
          ('MADHYA PRADESH                                                               ',
           70.63),
          ('GUJARAT                                                                      ',
           79.31),
          ('MAHARASHTRA                                                                  ',
           82.91),
          ('ANDHRA PRADESH                                                               ',
           67.66),
          ('KARNATAKA                                                                    ',
           75.6),
          ('TAMIL NADU                                                                   ',
           80.33)]

In [52]:     1   c2

Out[52]:  [('JAMMU AND KASHMIR                                                                          ',
           68.74),
          ('HIMACHAL PRADESH                                                                           ',
           83.78),
          ('PUNJAB                                                                                     ',
           76.68),
          ('CHANDIGARH                                                                                 ',
           86.43),
          ('UTTARANCHAL                                                                                ',
           79.63),
          ('HARYANA                                                                                    ',
           76.64),
          ('DELHI                                                                                      ',
           86.34),
          ('SIKKIM                                                                                     ',
           82.2),
          ('ARUNACHAL PRADESH                                                                          ',
           66.95),
          ('NAGALAND                                                                                   ',
           80.11),
          ('MANIPUR                                                                                    ',
           79.85),
          ('MIZORAM                                                                                    ',
           91.58),
          ('TRIPURA                                                                                    ',
           87.75),
          ('MEGHALAYA                                                                                  ',
           75.48),
          ('ASSAM                                                                                      ',
           73.18),
          ('CHHATTISGARH                                                                               ',
           71.04),
          ('DAMAN & DIU                                                                                ',
           87.07),
          ('DADRA & NAGAR HAVELI                                                                       ',
           77.65),
          ('GOA                                                                                        ',
           87.4),
          ('LAKSHADWEEP                                                                                ',
           92.28),
          ('KERALA                                                                                     ',

```
      93.91),
     ('PONDICHERRY                                                       ',
      86.55),
     ('ANDAMAN & NICOBAR ISLANDS                                          ',
      86.27)]
```

In [53]:
```
1 c3
```

Out[53]:
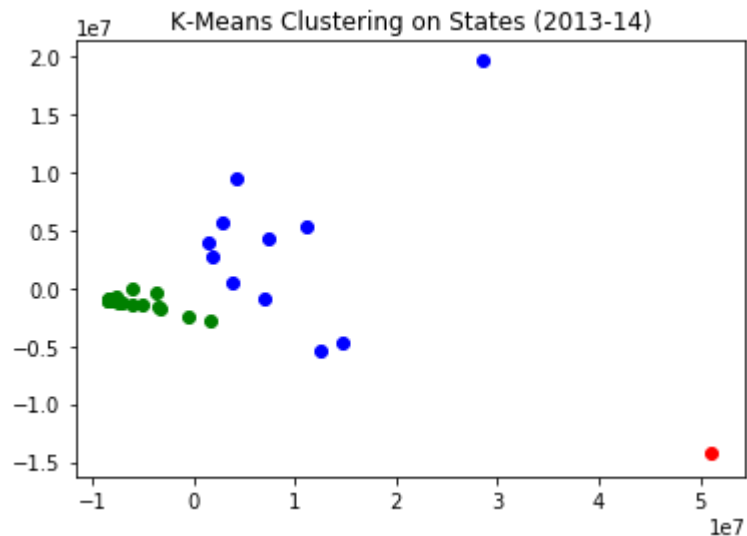```
[('UTTAR PRADESH                                                         ',
  69.72)]
```

In [54]:
```
1 plt.scatter(clus0[0], clus0[1], c=['blue'])
2 plt.scatter(clus1[0], clus1[1], c=['green'])
3 plt.scatter(clus2[0], clus2[1], c=['red'])
4 plt.title("K-Means Clustering on States (2013-14)")
5 plt.show()
```



## Observation of clustering

**K-Means clustering was applied to the states of India.**

**Clustering was performed for k=3 clusters.**

**Two clusters (represented by blue and green) have a good number of points in the third cluster. However, the third cluster (red) contains only a single state.**

**The state in the third cluster is Uttar Pradesh with a literacy rate of 69.72%.**

**Random points were sampled from the other two clusters and the details about those points were observed.**

# Visualization School

```
In [55]:    1  import numpy as np
            2  import pandas as pd
            3  import seaborn as sns
            4  import matplotlib.pyplot as plt
```

# Import Data

```
In [56]:    1  xls_06_07 = pd.ExcelFile('Data/processed/SRC_Rawdata_2006-07_mod.xls')
            2  facilities_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
            3  condition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
            4
            5  xls_07_08 = pd.ExcelFile('Data/processed/SRC_Rawdata_2007-08_mod.xls')
            6  facilities_07_08 = pd.read_excel(xls_07_08, 'School Facilities')
            7  condition_07_08 = pd.read_excel(xls_07_08, 'School Condition')
            8
            9  xls_08_09 = pd.ExcelFile('Data/processed/SRC_Rawdata_2008-09_mod.xls')
           10  facilities_08_09 = pd.read_excel(xls_08_09, 'School Facilities')
           11  condition_08_09 = pd.read_excel(xls_08_09, 'School Condition')
           12
           13  xls_09_10 = pd.ExcelFile('Data/processed/SRC_Rawdata_2009-10_mod.xls')
           14  facilities_09_10 = pd.read_excel(xls_09_10, 'School Facilities')
           15  condition_09_10 = pd.read_excel(xls_09_10, 'School Condition')
           16
           17  xls_10_11 = pd.ExcelFile('Data/processed/SRC_Rawdata_2010-11_mod.xls')
           18  facilities_10_11 = pd.read_excel(xls_10_11, 'School Facilities')
           19  condition_10_11 = pd.read_excel(xls_10_11, 'School Condition')
           20
           21  xls_11_12 = pd.ExcelFile('Data/processed/SRC_Rawdata_2011-12_mod.xls')
           22  facilities_11_12 = pd.read_excel(xls_11_12, 'School Facilities')
           23  condition_11_12 = pd.read_excel(xls_11_12, 'School Condition')
           24
           25  xls_12_13 = pd.ExcelFile('Data/processed/SRC_Rawdata_2012-13_mod.xls')
           26  facilities_12_13 = pd.read_excel(xls_12_13, 'School Facilities')
           27  condition_12_13 = pd.read_excel(xls_12_13, 'School Condition')
           28
           29  xls_13_14 = pd.ExcelFile('Data/processed/SRC_Rawdata_2013-14_mod.xls')
           30  facilities_13_14 = pd.read_excel(xls_13_14, 'School Facilities')
           31  condition_13_14 = pd.read_excel(xls_13_14, 'School Condition')
```

# Normalization

```
In [57]:    1  def normalize(df, f_range=(0,100)):
            2      tdf = f_range[0] + ((df - df.min())*(f_range[1]-f_range[0]))/(df.max()-df.min())
            3      return tdf
```

# Visualization

```
In [58]:  1  states = facilities_13_14['statname'][1:]
          2  states = [s[:3] for s in states]
          3  states_dict = dict()
          4  for i, s in enumerate(states):
          5      states_dict[i] = s
```

```
In [59]:  1  years = ['2006-07', '2007-08', '2008-09', '2009-10', '2010-11', '2011-12', '2012-13', '2013-14']
```
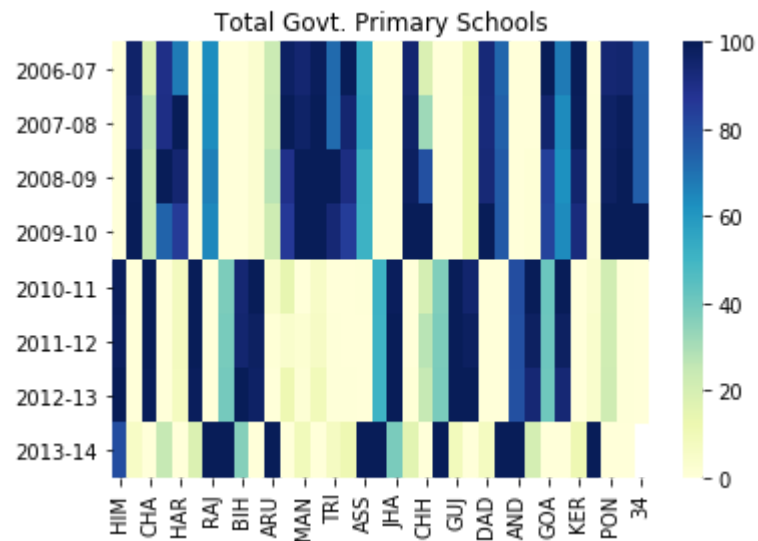
## Govt. Primary Schools

On observing the heatmaps below, we can observe that there has been a change in the number of schools in most states. States such as Uttarakhand, West Bengal, who have seen an increase in the number of schools have also recorded an increase in literacy rates. However, states such as Mizoram have seen a decline in the number of schools. This could be a result of shutting down of schools due to environmental crisis or other external factors. Despite the fall in the school count, the literacy rate has been stable on the higher side of the scale.

