# Words Apart
Mapping the Hidden Structure of Writing Style Through Multivariate Analysis

Sam Lee        Brigg Trendler        Daisy Harris

Department of Statistics, Brigham Young University

November 2025

## 1    Introduction

Writing styles exhibit systematic patterns in pronoun use, temporal framing, reasoning language, and narrative structure, and these patterns can reveal natural groupings of texts beyond their content. Using 1,000 texts scored on eighteen stylistic features (Collins 2003), we investigate whether such features form coherent style-based clusters and how these clusters relate to conventional genre labels. Our analysis proceeds in three steps. We first apply principal component analysis (PCA) to obtain an orthogonal representation of the stylistic variables and generate candidate clusterings for $k = 3, \ldots, 7$ using $k$–means, Gaussian mixture models, and Ward's method. We then evaluate each clustering using MANOVA separation, bootstrap stability, supervised misclassification, and prediction strength to assess distinctness and reproducibility. Finally, we interpret the selected clustering through canonical discriminant analysis and compare the resulting style groups to the original genre categories.

## 2    Cluster Analysis

We begin by conducting a principal component analysis (PCA) on the 18 linguistic features to obtain an orthogonal basis for clustering. Because the features are measured on heterogeneous scales, PCA is performed on the standardized variables which we define as $\widetilde{\boldsymbol{X}}$. Equivalently, our PCA uses the correlation matrix, rather than the raw covariance matrix[1], ensuring that each feature contributes equally. This is standard practice when variables differ substantially in scale. PCA is obtained by performing a spectral decomposition on the correlation matrix, $\boldsymbol{S}$, where $\boldsymbol{S} = \boldsymbol{C}\boldsymbol{D}\boldsymbol{C}'$. By definition, $\boldsymbol{C} = (\boldsymbol{c}_1, \ldots, \boldsymbol{c}_{18})$ is an orthonormal matrix of eigenvectors of $\boldsymbol{S}$ and $\boldsymbol{D} = \mathrm{diag}(d_1, \ldots, d_{18})$ is the diagonal matrix of $\boldsymbol{S}$ of eigenvalues ordered $d_1 \geq \cdots \geq d_{18} \geq 0$. The principal component score matrix is then, $\boldsymbol{Z} = \widetilde{\boldsymbol{X}}\boldsymbol{C}$, so that $\boldsymbol{Z}$ provides an orthogonal reparameterization of the standardized feature space. The $d$th column of $\boldsymbol{Z}$ is the projection of the data onto eigenvector $\boldsymbol{c}_d$ and has variance $d_d$. Figure 1 displays the cumulative proportion of variance explained by the eigenvalues, showing that the first thirteen components capture approximately 90% of the total variation[2]. Retaining only the principal components that capture the dominant patterns of covariance ensures that clustering operates on the subspace that yields the most separation in $\boldsymbol{X}$, rather than on directions with vanishing variance that inflate Euclidean distances[3]. All subsequent clustering analyses are therefore performed in this reduced PCA space[4].

To estimate a partition of the data for a given $k$ number of clusters that yields optimal cluster separation, we adopt the Calinski–Harabasz (CH) criterion, originally proposed by Caliński and Harabasz (1974), which formalizes cluster quality as a ratio of between-cluster to within-cluster dispersion. For an observation vector $\boldsymbol{x}_i$ with cluster labels $c(i) \in \{1, \ldots, k\}$, cluster means $\boldsymbol{\mu}_j$, cluster sizes $n_j$, and global mean $\boldsymbol{\mu}$, the Calinski–Harabasz index is defined as,

$$\mathrm{CH}(k) = \frac{B_k/(k-1)}{W_k/(n-k)}, \qquad B_k = \sum_{j=1}^{k} n_j \, \|\boldsymbol{\mu}_j - \boldsymbol{\mu}\|^2, \qquad W_k = \sum_{j=1}^{k} \sum_{i \in j} \|\boldsymbol{x}_i - \boldsymbol{\mu}_j\|^2. \qquad (1)$$

---

[1] Notationally, $\boldsymbol{S} = \mathrm{cor}(\boldsymbol{X}) = \frac{1}{n-1} \widetilde{\boldsymbol{X}}'\widetilde{\boldsymbol{X}}$, where where $\widetilde{X}_{ij} = \frac{X_{ij} - \bar{X}_{\cdot j}}{s_j}$, where $\bar{X}_{\cdot j}$ and $s_j$ are the sample mean and standard deviation of the $j$th variable.

[2] We also look at a scree plot for choosing the number of components to include. We find that five eigenvalues exceed one. Alternative analyses may wish to reduce the covariate space down to five dimensions for this reason. However, for the purposes of our analysis, we keep thirteen dimensions.

[3] Moreover, as shown by Ding and He (2004), the subspace spanned by the leading principal components is the one most aligned with the optimal $k$-means partition.

[4] That is, we consider the observation vector $\boldsymbol{x}_i \in \mathbb{R}^P$, where $P = 13$. We use the 13 eigenvectors corresponding to the 13 largest eigen values in $\boldsymbol{D}$. We then construct $\boldsymbol{Z} = \widetilde{\boldsymbol{X}}\boldsymbol{C}^{\star}$, where $\boldsymbol{C}^{\star}$ is the corresponding $n \times 13$ matrix of orthonormalized eigenvectors.

Intuitively, CH increases when clusters exhibit high between-cluster separation and low within-cluster dispersion. Because CH is a purely geometric criterion, it provides a model-agnostic benchmark across clustering procedures.

In this study, we propose three unsupervised clustering algorithms. We first implement $k$–means clustering (Hartigan and Wong 1979). This algorithm partitions the data into $k$ clusters by minimizing the within-cluster sum of squared Euclidean distances. For cluster assignments $c(i)$ and centroids $\boldsymbol{\mu}_j$, the objective function can be defined as,

$$\min_{\{c(i),\boldsymbol{\mu}_j\}} \sum_{j=1}^{k} \sum_{i:c(i)=j} \|\boldsymbol{x}_i - \boldsymbol{\mu}_j\|^2, \tag{2}$$

which is optimized iteratively by assigning each point to the nearest centroid and then repeatedly recomputing centroids as sample means and again until convergence. Because the $k$–means objective is nonconvex, we employ ten random initializations and record the best solution obtained.

Second, we fit a Gaussian mixture model (GMM) in which, motivated by our use of principal components, each component has a diagonal covariance matrix. This implies that, within each cluster, the coordinates are conditionally independent (under the Gaussian assumption) but may have cluster-specific variances. For the latent cluster label $\zeta_i \in \{1,\ldots,k\}$, we assume,

$$\boldsymbol{x}_i \mid \zeta_i = j \sim \mathcal{N}_P(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), \qquad \boldsymbol{\Sigma}_j = \mathrm{diag}(\sigma_{j1}^2, \ldots, \sigma_{jP}^2), \tag{3}$$

so that the component density factorizes as,

$$p(\boldsymbol{x}_i \mid \zeta_i = j) = \prod_{p=1}^{P} \mathcal{N}(x_{ip} \mid \mu_{jp}, \sigma_{jp}^2). \tag{4}$$

Defining $\pi_j$ as the mixing proportion for component $j$, with $\sum_{j=1}^{k} \pi_j = 1$. The GMM is fit through the expectation–maximization (EM) algorithm. In the E-step, we compute the posterior cluster probabilities

$$q_{ij} = \mathrm{Pr}(\zeta_i = j \mid \boldsymbol{x}_i) = \pi_j \prod_{p=1}^{P} \mathcal{N}(x_{ip} \mid \mu_{jp}, \sigma_{jp}^2) / \left( \sum_{\ell=1}^{k} \pi_\ell \prod_{p=1}^{P} \mathcal{N}(x_{ip} \mid \mu_{\ell p}, \sigma_{\ell p}^2) \right). \tag{5}$$

In the M-step, we update the parameters using these posterior weights:

$$\pi_j = \frac{1}{n} \sum_{i=1}^{n} q_{ij}, \qquad \mu_{jp} = \frac{\sum_{i=1}^{n} q_{ij} x_{ip}}{\sum_{i=1}^{n} q_{ij}}, \qquad \sigma_{jp}^2 = \frac{\sum_{i=1}^{n} q_{ij}\left(x_{ip} - \mu_{jp}\right)^2}{\sum_{i=1}^{n} q_{ij}}. \tag{6}$$

Final cluster assignments are obtained by the maximum a posteriori rule, $\widehat{\zeta}_i = \arg\max_{1 \leq j \leq k} q_{ij}$.

