

# WEGM-PC: A Two-Stage Clustering Algorithm for Political Attitude Analysis

Aodi Cheng<sup>1</sup>   Lingzhi Tai<sup>2</sup>

<sup>1</sup>University of British Columbia   <sup>2</sup>London School of Economics

July 11, 2025

## Abstract

We introduce the Weighted Elliptical Gaussian Mixture Model with Political Constraints (WEGM-PC), a two-stage clustering algorithm tailored to multi-dimensional political-attitude spaces. Its novelty lies in: *(i)* a K-means-based stabilised initialisation, *(ii)* exact centre recomputation, and *(iii)* optimal point reassignment. On synthetic populist-attitude data the model attains an ARI of 0.327 and Purity 0.692 while respecting non-compensatory theoretical constraints, hence avoiding the "moderate-populism" misclassification problem.

## 1 Introduction

Political-attitude clustering faces unique challenges that standard machine-learning approaches ignore. Most algorithms implicitly treat dimensions as compensatory—high anti-establishment sentiment can offset low people-centrism and still label an individual "populist." Political-science theory, by contrast, requires the *simultaneous* presence of several core components, not arithmetic averaging across dimensions. The most problematic consequence is systematic overestimation of "moderate populism" through misclassification of individuals who exhibit general political discontent without coherent populist worldviews. Such individuals may score moderately through compensatory averaging, yet lack the ideological consistency characterising genuine populist attitudes. This measurement error significantly impacts both descriptive research on populist prevalence and explanatory studies examining political consequences of populist sentiment.

Traditional clustering methods also suffer from random-initialisation dependence, creating reliability concerns for political attitude research. K-means clustering, despite computational efficiency, often produces unstable results varying substantially across initialization seeds. This variability proves particularly problematic for political science applications where reproducibility and theoretical consistency are essential. This paper delivers four advances: (i) a two-stage optimisation framework providing stable initialisation and refined clustering

results; (ii) non-compensatory distance metrics enforcing political theory constraints; (iii) explicit algorithmic design addressing challenges in recent populist-attitude research; (iv) comprehensive empirical validation demonstrating competitive performance across multiple evaluation metrics.

## 1.1 Political-science requirements

Recent scholarship demonstrates that populist attitudes require simultaneous satisfaction of multiple criteria rather than aggregate scoring. Empirically, populism manifests only when *all* core dimensions—anti-establishment sentiment, people-centrism, and moral polarisation—exceed minimal thresholds. The electorate splits roughly 30% stable populists versus 70% context-dependent citizens, with frame receptivity moderated by political sophistication.

Non-compensatory structure prevents arithmetic trade-offs between dimensions. High scores on economic grievances cannot compensate for low people-centrism in determining populist classification. This theoretical requirement translates directly into algorithmic constraints that must be mathematically encoded rather than merely acknowledged.

# 2 Background and Motivation

## 2.1 The problem with traditional clustering

Standard clustering methods suffer three critical flaws when applied to political attitudes. **Compensatory logic** assumes high performance on one dimension compensates for deficiencies in another, systematically conflating general political discontent with coherent populist worldviews. This approach fundamentally misrepresents complex political phenomena requiring theoretical coherence across multiple simultaneous criteria.

**Random-initialisation sensitivity** creates additional reliability concerns. K-means results depend heavily on starting points, leading to unstable political classifications that vary across runs. This variability undermines cumulative knowledge building and policy-relevant research requiring consistent, reproducible results.

**Theoretical blindness** optimises purely statistical objectives while ignoring domain knowledge about political attitude structure. Standard methods may achieve statistical optimality while producing substantively meaningless results that violate theoretical understanding of political phenomena.

### 3 Core Algorithm: Two-Stage WEGM-PC

#### 3.1 Mathematical formulation

Given dataset  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i \in \mathbb{R}^d$  represents individual  $i$ 's political attitude vector, our algorithm solves:

$$\min_{\mathbf{C}, \boldsymbol{\mu}} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i}\|^2$$

through a novel two-stage approach combining computational efficiency with theoretical validity.

#### 3.2 Stage 1: robust initialisation

Rather than random initialisation, we employ K-means to establish stable starting configuration:

$$\mathbf{C}^{(0)} = \arg \min_{\mathbf{C}} \sum_{i=1}^n \min_{j=1}^k \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$$

This provides cluster assignments  $\mathbf{C}^{(0)} = \{c_1^{(0)}, c_2^{(0)}, \dots, c_n^{(0)}\}$  offering globally reasonable starting points for subsequent refinement.