```
In [60]:   1  total_govt_schools_primary = pd.concat([facilities_06_07['govt_sch_1'], facilities_07_08['govt_sch_1'], fac
           2
           3  col_6 = pd.DataFrame(facilities_12_13['govt_sch_1'])
           4  col_6 = col_6.drop(col_6.index[0]).reset_index()
           5  del col_6['index']
           6
           7  col_7 = pd.DataFrame(facilities_13_14['govt_sch_1'])
           8  col_7 = col_7.drop(col_7.index[0]).reset_index()
           9  del col_7['index']
          10
          11  total_govt_schools_primary = pd.concat([total_govt_schools_primary, col_6, col_7], ignore_index=True, axis=
```

```
In [61]:   1  total_govt_schools_primary = total_govt_schools_primary.T
           2  total_govt_schools_primary.index = years
           3  total_govt_schools_primary = total_govt_schools_primary.rename(columns=states_dict)
```

```
In [62]:   1  trans_data = normalize(total_govt_schools_primary)
```

In [63]:
```python
ax = sns.heatmap(trans_data, cmap="YlGnBu")
ax.set_title('Total Govt. Primary Schools')
plt.show()
```
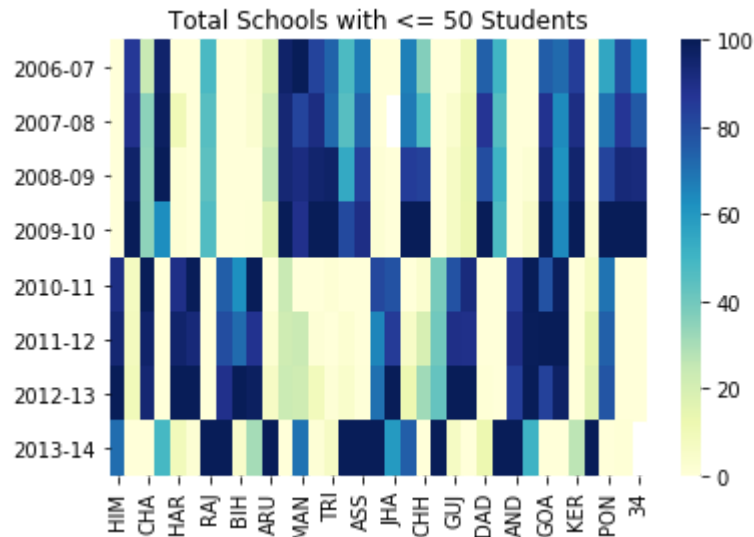
Total Govt. Primary Schools

In [64]:
```python
total_priv_schools_primary = pd.concat([facilities_06_07['pvt_sch_1'], facilities_07_08['pvt_sch_1'], facil

col_6 = pd.DataFrame(facilities_12_13['pvt_sch_1'])
col_6 = col_6.drop(col_6.index[0]).reset_index()
del col_6['index']

col_7 = pd.DataFrame(facilities_13_14['pvt_sch_1'])
col_7 = col_7.drop(col_7.index[0]).reset_index()
del col_7['index']

total_priv_schools_primary = pd.concat([total_priv_schools_primary, col_6, col_7], ignore_index=True, axis=
```

In [65]:
```python
total_priv_schools_primary = total_priv_schools_primary.T
total_priv_schools_primary.index = years
total_priv_schools_primary = total_priv_schools_primary.rename(columns=states_dict)
```

In [66]:
```python
trans_data = normalize(total_priv_schools_primary)
```

In [67]:
```python
1  ax = sns.heatmap(trans_data, cmap="YlGnBu")
2  ax.set_title('Total Private Primary Schools')
3  plt.show()
```



## Schools with less than 50 Students

In [68]:
```python
1  total_lt_50 = pd.concat([facilities_06_07['less_50_e_all'], facilities_07_08['less_50_e_all'], facilities_0
2
3  col_6 = pd.DataFrame(facilities_12_13['less_50_e_all'])
4  col_6 = col_6.drop(col_6.index[0]).reset_index()
5  del col_6['index']
6
7  col_7 = pd.DataFrame(facilities_13_14['less_50_e_all'])
8  col_7 = col_7.drop(col_7.index[0]).reset_index()
9  del col_7['index']
10
11  total_lt_50 = pd.concat([total_lt_50, col_6, col_7], ignore_index=True, axis=1)
```

In [69]:
```python
1  total_lt_50 = total_lt_50.T
2  total_lt_50.index = years
3  total_lt_50 = total_lt_50.rename(columns=states_dict)
```

```
In [70]:   1  trans_data = normalize(total_lt_50)
```

```
In [71]:   1  ax = sns.heatmap(trans_data, cmap="YlGnBu")
           2  ax.set_title('Total Schools with <= 50 Students')
           3  plt.show()
```



## Total Classrooms which need Repairs

```
In [72]:   1  total_cls_repairs = pd.concat([condition_06_07['cls_major5'], condition_07_08['cls_major5'], condition_08_0
           2
           3  col_6 = pd.DataFrame(condition_12_13['cls_major5'])
           4  col_6 = col_6.drop(col_6.index[0]).reset_index()
           5  del col_6['index']
           6
           7  col_7 = pd.DataFrame(condition_13_14['cls_major5'])
           8  col_7 = col_7.drop(col_7.index[0]).reset_index()
           9  del col_7['index']
          10
          11  total_cls_repairs = pd.concat([total_cls_repairs, col_6, col_7], ignore_index=True, axis=1)
```

```
In [73]:    1  total_cls_repairs = total_cls_repairs.T
            2  total_cls_repairs.index = years
            3  total_cls_repairs = total_cls_repairs.rename(columns=states_dict)
```

```
In [74]:    1  trans_data = normalize(total_cls_repairs)
```

```
In [75]:    1  ax = sns.heatmap(trans_data, cmap="YlGnBu")
            2  ax.set_title('Total Classrooms which need Repairs')
            3  plt.show()
```



## No. of Primary Schools Receiving Grants

In [76]:
```python
total_prim_sch_grants = pd.concat([condition_06_07['sdg_1'], condition_07_08['sdg_1'], condition_08_09['sdg

col_6 = pd.DataFrame(condition_12_13['sdg_1'])
col_6 = col_6.drop(col_6.index[0]).reset_index()
del col_6['index']

col_7 = pd.DataFrame(condition_13_14['sdg_1'])
col_7 = col_7.drop(col_7.index[0]).reset_index()
del col_7['index']

total_prim_sch_grants = pd.concat([total_prim_sch_grants, col_6, col_7], ignore_index=True, axis=1)
```

In [77]:
```python
total_prim_sch_grants = total_prim_sch_grants.T
total_prim_sch_grants.index = years
total_prim_sch_grants = total_prim_sch_grants.rename(columns=states_dict)
```

In [78]:
```python
trans_data = normalize(total_prim_sch_grants)
```

In [79]:
```python
ax = sns.heatmap(trans_data, cmap="YlGnBu")
ax.set_title('Total Primary School Grants Received')
plt.show()
```

In [ ]: `1`

## Imports

In [80]:
```python
#imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px

from sklearn.cluster import KMeans
import pandas as pd
import numpy as np
from sklearn import preprocessing
from sklearn.decomposition import PCA
from numpy import linalg as LA
import numpy.linalg as linalg
import copy

from collections import defaultdict
from sklearn.manifold import TSNE
```

In [81]:
```python
#Hashing corresponding to features.
metadata = pd.read_csv("data processed/2015_16_Statewise_Elementary_Metadata.csv")
metadata_hash = {}
for i in range(len(metadata)):
    metadata_hash[metadata.iat[i,0]] = metadata.iat[i,1]
```

**for read file**

In [82]:
```python
def make_dataframe(xls_06_07):
    basicdata_06_07 = pd.read_excel(xls_06_07, 'Basic Data')
    schoolfacility_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
    schoolfacility_06_07.drop(['statcd','ac_year','statname'], axis = 1, inplace=True)
    schoolcondition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
    schoolcondition_06_07.drop(['statcd','ac_year','statname'], axis = 1, inplace=True)
    enrolment_06_07 = pd.read_excel(xls_06_07, 'Enrolment')
    enrolment_06_07.drop(['statcd','ac_year','statname'], axis = 1, inplace=True)
    teacher_06_07 = pd.read_excel(xls_06_07, 'Teacher')
    teacher_06_07.drop(['statcd','ac_year','statname'], axis = 1, inplace=True)
    df = pd.concat([basicdata_06_07, schoolfacility_06_07, schoolcondition_06_07, enrolment_06_07, teacher_
    return df
```

In [83]:
```python
# basicdata_06_07_features = set(list(basicdata_06_07.columns))
# schoolfacility_06_07_features = set(list(schoolfacility_06_07.columns))
# schoolcondition_06_07_featuers = set(list(schoolcondition_06_07.columns))
# enrolment_06_07_features = set(list(enrolment_06_07.columns))
# teacher_06_07_features = set(list(teacher_06_07.columns))
xls_06_07 = pd.ExcelFile('data processed/SRC_Rawdata_2006-07_mod.xls')
basicdata_06_07 = pd.read_excel(xls_06_07, 'Basic Data')
colmn = ['statcd', 'ac_year', 'statname']
schoolfacility_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
schoolfacility_06_07.drop(colmn, inplace=True, axis=1)
schoolcondition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
schoolcondition_06_07.drop(colmn, inplace=True, axis=1)
enrolment_06_07 = pd.read_excel(xls_06_07, 'Enrolment')
enrolment_06_07.drop(colmn, inplace=True, axis=1)
teacher_06_07 = pd.read_excel(xls_06_07, 'Teacher')
teacher_06_07.drop(colmn, inplace=True, axis=1)

#Concate all excel into one
df06_07 = pd.concat([basicdata_06_07, schoolfacility_06_07,schoolcondition_06_07,enrolment_06_07,teacher_06

