

CSCA 5632 - Intro. to Machine Learning

Unsupervised Learning

Dyego Fernandes de Sousa

University of Colorado Boulder

November 28, 2025

This is a job for the Un**SUPER**vised Learning. Problem Statement

To showcase and apply techniques of **Unsupervised Learning** at the brand-level data to discover:

- **Hidden Patterns**

- **Natural groupings** Clusters of brands based on their multi-dimensional characteristics
- **Unusual patterns** and outliers that don't fit conventional categorizations
- **Relationships** between environmental performance, market positioning, and consumer targeting

- **Latent Features** that drive brand differentiation

The Hypothesis:

A company being compliant with ESG practices, doesn't mean that all of their brands also comply?

Primary Dataset: Comprehensive brand-level dataset, based on my **Supervised Learning project**. It contains:

- **Industry & Company Data:** ESG scores, emissions (scope1+2), revenues, environmental risk indicators
- **Brand Demographics:** Age group targeting, income levels, lifestyle segments
- **Operational Features:** Franchise model, online sales, fleet ownership, drive-through
- **Sustainability Metrics:** Electric vehicle %, ESG programs, sustainability awards, R&D spending

Dataset Shape:

- 3,600+ brands across multiple industries and parent companies
- 77+ features

- 1. **Data Loading and Preprocessing:** Loading, handling missing values, etc
- 2. **Feature Engineering:** Creating additional features, etc
- 3. **Feature Selection:** Select relevant features
- 4. **Feature Scaling:** Standardization
- 5. **Dimensionality Reduction:**
 - **PCA** for variance reduction and **latent feature discovery**
 - t-SNE for 2D visualization
- 6. **Hyperparameter Tuning:**
 - Grid search and multi-metric evaluation
- 7. **Clustering with Optimized Parameters:**
 - K-Means, Hierarchical, DBSCAN (outliers/noise)
- 8. **Cluster Interpretation:** Profile segments
- 9. **Validation:** Compare algorithms

Feature Engineering

Creating features that capture ESG risk and sustainability patterns.

```
1 greenwashing_combined = (  
    initial_greenwashing_risk + ... +  
    greenwashing_factor_4).clip(lower=0,  
    upper=1)  
2  
3 other_risks_combined = full_df['  
    deforestation_risk'].fillna(0) + ...  
    + (full_df['  
    product_portfolio_diversity'].fillna  
    (0) / 10)  
4  
5 env_risk_combined = greenwashing_combined  
    + other_risks_combined  
6 full_df['environmental_risk_score'] =  
    env_risk_combined
```

Figure: Feature Engineering

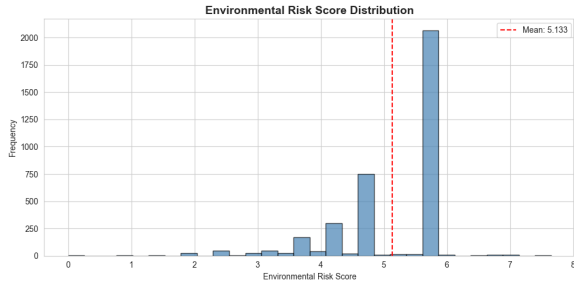


Figure: Environmental Risk Score plot

Features Selection for Clustering + Encoding

Focus on features that capture **Company-level ESG**, **Brand demographics**, **Operational characteristics**, **Sustainability metrics**.

```
1 # Define feature groups
2 company_esg_features = ['scope12_total', ..., '
   positive_innovation_risk']
3 demographic_features = ['income_middle', '
   income_high', 'income_premium', '
   income_snap_support']
4 lifestyle_features = []
5 operational_features = ['online_sales', ..., '
   supply_chain_localization_encoded']
6 sustainability_features = ['
   electric_vehicles_percent', ..., '
   sustainability_actions_encoded']
7 divergence_features = ['sustainability_positioning',
   ..., 'esg_premium_divergence']
8 market_features = ['customer_loyalty_index', ..., '
   year_of_foundation']
9 interaction_features = ['emissions_intensity', ...,
   'env_risk_concentration']
10
11 feature_columns = (company_esg_features +
   demographic_features + lifestyle_features +
   operational_features + sustainability_features
   + divergence_features + market_features +
   interaction_features)
```