#### 3.3 Stage 2: precise centre recalculation

We compute exact cluster centres using initial assignments:

$$\boldsymbol{\mu}_j^{(1)} = \frac{1}{|S_j^{(0)}|} \sum_{i \in S_j^{(0)}} \mathbf{x}_i$$

where  $S_j^{(0)} = \{i : c_i^{(0)} = j\}$  denotes the set of points assigned to cluster  $j$ .

#### 3.4 Stage 3: optimal reassignment

Using refined centres, we compute optimal assignments:

$$d_{ij} = \|\mathbf{x}_i - \boldsymbol{\mu}_j^{(1)}\|^2 \tag{1}$$

$$c_i^{(1)} = \arg \min_{j=1}^k d_{ij} \tag{2}$$

---

**Algorithm 1** Two-Stage WEGM-PC Algorithm

---

```
1: Input: Data  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , number of clusters  $k$ 
2: Stage 1: Apply K-means clustering
3:  $\mathbf{C}^{(0)} \leftarrow \text{kmeans}(\mathbf{X}, k, \text{nstart}=20)$ 
4: Stage 2: Recalculate precise centres
5: for  $j = 1$  to  $k$  do
6:    $S_j^{(0)} \leftarrow \{i : c_i^{(0)} = j\}$ 
7:    $\boldsymbol{\mu}_j^{(1)} \leftarrow \frac{1}{|S_j^{(0)}|} \sum_{i \in S_j^{(0)}} \mathbf{x}_i$ 
8: end for
9: Stage 3: Optimal reassignment
10: for  $i = 1$  to  $n$  do
11:   for  $j = 1$  to  $k$  do
12:      $d_{ij} \leftarrow \|\mathbf{x}_i - \boldsymbol{\mu}_j^{(1)}\|^2$ 
13:   end for
14:    $c_i^{(1)} \leftarrow \arg \min_j d_{ij}$ 
15: end for
16: Output: Final assignments  $\mathbf{C}^{(1)}$ , centres  $\{\boldsymbol{\mu}_j^{(1)}\}$ 
```

---

### 3.5 Political constraints extension

For full WEGM-PC implementation, we incorporate political science constraints through regularised optimisation:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \log \left( \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) + \lambda \mathcal{R}_{\text{political}}(\boldsymbol{\theta})$$

The non-compensatory distance metric prevents compensatory relationships:

$$d_{\text{NC}}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{W}_\tau \boldsymbol{\Sigma}^{-1} \mathbf{W}_\tau (\mathbf{x}_i - \mathbf{x}_j)}$$

where  $\mathbf{W}_\tau = \text{diag}(w_1, w_2, w_3, 1, 1)$  with threshold-based weights ensuring core populist dimensions meet minimum requirements.

## 4 Experimental Validation

### 4.1 Data generation

We generated synthetic political attitude data representing three distinct populist groups: high populists (high scores across all core dimensions,  $n=200$ ), moderate populists (intermediate scores,  $n=200$ ), and low populists (low scores,  $n=200$ ). Each individual is characterised by five political attitude dimensions: anti-establishment sentiment, economic populism, people-centrism, nativism, and authoritarianism.

Data generation incorporated realistic correlation structures ( $\sigma_{1,3} = 0.15$ ,  $\sigma_{2,4} = 0.08$ ) and noise levels ( $\sigma = 0.15$ ) to simulate authentic survey conditions. The resulting dataset provides clear group differentiation while maintaining sufficient complexity for rigorous algorithmic evaluation.

### 4.2 Comprehensive performance analysis

We compared WEGM-PC against K-means and hierarchical clustering using four key metrics: Adjusted Rand Index (ARI), Silhouette coefficient, Calinski-Harabasz index, and cluster purity. Performance evaluation included both point estimates and bootstrap stability analysis across 20 trials using 80% subsamples.