Finally, we apply Ward's hierarchical agglomerative method (Ward 1963), which constructs a dendrogram by successively merging the pair of clusters that yields the smallest increase in total within-cluster variance (a reparameterization of Eq. (1)), ensuring that each agglomeration step yields the least distortion to the overall within-cluster dispersion. This method differs fundamentally from the previous two in that it does not optimize over cluster assignments globally but instead builds a full hierarchy of nested partitions.

In addition to assessing separation quality, we also evaluate cluster stability. Namely, we estimate the extent to which the estimated partition is reproducible under perturbations of the dataset. Following a nonparametric bootstrap framework[5], we quantify stability using bootstrap resampling of the Jaccard coefficient. Formally, we define $\{\mathcal{G}_1, \ldots, \mathcal{G}_k\}$ as the $k$ clusters obtained from the original data, where $j \in \{1,\ldots,k\}$ indexes the original clusters. For a given bootstrap resample, we define $\mathcal{S} \subseteq \{1,\ldots,n\}$ as the set of original observations that appear at least once in the resampled dataset (the overlapping set), and let $\{\mathcal{H}_1^*, \ldots, \mathcal{H}_k^*\}$ denote the clusters obtained by reclustering the bootstrap sample, with $s \in \{1,\ldots,k\}$ indexing the bootstrap clusters. For each original cluster $\mathcal{G}_j$, we compute its agreement with each bootstrap cluster $\mathcal{H}_s^*$ on the restricted set $\mathcal{S}$ using the Jaccard index

$$J(\mathcal{G}_j, \mathcal{H}_s^*) = \frac{|(\mathcal{G}_j \cap \mathcal{S}) \cap (\mathcal{H}_s^* \cap \mathcal{S})|}{|(\mathcal{G}_j \cap \mathcal{S}) \cup (\mathcal{H}_s^* \cap \mathcal{S})|}. \tag{7}$$

---

[5]See Suzuki and Shimodaira (2006) and Hennig (2007) for formal treatments of cluster stability and uncertainty.

Table 1: Clustering quality (Calinski–Harabasz) and bootstrap Jaccard stability by method and number of clusters

| | $k = 3$ | | $k = 4$ | | $k = 5$ | | $k = 6$ | | $k = 7$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| Method | CH | $\bar{\gamma}$ (MCSE) | CH | $\bar{\gamma}$ (MCSE) | CH | $\bar{\gamma}$ (MCSE) | CH | $\bar{\gamma}$ (MCSE) | CH | $\bar{\gamma}$ (MCSE) |
| $k$-means | 187.69 | 0.906 (0.0024) | 151.85 | 0.688 (0.0053) | 134.31 | 0.647 (0.0049) | 124.27 | 0.678 (0.0050) | 115.27 | 0.639 (0.0047) |
| GMM | 65.27 | 0.560 (0.0052) | 65.08 | 0.588 (0.0066) | 87.18 | 0.622 (0.0043) | 85.53 | 0.548 (0.0041) | 73.02 | 0.465 (0.0044) |
| Ward | 157.72 | 0.608 (0.0031) | 128.49 | 0.509 (0.0026) | 112.32 | 0.477 (0.0040) | 103.56 | 0.487 (0.0042) | 94.68 | 0.459 (0.0033) |

Notes: CH denotes the Calinski–Harabasz index computed on the first thirteen principal components. $\bar{\gamma}$ is the mean bootstrap Jaccard stability coefficient, with Monte Carlo standard errors (MCSE) in parentheses, based on 500 nonparametric bootstrap resamples.

The bootstrap stability of $\mathcal{G}_j$ is then defined as the bootstrapped cluster most commensurate with the original clustering on the full data set for that bootstrap iteration: $\gamma_j = \max_{1 \leq s \leq k} J(\mathcal{G}_j, \mathcal{H}_s^*)$. We report the average stability together with Monte Carlo standard errors computed across 500 bootstrap replications in Table 1. Intuitively, larger values of $\gamma$ indicate clusters that replicate consistently under resampling, whereas smaller values signal unstable clusters that tend to merge, split, or disappear across bootstrap samples. $k$–means consistently attains the highest CH values[6] and the strongest bootstrap stability. Ward's method yields moderate separation but substantially weaker stability across all $k$, while the Gaussian mixture model shows lower CH values but competitive stability for $k = 4$ and $k = 5$. However, since $k$–means attains higher CH scores and more stability across all values of $k$, we consider clusters formed by this algorithm to be our "best" sets of clusters and will be used for the remainder of the analysis.
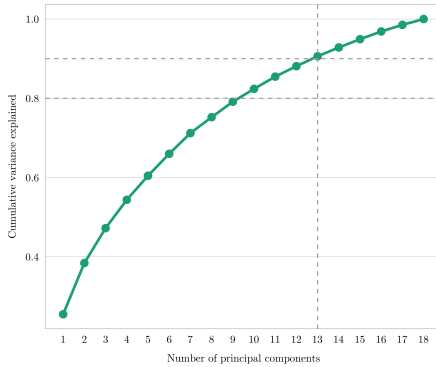


Figure 1: Cumulative variance explained by principal components used for clustering. The vertical dashed line at $k = 13$ denotes the point where $> 90\%$ of the total variation in $X$ is accounted for.
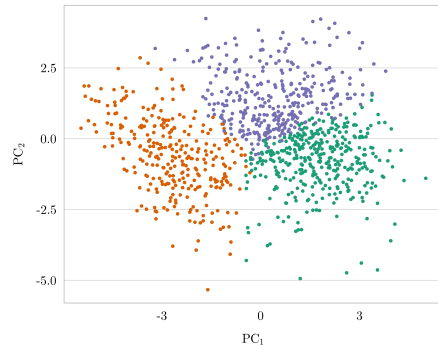


Figure 2: First two principal components plotted against the $k$-means clustering algorithm with $k = 3$ clusters (shown in different colors).

## 2.1 Optimal Number of Clusters

To assess how many clusters are supported by the data, we evaluated multivariate separation, supervised classification performance, and resampling-based stability for $k = 3, \ldots, 7$. We present these results in Table 2. We first conducted MANOVA tests to examine whether the multivariate mean vectors of the 13 retained principal components differed across the $k$–means partitions. For each $k$, we computed four standard test statistics: Wilks' $\Lambda$, Pillai's trace, the Hotelling–Lawley trace, and Roy's greatest root.

All tests were highly significant ($p < 0.0001$), indicating that at least one of the cluster centroids differ across multiple linear combinations of the principal components[7]. The corresponding canonical eigenstructure shows that for smaller $k$, most separation is concentrated in a single canonical dimension, whereas for larger $k$, the separation is distributed across several directions. These findings confirm that the PCA features capture meaningful between-group variation and that the resulting $k$–means partitions represent statistically distinct multivariate clusters. However, statistical significance alone does not determine an optimal $k$. We

---

[6]It should be noted that the particular values of the CH scores are generally irrelevant. Because the CH score is a function of total variance, the more principal components we use in our $\mathbf{Z}$ matrix, the smaller corresponding values of CH will be.

[7]Again, we stress by construction of the MANOVA tests, as the number of clusters (factors) grows, the probability that we find at least one significant different mean vector (even by chance) artificially increases.

