# College Scorecard: Cluster Analysis

*Michael L. Thompson*

*September 4, 2017*

## Contents

## Introduction

This is an exploratory analysis of the U.S. Dept. of Education College Scorecard database. My intent is to investigate patterns amongst the colleges as visualized using t-distributed Stochastic Neighbor Embedding (t-SNE). This method projects the high-dimensional data into two dimensions. From there, I can apply hierarchical clustering to identify clusters in the new 2-D space.

## Prepare Data

We read in the College Scorecard dataset and convert columns into Bayes factors, which accentuate differences amongst the colleges. Colleges having a disproportionately high number of students with a certain attribute – say, an SAT in excess of 1400 – will have highly positive Bayes factors for that attribute.

I strip out a lot of the variables that define the student body demographics. The idea is that I'd like to identify structure in the "outcome" variables – things like academic disciplines, completion rates, future earnings, credit default rates, etc. – and then later check if this structure is correlated to demographics – things like geographic location, campus setting, student ethnicity, etc.

```
glmdata_all <- DataSpec$studentBF %>%
  dplyr::select(
    c(-1, -(3:8)), -matches('_(WHITE|BLACK|ASIAN|OTHER|HISP|NRA|AIAN|UNKN)|2MOR|UNKN|NHPI|AIAN|BF_male|
    -matches('Challenge|_DEP_STAT_|notvet|le24y|OUTOFSTATE|prior|(^BF_[gl][et].+[0-9]+K$)|locale|FarWest
  ) %>%
  select_if( .predicate = function(x) any(x != x[[1]]) ) %>%
  filter( complete.cases(.) )
tsne_mat_all <- glmdata_all %>% select(-College) %>% as.matrix() %>% scale()
```

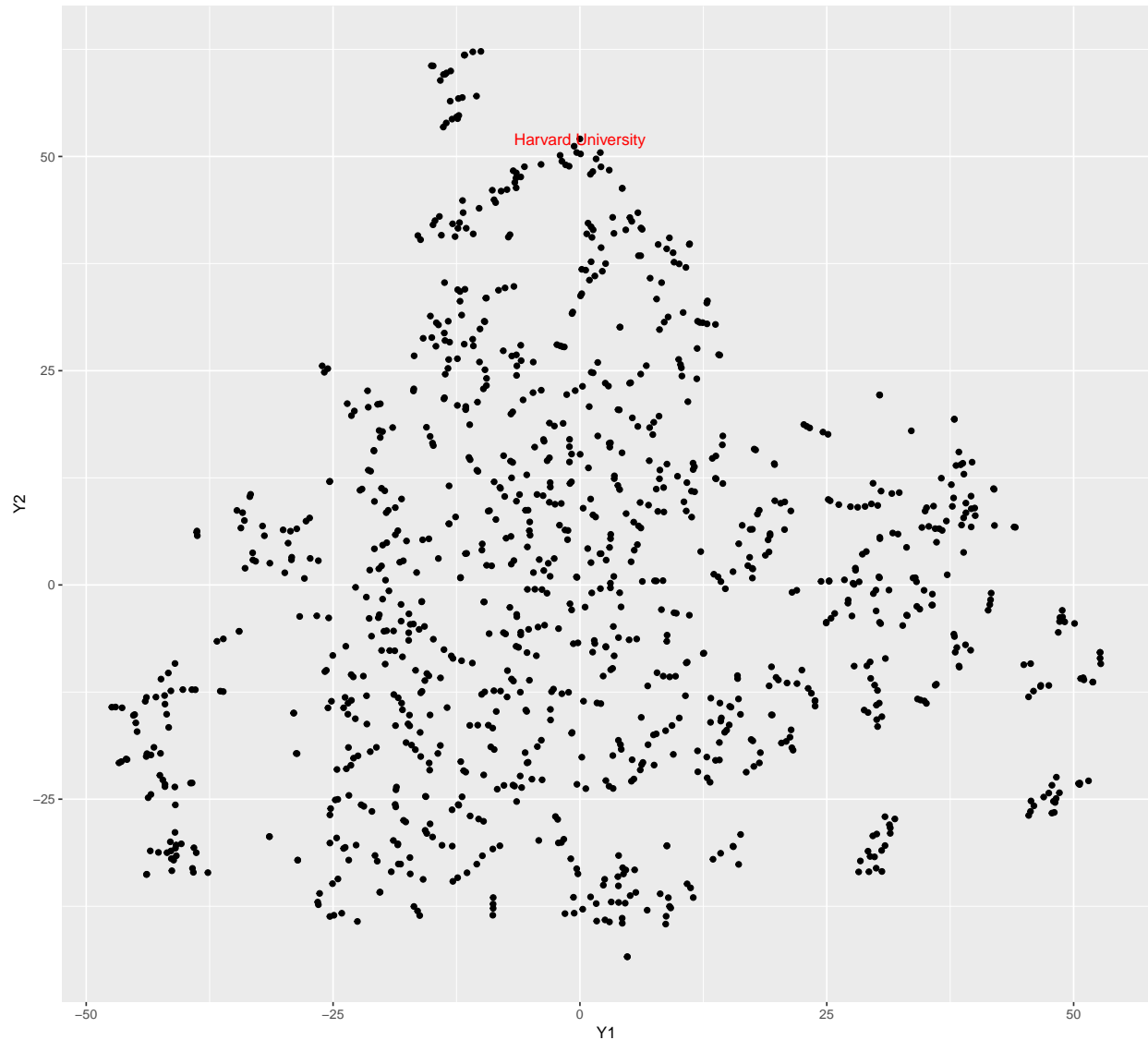## Perform t-Distributed Stochastic Neighbor Embedding (t-SNE)

Now, I'll map the data into a 2-D space using t-SNE.
Hopefully, it will be easy to see clusters of colleges.