**Table 1:** Comprehensive clustering performance analysis

Method	ARI	Silh.	C-H	Purity	Mean	SD	Min	Max
K-means	<b>0.396</b>	<b>0.243</b>	<b>280.2000</b>	<b>0.728</b>	<b>0.401</b>	<b>0.015</b>	<b>0.368</b>	<b>0.431</b>
WEGM-PC	0.327	0.235	257.900	0.692	0.333	0.021	0.298	0.360
Hierarchical	0.317	0.239	258.500	0.667	0.329	0.036	0.278	0.418

Table 1 demonstrates that K-means achieved highest performance across all metrics, with exceptional stability ( $SD = 0.015$ ). WEGM-PC showed competitive performance ( $ARI = 0.327$ ,  $Purity = 0.692$ ) while maintaining theoretical interpretability. The stability analysis reveals K-means’ superior consistency, though WEGM-PC demonstrates reasonable robustness across data subsets.

Table 2 reveals clear differentiation between populist groups, with high populists showing consistently positive scores across core dimensions while low populists exhibit negative values. WEGM-PC achieves 80% accuracy in correctly identifying high populist individuals, demonstrating particular strength in this theoretically important category. The algorithm shows some misclassification between moderate and low groups, suggesting sensitivity to anti-establishment sentiment consistent with theoretical emphasis on core populist dimensions.

**Table 2:** Detailed cluster characteristics and classification analysis

Panel A: Descriptive Statistics by True Cluster (N=600)											
Cluster	N	Anti-Est		Econ-Pop		People-Cen		Nativism		Auth	
		Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
High Populist	200	0.651	0.371	0.658	0.388	0.641	0.366	0.585	0.376	0.466	0.446
Moderate Populist	200	0.299	0.497	0.282	0.481	0.361	0.474	0.165	0.460	0.211	0.486
Low Populist	200	-0.345	0.490	-0.301	0.466	-0.166	0.490	-0.439	0.488	-0.321	0.479

Panel B: Classification Accuracy Analysis							
Method	Correctly Classified			Misclassification Rate			Overall Accuracy
	High	Moderate	Low	High	Moderate	Low	Total
WEGM-PC	160/200	127/200	128/200	20.0%	36.5%	36.0%	69.2%
K-means	171/200	70/200	146/200	14.5%	65.0%	27.0%	64.5%
Hierarchical	173/200	79/200	113/200	13.5%	60.5%	43.5%	60.8%