Table 2: Separation metrics for $k$-means clusters on principal component scores

| | $k = 3$ | $k = 4$ | $k = 5$ | $k = 6$ | $k = 7$ |
|---|---|---|---|---|---|
| *Panel A: MANOVA* | | | | | |
| Wilks' $\Lambda$ | 0.110 | 0.056 | 0.032 | 0.016 | 0.011 |
| Pillai's trace | 1.288 | 1.749 | 2.067 | 2.548 | 2.873 |
| Largest canonical root $\lambda_1$ | 3.436 | 3.380 | 4.486 | 4.936 | 5.074 |
| Percent separation in dim. 1 (%) | 76.5 | 60.9 | 63.7 | 59.5 | 56.8 |
| Canonical corr. $\rho_1$ | 0.880 | 0.878 | 0.904 | 0.912 | 0.914 |
| *Panel B: Supervised classification* | | | | | |
| Misclass. rate (RF / MN / XGB) | 0.23 / 0.48 / 0.10 | 0.30 / 0.49 / 0.14 | 0.41 / 0.68 / 0.20 | 0.41 / 0.86 / 0.21 | 0.45 / 0.64 / 0.27 |
| Prediction strength $ps(k)$ | 0.91 | 0.72 | 0.64 | 0.47 | 0.46 |

Notes: Entries report separation statistics for $k$-means clustering on the first thirteen principal components for $k = 3, \ldots, 7$. Panel A reports MANOVA-based measures of between-cluster separation. Panel B reports out-of-sample supervised classification performance for random forest (RF), multinomial logit (MN), and XGBoost (XGB) using 10-fold cross validation; entries in the misclassification row are ordered RF/MN/XGB for each $k$. All MANOVA tests reject equal mean vectors at $p < 0.0001$ (Wilks' $\Lambda$ $F$-tests not shown).

therefore complemented MANOVA with two external validation strategies: supervised misclassification and prediction strength. For each $k$, we treated the $k$–means labels as pseudo-outcomes and estimated out-of-sample misclassification rates using 10-fold cross-validation. Each fold held out one observation, trained three classifiers on the remaining data (random forest, multinomial logit, and XGBoost), and predicted the held-out label. This yields a measure of how well cluster structure generalizes. All three classifiers exhibited the same pattern: misclassification was lowest at $k = 3$ and increased monotonically with $k$. This reflects a well-known drawback of this approach, namely that misclassification necessarily increases as the number of classes grows, producing an implicit penalty on larger $k$.

To avoid this inflation, we also implemented the prediction strength criterion introduced by Tibshirani and Walther (2005), which evaluates how consistently pairs of points that co-cluster in the training data $(X_{\mathrm{tr}})$ also co-cluster in the test data $X_{\mathrm{te}}$. Formally,

$$\mathrm{ps}(k) = \min_{1 \leq j \leq k} \frac{1}{n_{k_j}(n_{k_j} - 1)} \sum_{i \neq i' \in A_{k_j}} D[C(X_{\mathrm{tr}}, k), X_{\mathrm{te}}]_{ii'}, \tag{8}$$

where $A_{k_j}$ denotes the set of points assigned to cluster $j$ in the training sample, and $D[\cdot]_{ii'}$ is an indicator of whether $i$ and $i'$ are assigned to the same cluster in both train and test partitions. Higher prediction strength indicates more reproducible clustering. Consistent with the misclassification results, $\mathrm{ps}(k)$ was highest at $k = 3$ and declined steadily thereafter.

Finally, following Ben-Hur et al. (2002), nonparametric bootstrap stability scores reported earlier (Table 1) can provide an additional assessment of robustness. These scores measure how reliably each cluster reappears under bootstrap resampling and re-clustering. Stability was strongest for $k = 3$ and weakened for larger $k$, aligning with both the supervised and resampling-based evidence. All criteria identify $k = 3$ as the most coherent and reproducible cluster structure. This value of $k$ therefore serves as the basis for our subsequent analysis.

## 3 Latent Cluster Evaluation

For $k = 3$, the MANOVA produced two nonzero characteristic roots (76.52% and 23.48%), the maximum possible for three groups[8]. The cluster centroids differ along a two–dimensional canonical subspace within the 13-dimensional PCA space, confirming that the three groups are not collinear.

To explore how these clusters differ in practice, we first examined the distribution of genre labels within each group (Figure 3). Cluster 1 is dominated by fiction (86%), forming a stylistically distinct group. Clusters 2 and 3 exhibit more mixed genre compositions: Cluster 2 contains more reporting, skills-and-hobbies, and official communications, while Cluster 3 includes more editorial, biography, and religious writing. To characterize the dimensions underlying this separation, we conducted a linear discriminant analysis (LDA) (Table 3) using the standardized linguistic variables. LDA identifies linear combinations of features that maximize separation among the three clusters. The first discriminant function captures roughly 75% of the discriminant trace and loads most strongly on first–person language, past events, thinking back,

---

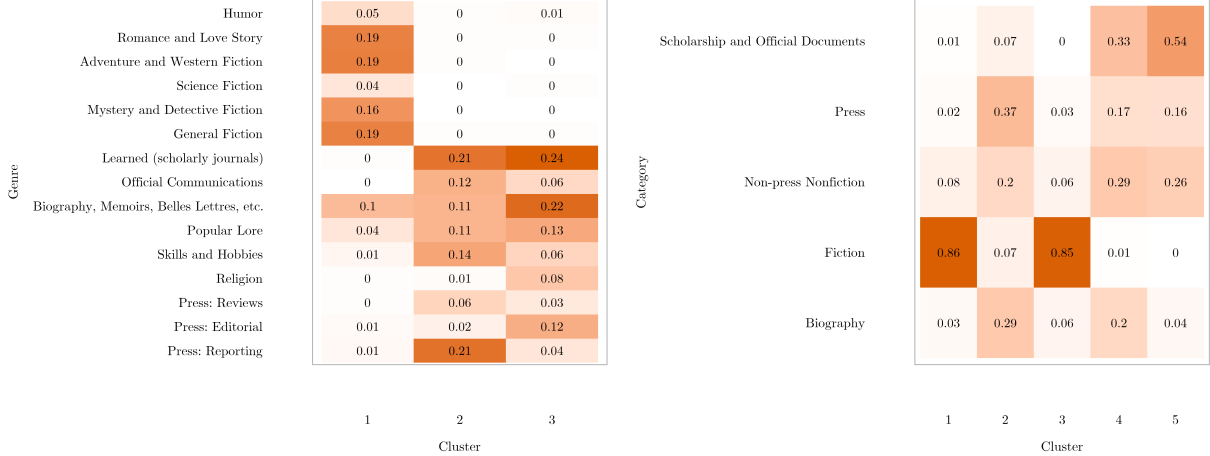[8]The essential dimensionality is computed as $s = \min(p, k - 1) = \min(13, 3 - 1) = 2$.

Figure 3: Heatmaps of the genres represented within three clusters (left) and the categories represented within five clusters (right). The intensity of the orange shade represents the proportion of observations within the column cluster belonging to the row (genre).

spatial framing, and linear guidance. These features are most prominent in Cluster 1, consistent with a more personal, narrative, and fiction-oriented style. The second discriminant function emphasizes reasoning, social ties, inner thinking, and negative thinking, with Cluster 3 showing higher scores than Cluster 2. This separation shows the distinction between reflective, introspective writing and more formal, action-oriented prose. We also examined the five-cluster solution and compared it to the broader "super-genre" categories

Table 3: LDA Loadings

| LD | FirstPerson | InnerThinking | ThinkPositive | ThinkNegative | ThinkAhead | ThinkBack | Reasoning | Share_SocTies | Direct_Activity |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| LD1 | -0.655 | 0.060 | 0.063 | 0.094 | -0.035 | -0.360 | 0.013 | 0.199 | 0.039 |
| LD2 | 0.123 | 0.492 | 0.002 | 0.398 | 0.082 | 0.031 | 0.553 | 0.622 | 0.093 |

| LD | Interacting | Notifying | LinearGuidance | WordPicture | SpaceInterval | Motion | PastEvents | TimeInterval | ShiftingEvents |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| LD1 | -0.324 | 0.137 | -0.485 | -0.088 | -0.411 | -0.122 | -0.551 | 0.046 | -0.013 |
| LD2 | 0.165 | 0.108 | 0.328 | 0.003 | 0.266 | -0.148 | 0.135 | 0.085 | -0.158 |

Notes: Entries report standardized linear discriminant loadings for LD1 and LD2 based on the 18 linguistic features. Larger absolute values indicate stronger contributions to the discriminant directions; signs indicate direction of association. LD1 accounts for roughly 75% of the total discriminant trace, with LD2 capturing the remaining orthogonal dimension.

(Figure 3, right panel). Two of the clusters are almost entirely fiction, while the remaining clusters mix genres in patterns that do not align with standard category boundaries. This inconsistency suggests that genre labels do not capture the dominant stylistic variation in the data and that clustering based on rhetorical feature may provided a more coherent representation of writing style.

## 4   Conclusion

In summary, our unsupervised clustering on stylistic features reveals coherent patterns in how texts are written that are only partially captured by conventional genre labels. Writing style, summarized through eighteen linguistic features, exhibits a clear and stable low-dimensional structure. After projecting the data onto principal components and evaluating $k = 3, \ldots, 7$ across Calinski–Harabasz indices, MANOVA separation, supervised misclassification, prediction strength, and bootstrap Jaccard stability, we consistently selected $k = 3$ as the most reproducible partition.