It takes a bit of trial and error (short of doing a formal hyperparameter optimization) to arrive at hyperparameters capable of generating discernible structure in a 2-D scatterplot.

```r
set.seed( 173 )
tsne_all <- Rtsne( tsne_mat_all, perplexity = 10, initial_dims = 50, theta = 0.5, max_iter = 2000 )

# Rotate coordinates so that high-prestige colleges appear at high Y2 coordinates,
# i.e. in the top center of the plot.
i_harvard <- grep( 'Harvard', glmdata_all$College )
harvard_coord <- tsne_all$Y[i_harvard,]
harvard_angle <- atan(harvard_coord[2]/harvard_coord[1])
rotate_angle  <- pi/2 - harvard_angle
rotation_matrix <- matrix(
  c(cos(rotate_angle),sin(rotate_angle),-sin(rotate_angle),cos(rotate_angle)),
  2,2, byrow = TRUE
)
tsne_all$Y %<>% { (.) %*% rotation_matrix }
if( abs(tsne_all$Y[i_harvard,2]) < abs(tsne_all$Y[i_harvard,1]) ){
  tmp <- tsne_all$Y[,1]
  tsne_all$Y[,1] <- tsne_all$Y[,2]
  tsne_all$Y[,2] <- tmp
}
if( tsne_all$Y[i_harvard,2] < 0 ){
  tsne_all$Y[,2] <- -tsne_all$Y[,2]
}
```

```r
tsne_all$Y %>%
  as_tibble() %>%
  setNames(c('Y1','Y2')) %>%
  {
    ggplot(.,aes(x=Y1,y=Y2)) +
      geom_point() +
      geom_text(
        data=(.) %>% mutate(College=glmdata_all$College) %>% filter(grepl('Harvard',College)),
        mapping=aes(label=College),
        color='red',
        size=4
      )
  } %>%
  print()
```
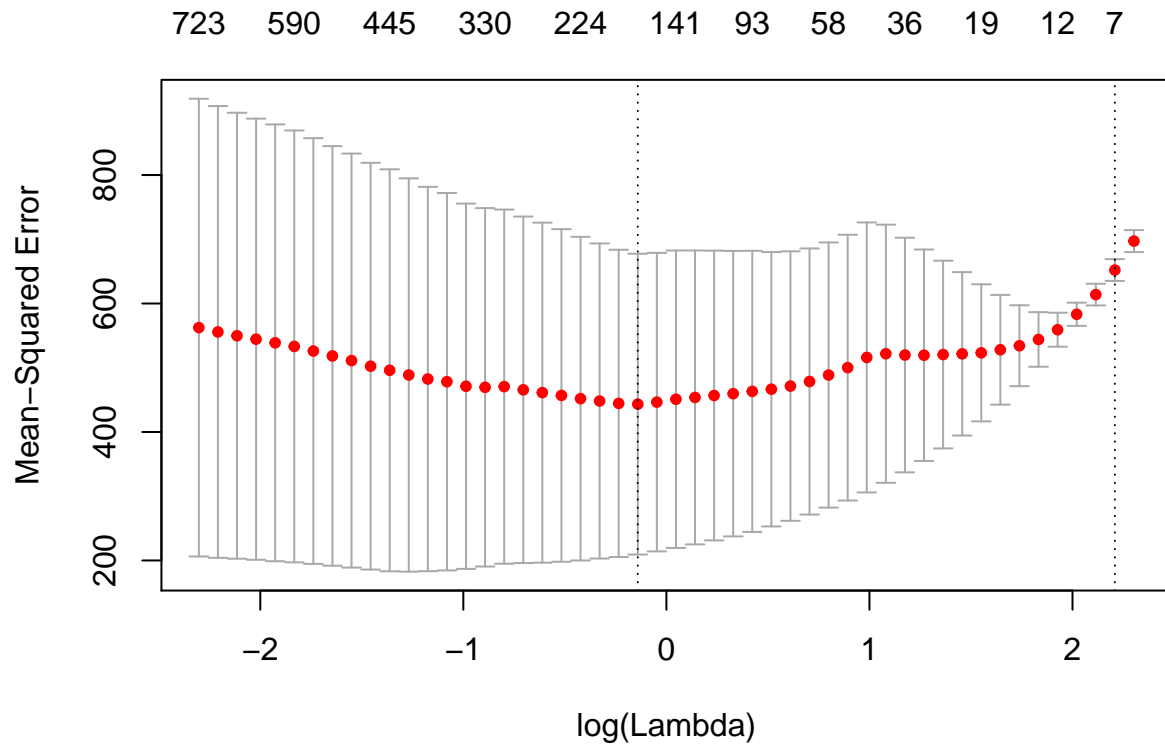
## Find Underlying Dimension Driving 2-D Structure

Using `glmnet`, I perform variable selection modeling of the 2-D t-SNE coordinates as responses vs. the original features from which the t-SNE coordinates were found. This way we'll have an approximate linear model showing which features contributed to which coordinate. As such, we'll have the basis for plotting a biplot of colleges overlayed on feature dimensions in 2-D, analogous to a PCA biplot.

```
mmat <- model.matrix( ~ .:. - 1, as.data.frame(tsne_mat_all))
# b <- eigen(cor(mmat))
# mmat <- mmat[,apply(b$vectors[,1:200],2,function(x) which.max(abs(x))) %>% unique() %>% sort()]

set.seed( 2393 )
tsne_glmnet_all <- cv.glmnet(
  x      = mmat,
  y      = tsne_all$Y,
  family = 'mgaussian',
  lambda = exp(seq(log(0.1),log(10),length.out = 50)))
```

```
)
plot( tsne_glmnet_all )
```



**Check the Predictions**

It can be tricky to find a subset of features and their interactions that both describe the t-SNE coordinates well *and* do not suffer from extreme collinearity, which can make the validation error at low `lambda` explode when applying `glmnet`.