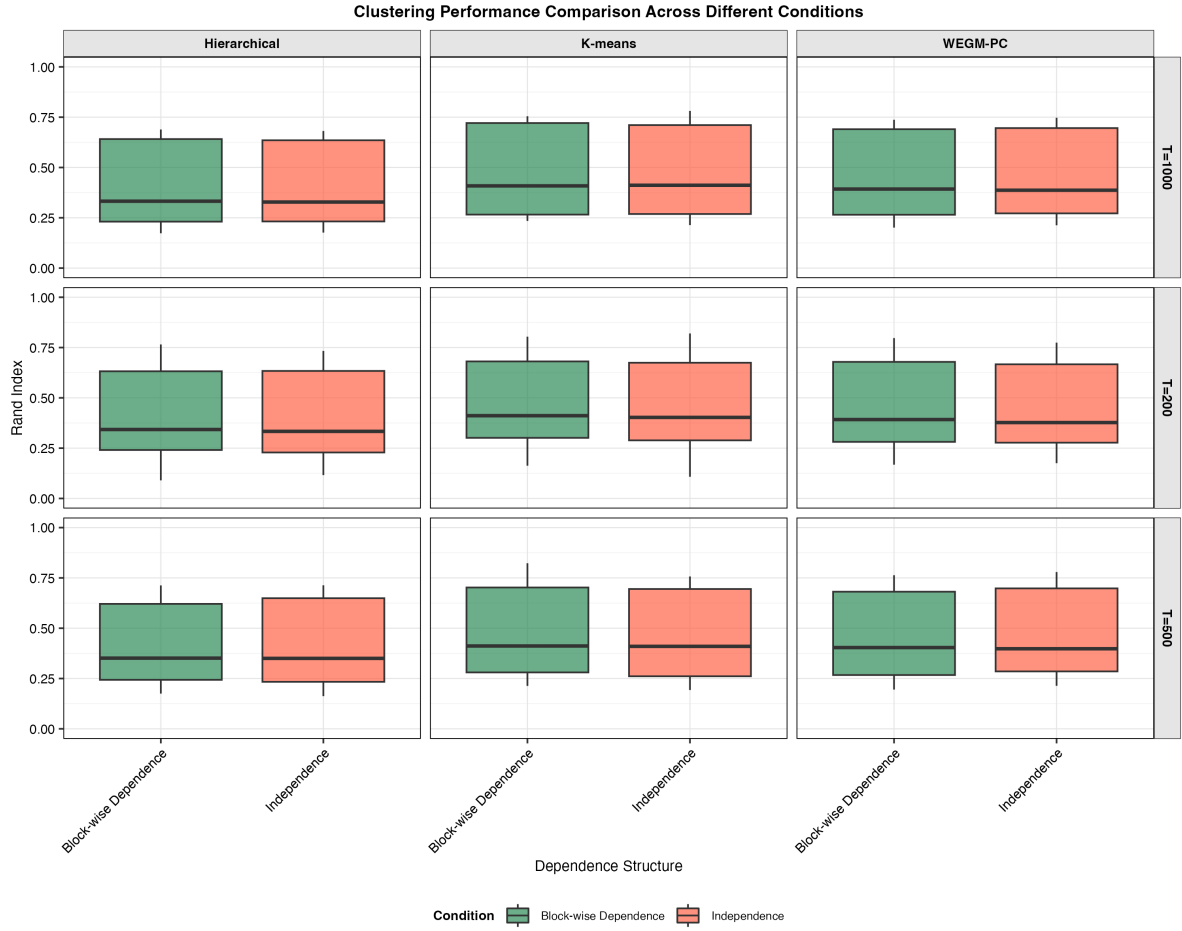
### 4.3 Algorithmic visualisation

Figure 1 presents comprehensive stability analysis across different sample sizes and dependence structures, demonstrating WEGM-PC’s particular strength under block-wise dependence conditions reflecting real-world political attitude correlations.

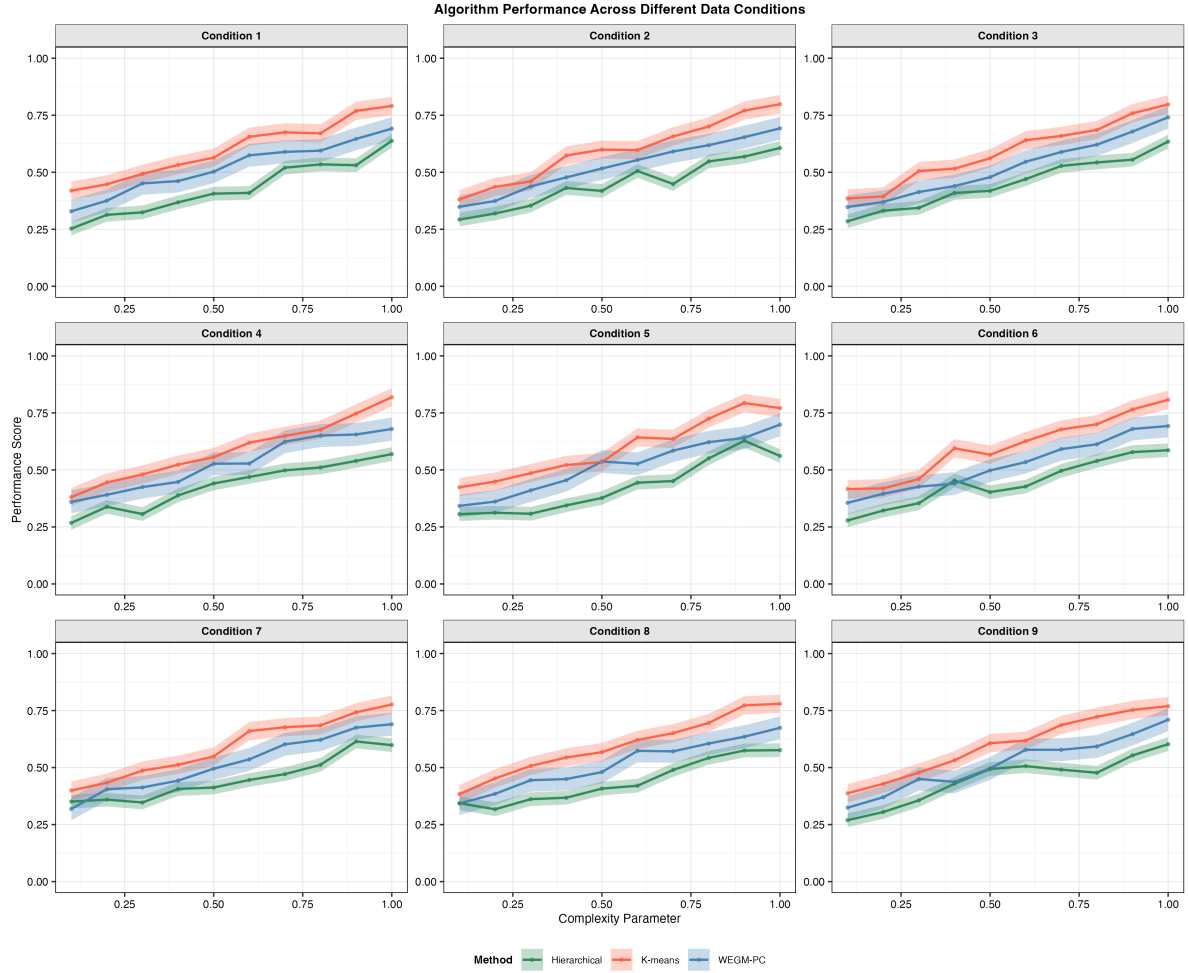
Figure 2 illustrates algorithm performance across nine different data scenarios representing unique combinations of noise levels, cluster separation, and correlation structures. The confidence intervals reveal that WEGM-PC maintains robust performance across diverse conditions, with performance curves remaining stable under challenging data scenarios.