Figure: Logical grouping of features

```
1 def encode_esg_programs(val):
2 def encode_sustainability_actions(val):
3 def encode_fossil_fuel(val):
4 def encode_diversity(val):
5 def encode_localization(val):
6
7 full_df['esg_programs_encoded'] = full_df['
   esg_programs'].apply(encode_esg_programs)
8 full_df['sustainability_actions_encoded'] = full_df[
   'sustainability_actions'].apply(
   encode_sustainability_actions)
9 full_df['fossil_fuel_reliance_encoded'] = full_df['
   fossil_fuel_reliance'].apply(encode_fossil_fuel)
10 full_df['product_portfolio_diversity_encoded'] =
   full_df['product_portfolio_diversity'].apply(
   encode_diversity)
11 full_df['supply_chain_localization_encoded'] =
   full_df['supply_chain_localization'].apply(
   encode_localization)
```

Figure: Encoding categorical features

Feature Engineering and Scaling

Examine alignment/misalignment between:

- Brand sustainability positioning vs. company environmental performance
- Premium pricing vs. actual ESG investment
- Marketing claims vs. operational reality
- Check the **Appendix** for more details

```
1 full_df['sustainability_positioning'] = ( full_df['electric_vehicles_percent'] * 0.4 + (0.3 if
    full_df['esg_programs'].count() > 0 else 0) + full_df['major_sustainability_award_last5y'].astype
    (float) * 0.3)
2 emissions_normalized = (full_df['scope12_total'] - full_df['scope12_total'].min()) / (full_df['
    scope12_total'].max() - full_df['scope12_total'].min())
3 full_df['company_env_risk'] = (emissions_normalized * 0.5 + full_df['environmental_risk_score'] *
    0.5)
4 full_df['sustainability_divergence'] = full_df['sustainability_positioning'] - (1 - full_df['
    company_env_risk'])
5 full_df['premium_positioning'] = (full_df['age_seniors'].astype(float) * 0.4 + full_df['
    income_premium'].astype(float) * 0.6)
6 full_df['esg_premium_divergence'] = full_df['premium_positioning'] - (1 - full_df['company_env_risk'
    ])
```

Figure: Feature Engineering and Scaling

Divergence and Hipocresy Analysis Result

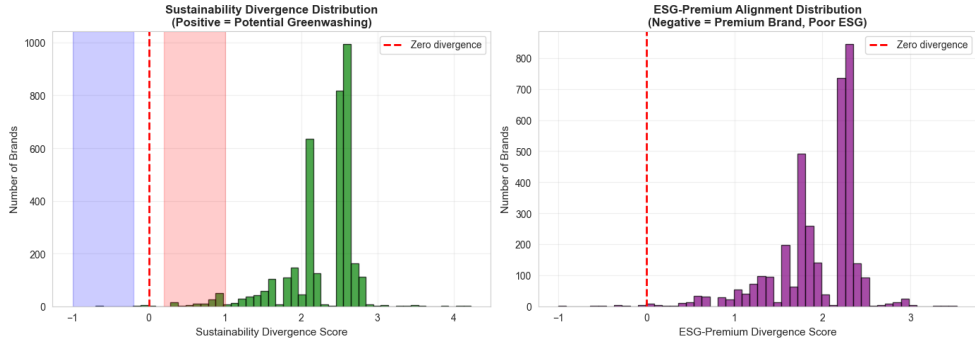


Figure: Divergence plot

Dimensionality Reduction

Applying PCA to reduce dimensionality

```
1 X_pca, pca_model = reducer.apply_pca(  
    X_scaled, n_components=0.95)
```

```
1 variance_df = reducer.  
    get_variance_dataframe()  
2 fig = viz.plot_variance_explained(  
    variance_df, n_components=min(15,  
    X_pca.shape[1]))  
3 plt.show()
```

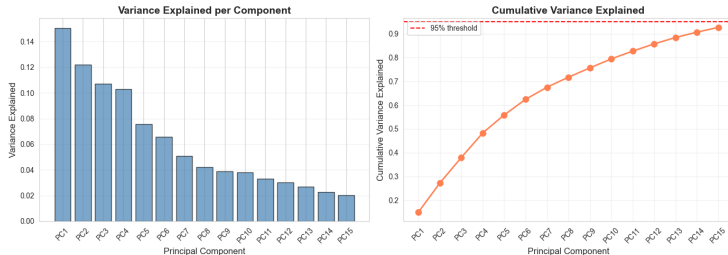


Figure: Variance Explained after PCA plots

Top Features per Component

```
1 # Analyze top features per component
2 top_features = reducer.
   get_top_features_per_component(
   feature_columns, n_top=5)
3
4 print("Top 5 features for first 3
   principal components:")
5 print(top_features[top_features['
   Component'].isin(['PC1', 'PC2',
   'PC3'])])
6
7 # Plot feature importance for PC1
8 fig = viz.plot_feature_importance(
   top_features, component='PC1',
   n_features=10)
9 plt.show()
```