Judging from the cross-validation curve above and the observed vs. predicted plots below, it looks like we've got a pretty good model.

```
lambda <- tsne_glmnet_all %$% lambda.min #{ exp( mean(log(c(lambda.min,lambda.1se))) ) }
pred <- tsne_glmnet_all %>% predict( newx = mmat, s = lambda ) %>% drop()
plot(pred[,1],tsne_all$Y[,1])
```

4

```r
plot(pred[,2],tsne_all$Y[,2])
```

## Visualize the Colleges in 2-D

```
tsne_glmnet_coef_all <- tsne_glmnet_all %>% coef( s = lambda )
# tsne_glmnet_coef_all$y1[-1] %>%
# { (.)[abs((.)[,1])>0,1] } %>%
# { data_frame(Coefficient = names(.), value = round(.,2)) } %>%
#   print()
# tsne_glmnet_coef_all$y2[-1] %>%
# { (.)[abs((.)[,1])>0,1] } %>%
# { data_frame(Coefficient = names(.), value = round(.,2)) } %>%
#   print()

tsne_coef_df_all <-
  tsne_glmnet_coef_all$y1 %>%
  as.matrix() %>%
  as.data.frame() %>%
  as_tibble() %>%
  rownames_to_column() %>%
  setNames(c("Coefficient","Y1")) %>%
  full_join(
    tsne_glmnet_coef_all$y2 %>%
      as.matrix() %>%
      as.data.frame() %>%
      as_tibble() %>%
```

```
        rownames_to_column() %>%
        setNames(c("Coefficient","Y2")),
      by = "Coefficient"
  ) %>%
  filter( abs(Y1) > 1.0E-9 | abs(Y2) > 1.0E-9 ) %>% slice(-1)

tsne_coef_df_all %>% mutate(mag = sqrt(Y1^2+Y2^2)) %>% arrange(desc(mag)) %>% print(n = 30)
```

```
## # A tibble: 164 x 4
##                                Coefficient          Y1          Y2
##                                      <chr>       <dbl>       <dbl>
## 1                  BF_ScienceTechnologies  4.27786035 -1.83442089
## 2                     BF_ForeignLanguages  1.50917251  3.11063075
## 3     BF_SAT_gt800le1000:BF_MechanicRepair  2.15702170  1.82800909
## 4                         BF_discBreadth   2.59109874  0.95726708
## 5               BF_AgricultureAgriculture   2.71810937 -0.07382335
## 6                           BF_SAT_gt1400  -0.39459961  2.44873066
## 7                      BF_PersonalCulinary  2.07235738 -1.12785793
## 8                         BF_fsend_5_2005  -0.74049350  2.11707889
## 9                     BF_VisualPerforming   2.18231155 -0.05955655
## 10                            BF_veteran   1.36254811 -1.60867889
## 11            BF_EngineeringTechnologies   1.58880282 -1.37435382
## 12                      BF_pell_ever_2005  1.21985948 -1.68407873
## 13               BF_PhilosophyReligious   0.73841205  1.90518852
## 14                    BF_PhysicalSciences  0.97916719  1.70908410
## 15                 BF_TheologyReligious  -1.93590613  0.03208227
## 16        BF_CommunicationsTechnologies   1.68233696  0.18596031
## 17                             BF_CDR3est  0.26320881 -1.67190922
## 18                           BF_AreaEthnic  0.96254797  1.33198095
## 19                       BF_MechanicRepair  1.55889711  0.02948046
## 20                             BF_History   0.89410980  1.20355083
## 21                          BF_p_gt48Kle75K  0.07718334  1.49090462
## 22                        BF_FamilyConsumer  1.36361712 -0.10300407
## 23                     BF_SAT_gt800le1000  -0.20755623 -1.32293371
## 24                       BF_EnglishLanguage  1.15850950  0.59568459
## 25               BF_TransportationMaterials  0.92382710 -0.90084071
## 26 BF_PersonalCulinary:BF_ForeignLanguages -1.27268893 -0.06349356
## 27                     BF_NaturalResources  0.88801773  0.89911025
## 28                           BF_not1stgen  -0.79841373  0.93536976
## 29                    BF_HealthProfessions  1.13516055 -0.41607078
## 30                       BF_SocialSciences  0.76216239  0.91547352
## # ... with 134 more rows, and 1 more variables: mag <dbl>
```

```
# tsne_coef_df %>%
# {
#   ggplot(., aes(x=Y1,y=Y2,label=Coefficient)) +
#     geom_point() +
#     geom_text( check_overlap = TRUE )
# } %>%
#   print()

key_terms <- tsne_coef_df_all %>%
  mutate(mag= sqrt(Y1^2+Y2^2)) %>%
  filter(abs(mag)>quantile(abs(mag),0.9)) %>%
```

```r
  arrange(desc(mag)) %$% Coefficient %>% setdiff("(Intercept)")


college_names <- glmdata_all %$%
  College %>%
  { gsub('^[0-9_]+','',. ) } %>%
  { gsub('Northwestern University','NU',.) } %>%
  { gsub('University of Notre Dame','Notre Dame U.',.) } %>%
  { gsub('Cornell University','Cornell U.',.) } %>%
  { gsub('California','Cal',. ) } %>%
  { gsub('Mass.+Inst.+Tech.+','MIT',. ) } %>%
  { gsub('(Mass|Penn|Wash)[^ ]+ *','\\1',.) } %>%
  { gsub('Polytechnic','Poly',. ) } %>%
  { gsub('Institute of Tech[^ ]+','IT',. ) } %>%
  { gsub('Tech.+Inst.+','Tech',. ) } %>%
  { gsub('State','St',. ) } %>%
  { gsub('University','U',. ) } %>%
  { gsub('(U of )|( U$)','',. ) } %>%
  { gsub('College','Col',. ) } %>%
  { gsub('New York','NY',.)} %>%
  { gsub('International','Intl',.) } %>%
  { gsub('North[^ ]+','N',.)} %>%
  { gsub('South[^ ]+','S',.)} %>%
  { gsub('West[^ ]+','W',.)} %>%
  { gsub('East[^ ]+','E',.)} %>%
  { gsub(' U-','-',.)} %>%
  { gsub('-Penn St ','',.)} %>%
  { gsub(' Col *$','',.)} %>%
  { gsub('-(Main)* Campus','',.)} %>%
  { gsub('^PennSt([^-]+)$','Penn St-\\1',.)} %>%
  { gsub(' and ','&',.)} %>%
  { gsub('Agricultural & Mechanical','A&M',.)}

st_abb <- state.abb %>% setNames( state.name )
for( st_nm in names(st_abb) ){
  college_names %<>% { gsub(st_nm,st_abb[st_nm],.) }
}

categories <- {
  mmat[,key_terms] %*%
    (tsne_coef_df_all %>% filter(Coefficient %in% key_terms) %$% Y2)
} %>%
  sapply(
    function(x,q){ length(q) - sum(x>q) + 1 },
    q=quantile(.,c(0.1,0.25,0.75,0.9))
  ) %>%
  factor()

tsne_df_all <- tsne_all$Y %>%
  as_tibble() %>%
  setNames(c("Y1","Y2")) %>%
  mutate(
    College = college_names,
```

```
    category = categories,
    BF_Income_gt110K = glmdata_all %$% {10.0^BF_p_gt110K}
  ) %>%
  dplyr::select( College, category, BF_Income_gt110K, everything() ) %>%
  mutate_at(funs(1.7*scale(.)),.vars=vars(Y1,Y2))
```