Canonical discriminant analysis further showed that these three clusters differ primarily along two stylistic dimensions. The first corresponds to a personal, narrative style marked by first-person language, past-oriented framing, and spatial or event description, aligning with fiction-heavy texts. The second separates reflective, reasoning-focused writing from more formal, action- or report-oriented prose. Genre labels do not consistently track these stylistic dimensions, particularly in higher-$k$ solutions where clusters mix super-genres. This suggests that stylistic features provide a more coherent representation of writing variation than genre alone.

5

# References

Ben-Hur, A. et al. (2002). "A Stability Based Method for Discovering Structure in Clustered Data". In: *Proceedings of the Pacific Symposium on Biocomputing*, pp. 6–17.

Caliński, T. and J Harabasz (1974). "A dendrite method for cluster analysis". In: *Communications in Statistics* 3.1, pp. 1–27. DOI: 10.1080/03610927408827101. eprint: https://www.tandfonline.com/doi/pdf/10.1080/03610927408827101. URL: https://www.tandfonline.com/doi/abs/10.1080/03610927408827101.

Collins, Jeff (2003). "Variations in Written English: Characterizing the Rhetorical Language Choices in the Brown Corpus of Texts". Department of English, Dietrich College of Humanities and Social Sciences. PhD dissertation. Carnegie Mellon University. URL: https://www.cmu.edu/dietrich/english/academic-programs/phd-rhetoric/phd-rhetoric-dissertation_abstracts/Jeff%20Collins.html.

Ding, Chris and Xiaofeng He (2004). "Cluster Structure of K-means Clustering via Principal Component Analysis". In: ed. by Honghua Dai et al., pp. 414–418.

Hartigan, J. A. and M. A. Wong (1979). "Algorithm AS 136: A K-Means Clustering Algorithm". In: *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 28.1, pp. 100–108. ISSN: 00359254, 14679876. URL: http://www.jstor.org/stable/2346830 (visited on 11/17/2025).

Hennig, Christian (2007). "Cluster-wise assessment of cluster stability". In: *Computational Statistics & Data Analysis* 52.1, pp. 258–271. ISSN: 0167-9473. DOI: https://doi.org/10.1016/j.csda.2006.11.025. URL: https://www.sciencedirect.com/science/article/pii/S0167947306004622.

Suzuki, R. and H. Shimodaira (2006). "Pvclust: An R Package for Assessing the Uncertainty in Hierarchical Clustering". In: *Bioinformatics* 22.12, pp. 1540–1542. DOI: 10.1093/bioinformatics/btl117.

Tibshirani, Robert and Guenther Walther (2005). "Cluster Validation by Prediction Strength". In: *Journal of Computational and Graphical Statistics* 14.3, pp. 511–528. ISSN: 10618600. URL: http://www.jstor.org/stable/27594130 (visited on 11/17/2025).

Ward, Joe H. (1963). "Hierarchical Grouping to Optimize an Objective Function". In: *Journal of the American Statistical Association* 58.301, pp. 236–244. ISSN: 01621459, 1537274X. URL: http://www.jstor.org/stable/2282967 (visited on 11/17/2025).

# 5  Appendix - R Code

## 5.1  Data Setup

```r
source("pca.R")

# Read in the data
# Credit: Patric Platts
collins <- read.table(
  "https://tofu.byu.edu/docs/files/stat666/datasets/collins.txt",
  header = TRUE,
  sep = "\t",
  skip = 10, # skip the 10 lines of description before the header
  check.names = FALSE
)

X <- collins[, -c(1, 20:24)]

#PCA on the correlation matrix (this is done by scaling)
X.pca <- pca(X, center = T, scale. = T)
```

## 5.2  PCA

```r
pca <- function(X, k = NULL, center = TRUE, scale. = FALSE) {
  X <- as.matrix(X)
  n <- nrow(X)
  p <- ncol(X)

  if (is.null(k)) {
    k = min(n - 1, p)
  }

  if (k > min(n - 1L, p)) {
    stop("k must be <= min(n - 1, p).")
  }

  ## center/scale columns
  Xc <- if (center || scale.) {
    scale(X, center = center, scale = scale.)
  } else {
    X
  }

  ## covariance matrix and spectral decomposition
  S <- crossprod(Xc) / (n - 1L) # p x p covariance
  eig <- eigen(S, symmetric = TRUE)

  ## loadings (eigenvectors) and scores (principal components)
  V <- eig$vectors[, 1:k, drop = FALSE] # p x k
  scores <- Xc %*% V # n x k

  colnames(scores) <- paste0("PC", seq_len(k))
  attr(scores, "eigenvalues") <- eig$values[1:k]

  scores # n x k matrix of principal components
}

pca_elbow_plot <- function(
  X,
```

```r
  k.star = 13,
  max_k = NULL,
  center = TRUE,
  scale. = FALSE,
  target_lines = c(0.80, 0.90)
) {
  X <- as.matrix(X)
  n <- nrow(X)
  p <- ncol(X)

  k_full <- min(n - 1L, p)
  if (is.null(max_k)) {
    max_k <- k_full
  }
  max_k <- min(max_k, k_full)

  pcs <- pca(X, k = k_full, center = center, scale. = scale.)
  eigvals <- attr(pcs, "eigenvalues")

  var_prop <- eigvals / sum(eigvals)
  cum_var_prop <- cumsum(var_prop)

  df <- data.frame(
    k = 1:max_k,
    cum_var = cum_var_prop[1:max_k]
  )

  ggplot(df, aes(k, cum_var, color = "PC")) +
    geom_line(linewidth = 1.2) +
    geom_point(size = 3) +
    geom_hline(
      data = data.frame(y = target_lines),
      aes(yintercept = y),
      linetype = "dashed",
      linewidth = 0.5,
      color = "grey60"
    ) +
    geom_vline(
      aes(xintercept = k.star),
      linetype = "dashed",
      linewidth = 0.5,
      color = "grey60"
    )+
    scale_color_paper(guide = "none") +
    scale_x_continuous(
      breaks = seq_len(max_k),
      minor_breaks = NULL
    ) +
    labs(
      x = "Number of principal components",
      y = "Cumulative variance explained",
      # title = "Elbow Plot for Principal Components",
      # subtitle = "Cumulative proportion of variance explained"
    ) +
    theme_paper(text_size = 64)
}
```

### 5.3 Calinkski-Harabasz Scores

```r
#Computes the Calinski-Harabasz Score
ch.score <- function(X, G) {
  #X is nxk data matrix
  #G is a vector of group indices corresponding to which groups
  #each row belongs in

  K <- unique(G)

  n <- nrow(X)
  n.j <- tapply(as.data.frame(X), G, nrow)
  mu.j <- tapply(as.data.frame(X), G, colMeans)
  mu <- colMeans(X)

  B.k <- sapply(K, function(k) {
    n.j[k] * crossprod(mu - mu.j[[k]])
  }) |>
    sum()

  W.k <- sapply(K, function(k) {
    m.j <- mu.j[[k]]
    X.g <- X[G == k, ]
    s.g <- apply(X.g, 1, function(x) {
      crossprod(x - m.j)
    }) |>
      sum()
  }) |>
    sum()

  ch.k <- (B.k / (length(K) - 1)) / (W.k / (n - length(K)))
  ch.k
}
```