statename = list(df06_07['statname'])
```

## Read and combine all year files

In [84]:
```python
1  #read 2006-07
2  xls_06_07 = pd.ExcelFile('data processed/SRC_Rawdata_2006-07_mod.xls')
3  df_06_07 = make_dataframe(xls_06_07)
4  df_06_07.drop(['ner_p','ner_up','retentionrate'], axis = 1, inplace=True)
5  df_06_07 = df_06_07.sort_index(axis=1)
```

In [85]:
```python
1   f0 = set(list(df_06_07.sort_index(axis=1).columns))
2
3   #read 2007-08
4   xls_07_08 = pd.ExcelFile('data processed/SRC_Rawdata_2007-08_mod.xls')
5   df_07_08 = make_dataframe(xls_07_08)
6   df_07_08.drop(['attendance_b_p','attendance_g_p','cls_minor','kitshed_py_1','kitshed_py_2','kitshed_py_3','
7   df_07_08 = df_07_08.sort_index(axis=1)
8
9   f2 = set(list(df_07_08.sort_index(axis=1).columns))
10
11  # concate 2006_08
12  frames = [df_06_07, df_07_08]
13  df_06_08 = pd.concat(frames,ignore_index=True)
14  df_06_08.drop(['ptr_py_1','ptr_py_2','ptr_py_3','ptr_py_4','ptr_py_5','ptr_py_all','scr_py_1','scr_py_2','s
15  print(df_06_08.shape)
```

(70, 630)

In [86]:
```python
1  f0 = set(list(df_06_08.sort_index(axis=1).columns))
2
3  #read 2008-09
4  xls_08_09 = pd.ExcelFile('data processed/SRC_Rawdata_2008-09_mod.xls')
5  df_08_09 = make_dataframe(xls_08_09)
6  df_08_09.drop(['attendance_b_p','attendance_g_p','cls_minor','kitshed_py_1','kitshed_py_2','kitshed_py_3','
7  df_08_09 = df_08_09.sort_index(axis=1)
8
9  f3 = set(list(df_08_09.sort_index(axis=1).columns))
10
11 # concate 2006_09
12 frames = [df_06_08, df_08_09]
13 df_06_09 = pd.concat(frames,ignore_index=True)
14 df_06_09.drop(['avg_tch_py_1','avg_tch_py_2','avg_tch_py_3','avg_tch_py_4','avg_tch_py_5','avg_tch_py_all',
15 print(df_06_09.shape)
```

(105, 612)

In [ ]:
```python
1
```

In [87]:
```python
1  f0 = set(list(df_06_09.sort_index(axis=1).columns))
2
3  #read 2009-10
4  xls_09_10 = pd.ExcelFile('data processed/SRC_Rawdata_2009-10_mod.xls')
5  df_09_10 = make_dataframe(xls_09_10)
6  df_09_10.drop(['attendance_b_p','attendance_g_p','cls_minor','kitshed_py_1','kitshed_py_2','kitshed_py_3','
7  df_09_10 = df_09_10.sort_index(axis=1)
8
9  f4 = set(list(df_09_10.sort_index(axis=1).columns))
10
11 # concate 2006_10
12 frames = [df_06_09, df_09_10]
13 df_06_10 = pd.concat(frames,ignore_index=True)
14 print(df_06_10.shape)
```

(140, 612)

In [88]:

```python
f0 = set(list(df_06_10.sort_index(axis=1).columns))

#read 2010-11
xls_10_11 = pd.ExcelFile('data processed/SRC_Rawdata_2010-11_mod.xls')
df_10_11 = make_dataframe(xls_10_11)
df_10_11.columns = map(str.lower, df_10_11.columns)
df_10_11.columns = df_10_11.columns.str.replace(' ', '_')
df_10_11 = df_10_11.sort_index(axis=1)

f5 = set(list(df_10_11.sort_index(axis=1).columns))
d5_1 = list(f5.difference(f0))

df_10_11.drop(d5_1, axis = 1, inplace=True)
f5 = set(list(df_10_11.sort_index(axis=1).columns))
d5_2 = list(f0.difference(f5))

df_06_10.drop(d5_2, axis = 1, inplace=True)
f0 = set(list(df_06_10.sort_index(axis=1).columns))

df_10_11 = df_10_11.loc[:, ~df_10_11.columns.duplicated()]

# concate 2006_11
frames = [df_06_10, df_10_11]
df_06_11 = pd.concat(frames,ignore_index=True)
print(df_06_11.shape)
```

(175, 604)

```
In [89]:   1  f0 = set(list(df_06_11.columns))
           2
           3  #read 2011-12
           4  xls_11_12 = pd.ExcelFile('data processed/SRC_Rawdata_2011-12_mod.xls')
           5  df_11_12 = make_dataframe(xls_11_12)
           6  df_11_12.columns = map(str.lower, df_11_12.columns)
           7  df_11_12.columns = df_11_12.columns.str.replace(' ', '_')
           8  df_11_12 = df_11_12.sort_index(axis=1)
           9
          10  f6 = set(list(df_11_12.columns))
          11  d6_1 = list(f6.difference(f0))
          12  df_11_12.drop(d6_1, axis = 1, inplace=True)
          13
          14  f6 = set(list(df_11_12.columns))
          15  d6_2 = list(f0.difference(f6))
          16  df_06_11.drop(d6_2, axis = 1, inplace=True)
          17  f0 = set(list(df_06_11.columns))
          18
          19  # concate 2006_12
          20  frames = [df_06_11, df_11_12]
          21  df_06_12 = pd.concat(frames,ignore_index=True)
          22  print(df_06_12.shape)
          23
```

(211, 512)

```
In [90]:   1  f0 = set(list(df_06_12.columns))
           2
           3  #read 2012-13
           4  xls_12_13 = pd.ExcelFile('data processed/SRC_Rawdata_2012-13_mod.xls')
           5  df_12_13 = make_dataframe(xls_12_13)
           6  df_12_13.columns = map(str.lower, df_12_13.columns)
           7  df_12_13.columns = df_12_13.columns.str.replace(' ', '_')
           8  df_12_13 = df_12_13.sort_index(axis=1)
           9
          10  f7 = set(list(df_12_13.columns))
          11  d7_1 = list(f0.difference(f7))
          12  df_06_12.drop(d7_1, axis = 1, inplace=True)
          13  f0 = set(list(df_06_12.columns))
          14
          15  d7_2 = list(f7.difference(f0))
          16  df_12_13.drop(d7_2, axis = 1, inplace=True)
          17  f7 = set(list(df_12_13.columns))
          18
          19  # concate 2006_13
          20  frames = [df_06_12, df_12_13]
          21  df_06_13 = pd.concat(frames,ignore_index=True)
          22  print(df_06_13.shape)
```

(247, 464)

In [91]:
```python
1  f0 = set(list(df_06_13.columns))
2
3  #read 2013-14
4  xls_13_14 = pd.ExcelFile('data processed/SRC_Rawdata_2013-14_mod.xls')
5  df_13_14 = make_dataframe(xls_13_14)
6  df_13_14.columns = map(str.lower, df_13_14.columns)
7  df_13_14.columns = df_13_14.columns.str.replace(' ', '_')
8  df_13_14 = df_13_14.sort_index(axis=1)
9
10 f8 = set(list(df_13_14.columns))
11 d8_1 = list(f0.difference(f8))
12 df_06_13.drop(d8_1, axis = 1, inplace=True)
13 f0 = set(list(df_06_13.columns))
14
15 d8_2 = list(f8.difference(f0))
16 df_13_14.drop(d8_2, axis = 1, inplace=True)
17 f8 = set(list(df_13_14.columns))
18
19 # concate 2006_14
20 frames = [df_06_13, df_13_14]
21 df_06_14 = pd.concat(frames,ignore_index=True)
22 print(df_06_14.shape)
23 f0 = set(list(df_06_14.columns))
```

(283, 435)

In [92]:
```python
1  #file combined from 2006-2014 and stored in allyear.csv
2  print(df_06_14.shape)
3  df_06_14.to_csv('allyear.csv')
```

(283, 435)

In [93]:
```python
1  all_feature = list(df06_07.columns)
```

In [94]:

```python
#remove unnecessary data field
df06_07 = df06_07.drop('statcd', axis=1)
df06_07 = df06_07.drop('ac_year', axis=1)
df06_07 = df06_07.drop('statname', axis=1)
df06_07 = df06_07.drop('area_sqkm', axis=1)

all_features_df06_07 = list(df06_07.columns)
df06_07  = df06_07.fillna(0)

#Normalization
x = df06_07.values
min_max_scaler = preprocessing.Normalizer()
x_scaled = min_max_scaler.fit_transform(x)

x_scaled= x_scaled.tolist()
dataset=copy.deepcopy(x_scaled)
```

## Finding Value of K for performing K-Means Clustering

```
In [95]:    1  #Here we find out best k value to perform kmeans clustering.
            2  dataset=copy.deepcopy(x_scaled)
            3  cost =[]
            4  for i in range(1, 20):
            5      KM = KMeans(n_clusters = i, max_iter = 500)
            6      KM.fit(dataset)
            7      cost.append(KM.inertia_)
            8
            9  plt.plot(range(1, 20), cost, color ='g', linewidth ='3')
           10  plt.xlabel("Value of K")
           11  plt.ylabel("Sqaured Error (Cost)")
           12  plt.show()
```



## Data visualization

```
In [96]:    1  from sklearn.decomposition import TruncatedSVD
```

In [97]:

```python
final_total_feature = list(df_06_14.columns)
feature_drop4clustering = ['statcd','area_sqkm','blocks','clusters','villages']
df_06_14.drop(feature_drop4clustering, axis = 1, inplace=True)

update_statename= []
for i in list(df_06_14['statname']):
    update_statename.append(str(i).split('    ')[0])

new_df = pd.DataFrame({'statname': update_statename})
df_06_14.update(new_df)

# df_06_14['year_state'] = df_06_14[['ac_year', 'statname']].apply(lambda x: ''.join(str(x)), axis=1)

labels = []
for i,j in zip((list(df_06_14['ac_year'])),(list(df_06_14['statname']))):
    labels.append(str(str(i)+" "+str(j)))