**Show Biplot for Structure Interpretaion**

we can overlay the feature dimensions on the college scatterplot in the 2-D t-SNE coordinate space. This allows us to more easily interpret the structure we're seeing. However, some of the interaction terms, in particular, are tricky to interpret because they have a positive value for a college if both of the features in the product making up the interaction have the same sign. So it could be that the college has a disproportionately higher *or* lower number of students having the attributes of both of the corresponding features.
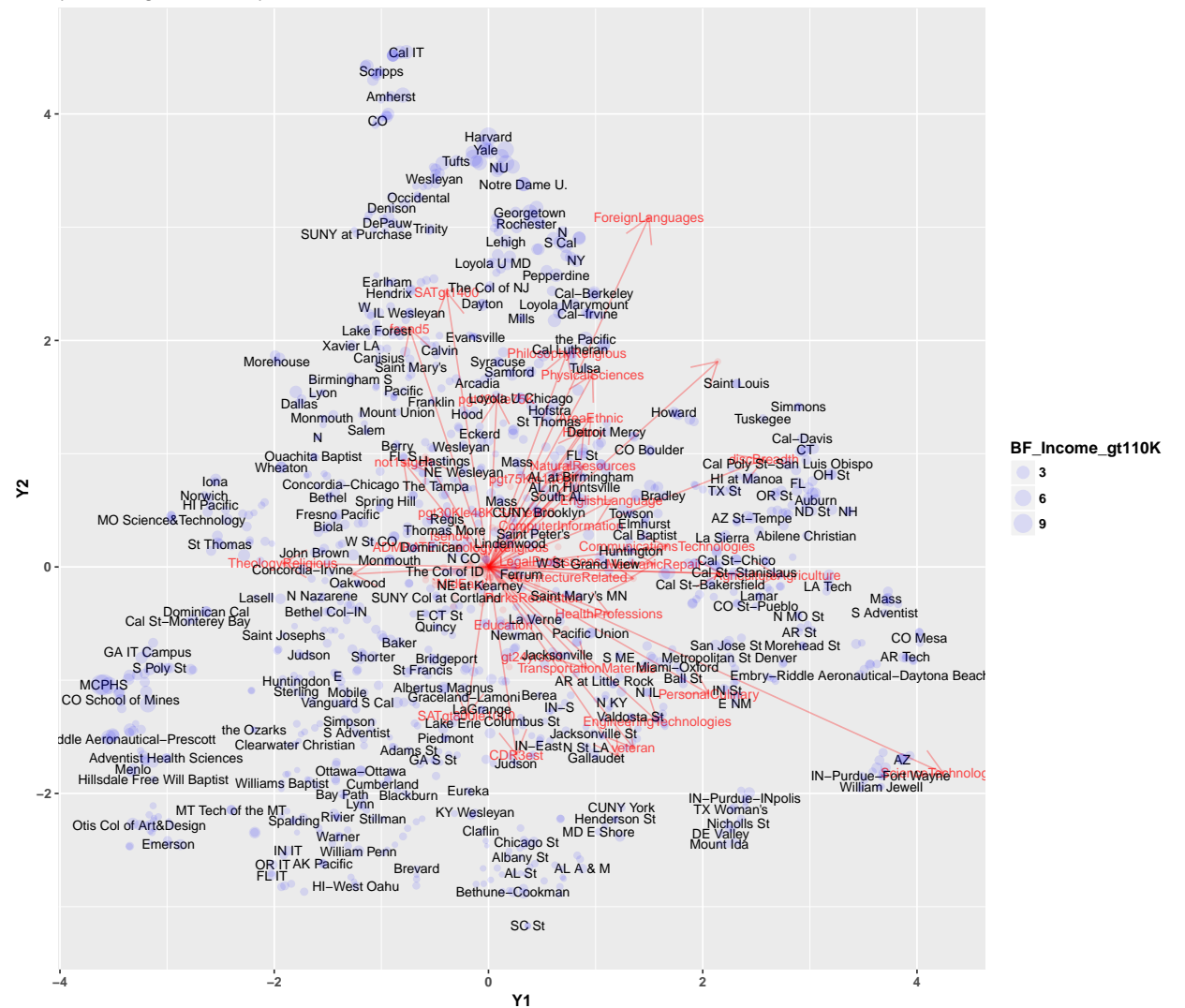
```
f_mult <- max(sqrt(tsne_df_all$Y1^2 + tsne_df_all$Y2^2))/max(sqrt(tsne_coef_df_all$Y1[-1]^2 + tsne_coef_
y2_min <- -3.5
tsne_coef_df_all %>%
  mutate(
    mag = sqrt(Y1^2 + Y2^2),
    Y2 = pmax(y2_min,Y2*f_mult),
    Y1 = Y1*f_mult,
    Coefficient = gsub('\\([^)]+\\)|(_*2005)|_','',gsub('BF_','',Coefficient))
  ) %>%
  {
    ggplot(., aes( x = Y1, y = Y2 ) ) +
      geom_point( color = 'red', alpha = 0.1 ) +
      geom_text(
        aes( label = Coefficient),
        color = 'red',
        alpha = 0.7,
        size = 3,
        check_overlap = TRUE
      ) +
      geom_segment(
        inherit.aes = FALSE,
        data = (.) %>% filter(mag>1),
        aes( x=0, y=0, xend=Y1, yend=Y2 ),
        color = 'red',
        alpha = 0.3,
        arrow = arrow(length = unit(0.03, "npc"))
      ) +
      geom_text(
        inherit.aes = FALSE,
        data = tsne_df_all,
        aes( x=Y1, y=Y2, label=College ),
        mapping=,
        color = 'black',
        size=3,
        check_overlap = TRUE
      ) +
      geom_point( data=tsne_df_all, aes(x=Y1,y=Y2, size = BF_Income_gt110K ), color='blue',alpha=0.1) +