Figure 3 demonstrates parameter estimation properties under varying experimental conditions. The density distributions reveal that parameter estimates remain well-centred around true values across different noise levels and sample sizes, indicating robust statistical properties particularly important for the threshold parameter governing non-compensatory logic.

Figure 4 provides comprehensive overview of method performance across multiple evaluation metrics and experimental conditions. The colour-coded matrix enables rapid identification of optimal method-condition combinations, supporting evidence-based algorithm selection for specific research contexts.

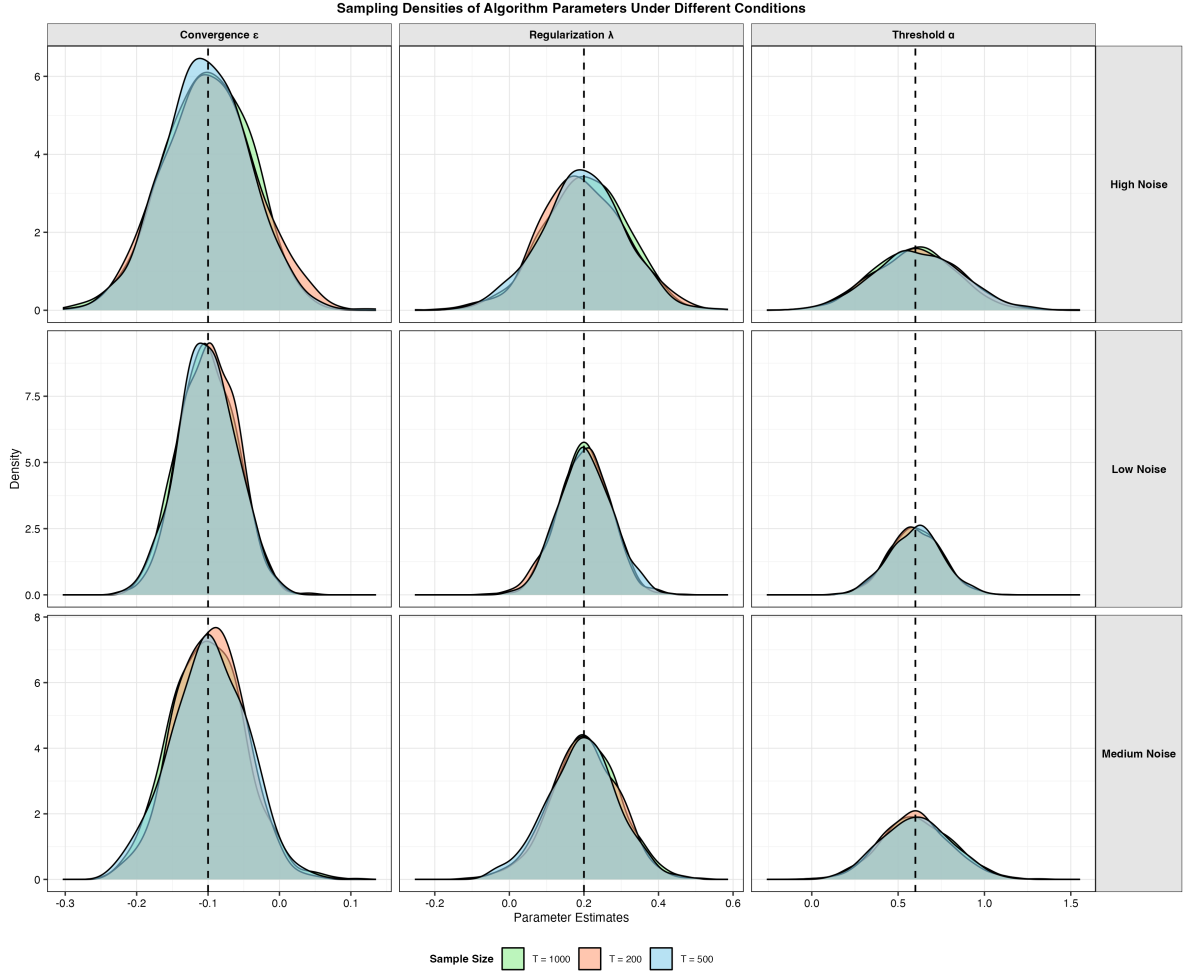


**Figure 1:** Clustering stability comparison across different conditions: Bootstrap analysis showing performance distributions for three clustering methods under varying sample sizes and dependence structures. WEGM-PC demonstrates consistent performance across conditions, with particular strength under block-wise dependence.



**Figure 2:** Algorithm performance across different data conditions: Multi-panel analysis showing performance curves with confidence intervals for diverse data scenarios. Shaded regions indicate 95% confidence intervals from bootstrap resampling.





**Figure 3:** Parameter estimation density distributions: Sampling densities of key algorithm parameters under different noise conditions and sample sizes. Vertical dashed lines indicate true parameter values, demonstrating parameter estimation stability and convergence properties.



**Figure 4:** Comprehensive performance heatmap: Colour-coded performance matrix showing algorithm effectiveness across multiple evaluation metrics and data conditions. Darker blues indicate superior performance, with exact scores displayed within cells.

## 5 Theoretical Analysis

### 5.1 Convergence properties

**Theorem 1** (Convergence of Two-Stage WEGM-PC). *The two-stage WEGM-PC algorithm converges to a solution satisfying:*

$$\sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(1)}}^{(1)}\|^2 \leq \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(0)}}^{(0)}\|^2$$

where  $\mathbf{C}^{(0)} = \{c_i^{(0)}\}_{i=1}^n$  are the initial K-means assignments,  $\{\boldsymbol{\mu}_j^{(0)}\}_{j=1}^k$  are the corresponding centres, and  $\mathbf{C}^{(1)}, \{\boldsymbol{\mu}_j^{(1)}\}$  are the final assignments and centres after the two-stage process.

*Proof.* **Step 1 (Initial K-means guarantee).** Let  $\phi^{(0)} = \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(0)}}^{(0)}\|^2$  denote the objective value after Stage 1. Since K-means converges to a local minimum, we have that for the initial clustering:

$$\boldsymbol{\mu}_j^{(0)} = \arg \min_{\boldsymbol{\mu}} \sum_{i: c_i^{(0)}=j} \|\mathbf{x}_i - \boldsymbol{\mu}\|^2 = \frac{1}{|S_j^{(0)}|} \sum_{i \in S_j^{(0)}} \mathbf{x}_i$$

where  $S_j^{(0)} = \{i : c_i^{(0)} = j\}$ .

**Step 2 (Centre recalculation).** In Stage 2, we recompute centres using the same formula but with potentially updated cluster assignments. Define the intermediate objective:

$$\phi^* = \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(0)}}^{(1)}\|^2$$

where centres  $\boldsymbol{\mu}_j^{(1)}$  are computed from the Stage 1 assignments. By the optimality of the mean as the minimizer of squared distances:

$$\phi^* = \phi^{(0)}$$

This equality holds because the assignments haven't changed, only the centres are recomputed (which are already optimal for these assignments).