## 5.4 $K-$means Clustering

```r
source("data.R")
source("calinski.R")

library(boot)

kmeans_clusters <- function(X, k, Bstart = 50, seed = NULL) {
  if (!is.null(seed)) {
    set.seed(seed)
  }
  fit <- kmeans(X, centers = k, nstart = Bstart, iter.max = 100)
  fit$cluster
}

kmeans_ch <- function(X, k.grid = 3:7, Bstart = 50, seed = NULL) {
  if (!is.null(seed)) {
    set.seed(seed)
  }

  ch_vec <- sapply(k.grid, function(k) {
    g <- kmeans_clusters(X, k, Bstart = Bstart)
    ch.score(X, g)
  })

  data.frame(k = k.grid, CH = ch_vec)
```

```r
}

kmeans_jaccard_boot <- function(
  X,
  k,
  R = 200,
  Bstart = 50,
  seed = NULL
) {
  if (!is.null(seed)) {
    set.seed(seed)
  }

  n <- nrow(X)
  cl0 <- kmeans_clusters(X, k, Bstart = Bstart)
  boot_fun <- function(data, indices) {
    idx <- indices
    Xb <- data[idx, , drop = FALSE]

    cl_boot <- kmeans_clusters(Xb, k, Bstart = Bstart)

    idx_star <- sort(unique(idx))
    cl0_star <- cl0[idx_star]

    first_pos <- tapply(seq_along(idx), idx, `[`, 1)
    pos_star <- as.integer(first_pos[as.character(idx_star)])
    clb_star <- cl_boot[pos_star]

    gamma_C <- numeric(k)

    for (c in seq_len(k)) {
      C_pts <- idx_star[cl0_star == c]
      if (length(C_pts) == 0L) {
        gamma_C[c] <- 0
      } else {
        jacc_c <- sapply(seq_len(k), function(d) {
          D_pts <- idx_star[clb_star == d]
          if (length(D_pts) == 0L) {
            return(0)
          }
          inter <- length(intersect(C_pts, D_pts))
          union <- length(union(C_pts, D_pts))
          if (union == 0L) 0 else inter / union
        })
        gamma_C[c] <- max(jacc_c)
      }
    }

    gamma_C
  }

  b <- boot(
    data = X,
    statistic = boot_fun,
    R = R,
    parallel = "multicore"
  )

  gamma_bar <- colMeans(b$t)
```

```
84     gamma_mcse <- apply(b$t, 2, sd) / sqrt(R)

85

86     mean_gamma <- mean(gamma_bar)
87     mean_gamma_mcse <- sd(rowMeans(b$t)) / sqrt(R)

88

89     list(
90       gamma_bar = gamma_bar,
91       gamma_mcse = gamma_mcse,
92       mean_gamma = mean_gamma,
93       mean_gamma_mcse = mean_gamma_mcse,
94       boot_mat = b$t,
95       boot_obj = b
96     )
97   }

98

99   k.grid <- 3:7

100

101  kmeans.summary <- function(X, iters = 500, seed = 666) {
102    kmeans_ch_summary <- kmeans_ch(X, k.grid, Bstart = 10, seed = 666)
103    print(kmeans_ch_summary)

104

105    kmeans_stab <- lapply(k.grid, function(k) {
106      kmeans_jaccard_boot(
107        X,
108        k = k,
109        R = iters,
110        Bstart = 10,
111        seed = seed
112      )
113    })
114    names(kmeans_stab) <- paste0("k", k.grid)
115    print(kmeans_stab)

116

117    kmeans_stab_summary <- data.frame(
118      k = k.grid,
119      mean_gamma = sapply(kmeans_stab, \(z) z$mean_gamma),
120      mean_gamma_mcse = sapply(kmeans_stab, \(z) z$mean_gamma_mcse)
121    )

122

123    print(kmeans_stab_summary)
124  }

125

126  kmeans.summary(X = X.pca[, 1:13], iters = 500)

127

128  # k        CH
129  # 1 3 187.6944
130  # 2 4 151.8520
131  # 3 5 134.3061
132  # 4 6 124.2679
133  # 5 7 115.2718
134  #    k mean_gamma mean_gamma_mcse
135  # k3 3  0.9059998     0.002375117
136  # k4 4  0.6883734     0.005286859
137  # k5 5  0.6472300     0.004920278
138  # k6 6  0.6775243     0.005006445
139  # k7 7  0.6394362     0.004726481

140

141  get.data <- function(X, k, seed) {
142    as.data.frame(cbind(
```

```
143      X,
144      "label" = kmeans_clusters(X, k, Bstart = 10, seed = seed)
145    ))
146  }
147
148  create.data <- function() {
149    for (k in k.grid) {
150      dat.k <- get.data(X.pca[, 1:13], k, seed = 666)
151      readr::write_csv(dat.k, stringr::str_c("Data/dat-", k, ".csv"))
152    }
153  }
```

## 5.5  Gaussian Mixture Model Clustering

```
1   source("data.R")
2   source("calinski.R")
3
4   library(mclust)
5
6   gmm_clusters <- function(X, k) {
7     fit <- Mclust(X, G = k, modelNames = "VVI")
8     fit$classification
9   }
10
11  gmm_ch <- function(X, k.grid = 3:7, seed = NULL) {
12    if (!is.null(seed)) {
13      set.seed(seed)
14    }
15
16    ch_vec <- sapply(k.grid, function(k) {
17      g <- gmm_clusters(X, k) #sample(1:k, nrow(X), replace = T)
18      ch.score(X, g)
19    })
20
21    data.frame(k = k.grid, CH = ch_vec)
22  }
23
24  gmm_jaccard_boot <- function(X, k, R = 200, seed = NULL) {
25    if (!is.null(seed)) {
26      set.seed(seed)
27    }
28
29    n <- nrow(X)
30    cl0 <- gmm_clusters(X, k)
31
32    boot_fun <- function(data, indices) {
33      idx <- indices
34      Xb <- data[idx, , drop = FALSE]
35
36      cl_boot <- gmm_clusters(Xb, k)
37
38      idx_star <- sort(unique(idx))
39      cl0_star <- cl0[idx_star]
40
41      first_pos <- tapply(seq_along(idx), idx, '[', 1)
42      pos_star <- as.integer(first_pos[as.character(idx_star)])
43      clb_star <- cl_boot[pos_star]
44
```

```r
      gamma_C <- numeric(k)

      for (c in seq_len(k)) {
        C_pts <- idx_star[cl0_star == c]
        if (length(C_pts) == 0L) {
          gamma_C[c] <- 0
        } else {
          jacc_c <- sapply(seq_len(k), function(d) {
            D_pts <- idx_star[clb_star == d]
            if (length(D_pts) == 0L) {
              return(0)
            }
            inter <- length(intersect(C_pts, D_pts))
            union <- length(union(C_pts, D_pts))
            if (union == 0L) 0 else inter / union
          })
          gamma_C[c] <- max(jacc_c)
        }
      }

      gamma_C
    }

  b <- boot(
    data = X,
    statistic = boot_fun,
    R = R,
    parallel = "multicore"
  )

  gamma_bar <- colMeans(b$t)
  gamma_mcse <- apply(b$t, 2, sd) / sqrt(R)

  mean_gamma <- mean(gamma_bar)
  mean_gamma_mcse <- sd(rowMeans(b$t)) / sqrt(R)

  list(
    gamma_bar = gamma_bar,
    gamma_mcse = gamma_mcse,
    mean_gamma = mean_gamma,
    mean_gamma_mcse = mean_gamma_mcse,
    boot_mat = b$t,
    boot_obj = b
  )
}

k.grid <- 3:7

gmm.summary <- function(X, iters = 500, seed = 666) {
  gmm_ch_summary <- gmm_ch(X, k.grid, seed = seed)
  print(gmm_ch_summary)
  gmm_stab <- lapply(k.grid, function(k) {
    gmm_jaccard_boot(X, k = k, R = iters, seed = seed)
  })
  names(gmm_stab) <- paste0("k", k.grid)

  gmm_stab_summary <- data.frame(
    k = k.grid,
    mean_gamma = sapply(gmm_stab, \(z) z$mean_gamma),
```

```
104      mean_gamma_mcse = sapply(gmm_stab, \(z) z$mean_gamma_mcse)
105    )
106    print(gmm_stab_summary)
107 }
108
109 gmm.summary(X = X.pca[, 1:13], 500)
110 #  k       CH
111 # 1 3 65.27188
112 # 2 4 65.08038
113 # 3 5 87.18404
114 # 4 6 85.52907
115 # 5 7 73.02151
116 #   k mean_gamma mean_gamma_mcse
117 # k3 3  0.5604275     0.005184017
118 # k4 4  0.5876058     0.006601812
119 # k5 5  0.6217609     0.004314503
120 # k6 6  0.5483159     0.004132778
121 # k7 7  0.4651430     0.004417613
```