      ggtitle( "t-SNE Biplot" , subtitle = "(blue = college, red = feature)") +
      theme( text = element_text( face = 'bold' ) ) #+
```

9

```
    #scale_y_continuous(limits = c(y2_min,4))
    #scale_y_continuous(limits = c(y2_min,4))
} %>%
print()
```

**t−SNE Biplot**

**(blue = college, red = feature)**



```
select_colleges <- c(
  '^OH St', '^MI-Ann Arbor', '^Purdue$', '^NU$',
  'Harvard', 'Yale', 'Princeton','^Penn$','^Cornell U\\.$','^Brown$',
  '^Howard$','Tuskegee','Hampton','Morehouse','Grambling','Bethune-Cookman',
  'Stanford','Johns Hopkins','Duke','Vanderbilt','Rice','Wash.+St Louis',
  'Notre Dame U\\.','^Pomona$','Harvey Mudd','Swarthmore',
  'MIT','Cal *IT'
)
tsne_select <- tsne_df_all %>%
  slice( sapply( select_colleges, function(nm_regex) grep(nm_regex,(.)$College) ) ) %$%
  set_rownames(as.matrix(select(.,Y1,Y2)),College) %>%
  round(1)
```

Here are the t-SNE 2-D coordinates for some notable universities:

- **Big 10**
  - Ohio State: 3.2, 0.8
  - Michigan: 0.7, 2.8
  - Purdue: 3.6, -0.9
  - Northwestern: 0.1, 3.5
- **Ivy League**
  - Harvard: 0, 3.8
  - Yale: 0, 3.7
  - Princeton: -0.8, 4.2
  - Penn: 0.2, 3.6
  - Cornell: 0.2, 3.5
  - Brown: 0.1, 3.5
- **HBCUs**
  - Howard: 1.7, 1.4
  - Tuskegee: 2.6, 1.3
  - Hampton Institute: 0.6, -0.7
  - Morehouse: -2, 1.8
  - Grambling: 0.5, -2.8
  - Bethune-Cookman: 0.2, -2.9
- **Others**
  - Stanford: -0.1, 3.6
  - Johns Hopkins: 0.4, 3.2
  - Duke: -0.1, 3.6
  - Vanderbilt: -0.1, 3.6
  - Rice: -0.2, 3.7
  - Washington U.-St. Louis: 0.1, 3.6
  - Notre Dame: 0.3, 3.4
  - Pomona: -0.9, 4.5
  - Harvey Mudd: -0.8, 4.5
  - Swarthmore: -0.9, 4.5
  - MIT: 0.2, 3.7
  - CalTech: -0.8, 4.5