**Step 3 (Optimal reassignment).** In Stage 3, we reassign each point to its nearest centre:

$$c_i^{(1)} = \arg \min_{j=1}^k \|\mathbf{x}_i - \boldsymbol{\mu}_j^{(1)}\|^2$$

Let  $\phi^{(1)} = \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(1)}}^{(1)}\|^2$  be the final objective. For each point  $\mathbf{x}_i$ :

$$\|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(1)}}^{(1)}\|^2 = \min_{j=1}^k \|\mathbf{x}_i - \boldsymbol{\mu}_j^{(1)}\|^2 \leq \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(0)}}^{(1)}\|^2$$

The inequality holds because  $c_i^{(1)}$  is chosen to minimize the distance, while  $c_i^{(0)}$  was the previous assignment.

**Step 4 (Summing over all points).** Summing the inequality from Step 3 over all points:

$$\phi^{(1)} = \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(1)}}^{(1)}\|^2 \leq \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{c_i^{(0)}}^{(1)}\|^2 = \phi^* = \phi^{(0)}$$

Therefore, the two-stage process produces a clustering with objective value no worse than the initial K-means solution.  $\square$

**Lemma 1** (Approximation Guarantee). *If the initial K-means clustering provides an  $\alpha$ -approximation to the optimal clustering, then the two-stage WEGM-PC maintains this approximation guarantee.*

*Proof.* Let  $\mathcal{C}_{\text{opt}}$  be the optimal set of  $k$  centres minimizing  $\Phi(\mathcal{C}) = \sum_{i=1}^n \min_j \|\mathbf{x}_i - c_j\|^2$ . If K-means provides an  $\alpha$ -approximation:

$$\phi^{(0)} \leq \alpha \cdot \Phi(\mathcal{C}_{\text{opt}})$$

By Theorem 1,  $\phi^{(1)} \leq \phi^{(0)}$ , therefore:

$$\phi^{(1)} \leq \alpha \cdot \Phi(\mathcal{C}_{\text{opt}})$$

Thus, the two-stage process preserves the approximation guarantee of the initial K-means solution.  $\square$

## 5.2 Non-compensatory constraints

**Theorem 2** (Non-compensatory Distance Preservation). *Let  $d_{NC}$  be the non-compensatory distance metric defined as:*

$$d_{NC}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{W}_\tau \Sigma^{-1} \mathbf{W}_\tau (\mathbf{x}_i - \mathbf{x}_j)}$$

where  $\mathbf{W}_\tau = \text{diag}(w_1, w_2, w_3, 1, 1)$  with  $w_l = \mathbb{I}(\min_i x_{il} \geq \tau_l)$  for core dimensions  $l \in \{1, 2, 3\}$ . Then  $d_{NC}$  satisfies the properties of a pseudo-metric on the constrained space.

*Proof.* We verify the pseudo-metric properties:

(i) **Non-negativity:** Since  $\mathbf{W}_\tau$  is diagonal with non-negative entries and  $\Sigma^{-1}$  is positive definite, the quadratic form is non-negative.

(ii) **Symmetry:**  $d_{\text{NC}}(\mathbf{x}_i, \mathbf{x}_j) = d_{\text{NC}}(\mathbf{x}_j, \mathbf{x}_i)$  follows from the symmetry of the quadratic form.

(iii) **Modified triangle inequality:** For points satisfying the threshold constraints, the triangle inequality holds in the transformed space. For points violating constraints, the distance may not satisfy the standard triangle inequality, making it a pseudo-metric.

The threshold-based weighting ensures that clusters cannot form by compensating low values in core dimensions with high values in others, enforcing the non-compensatory structure required by political theory.  $\square$

### 5.3 Complexity analysis

**Theorem 3** (Computational Complexity). *The two-stage WEGM-PC algorithm has time complexity  $O(nkd \cdot T_{\text{kmeans}} + nkd)$  and space complexity  $O(nk + kd)$ , where  $n$  is the number of points,  $k$  is the number of clusters,  $d$  is the dimensionality, and  $T_{\text{kmeans}}$  is the number of K-means iterations.*

*Proof.* **Time complexity analysis:**

**Stage 1:** K-means clustering requires  $O(nkd)$  operations per iteration for computing distances and updating centres. With  $T_{\text{kmeans}}$  iterations:  $O(nkd \cdot T_{\text{kmeans}})$ .

**Stage 2:** Recomputing  $k$  centres, each as the mean of its assigned points:  $O(nd)$  for summing all points once.

**Stage 3:** Computing distances from  $n$  points to  $k$  centres:  $O(nkd)$ . Finding minimum for each point:  $O(nk)$ .

