# heatmap of topics over features (from mixEHR outputs)

October 27, 2019

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
%matplotlib inline
qn = pd.read_csv("../sorted_newpheno.csv")

cols1 =␣
 →["typeid","featureid","topic0","topic1","topic2","topic3","topic4","topic5","topic6","topic7"
# cols2 =␣
 →["varid","stateid","typeid","varid","topic0","topic1","topic2","topic3","topic4","topic5","to
phi_norm = pd.read_csv("/home/mcb/li_lab/wliu92/data3/Oct_25/
 →oct25_5col_JCVB0_nmar_K10_iter300_phi_normalized.csv",names = cols1)
# eta = pd.read_csv("/home/mcb/li_lab/wliu92/data3/tr_out/
 →training_5col_JCVB0_nmar_K10_iter300_eta_normalized.csv",names = cols1)
# psi_norm = pd.read_csv("/home/mcb/li_lab/wliu92/data3/tr_out/
 →training_5col_JCVB0_nmar_K10_iter300_psi.csv",names = cols1)
# phi_norm["typeid"].unique()
```

```python
lst = []
for i in range (0,476,4):#loop each var
    large = -1
    largeid = i
    for l in range (i, i+4,1):
        tmp = eta.iloc[l,3]
        if (tmp>large):
            large = tmp
            largeid = l
        else:
            continue
    lst.append(largeid)#the largest row indexes
```

```python
#lst is the lst of row indexes with persumably max prob for each topic
data = pd.DataFrame(eta.iloc[i,2:] for i in lst)
x = np.array([])
newlst = []
for i in range(0,119,1):
```

1

```
        count = 0
        for j in range(0,10,1):
            if data.iloc[i,j]==1:
                count +=1
#           else:
#               data.iloc[i,j] = 1/data.iloc[i,j]
        if count !=10:
            if x.size ==0:
                x = np.array(data.iloc[i])
                newlst.append(i+1)
            else:
                x = np.vstack((x,np.array(data.iloc[i])))
                newlst.append(i+1)
        else:
            continue
```

[396]:
```
y_labels = []
y_label_states = []
for i in newlst:
    y_labels.append(qn.iloc[:,i].name)
for l in lst:
    y_label_states.append(int (eta.iloc[l,1]))
# y_label_states
```

[475]:
```
cols =␣
 ↪["topic0","topic1","topic2","topic3","topic4","topic5","topic6","topic7","topic8","topic9"]
# features = [i for i in range(0,23,1)]
fig,ax = plt.subplots(figsize = (10,10))
im = ax.imshow(x,cmap = "Reds")
plt.colorbar(im)
ax.set_xticks(np.arange(0,10,step=1.0))
ax.set_yticks(np.arange(0,23,step=1.0))
ax.set_xticklabels(cols)
ax.set_yticklabels(y_labels)
plt.setp(ax.get_xticklabels(), rotation=45, ha = "right", rotation_mode =␣
 ↪"anchor")
ax.set_title("topic probabilities over features")
plt.show()
```

## topic probabilities over features