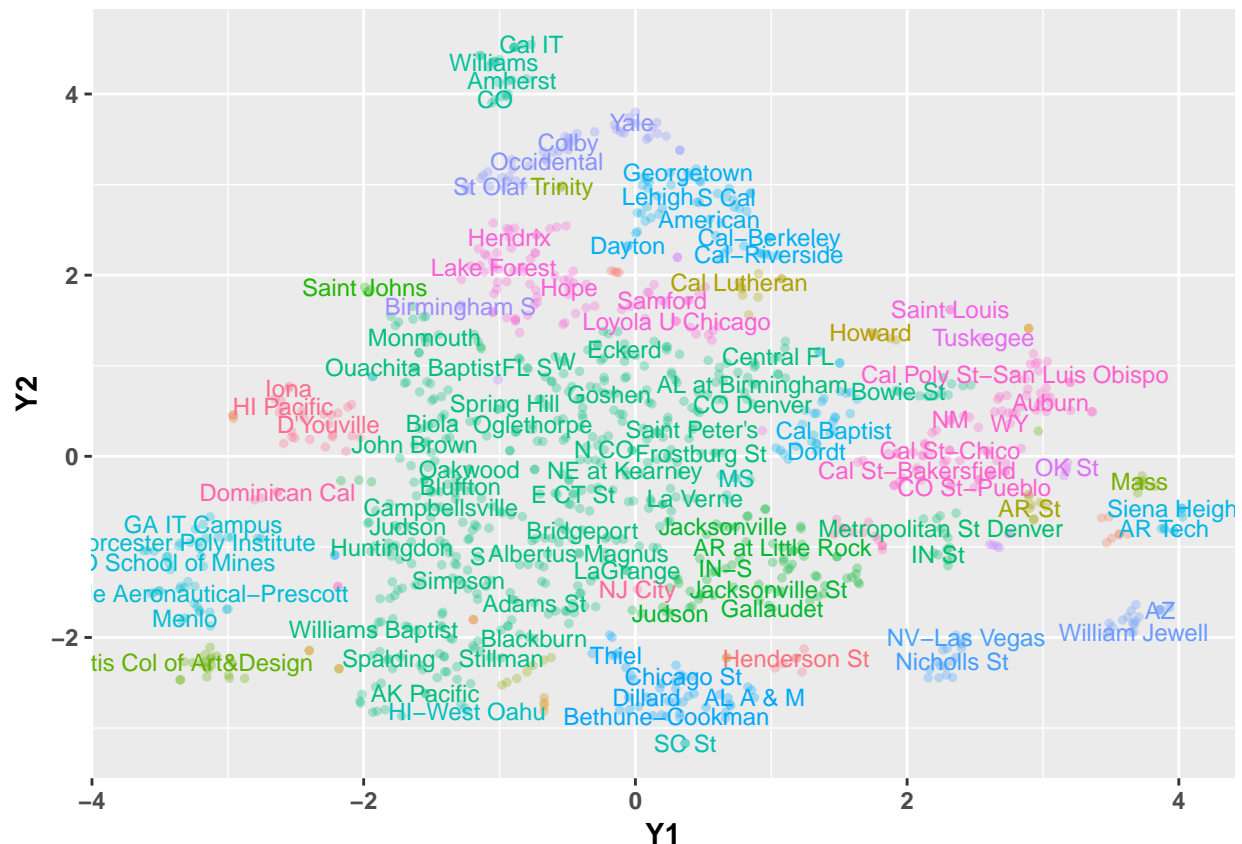**Perform Hierarchical Clustering**

Now, I perform cluster analysis. Hierarchical clustering is a quick way to identify clusters in the 2-D t-SNE space. We can then color the clusterings in a scatterplot to more easily visualize the structure.

```
tsne_mat_hc_all <- tsne_df_all %>% select(Y1,Y2) %>% as.matrix() %>% set_rownames(tsne_df_all$College)
hc_all <- hclust( d = dist( tsne_mat_hc_all ), method = 'single' )
n_cluster <- 55
cluster_id_all <- cutree( hc_all, k = n_cluster )

# plot( tsne_mat_hc, pch=20, cex=0.5 )
# for(j in seq_along(cl)){
#   points( tsne_mat_hc[ cl[[j]], ], pch=20, col=j, cex=1)
# }

# randomize so adjacent clusters are more likely to have very different colors.
set.seed(137)
cluster_id_all <- setNames( sample.int(n_cluster)[cluster_id_all], names(cluster_id_all) )
```

```
tsne_mat_hc_all %>%
  as_tibble() %>%
  mutate( College = names(cluster_id_all), cluster = factor( cluster_id_all ) ) %>%
  {
    ggplot(.,aes( x = Y1, y = Y2, color = cluster ) ) +
      geom_point( size = 1, alpha = 0.3 ) +
      geom_text( aes(label = College ), size = 3, check_overlap = TRUE ) +
      theme(
        text = element_text( face = 'bold' ),
        legend.position = 'none'
      )
  } %>%
  print()
```