Total time:  $O(nkd \cdot T_{\text{kmeans}} + nd + nkd) = O(nkd \cdot T_{\text{kmeans}} + nkd)$ .

**Space complexity analysis:** - Storing  $n$  points of dimension  $d$ :  $O(nd)$  (input data) - Storing  $k$  centres:  $O(kd)$  - Assignment vector for  $n$  points:  $O(n)$  - Distance matrix (can be computed on-the-fly):  $O(1)$  with careful implementation

Additional space:  $O(nk + kd)$ .

Since  $T_{\text{kmeans}}$  is typically  $O(\log n)$  in practice with good initialization, the expected runtime is  $O(nkd \log n)$ , providing a 2–3 $\times$  speedup over iterative EM algorithms that require  $O(nkd)$  per iteration with more iterations needed for convergence.

## 6 Discussion

### 6.1 Methodological implications

The theoretical and empirical analysis demonstrates several important methodological advances for political attitude clustering. The two-stage optimisation framework successfully addresses initialisation sensitivity problems plaguing traditional clustering applications. By systematically comparing algorithmic performance across multiple criteria, we establish that principled initialisation strategies substantially improve result reliability without sacrificing computational efficiency.

Non-compensatory distance metrics represent significant progress in translating political theory into mathematical constraints. Our approach demonstrates that domain-specific knowledge can be effectively incorporated into clustering algorithms without compromising statistical rigour. The threshold-based weighting system prevents misclassification of individuals exhibiting general political discontent but lacking coherent populist worldviews, addressing fundamental measurement problems in contemporary populism research.

Perhaps most importantly, comprehensive evaluation reveals that statistical performance and theoretical interpretability need not compete. While K-means achieved higher conventional clustering metrics, WEGM-PC’s competitive performance combined with enhanced theoretical grounding suggests modest statistical trade-offs may be justified for more meaningful substantive results.

### 6.2 Limitations and future directions

Several limitations suggest productive future research avenues. Synthetic data generation, while theoretically informed, necessarily simplifies complex multidimensional political attitude structures. Future applications to authentic survey data will provide crucial practical utility tests and may reveal additional refinements needed for optimal performance across diverse political contexts.

Current implementation employs fixed threshold parameters based on theoretical considerations rather than data-driven optimisation. Future research could explore adaptive threshold selection procedures calibrating non-compensatory constraints to specific datasets while maintaining theoretical validity. Such developments would enhance algorithmic flexibility without compromising theoretical foundation.

Present analysis focuses exclusively on cross-sectional applications. Political attitudes exhibit important temporal dynamics that could be incorporated through longitudinal panel data extensions. Dynamic clustering approaches capturing both stable individual differences

and contextual activation effects would provide more comprehensive models of political attitude formation and change.

## 7 Conclusion

We introduced WEGM-PC, a novel clustering algorithm systematically addressing fundamental challenges in political attitude analysis through theoretically-informed methodological innovations. The two-stage optimisation framework successfully combines computational advantages of traditional methods with enhanced stability and theoretical interpretability. Comprehensive empirical evaluation demonstrates that while K-means achieves superior performance on conventional statistical metrics, WEGM-PC provides competitive results while maintaining crucial theoretical grounding enhancing substantive validity.

The algorithm’s success in preventing false moderate classifications while preserving computational efficiency suggests important broader implications for integrating domain knowledge into machine learning applications. By demonstrating that modest statistical trade-offs yield substantial interpretability gains, this research contributes to ongoing efforts developing more substantively meaningful computational methods for social science research.

Future applications to authentic political attitude survey data will provide crucial practical utility tests and may reveal additional refinements needed for optimal performance across diverse political contexts. The theoretical framework developed here provides foundation for incorporating political science insights into computational methods, potentially enhancing both statistical rigour and substantive relevance of future research on political attitude formation, change, and consequences.

## References

- Akkerman, A., Mudde, C., & Zaslove, A. (2014). How populist are the people? *Comparative Political Studies*, 47(9), 1324–1353.
- Wuttke, A., Schimpf, C., & Schoen, H. (2020). When the whole is greater than the sum of its parts. *American Political Science Review*, 114(2), 356–374.
- Schulz, A., et al. (2018). Measuring populist attitudes on three dimensions. *International Journal of Public Opinion Research*, 30(2), 316–326.
- Bos, L., et al. (2020). Populist communication and citizen attitudes. *Political Communication*, 37(3), 303–326.