# new_df = pd.DataFrame({'year_state': year_state})
# df_06_14.update(new_df)

f12= ['ac_year','statname']
df_06_14.drop(f12, axis = 1, inplace=True)

all_features = list(df_06_14.columns)
df_06_14  = df_06_14.fillna(0)

for i in list(df_06_14.columns):
    try:
        df_06_14[i].replace('-', 0, inplace=True)
        df_06_14[i].replace(' ', 0, inplace=True)
    except:
        pass
```

In [98]:
```python
1  # Normalization of dataset
2  x = df_06_14.values.T
3  min_max_scaler = preprocessing.Normalizer()
4  x_scaled = min_max_scaler.fit_transform(x)
5
6  # df = pd.DataFrame(x_scaled)
7
8  #TSNE PLOTS
9  x_scaled= x_scaled.tolist()
10 dataset=copy.deepcopy(x_scaled)
11
12 X_embedded = TSNE(n_components=2).fit_transform(dataset)
13
14 x= [i[0] for i in X_embedded]
15 y= [i[1] for i in X_embedded]
16 df123 = pd.DataFrame(list(zip(x, y)), columns =['X', 'Y'])
17 fig = px.scatter(df123,x='X', y='Y',title="TSNE PLOT OF COMBINED DATA-SET")
18 fig.show()
19
```

## TSNE PLOT OF COMBINED DATA-SET

**K_Means**

In [99]:
```python
#APPLY KMEANS AND PLOT CLUSTERS
kmeans = KMeans(n_clusters=10, random_state=0).fit(dataset)
# kmeans.labels_

knn_labels = list(kmeans.labels_)

cluster_name = defaultdict(list)
for i,j in zip(knn_labels,labels):
    try:
        cluster_name[i].append(j)
    except:
        cluster_name[i] = j

df = pd.DataFrame(list(zip(x, y,knn_labels,labels)), columns =['X', 'Y','Label','distname'])
fig = px.scatter(df, x='X', y='Y', color='Label',title = "K-MEANS CLUSTERING")
fig.show()
```

## K-MEANS CLUSTERING

In [100]:   1  cluster_name

Out[100]: defaultdict(list,
            {6: ['2006-07 A & N Islands',
              '2006-07 Andhra Pradesh',
              '2006-07 Arunachal Pradesh',
              '2006-07 Assam',
              '2006-07 Chhattisgarh',
              '2006-07 D & N Haveli',
              '2006-07 Daman & Diu',
              '2006-07 Haryana',
              '2006-07 Himachal Pradesh',
              '2006-07 Kerala',
              '2006-07 Lakshadweep',
              '2006-07 Madhya Pradesh',
              '2006-07 Maharashtra',
              '2006-07 Nagaland',
              '2006-07 Orissa',
              '2006-07 Puducherry',
              '2006-07 Punjab',
              '2006-07 Rajasthan',
              '2006 07 Sikkim'

**Clustering performed on the whole dataset : After combining the data for all the years and then perform clustering over the dataset, we find the value of k by assigning the value of k and find squared error between the centroid and after a threshold select the value of k which 10.**

In [101]:
```python
1  # Normalization of dataset
2  x = df_06_14.values.T
3  min_max_scaler = preprocessing.Normalizer()
4  x_scaled = min_max_scaler.fit_transform(x)
5  # df = pd.DataFrame(x_scaled)
```

## SVD performed

In [102]:

```python
svd = TruncatedSVD(n_components=282, n_iter=7, random_state=42)
x_scaled = svd.fit_transform(x_scaled)
dataset=copy.deepcopy(x_scaled)
```

In [103]:
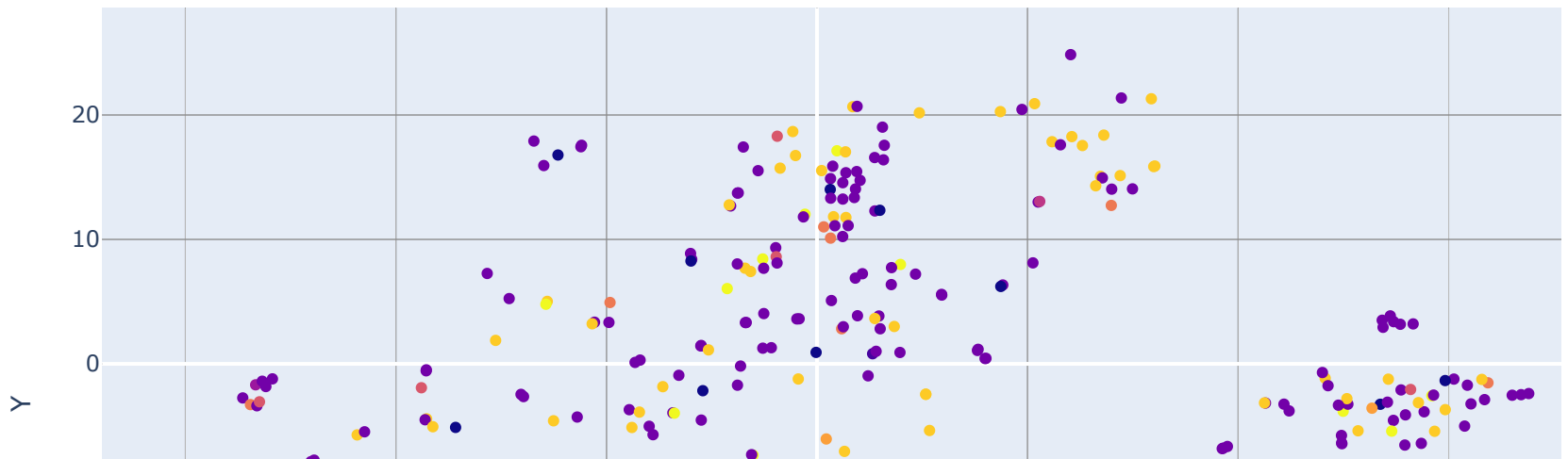
```python
x = df_06_14.values.T
min_max_scaler = preprocessing.Normalizer()
x_scaled = min_max_scaler.fit_transform(x)
x_scaled= x_scaled.tolist()
dataset=copy.deepcopy(x_scaled)

kmeans = KMeans(n_clusters=10, random_state=0).fit(np.array(dataset).T)
# kmeans.labels_
knn_labels = list(kmeans.labels_)
cluster_name = defaultdict(list)
for i,j in zip(knn_labels,labels):
    try:
        cluster_name[i].append(j)
    except:
        cluster_name[i] = j
x= [i[0] for i in X_embedded]
y= [i[1] for i in X_embedded]

df = pd.DataFrame(list(zip(x, y,knn_labels,labels)), columns =['X', 'Y','Label','distname'])
fig = px.scatter(df, x='X', y='Y', color='Label',title = "Perform Kmeans of SVD output data")
fig.show()

```

## Perform Kmeans of SVD output data

*conclusion - After SVD, K-Means clustering is not performed well*

## Visualization

**Read file of 2006-07 and 2016-17(For Comparision between Goverment vs private school ratio through visualization)**
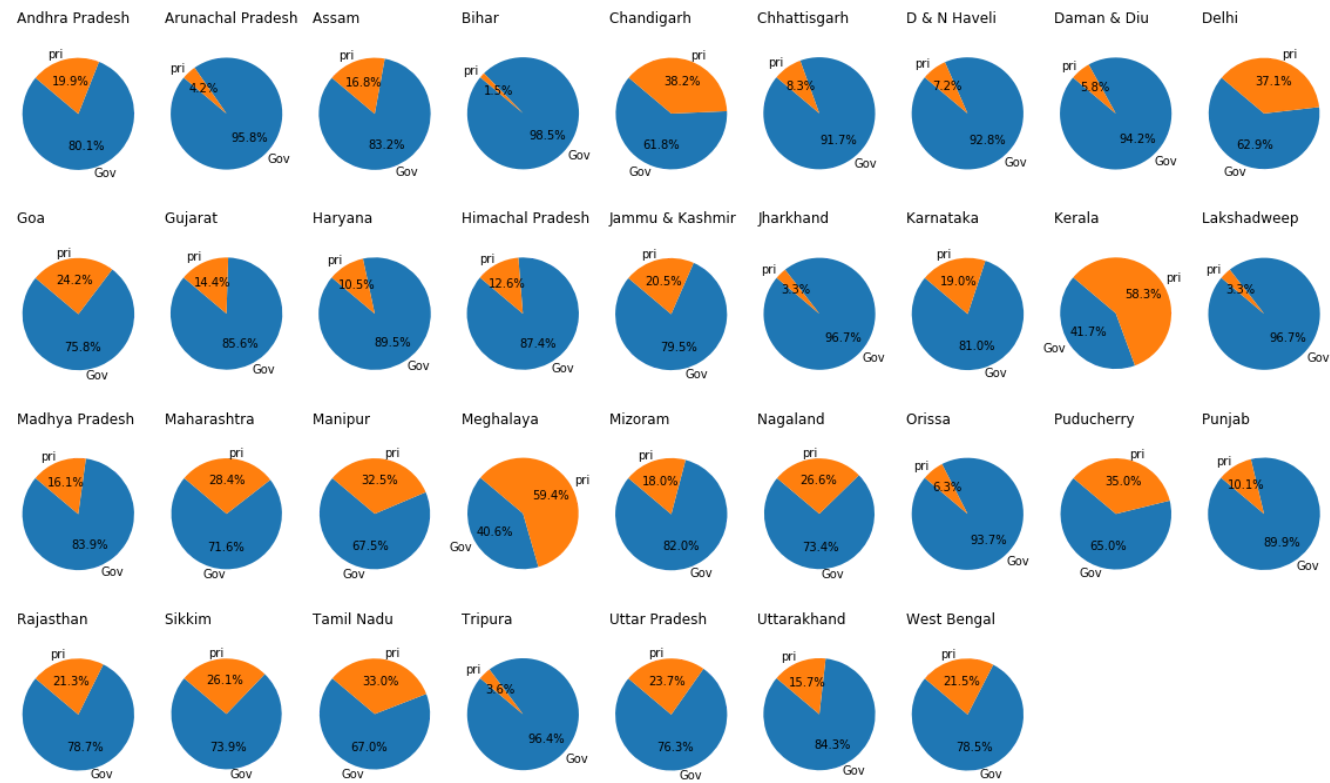
In [104]:

```python
1   #read
2   xls_06_07 = pd.ExcelFile('data processed/SRC_Rawdata_2006-07_mod.xls')
3   basicdata_06_07 = pd.read_excel(xls_06_07, 'Basic Data')
4   schoolfacility_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
5   schoolcondition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
6   enrolment_06_07 = pd.read_excel(xls_06_07, 'Enrolment')
7   teacher_06_07 = pd.read_excel(xls_06_07, 'Teacher')
8
9   #drop unnecessary features
10  basicdata_06_07.drop(['statcd','ac_year','area_sqkm','districts_covered','blocks','clusters'], axis = 1, in
11
12  #prepare data frame for plot
13  df1 = basicdata_06_07.filter(['statname','schools','literacy_rate'], axis=1)
14  df2 = schoolfacility_06_07.govt_sch_1 + schoolfacility_06_07.govt_sch_2 + schoolfacility_06_07.govt_sch_3 +
15  df3 = schoolfacility_06_07.pvt_sch_1 + schoolfacility_06_07.pvt_sch_2 + schoolfacility_06_07.pvt_sch_3 + sc
16
17  govpercentage = (df2/df1.schools)*100
18  privpercentage = (df3/df1.schools)*100
19
20  df = pd.concat([df1, df2, df3,govpercentage,privpercentage], axis=1, sort=False)
21  df = df.rename(columns = {"statname": "STATNAME", "schools":"SCHTOT", 0:"SCHTOTG", 1:"SCHTOTPR", 2:"GOVPER"
22
23  plt.figure(figsize=(20,12))
24  print("Government vs Private School Ratio In 2006-07")
25  for i in range(1,len(df)):
26      plt.subplot(4,9,i)
27      plt.title(df['STATNAME'][i],loc='left')
28      top = ['Gov','pri']
29      uttar = df.loc[df['STATNAME'] == df['STATNAME'][i],:]
30      value =[float(uttar['SCHTOTG']/uttar['SCHTOT'])*100,float(uttar['SCHTOTPR']/uttar['SCHTOT'])*100]
31      plt.pie(value, labels=top, autopct='%1.1f%%',startangle=140)
32      plt.axis('equal')
33  plt.show()
34
```
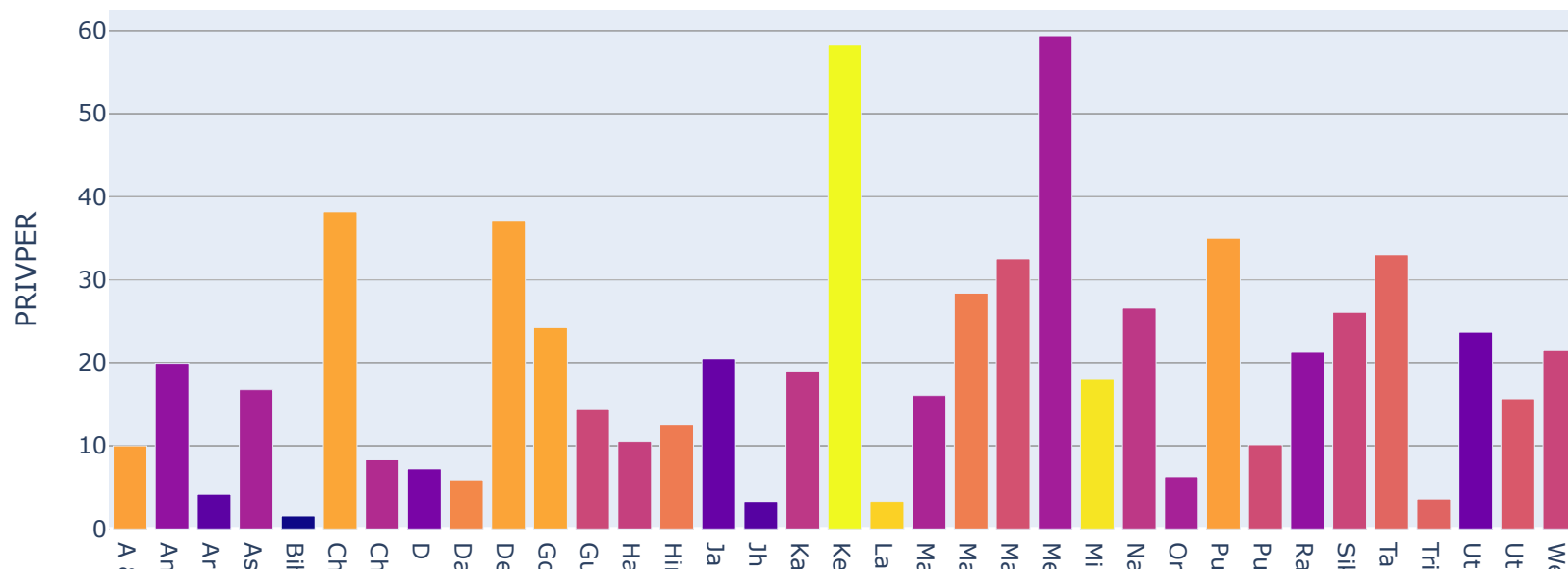
Government vs Private School Ratio In 2006-07
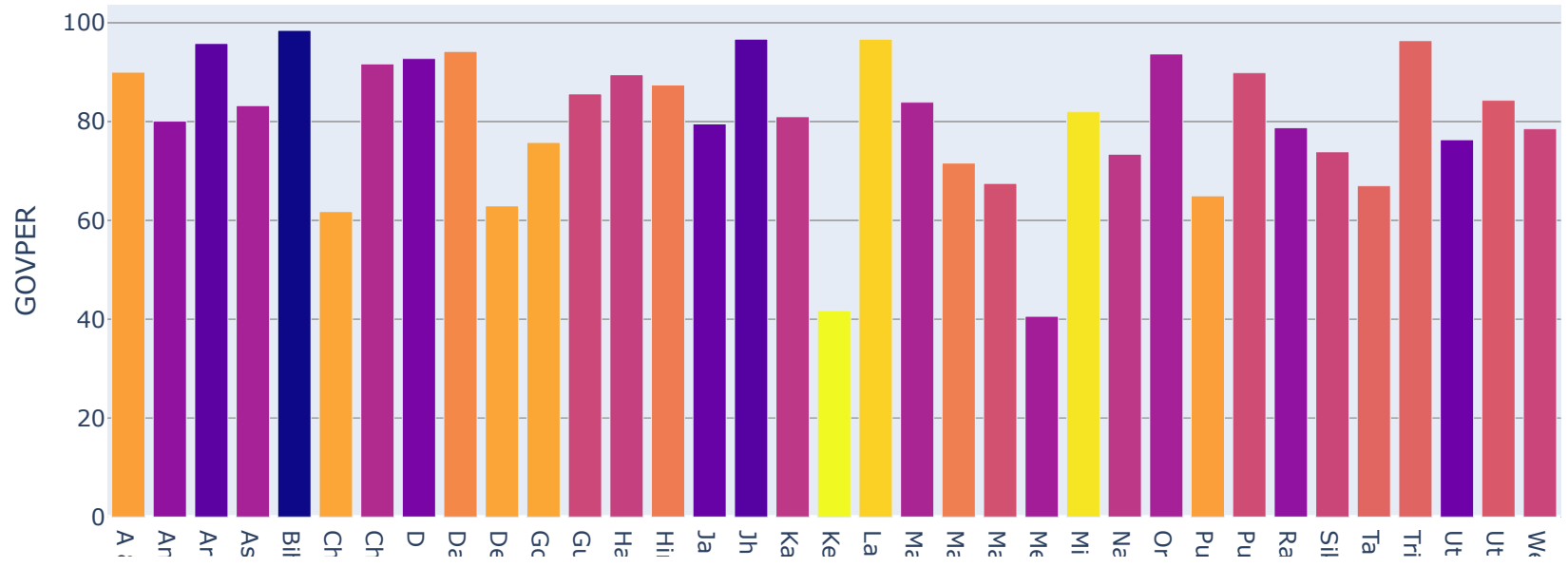
**Literacy rate of gov. and private school**

In [105]:
```python
fig = px.bar(df, x='STATNAME', y='PRIVPER',color='literacy_rate', title = "LITERACY RATE OF PRIVATE SCHOOL
fig.show()
```

## LITERACY RATE OF PRIVATE SCHOOL DIFFERENT STATE IN 2006-07

In [106]:
```python
fig = px.bar(df, x='STATNAME', y='GOVPER',color='literacy_rate', title = "LITERACY RATE OF GOVERNMENT SCHOO
fig.show()
```

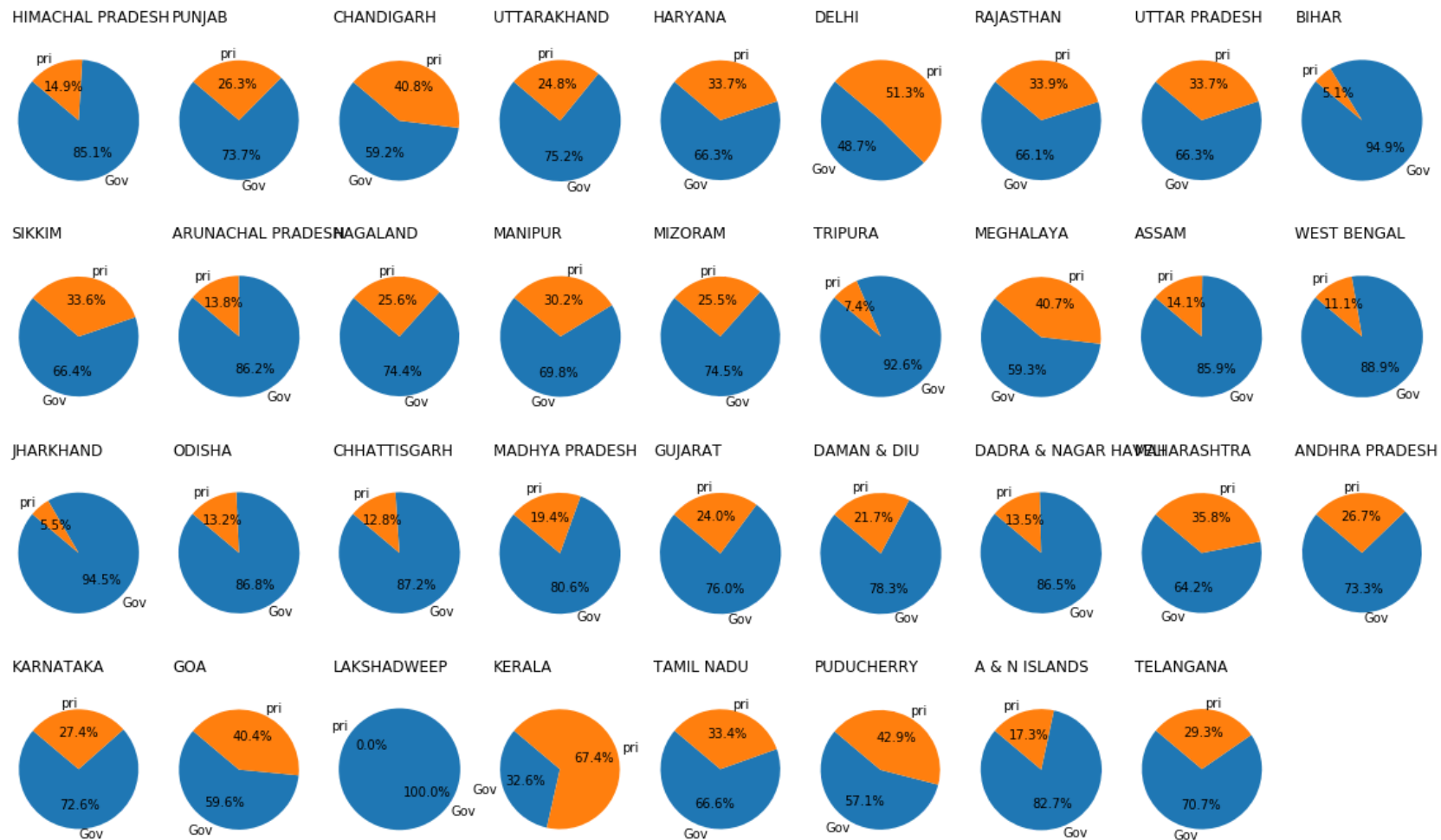LITERACY RATE OF GOVERNMENT SCHOOL DIFFERENT STATE IN 2006-07

In [107]:

```python
data = pd.read_csv("data processed/2016-17_Elementary.csv")
features = list(data.columns)
unique = list(set(list(metadata_hash.keys())).intersection(set(features)))

x = []
for i in features:
    if i not in list(metadata_hash.keys()):
        x.append(i)
plt.figure(figsize=(20,12))
print("Government vs Private School Ratio In 2016-17")
for i in range(1,len(data)):
    plt.subplot(4,9,i)
    plt.title(data['STATNAME'][i],loc='left')
    top = ['Gov','pri']
    uttar = data.loc[data['STATNAME'] == data['STATNAME'][i],:]
    value =[float(uttar['SCHTOTG']/uttar['SCHTOT'])*100,float(uttar['SCHTOTP']/uttar['SCHTOT'])*100]
    plt.pie(value, labels=top, autopct='%1.1f%%',startangle=140)
    plt.axis('equal')
plt.show()
```

Government vs Private School Ratio In 2016-17
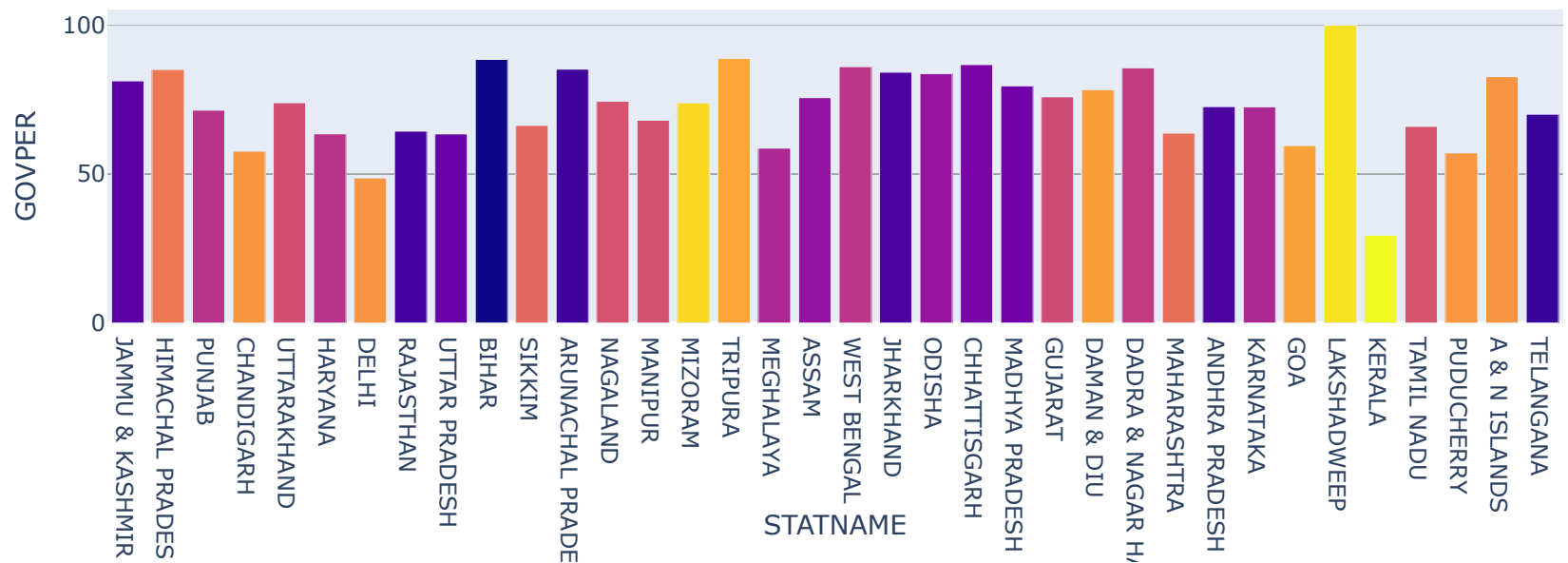
**HIMACHAL PRADESH PUNJAB** — pri 14.9% / Gov 85.1%

**PUNJAB** — pri 26.3% / Gov 73.7%

**CHANDIGARH** — pri 40.8% / Gov 59.2%

**UTTARAKHAND** — pri 24.8% / Gov 75.2%

**HARYANA** — pri 33.7% / Gov 66.3%

**DELHI** — pri 51.3% / Gov 48.7%

**RAJASTHAN** — pri 33.9% / Gov 66.1%

**UTTAR PRADESH** — pri 33.7% / Gov 66.3%

**BIHAR** — pri 5.1% / Gov 94.9%

**SIKKIM** — pri 33.6% / Gov 66.4%

**ARUNACHAL PRADESH** — pri 13.8% / Gov 86.2%

**NAGALAND** — pri 25.6% / Gov 74.4%

**MANIPUR** — pri 30.2% / Gov 69.8%

**MIZORAM** — pri 25.5% / Gov 74.5%

**TRIPURA** — pri 7.4% / Gov 92.6%

**MEGHALAYA** — pri 40.7% / Gov 59.3%

**ASSAM** — pri 14.1% / Gov 85.9%

**WEST BENGAL** — pri 11.1% / Gov 88.9%

**JHARKHAND** — pri 5.5% / Gov 94.5%

**ODISHA** — pri 13.2% / Gov 86.8%

**CHHATTISGARH** — pri 12.8% / Gov 87.2%

**MADHYA PRADESH** — pri 19.4% / Gov 80.6%

**GUJARAT** — pri 24.0% / Gov 76.0%

**DAMAN & DIU** — pri 21.7% / Gov 78.3%

**DADRA & NAGAR HAVELI** — pri 13.5% / Gov 86.5%

**MAHARASHTRA** — pri 35.8% / Gov 64.2%

**ANDHRA PRADESH** — pri 26.7% / Gov 73.3%

**KARNATAKA** — pri 27.4% / Gov 72.6%

**GOA** — pri 40.4% / Gov 59.6%

**LAKSHADWEEP** — pri 0.0% / Gov 100.0%

**KERALA** — pri 67.4% / Gov 32.6%

**TAMIL NADU** — pri 33.4% / Gov 66.6%

**PUDUCHERRY** — pri 42.9% / Gov 57.1%

**A & N ISLANDS** — pri 17.3% / Gov 82.7%

**TELANGANA** — pri 29.3% / Gov 70.7%

In [108]:
```python
1  GOVPER = (data['SCHTOTG']/data['SCHTOT'])*100
2  PRIVPER = (data['SCHTOTP']/data['SCHTOT'])*100
3  data = pd.concat([data,GOVPER,PRIVPER], axis=1, sort=False)
4  data = data.rename(columns = {0:"GOVPER", 1:"PRIVPER"})
5
6  fig = px.bar(data, x='STATNAME', y='GOVPER',color='OVERALL_LI', title = "LITERACY RATE OF GOVERNMENT SCHOOL
7  fig.show()
```
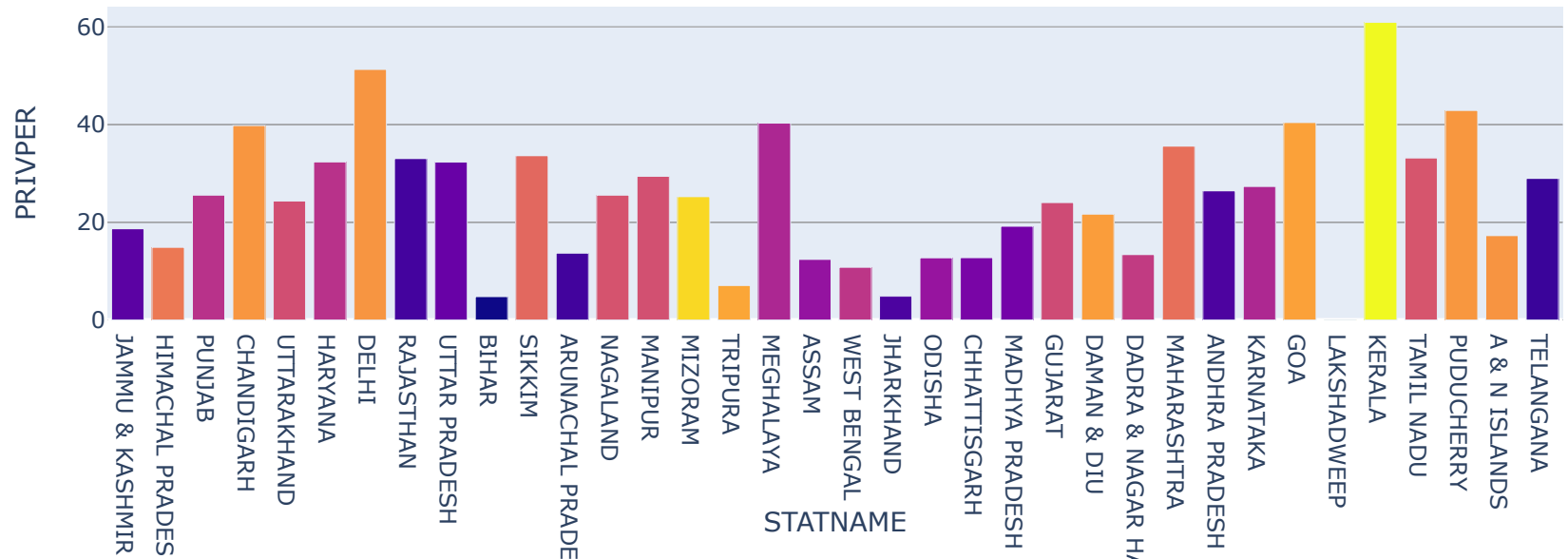
## LITERACY RATE OF GOVERNMENT SCHOOL DIFFERENT STATE IN 2016-17

In [109]:
```
1  fig = px.bar(data, x='STATNAME', y='PRIVPER',color='OVERALL_LI', title = "LITERACY RATE OF PRIVATE SCHOOL D
2  fig.show()
```

## LITERACY RATE OF PRIVATE SCHOOL DIFFERENT STATE IN 2016-17



## Observations by comparing literacy rate plots of 2006-07 and 2016-17:

We can observe that if the number of private schools is more in a state, then the literacy rate is also high, with the exception of Meghalaya.

A trend similar to private schools is not observed in the case of government schools.

This means that the resources and quality of education provided by government schools need to be improved instead of building new ones. And, more private schools must be built as they are improving the literacy rate.

**The quality of private schools in Meghalaya is not very good.**

## HOW SINGLE ROOM/TEACHER IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 2006-07