**Show Biplot with Cluster Coloring**

Finally, we can overlay the feature dimensions on the 2-D

```
cluster_id_all <- cutree( hc_all, k = n_cluster )
y2_min <- -4
y2_max <- 3.49
y1 <- range(tsne_mat_hc_all[,1])
y1[1] <- 0.5*floor(y1[1]/0.5)
y1[2] <- 0.5*ceiling(y1[2]/0.5)
y2 <- range(tsne_mat_hc_all[,2])
```

```r
y2[1] <- 0.5*floor(y2[1]/0.5)
y2[2] <- 0.5*ceiling(y2[2]/0.5)

is_out_of_bounds <- function(x,bounds){ x<bounds[1] | x>bounds[2] }
# Assumes that value violating bounds is of same sign as bound violated AND that bounds are of opposite
bound_factor <- function(x,bounds){
  f1 <- ifelse(x<bounds[1],x/bounds[1],0)
  f2 <- ifelse(x>bounds[2],x/bounds[2],0)
  mapply(function(b1,b2) if(b1>b2) c(1,b1) else c(2,b2),f1,f2)
}
tsne_modified <- tsne_coef_df_all %>%
  mutate(
    mag = sqrt(Y1^2 + Y2^2) ,
    Coefficient = gsub('\\([^)]+\\)|(_*2005)|_','',gsub('BF_','',Coefficient)),
    Y1  = f_mult*Y1,
    Y2  = f_mult*Y2
  )

# check bounds to find if any violated
bchk1 <- bound_factor(tsne_modified$Y1,y1)
bchk2 <- bound_factor(tsne_modified$Y2,y2)
# bound on Y1 violated
w1 <- which(bchk1[2,] != 0)
# bound on Y2 violated
w2 <- which(bchk2[2,] != 0)
# Keep only coord Y1 or Y2 violated the most by each violating pt.
for( i in intersect(w1,w2)) { if(bchk1[2,i]>bchk2[2,i]) w2<-setdiff(w2,i) else w1<- setdiff(w1,i) }
# bound on Y1 violated: fix it
for( i in w1 ){
  tsne_modified$Y2[i] <- tsne_modified$Y2[i]*y1[bchk1[1,i]]/tsne_modified$Y1[i]
  tsne_modified$Y1[i] <- y1[bchk1[1,i]]
}
# bound on Y2 violated: fix it
for( i in w2 ){
  tsne_modified$Y1[i] <- tsne_modified$Y1[i]*y2[bchk2[1,i]]/tsne_modified$Y2[i]
  tsne_modified$Y2[i] <- y2[bchk2[1,i]]
}

tsne_modified %>%
  {
    ggplot(., aes( x = Y1, y = Y2 ) ) +
      geom_point( color = 'red', alpha = 0.1 ) +
      geom_segment(
        inherit.aes = FALSE,
        data = (.) %>% filter(mag>1),
        aes( x=0, y=0, xend=Y1, yend=Y2 ),
        color = 'red',
        alpha = 0.3,
        arrow = arrow(length = unit(0.03, "npc"))
      ) +
      geom_text(
        inherit.aes = FALSE,
        data = tsne_mat_hc_all %>%
```
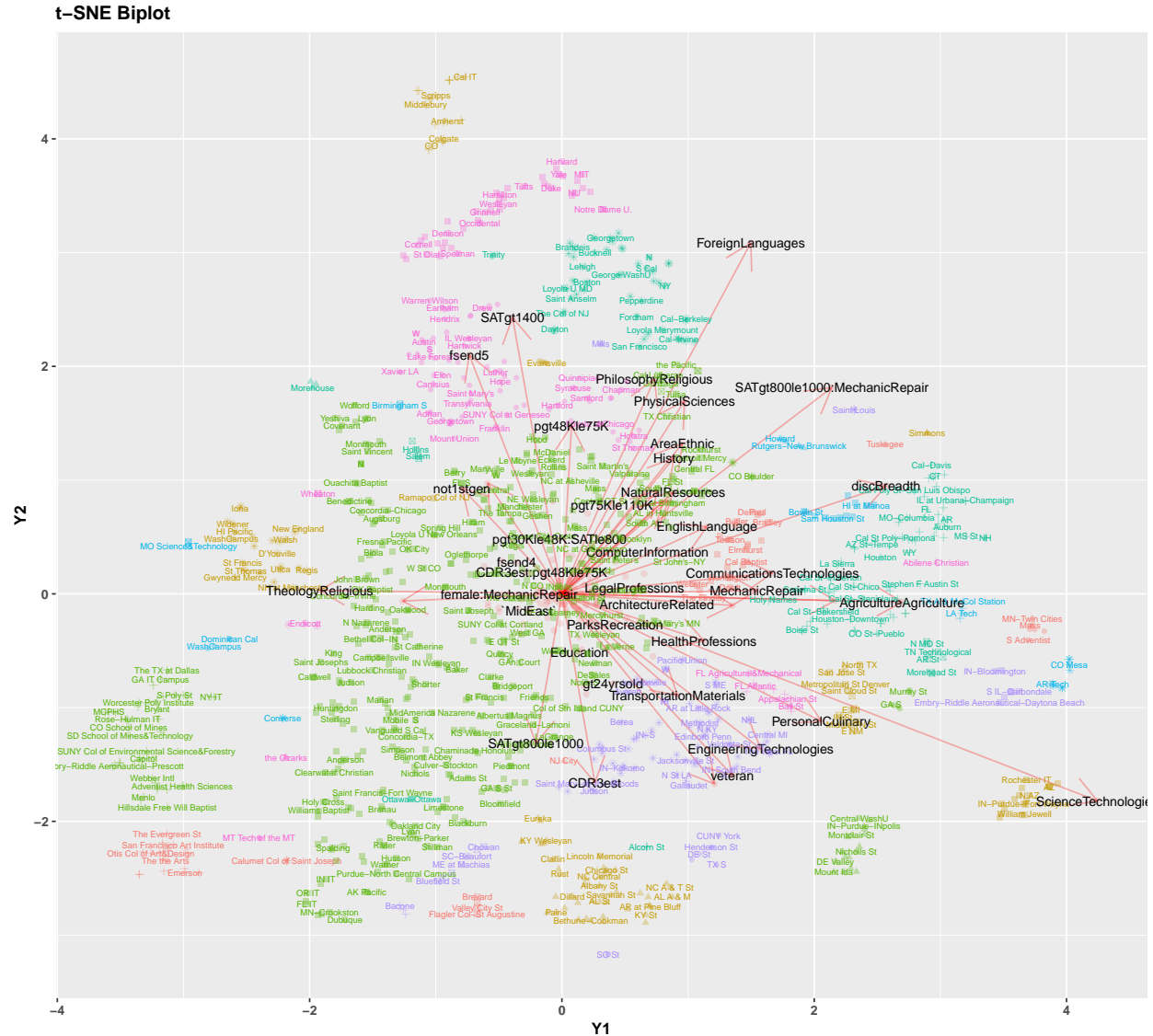
```
      as_tibble() %>%
      mutate(
        College = names(cluster_id_all),
        cluster = factor( (cluster_id_all %% 7) + 1 )
      ),
    aes( x=Y1, y=Y2, label=College, color = cluster ),
    mapping=,
    show.legend = FALSE,
    size=2,
    check_overlap = TRUE
  ) +
  geom_text(
    aes( label = Coefficient ),
    color = 'black',
    size = 3,
    check_overlap = TRUE
  )  +
  geom_point(
    data = tsne_mat_hc_all %>%
      as_tibble() %>%
      mutate(
        College = names(cluster_id_all),
        cluster = factor( (cluster_id_all %% 7) + 1 ),
        cluster_shape = factor( (cluster_id_all %% 6) + 1 )
      ),
    aes(x=Y1,y=Y2, color = cluster, shape = cluster_shape ),
    show.legend = FALSE,
    alpha=0.3
  ) +
  ggtitle( "t-SNE Biplot" ) +
  theme( text = element_text( face = 'bold' ) ) #+
  #scale_y_continuous(limits = c(y2_min,5))
} %>%
print()
```