## 5.6   Ward Aglomerative Clustering

```r
1  source("data.R")
2  source("calinski.R")
3
4  library(boot)
5
6  ward_clusters <- function(X, k) {
7    hc <- hclust(dist(X), method = "ward.D2")
8    cutree(hc, k = k)
9  }
10
11 ward_ch <- function(X, k.grid = 3:7) {
12   hc <- hclust(dist(X), method = "ward.D2")
13   ch_vec <- sapply(k.grid, function(k) {
14     groups <- cutree(hc, k = k)
15     ch.score(X, groups)
16   })
17   data.frame(k = k.grid, CH = ch_vec)
18 }
19
20 ward_jaccard_boot <- function(X, k, R = 200, seed = NULL) {
21   if (!is.null(seed)) {
22     set.seed(seed)
23   }
24
25   n <- nrow(X)
26   cl0 <- ward_clusters(X, k)
27
28   boot_fun <- function(data, indices) {
29     idx <- indices
30     Xb <- data[idx, , drop = FALSE]
31
32     cl_boot <- ward_clusters(Xb, k)
33
34     #original points that appear in bootstrap sample
35     idx_star <- sort(unique(idx))
36
37     #labels of those points in original clustering
```

```r
    cl0_star <- cl0[idx_star]

    #align bootstrap labels to same original indices
    first_pos <- tapply(seq_along(idx), idx, '[', 1)
    pos_star <- as.integer(first_pos[as.character(idx_star)])
    clb_star <- cl_boot[pos_star]

    #Jaccard per original cluster
    gamma_C <- numeric(k)

    for (c in seq_len(k)) {
      C_pts <- idx_star[cl0_star == c]
      if (length(C_pts) == 0L) {
        gamma_C[c] <- 0
      } else {
        jacc_c <- sapply(seq_len(k), function(d) {
          D_pts <- idx_star[clb_star == d]
          if (length(D_pts) == 0L) {
            return(0)
          }
          inter <- length(intersect(C_pts, D_pts))
          union <- length(union(C_pts, D_pts))
          if (union == 0L) 0 else inter / union
        })
        gamma_C[c] <- max(jacc_c)
      }
    }

    gamma_C
  }

  b <- boot(
    data = X,
    statistic = boot_fun,
    R = R,
    parallel = "multicore"
  )

  gamma_bar <- colMeans(b$t)
  gamma_mcse <- apply(b$t, 2, sd) / sqrt(R)

  mean_gamma <- mean(gamma_bar)
  mean_gamma_mcse <- sd(rowMeans(b$t)) / sqrt(R)

  list(
    gamma_bar = gamma_bar,
    gamma_mcse = gamma_mcse,
    mean_gamma = mean_gamma,
    mean_gamma_mcse = mean_gamma_mcse,
    boot_mat = b$t,
    boot_obj = b
  )
}

k.grid <- 3:7

ward.summary <- function(X, iters = 500, seed = 666) {
  ward_ch_summary <- ward_ch(X, k.grid)
  print(ward_ch_summary)
```

```
 97    ward_stab <- lapply(k.grid, function(k) {
 98      ward_jaccard_boot(X, k = k, R = iters, seed = seed)
 99    })
100    names(ward_stab) <- paste0("k", k.grid)
101
102    stab_summary <- data.frame(
103      k = k.grid,
104      mean_gamma = sapply(ward_stab, \(z) z$mean_gamma),
105      mean_gamma_mcse = sapply(ward_stab, \(z) z$mean_gamma_mcse)
106    )
107    print(stab_summary)
108  }
109
110  ward.summary(X.pca[, 1:13], 500)
111  # k        CH
112  # 1 3  157.71526
113  # 2 4  128.48701
114  # 3 5  112.32357
115  # 4 6  103.55704
116  # 5 7   94.68403
117  #     k mean_gamma  mean_gamma_mcse
118  # k3 3   0.6076719     0.003099954
119  # k4 4   0.5093733     0.002553549
120  # k5 5   0.4772990     0.003997354
121  # k6 6   0.4865560     0.004196353
122  # k7 7   0.4591030     0.003292549
```

### 5.7 Power Score Calculations

```
 1  source("kmeans.R")
 2
 3  X <- X.pca[,1:13]
 4
 5  power_score <- function(datasets, label_col = "label", n_folds = 10, seed = 666,
        Bstart = 100) {
 6    set.seed(seed)
 7    ks <- names(datasets)
 8
 9    results <- data.frame(
10      k = as.numeric(gsub("k", "", ks)),
11      PS = NA
12    )
13
14    power_score <- function(X_train, X_test, y_test, k, Bstart) {
15      km_fit <- kmeans(X_train, centers = k, nstart = Bstart)
16      centroids <- km_fit$centers
17      train_clusters <- km_fit$cluster
18
19      predict_clusters <- function(x) {
20        dists <- apply(centroids, 1, function(centroid) sum((x - centroid)^2))
21        which.min(dists)
22      }
23
24      test_clusters <- apply(X_test, 1, predict_clusters)
25
26      ps_per_cluster <- numeric(k)
27      for (j in 1:k) {
28        idx <- which(test_clusters == j)
```

```r
        if (length(idx) < 2) {
          ps_per_cluster[j] <- 0
        } else {
          D <- outer(as.numeric(y_test[idx]), as.numeric(y_test[idx]), FUN =
              function(a, b) as.numeric(a == b))
          ps_per_cluster[j] <- sum(D) / (length(idx) * (length(idx) - 1))
        }
      }

    PS <- min(ps_per_cluster)
    return(PS)
  }

  for (ki in seq_along(ks)) {
    data <- datasets[[ki]]
    y <- factor(data[[label_col]])
    X <- data[, setdiff(names(data), label_col)]
    n <- nrow(X)
    folds <- sample(rep(1:n_folds, length.out = n))

    ps_folds <- numeric(n_folds)

    for (fold in 1:n_folds) {
      train_idx <- which(folds != fold)
      test_idx  <- which(folds == fold)

      X_train <- X[train_idx, ]
      X_test  <- X[test_idx, ]
      y_train <- y[train_idx]
      y_test  <- y[test_idx]

      k <- results$k[ki]
      ps_folds[fold] <- power_score(X_train, X_test, y_test, k, Bstart)
    }

    results$PS[ki] <- mean(ps_folds)
  }

  return(results)
}

datasets <- list(
  k3 = data_k3,
  k4 = data_k4,
  k5 = data_k5,
  k6 = data_k6,
  k7 = data_k7
)

results_power_score <- power_score(datasets, label_col = "label", n_folds = 10)
results_power_score
```