```python
In [110]:
    1  xls_06_07 = pd.ExcelFile('data processed/SRC_Rawdata_2006-07_mod.xls')
    2  basicdata_06_07 = pd.read_excel(xls_06_07, 'Basic Data')
    3  colmn = ['statcd', 'ac_year', 'statname']
    4  schoolfacility_06_07 = pd.read_excel(xls_06_07, 'School Facilities')
    5  schoolfacility_06_07.drop(colmn, inplace=True, axis=1)
    6  schoolcondition_06_07 = pd.read_excel(xls_06_07, 'School Condition')
    7  schoolcondition_06_07.drop(colmn, inplace=True, axis=1)
    8  enrolment_06_07 = pd.read_excel(xls_06_07, 'Enrolment')
    9  enrolment_06_07.drop(colmn, inplace=True, axis=1)
   10  teacher_06_07 = pd.read_excel(xls_06_07, 'Teacher')
   11  teacher_06_07.drop(colmn, inplace=True, axis=1)
   12
   13  #Concate all excel into one
   14  df06_07 = pd.concat([basicdata_06_07, schoolfacility_06_07,schoolcondition_06_07,enrolment_06_07,teacher_06
   15
   16  statename = list(df06_07['statname'])
   17  all_feature = list(df06_07.columns)
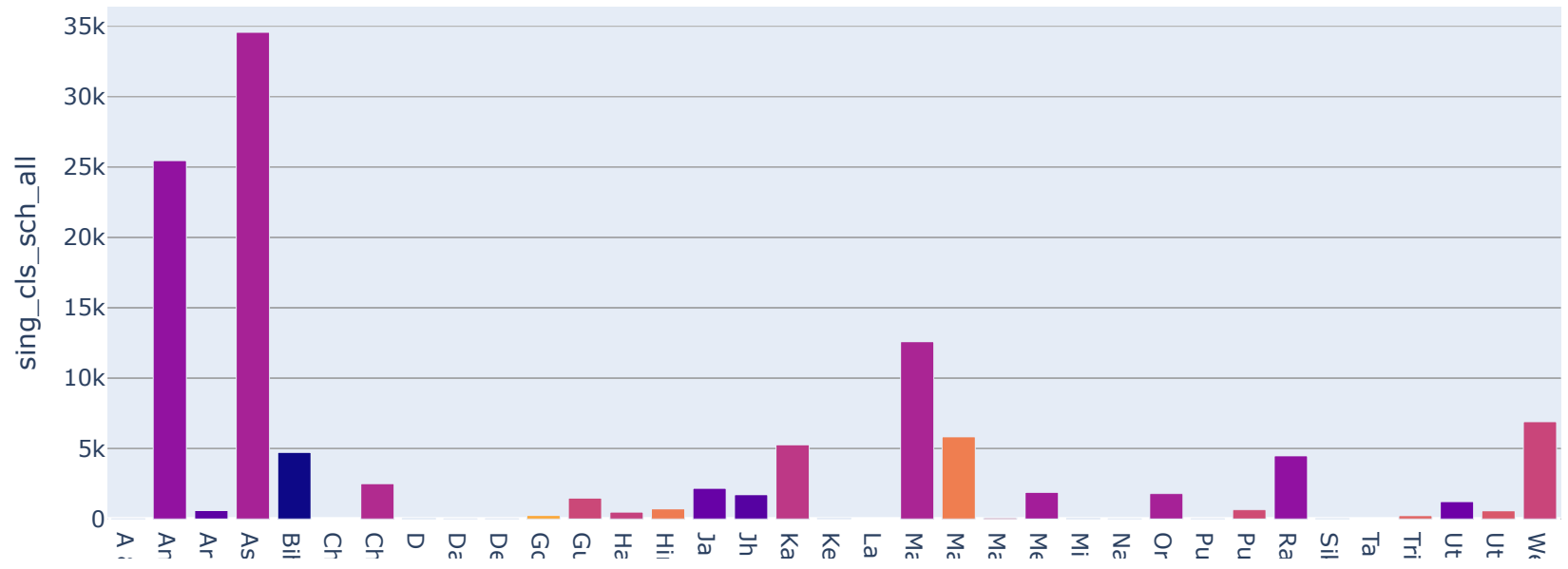```

In [111]:
```python
fig = px.bar(df06_07, x='statname', y='sing_tch_sch_all',color='literacy_rate', title="HOW SINGLE TEACHER I
fig.show()
```

HOW SINGLE TEACHER IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 2

In [112]:
```python
fig = px.bar(df06_07, x='statname', y='sing_cls_sch_all',color='literacy_rate', title="HOW SINGLE CLASS IN
fig.show()
```
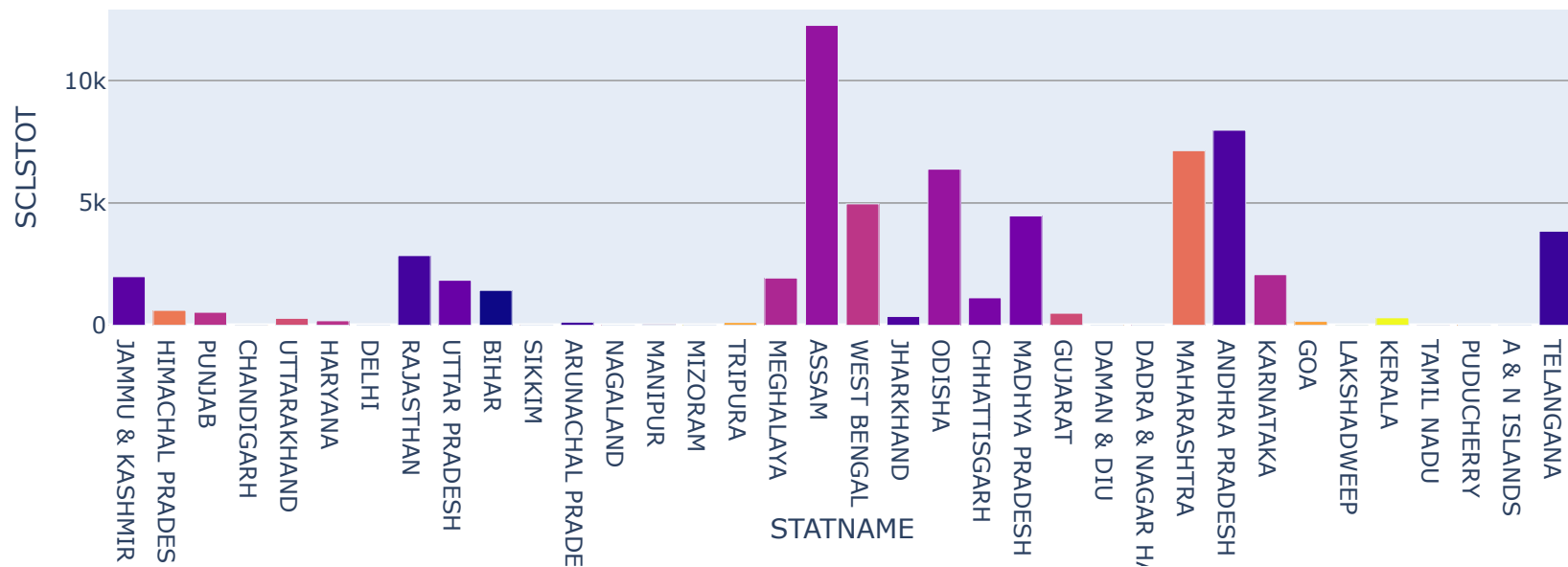
## HOW SINGLE CLASS IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 20(