Figure: Visualization code

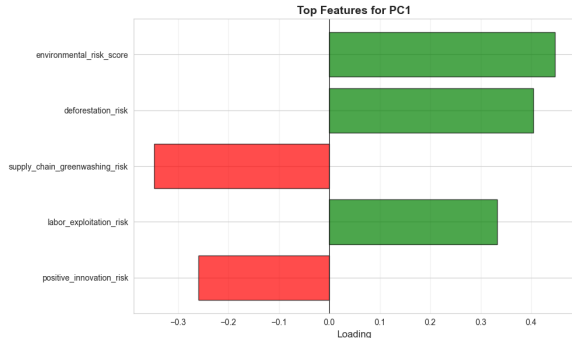


Figure: Visualization

Latent Feature Interpretation

Latent features discovered. Check the **Appendix** for the code.

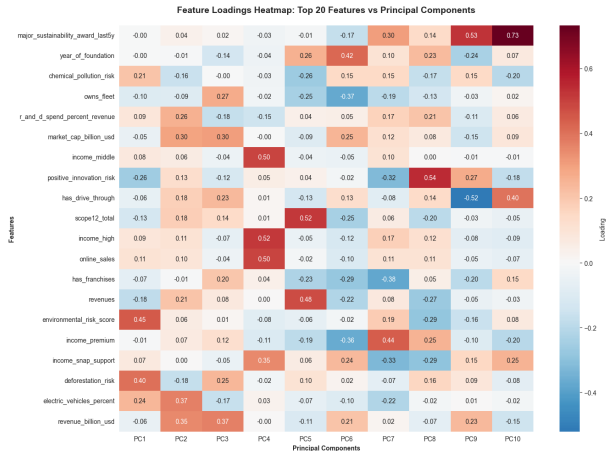


Figure: Latent Features

PCA Visualizing the Latent Space

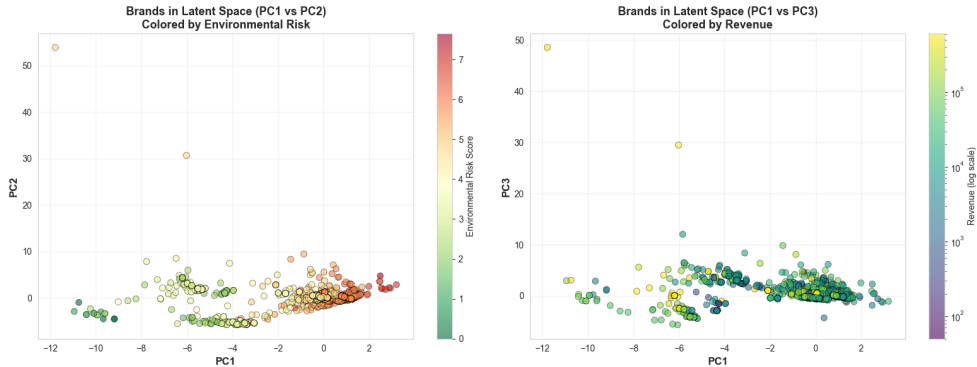


Figure: Plots in the Latent Space

Hyperparameter Tuning - K-Means

Systematic hyperparameter tuning on K-Means, Hierarchical, and DBSCAN Best parameters used in final clustering and analysis.

```
1 tuner = HyperparameterTuner(verbose=True)
```

Figure: Instantiating the HP tuner, check the Appendix and the git repo for more details

```
1 kmeans_tuning_results = tuner.tune_kmeans(  
2     X_pca,  
3     param_grid=kmeans_param_grid,  
4     metric=OPTIMIZATION_METRIC  
5 )  
6  
7 fig = viz.plot_hyperparameter_tuning(  
8     kmeans_tuning_results,  
9     algorithm='kmeans',  
10    metrics=['silhouette', 'calinski_harabasz', '  
11            'davies_bouldin', 'inertia']  
12 )  
13 plt.show()
```