## 5.8 Misclassification Rates

```r
library(randomForest)
library(nnet)
library(xgboost)
library(tidyverse)
```

```r
supervised_cv_all <- function(datasets, label_col = "label", n_folds = 10, seed =
    666) {
  set.seed(seed)
  ks <- names(datasets)

  results <- data.frame(
    k = as.numeric(gsub("k", "", ks)),
    Random_Forest = NA,
    Multinomial_Logistic = NA,
    XGBoost = NA
  )

  for (ki in seq_along(ks)) {
    data <- datasets[[ki]]
    y <- factor(data[[label_col]])
    X <- data[, setdiff(names(data), label_col)]
    n <- nrow(X)
    folds <- sample(rep(1:n_folds, length.out = n))

    rf_mis <- numeric(n_folds)
    multinom_mis <- numeric(n_folds)
    xgb_mis <- numeric(n_folds)

    for (fold in 1:n_folds) {
      train_idx <- which(folds != fold)
      test_idx  <- which(folds == fold)

      X_train <- X[train_idx, ]
      X_test  <- X[test_idx, ]
      y_train <- y[train_idx]
      y_test  <- y[test_idx]

      # Random Forest
      rf_model <- randomForest(x = X_train, y = y_train, ntree = 500)
      rf_pred <- predict(rf_model, X_test)
      rf_mis[fold] <- mean(rf_pred != y_test)

      # Multinomial Logistic Regression
      multinom_model <- multinom(y_train ~ ., data = X_train, trace = FALSE)
      multinom_pred <- predict(multinom_model, X_test, type = "class")
      multinom_mis[fold] <- mean(multinom_pred != y_test)

      # XGBoost
      y_train_num <- as.numeric(y_train) - 1
      y_test_num <- as.numeric(y_test) - 1
      dtrain <- xgb.DMatrix(data = as.matrix(X_train), label = y_train_num)
      dtest  <- xgb.DMatrix(data = as.matrix(X_test), label = y_test_num)
      params <- list(objective = "multi:softmax", num_class = length(unique(y)),
          eval_metric = "merror")
      xgb_model <- xgb.train(params = params, data = dtrain, nrounds = 100,
          verbose = 0)
      xgb_pred <- predict(xgb_model, dtest)
      xgb_mis[fold] <- mean(xgb_pred != y_test_num)
    }

    results$Random_Forest[ki] <- mean(rf_mis)
    results$Multinomial_Logistic[ki] <- mean(multinom_mis)
    results$XGBoost[ki] <- mean(xgb_mis)
```

```
61        }
62
63      return(results)
64    }
65
66    datasets <- list(
67      k3 = data_k3,
68      k4 = data_k4,
69      k5 = data_k5,
70      k6 = data_k6,
71      k7 = data_k7
72    )
73
74    results_all <- supervised_cv_all(datasets, label_col = "label", n_folds = 10)
75    results_all
76
77
78
79    supervised_cv_confidence <- function(datasets, label_col = "label", n_folds = 10,
        seed = 666) {
80      set.seed(seed)
81      ks <- names(datasets)
82
83      results <- data.frame(
84        k = as.numeric(gsub("k", "", ks)),
85        Random_Forest = NA,
86        Multinomial_Logistic = NA,
87        XGBoost = NA
88      )
89
90      for (ki in seq_along(ks)) {
91        data <- datasets[[ki]]
92        y <- factor(data[[label_col]])
93        X <- data[, setdiff(names(data), label_col)]
94        n <- nrow(X)
95        folds <- sample(rep(1:n_folds, length.out = n))
96
97        rf_conf <- numeric(n_folds)
98        multinom_conf <- numeric(n_folds)
99        xgb_conf <- numeric(n_folds)
100
101       for (fold in 1:n_folds) {
102         train_idx <- which(folds != fold)
103         test_idx  <- which(folds == fold)
104
105         X_train <- X[train_idx, ]
106         X_test  <- X[test_idx, ]
107         y_train <- y[train_idx]
108         y_test  <- y[test_idx]
109
110         # Random Forest
111         rf_model <- randomForest(x = X_train, y = y_train, ntree = 500)
112         rf_pred_prob <- predict(rf_model, X_test, type = "prob")
113         rf_conf[fold] <- mean(apply(rf_pred_prob, 1, max))
114
115         # Multinomial Logistic Regression
116         multinom_model <- nnet::multinom(y_train ~ ., data = X_train, trace = FALSE)
117         multinom_pred_prob <- predict(multinom_model, X_test, type = "probs")
118         multinom_conf[fold] <- mean(apply(multinom_pred_prob, 1, max))
```

```
119
120        # XGBoost
121        y_train_num <- as.numeric(y_train) - 1
122        y_test_num <- as.numeric(y_test) - 1
123        dtrain <- xgb.DMatrix(data = as.matrix(X_train), label = y_train_num)
124        dtest  <- xgb.DMatrix(data = as.matrix(X_test), label = y_test_num)
125        params <- list(objective = "multi:softprob", num_class = length(unique(y)),
              eval_metric = "mlogloss")
126        xgb_model <- xgb.train(params = params, data = dtrain, nrounds = 100,
              verbose = 0)
127        xgb_pred_prob <- matrix(predict(xgb_model, dtest), ncol = length(unique(y)),
              byrow = TRUE)
128        xgb_conf[fold] <- mean(apply(xgb_pred_prob, 1, max))
129      }
130
131      results$Random_Forest[ki] <- mean(rf_conf)
132      results$Multinomial_Logistic[ki] <- mean(multinom_conf)
133      results$XGBoost[ki] <- mean(xgb_conf)
134    }
135
136    return(results)
137 }
138
139 # Example usage
140 results_all_conf <- supervised_cv_confidence(datasets, label_col = "label",
       n_folds = 5)
141 results_all_conf
```

## 5.9   Heatmaps by Genre

```
1
2  source("themes.R")
3  source("data.R")
4
5  library(tidyverse)
6
7  k3 <- read.csv("Data/dat-3.csv")
8
9  all_data <- bind_cols(k3, collins)
10
11 mean_mat <- all_data %>%
12    group_by(label) %>%
13    summarize(across(
14      c(
15        "FirstPerson",
16        "InnerThinking",
17        "ThinkPositive",
18        "ThinkNegative",
19        "ThinkAhead",
20        "ThinkBack",
21        "Reasoning",
22        "Share_SocTies",
23        "Direct_Activity",
24        "Interacting",
25        "Notifying",
26        "LinearGuidance",
27        "WordPicture",
28        "SpaceInterval",
```

```r
        "Motion",
        "PastEvents",
        "TimeInterval",
        "ShiftingEvents",
        "Text_Coverage"
      ),
      mean
  )) %>%
  column_to_rownames("label") %>%
  as.matrix()

tbl <- table(all_data$Genre, all_data$label)
tbl_prop <- prop.table(tbl, margin = 2)
df_heat_prop <- as.data.frame(tbl_prop)
colnames(df_heat_prop) <- c("genre", "cluster", "proportion")
df_heat_prop <- df_heat_prop %>%
  mutate(
    genre_names = case_when(
      genre == 1 ~ "Press: Reporting",
      genre == 2 ~ "Press: Editorial",
      genre == 3 ~ "Press: Reviews",
      genre == 4 ~ "Religion",
      genre == 5 ~ "Skills and Hobbies",
      genre == 6 ~ "Popular Lore",
      genre == 7 ~ "Biography, Memoirs, Belles Lettres, etc.",
      genre == 8 ~ "Official Communications",
      genre == 9 ~ "Learned (scholarly journals)",
      genre == 10 ~ "General Fiction",
      genre == 11 ~ "Mystery and Detective Fiction",
      genre == 12 ~ "Science Fiction",
      genre == 13 ~ "Adventure and Western Fiction",
      genre == 14 ~ "Romance and Love Story",
      genre == 15 ~ "Humor"
    ),
    genre_names = factor(
      genre_names,
      levels = c(
        "Press: Reporting",
        "Press: Editorial",
        "Press: Reviews",
        "Religion",
        "Skills and Hobbies",
        "Popular Lore",
        "Biography, Memoirs, Belles Lettres, etc.",
        "Official Communications",
        "Learned (scholarly journals)",
        "General Fiction",
        "Mystery and Detective Fiction",
        "Science Fiction",
        "Adventure and Western Fiction",
        "Romance and Love Story",
        "Humor"
      )
    ),
    category = case_when(
      genre %in% c(1, 2, 3) ~ "Press",
      genre %in% c(4, 5, 6) ~ "Non-press Nonfiction",
      genre == 7 ~ "Biography",
      genre %in% c(8, 9) ~ "Scholarship and Official Documents",
```

```
 88         genre %in% c(10, 11, 12, 13, 14, 15) ~ "Fiction"
 89       ),
 90       category = factor(
 91         category,
 92         levels = c(
 93           "Press",
 94           "Non-press Nonfiction",
 95           "Biography",
 96           "Scholarship and Official Documents",
 97           "Fiction"
 98         )
 99       )
100     )
101
102
103  # Heatmap by genre
104  p1 <- ggplot(
105     df_heat_prop,
106     aes(x = cluster, y = genre_names, fill = proportion)
107  ) +
108     geom_tile() +
109     geom_text(aes(label = round(proportion, 2)), color = "black", size = 16) + #
            larger numbers on tiles
110     scale_fill_gradient2_paper() +
111     theme_paper(base_size = 64, text_size = 48) +
112     labs(x = "Cluster", y = "Genre", fill = "Proportion") +
113     theme(
114       legend.position = "none",
115       panel.grid = element_blank()
116     )
117
118  k5 <- read.csv("Data/dat-5.csv")
119  all_data2 <- bind_cols(k5, collins)
120
121  all_data2$category <- case_when(
122     all_data2$Genre %in% c(1, 2, 3) ~ "Press",
123     all_data2$Genre %in% c(4, 5, 6) ~ "Non-press Nonfiction",
124     all_data2$Genre == 7 ~ "Biography",
125     all_data2$Genre %in% c(8, 9) ~ "Scholarship and Official Documents",
126     all_data2$Genre %in% c(10, 11, 12, 13, 14, 15) ~ "Fiction"
127  )
128
129  tbl <- table(all_data2$category, all_data2$label)
130  tbl_prop <- prop.table(tbl, margin = 2)
131  df_heat_prop <- as.data.frame(tbl_prop)
132  colnames(df_heat_prop) <- c("category", "cluster", "proportion")
133
134  df_heat_prop_cat <- df_heat_prop %>%
135     group_by(category, cluster) %>%
136     summarize(proportion = sum(proportion)) %>%
137     ungroup()
138
139  # Heatmap by category
140  p3 <- ggplot(
141     df_heat_prop_cat,
142     aes(x = cluster, y = category, fill = proportion)
143  ) +
144     geom_tile() +
145     geom_text(aes(label = round(proportion, 2)), color = "black", size = 16) +
```

```
146    scale_fill_gradient2_paper () +
147    theme_paper ( base_size = 64 , text_size = 48) +
148    labs (x = "Cluster", y = "Category", fill = "Proportion") +
149    theme (
150      legend.position = "none",
151      panel.grid = element_blank ()
152    )
```