In [113]:
```python
#READ 2016-17
data16_17 = pd.read_csv("data processed/2016-17_Elementary.csv")
data16_17.head()
enrol_per = (data16_17['ENRTOT']/(data16_17['TOTPOPULAT']*1000))*100
df16_17 = pd.concat([data16_17, enrol_per], axis=1, sort=False)
df16_17 = df16_17.rename(columns = {0:"enrolment_score"})
```
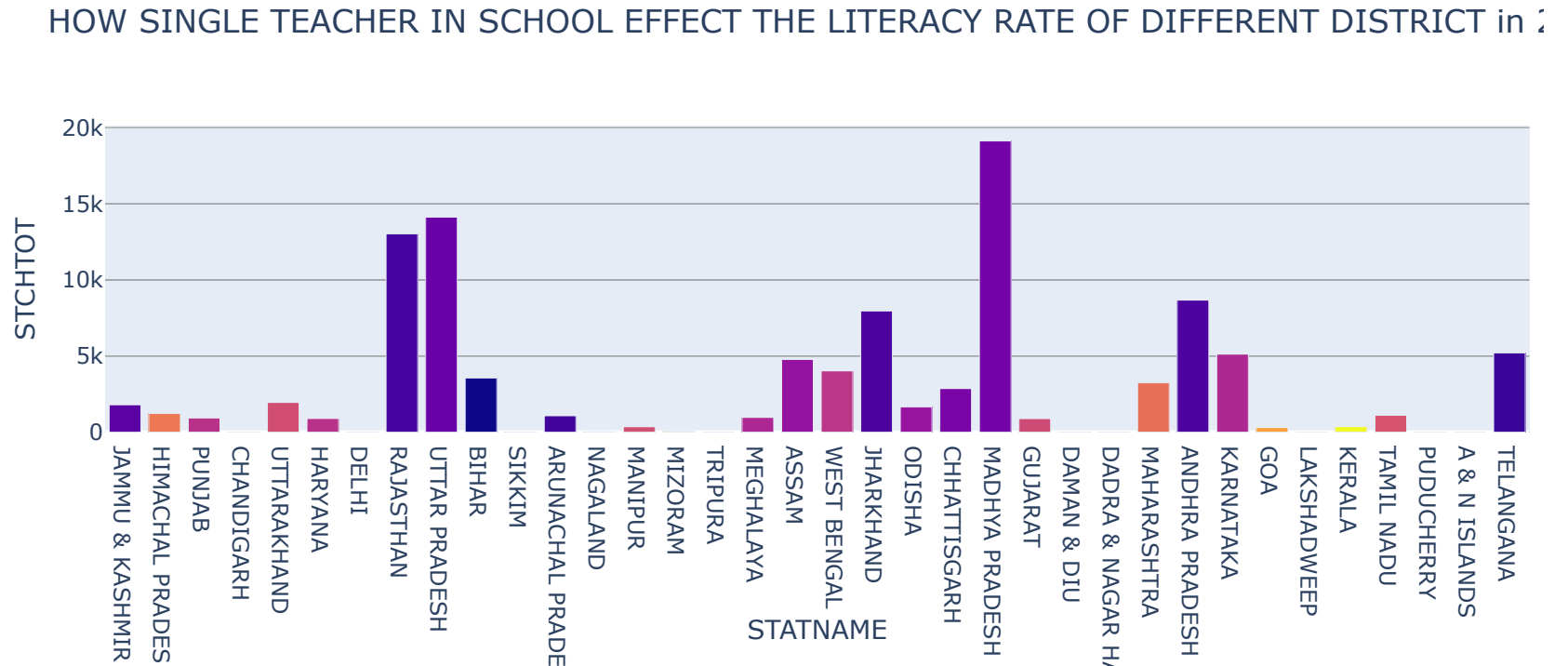
In [114]:
```python
fig = px.bar(data16_17, x='STATNAME', y='SCLSTOT',color='OVERALL_LI', title="HOW SINGLE CLASS IN SCHOOL EFF
fig.show()
```

HOW SINGLE CLASS IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 201

In [115]:
```python
fig = px.bar(data16_17, x='STATNAME', y='STCHTOT',color='OVERALL_LI', title="HOW SINGLE TEACHER IN SCHOOL E
fig.show()
```

HOW SINGLE TEACHER IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 2



## Observations by SINGLE ROOM/TEACHER IN SCHOOL EFFECT THE LITERACY RATE OF DIFFERENT DISTRICT in 2006-07 and 2016-17:

The literacy rate is generally poor for a single class or single teacher schools.

It is clear that the literacy rate decreases if we keep increasing the number of a single class or single teacher schools.

Single class and single teacher schools of Maharashtra have a decent literacy rate. It should be worked out why they are working so well in Maharashtra and then applied to other states as well.
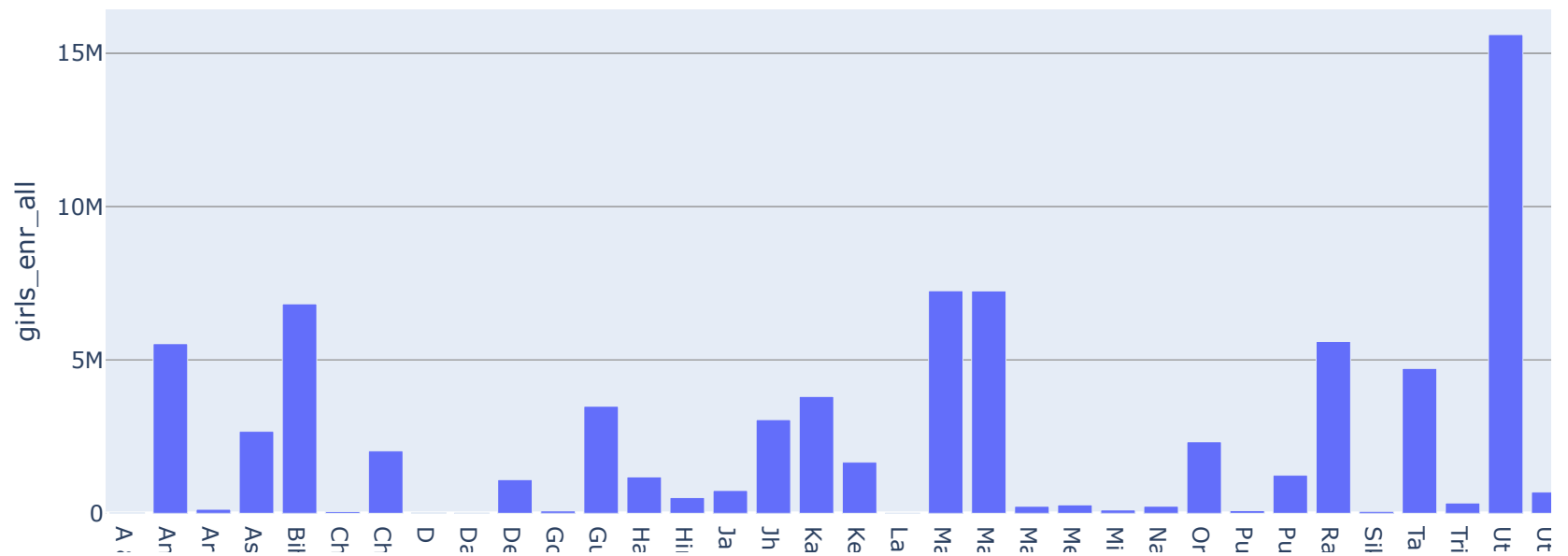
```
In [116]:   1  total_enrol = ['govt_enr_2','govt_enr_3','govt_enr_4','govt_enr_5','govt_enr_9','pvt_enr_1','pvt_enr_2','pv
```

```
In [117]:   1  val = df06_07['govt_enr_1']
            2  for i in total_enrol:
            3      df06_07['govt_enr_1'] += df06_07[i]
            4  val = pd.DataFrame(val)
```

```
In [118]:   1  # girls_enr_all
            2  fig = px.bar(df06_07, x='statname', y='girls_enr_all',title="GIRLS ENROLLMENT OF DIFFERENT STATE IN 2006-07
            3  fig.show()
```
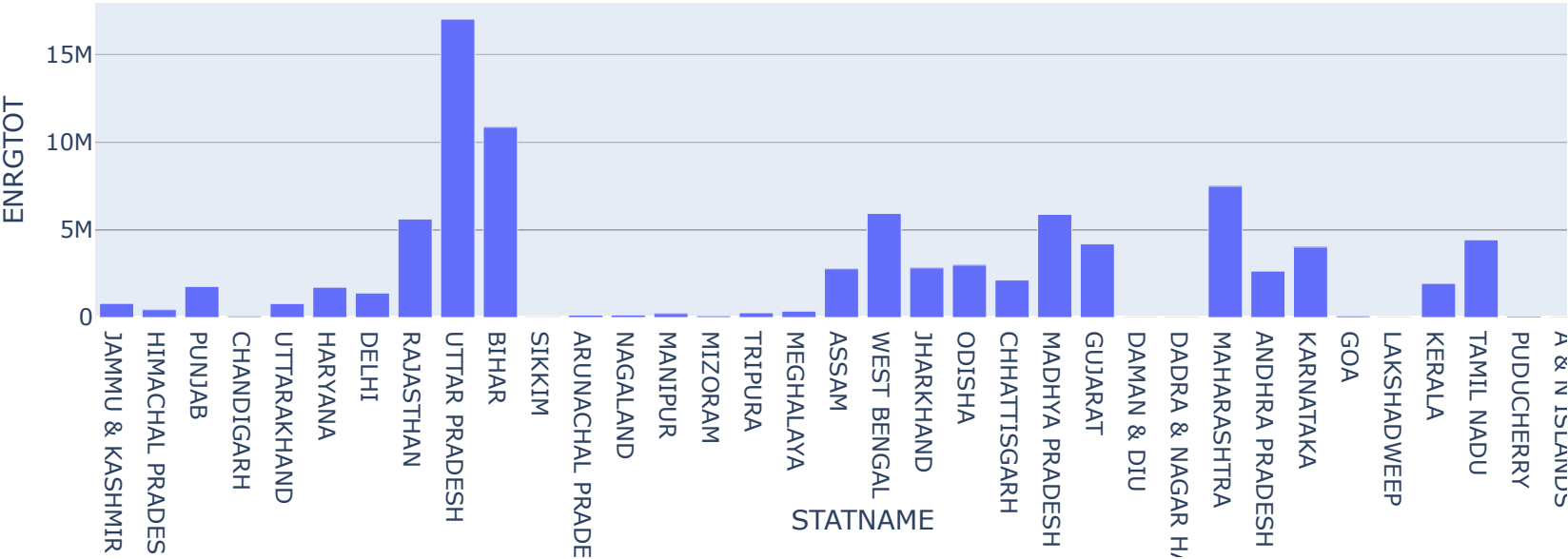
## GIRLS ENROLLMENT OF DIFFERENT STATE IN 2006-07

In [119]:
```python
data16_17['ENRTOTG'] += data16_17['ENRTOTP']
```

In [120]:
```python
# girls_enr_all
fig = px.bar(data16_17, x='STATNAME', y='ENRGTOT',title="GIRLS ENROLLMENT OF DIFFERENT STATE IN 2016-17", h
fig.show()
```

# GIRLS ENROLLMENT OF DIFFERENT STATE IN 2016-17



In [ ]:
```

```

In [ ]:
```

```