```
import seaborn as sns
g = sns.clustermap(x, cmap = "Reds",row_cluster = True,col_cluster = True,
 →yticklabels = y_labels,xticklabels = cols,method = 'single',annot = True)
```

The heatmap displays clustered values with rows and columns. Row labels (top to bottom): sex, education, employment, income, age, EDS_short_total, CTQ_PHYS_AB, CTQTOT, CTQ_EMOT_AB, BDItotalscore, CTQ_SEX_AB, BDI_CAT, TEI_sexandphys_01yn, Childabphyssexemot_ctq_01modandsev, CTQ_SEXEMOTANDPHYS_AB_012_modandsev, SUICIDE_ATTEMPT_LIFETIME, PSYCHIATRIC_HOSP_LIFETIME, CTQ_physab_cat_modandsev, CTQ_EMOT_cat_modandsev, subs_abuse_past, CTQ_sexab_cat_modandsev, race_ethnic, subs_abuse_now. Column labels (left to right): topic7, topic0, topic2, topic9, topic1, topic4, topic5, topic8, topic3, topic6.

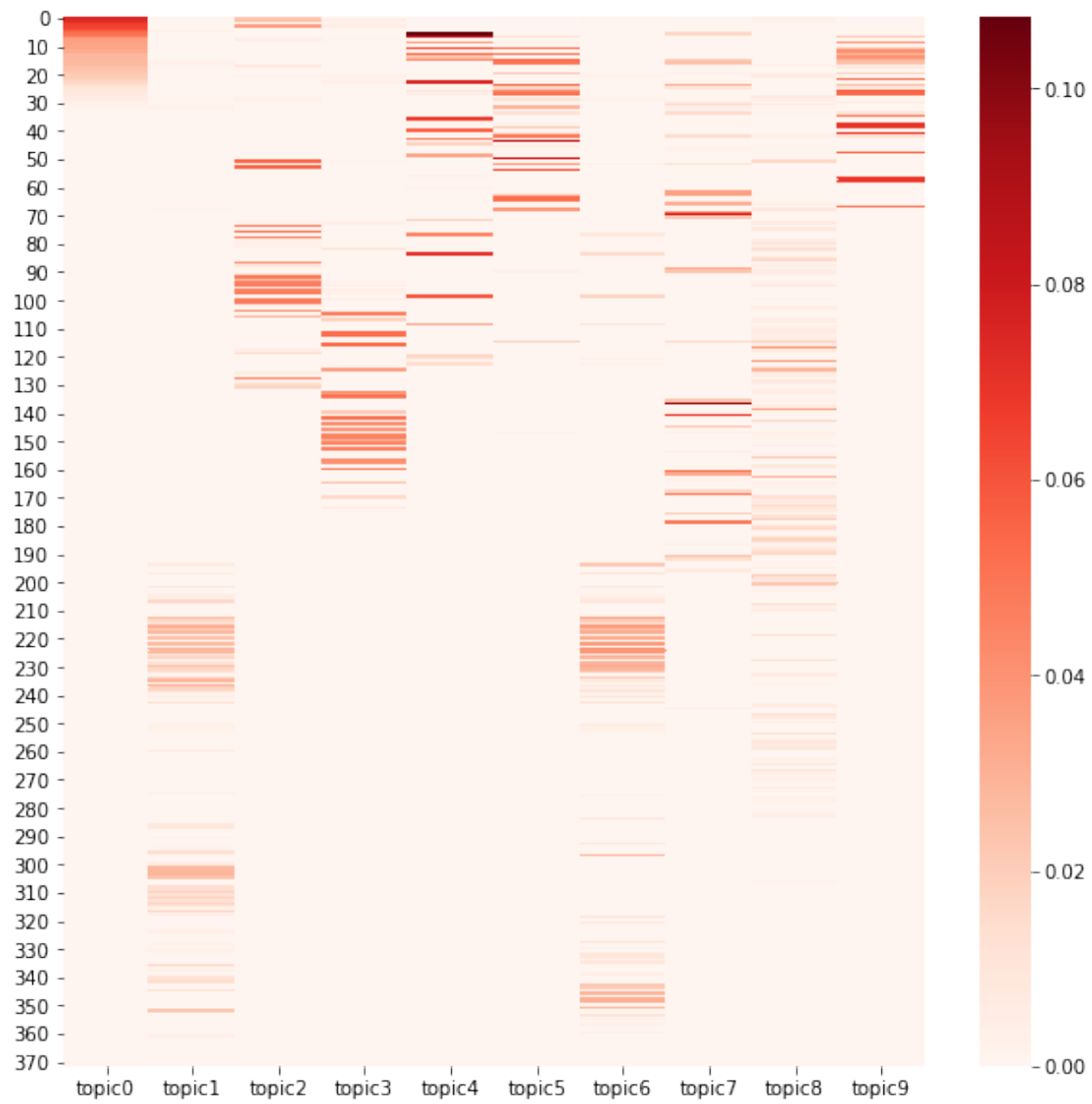| | topic7 | topic0 | topic2 | topic9 | topic1 | topic4 | topic5 | topic8 | topic3 | topic6 |
|---|---|---|---|---|---|---|---|---|---|---|
| sex | 0.67 | 0.73 | 0.73 | 0.74 | 0.57 | 0.54 | 0.86 | 0.84 | 0.8 | 0.78 |
| education | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 | 0.57 | 0.57 | 0.58 | 0.58 |
| employment | 0.67 | 0.66 | 0.64 | 0.66 | 0.67 | 0.67 | 0.69 | 0.69 | 0.69 | 0.67 |
| income | 0.31 | 0.31 | 0.3 | 0.3 | 0.3 | 0.3 | 0.31 | 0.31 | 0.31 | 0.31 |
| age | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 |
| EDS_short_total | 0.4 | 0.39 | 0.37 | 0.37 | 0.38 | 0.39 | 0.61 | 0.57 | 0.48 | 0.47 |
| CTQ_PHYS_AB | 0.18 | 0.56 | 0.81 | 0.82 | 0.6 | 0.76 | 0.051 | 0.063 | 0.085 | 0.1 |
| CTQTOT | 0.098 | 0.46 | 0.87 | 0.88 | 0.79 | 0.86 | 0.024 | 0.03 | 0.035 | 0.043 |
| CTQ_EMOT_AB | 0.12 | 0.65 | 0.9 | 0.9 | 0.87 | 0.87 | 0.021 | 0.019 | 0.048 | 0.057 |
| BDItotalscore | 0.65 | 0.78 | 0.85 | 0.85 | 0.84 | 0.85 | 0.058 | 0.082 | 0.13 | 0.21 |