## 5.10   Linear Discriminant Analysis

```
1  library (MASS)
2
3  source ("data.R")
4
5  X$label <- as.factor (k3$label)
6  X_scaled <- scale (X[, -which (names (X) == "label")])
7  lda_model <- lda (X$label ~ ., data = as.data.frame (X_scaled))
8
9  lda_model$scaling
```

## 5.11   Figures

### 5.11.1   Theme for Paper

```
1  ##Themes.R
2
3  library (ggplot2)
4  library (sysfonts)
5
6  font_add ("cm", regular = "fonts/cmunrm.ttf")
7  showtext :: showtext_auto ()
8
9  theme_paper <- function (base_size = 14 , text_size = 14 , base_family = "cm") {
10   theme_minimal ( base_size = base_size , base_family = base_family ) %+replace%
11     theme (
12       # Text
13       plot.title = element_text (
14         face = "bold",
15         size = text_size + 2,
16         hjust = 0.5,
17         margin = margin (b = 6)
18       ),
19       plot.subtitle = element_text (
20         size = text_size ,
21         hjust = 0.5,
22         margin = margin (b = 8)
23       ),
24       plot.caption = element_text (
25         size = text_size - 2,
26         hjust = 1,
27         margin = margin (t = 6)
28       ),
29       axis.title.x = element_text (size = text_size , margin = margin (t = 8)),
30       axis.title.y = element_text (
31         size = text_size ,
32         angle = 90 ,
33         margin = margin (r = 8)
34       ),
```

```r
        axis.text = element_text(size = text_size - 1),

        # Panel & grid
        panel.grid.major.x = element_blank(),
        panel.grid.minor.x = element_blank(),
        panel.grid.minor.y = element_blank(),
        #panel.grid.major.y = element_line(linewidth = 0.3, color = "grey85"),
        panel.border = element_rect(fill = NA, color = "grey60", linewidth = 0.4),

        # Legend
        legend.position = "bottom",
        legend.title = element_text(size = base_size - 1, face = "bold"),
        legend.text = element_text(size = base_size - 2),
        legend.key.width = unit(1.2, "lines"),

        # Strips (for facets)
        strip.background = element_rect(fill = "grey90", color = NA),
        strip.text = element_text(face = "bold", size = base_size - 1),

        # Margins
        plot.margin = margin(t = 8, r = 8, b = 8, l = 8)
    )
}

paper_palette <- c(
  "#1b9e77",
  "#d95f02",
  "#7570b3",
  "#e7298a",
  "#66a61e",
  "#e6ab02"
)

scale_color_paper <- function(...) {
  scale_color_manual(values = paper_palette, ...)
}

scale_fill_paper <- function(...) {
  scale_fill_manual(values = paper_palette, ...)
}

scale_fill_gradient2_paper <- function(
  low = "#7570b3",
  mid = "white",
  high = "#d95f02",
  midpoint = 0,
  ...
) {
  scale_fill_gradient2(
    low = low,
    mid = mid,
    high = high,
    midpoint = midpoint,
    ...
  )
}
```

### 5.11.2 Reproduce Figures

```r
source("themes.R")
source("data.R")

library(patchwork)

#PCA ######################################################

pca_elbow_plot(X, center = T, scale. = T)
ggsave("Figures/PCA-elbow.png", width = 6, height = 5, dpi = 600)

#How many components should we use according to a scree plot?
max(which(eigen(cor(X))$values >= 1))

#KMEANS ######################################################

source("kmeans.R")

kg <- kmeans_clusters(X.pca[, 1:13], 3, Bstart = 10)
kg.df <- data.frame(
  pc1 = X.pca[, 1],
  pc2 = X.pca[, 2],
  labels = factor(kg)
)

kmeans.plot <- ggplot(kg.df, aes(x = pc1, y = pc2, color = labels)) +
  geom_point(size = 0.5) +
  scale_color_paper(guide = "none") +
  labs(
    x = expression(PC[1]),
    y = expression(PC[2])
  ) +
  theme_paper(text_size = 64)
ggsave("Figures/kmeans.png", kmeans.plot, width = 5, height = 4, dpi = 600)

# Genre Classification ######################################

source("cluster_genres.R")

genre.cat <- p1 + p3

ggsave(
  "Figures/cluster_genre_heatmaps.png",
  genre.cat,
  width = 10,
  height = 4,
  dpi = 600
)
```

## 6 Appendix - SAS Code

### 6.1 MANOVA

```sas
/* Code generated on: 11/14/25, 1:43 PM */

/*************************************
```

```sas
                    K = 3
*******************************************/

FILENAME REFFILE '/home/u64365443/sasuser.v94/dat-3.csv';

PROC IMPORT DATAFILE=REFFILE
        DBMS=CSV
        OUT=WORK.PCA;
        GETNAMES=YES;
RUN;


ods select CharStruct MultStat;

proc glm data = PCA;
    class label;
    model PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 =
            label;

    /* Request Wilks', Pillai's, Hotelling-Lawley, Roy's tests */
    manova h = label / printe printh;
run;
quit;

/*****************************************
                    K = 4
*******************************************/

%web_drop_table(WORK.PCA);

FILENAME REFFILE '/home/u64365443/sasuser.v94/dat-4.csv';

PROC IMPORT DATAFILE=REFFILE
        DBMS=CSV
        OUT=WORK.PCA;
        GETNAMES=YES;
RUN;

%web_open_table(WORK.PCA);

ods select CharStruct MultStat;
proc glm data = PCA;
    class label;
    model PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 =
            label;

    /* Request Wilks', Pillai's, Hotelling-Lawley, Roy's tests */
    manova h = label / printe printh;
run;
quit;

/*****************************************
                    K = 5
*******************************************/

%web_drop_table(WORK.PCA);

FILENAME REFFILE '/home/u64365443/sasuser.v94/dat-5.csv';
```

```sas
PROC IMPORT DATAFILE=REFFILE
        DBMS=CSV
        OUT=WORK.PCA;
        GETNAMES=YES;
RUN;

%web_open_table(WORK.PCA);

ods select CharStruct MultStat;
proc glm data = PCA;
    class label;
    model PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 =
            label;

    /* Request Wilks', Pillai's, Hotelling-Lawley, Roy's tests */
    manova h = label / printe printh;
run;
quit;

/****************************************
                  K = 6
****************************************/

%web_drop_table(WORK.PCA);

FILENAME REFFILE '/home/u64365443/sasuser.v94/dat-6.csv';

PROC IMPORT DATAFILE=REFFILE
        DBMS=CSV
        OUT=WORK.PCA;
        GETNAMES=YES;
RUN;

%web_open_table(WORK.PCA);

ods select CharStruct MultStat;
proc glm data = PCA;
    class label;
    model PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 =
            label;

    /* Request Wilks', Pillai's, Hotelling-Lawley, Roy's tests */
    manova h = label / printe printh;
run;
quit;

/****************************************
                  K = 7
****************************************/

%web_drop_table(WORK.PCA);

FILENAME REFFILE '/home/u64365443/sasuser.v94/dat-7.csv';

PROC IMPORT DATAFILE=REFFILE
        DBMS=CSV
        OUT=WORK.PCA;
        GETNAMES=YES;
RUN;
```

```
123
124 %web_open_table ( WORK.PCA );
125
126 ods select CharStruct MultStat ;
127 proc glm data = PCA ;
128     class label ;
129     model PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 =
130         label ;
131
132     /* Request Wilks', Pillai's, Hotelling-Lawley, Roy's tests */
133     manova h = label / printe printh ;
134 run ;
135 quit ;
```