**t−SNE Biplot**

## Conclusions

We do find some structure in the plot. And, the rotation of the axis to put Harvard University at the top-center helps us to interpret the axes and give meaning to that structure.

**Notable clusters**

Clusters are colored with repeating colors and marked with repeating symbols, reflecting a limit of **ggplot2**. But each cluster should have an unique color-symbol combination.

- **Ivy Leagues** (orange pluses @ {0,3})
- **Big Publics (Land-Grants)** (lavender asterisks @ {-2,1.5})
    - **Cal-State System**
- **HCBU: haves & have-nots** (gold triangles @ {-2.5,-2})
- **Techies** (Orange dots @ {3.3,-1.7})
- **Artsies** (Green dpts @ {2,-3.5})

- **Religious** (green boxed X's @ {2.5,0})
- **Back-to-Schools & Late-Bloomers** (green pluses @ {-1,-1}) (>24yrs-old, veterans)

## Interpretation of Quadrants

The combination of cluster locations and Bayes factors feature rays helps us assign meaning to each quadrant of the biplot.

### Elite private & top-academic public, wealthy & smart

The vertical `Y2` axis is now almost perfectly aligned with the ray `pgt110K`, which is the ($\log_{10}$) Bayes factor capturing the prevalance of students from families with annual incomes greater than \$110,000. All the Ivy League, "Ivy wannabes", and top-academic public universities (e.g., Cal-Berkeley, U. Michigan-Ann Arbor) are aligned along the positive vertical axis. That axis is almost perfectly countered by the downward-pointed ray `SATle800`, which is the Bayes factor capturing the prevalance of students with combined Verbal & Math SAT scores less than or equal to 800, i.e., the lowest tail of SAT scores.

### Breadth versus specialization

The horizontal `Y1` axis isn't so readily interpretable. However, we see the ray `discBreadth`, whcih is the feature capturing the entropy (variety) in academic disciplines in which degrees are offered from the college, is pointing into the upper-left corner of the plot. So colleges aligned along this ray in the upper-right quadrant are the big public state universities that offer a broad range of degrees. On the other hand, the narrowly, highly specialized colleges appear in the lower-right quadrant of the plot.

### Pell grants & high 3-yr credit default rates

The colleges in the lower-left quadrant are the colleges most strongly aligned with rays `pellever`, which captures prevalence of students having ever received a federal Pell grant, and `CDR3est`, which captures prevalence of students defaulting on student loans within 3 years of leaving the college.

### More privates, but less elite

The upper-right quadrant is aligned with `SAT1400` (highest SAT students), `fsend5` (applied to many colleges), and `pgt48Kle75K` (mid-income families).

## Summary

This was an exploratory analysis investigating structure in the U.S. Dept. of Education College Scorecard dataset.