| CTQ_SEX_AB | 0.95 | 0.95 | 0.96 | 0.96 | 0.96 | 0.96 | 0.058 | 0.1 | 0.28 | 0.51 |
| BDI_CAT | 0.95 | 0.95 | 0.97 | 0.97 | 0.96 | 0.96 | 0.022 | 0.029 | 0.064 | 0.26 |
| TEI_sexandphys_01yn | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.015 | 0.019 | 0.043 | 0.12 |
| Childabphyssexemot_ctq_01modandsev | 0.96 | 0.97 | 0.98 | 0.98 | 0.97 | 0.98 | 0.0096 | 0.011 | 0.02 | 0.032 |
| CTQ_SEXEMOTANDPHYS_AB_012_modandsev | 0.96 | 0.96 | 0.97 | 0.97 | 0.97 | 0.97 | 0.011 | 0.014 | 0.034 | 0.05 |
| SUICIDE_ATTEMPT_LIFETIME | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.083 | 0.44 | 0.98 | 0.98 |
| PSYCHIATRIC_HOSP_LIFETIME | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.09 | 0.52 | 0.97 | 0.98 |
| CTQ_physab_cat_modandsev | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.017 | 0.024 | 0.8 | 0.97 |
| CTQ_EMOT_cat_modandsev | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.015 | 0.021 | 0.93 | 0.98 |
| subs_abuse_past | 0.95 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.073 | 0.18 | 0.42 | 0.91 |
| CTQ_sexab_cat_modandsev | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.02 | 0.028 | 0.11 | 0.94 |
| race_ethnic | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.88 | 0.9 | 0.92 | 0.92 |
| subs_abuse_now | 0.88 | 0.87 | 0.88 | 0.88 | 0.87 | 0.88 | 0.43 | 0.79 | 0.86 | 0.87 |

[609]:
```python
# phi_norm.iloc[0:5]
```

[622]:
```python
new_arr = np.array(phi_norm.iloc[:,2:])

g1,ax1 = plt.subplots (1,1,figsize = (10,10))
g1 = sns.heatmap(new_arr, ax =ax1, xticklabels = cols, yticklabels=10, cmap =
 ↪"Reds",annot = False)
```

```
[669]: g1 = sns.clustermap(new_arr, xticklabels = cols, yticklabels=10, cmap = "Reds",␣
       ↪z_score = 0,method = "weighted",annot = False)
```