Figure: HP tuning K-Means

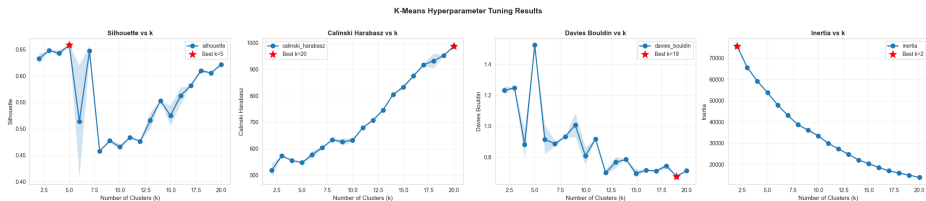


Figure: HP tuning K-Means - plot

Hyperparameter Tuning - Hierarchical

```
1 hierarchical_tuning_results = tuner.tune_hierarchical(  
2     X_pca,  
3     param_grid=hierarchical_param_grid,  
4     metric=OPTIMIZATION_METRIC  
5 )  
6 fig = viz.plot_hyperparameter_tuning(  
7     hierarchical_tuning_results,  
8     algorithm='hierarchical',  
9     metrics=['silhouette', 'calinski_harabasz', 'davies_bouldin']  
10 )  
11 plt.show()
```

Figure: HP Tuning Hierarchical

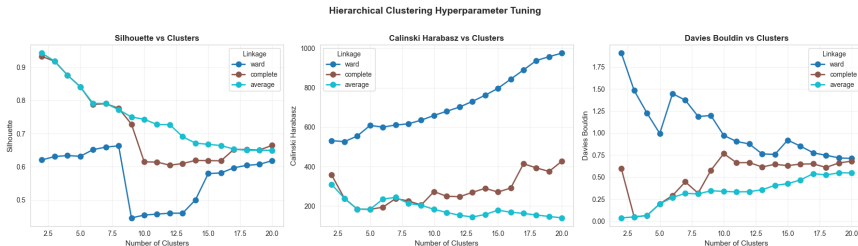


Figure: HP Tuning Hierarchical - plot

Hyperparameter Tuning - DBSCAN

```
1 dbscan_tuning_results = tuner.tune_dbscan(  
2     X_tsne,  
3     param_grid=dbscan_param_grid,  
4     metric=OPTIMIZATION_METRIC,  
5     min_cluster_size=10  
6 )  
7 fig = viz.plot_hyperparameter_tuning(  
8     dbscan_tuning_results,  
9     algorithm='dbscan'  
10 )  
11 plt.show()
```

Figure: HP Tuning DBSCAN

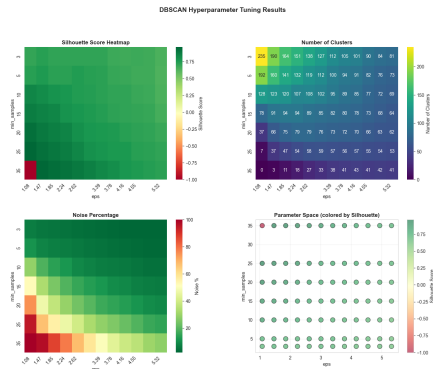


Figure: HP Tuning DBSCAN - plot

Model Selection Based on Hyperparameter Tuning Results

Models and best hyperparameters comparison:

Algorithm	Silhouette	Calinski-Harabasz	Davies-Bouldin	n_clusters	Noise %	Best Hyperparameters
Hierarchical	0.9429	310.03	0.0382	2	0%	linkage=average, metric=euclidean
DBSCAN	0.6829	3000.96	0.4708	10	13.3%	eps=1.4671, min_samples=20
K-Means	0.6581	547.22	1.5229	5	0%	n_clusters=5, n_init=10, max_iter=100

Model Selection Rationale:

- **Hierarchical:** Best overall with highest silhouette (0.94) and lowest Davies-Bouldin (0.04).
- **DBSCAN:** Highest Calinski-Harabasz (3000.96), useful for outlier detection.
- **K-Means:** Balanced baseline with moderate cluster count (k=5).

Decision: Primary focus on Hierarchical, with DBSCAN for outlier analysis.

Best Hierarchical

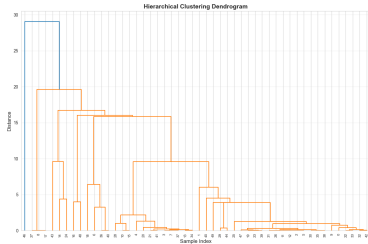


Figure: Dendrogram



Figure: Segmentation

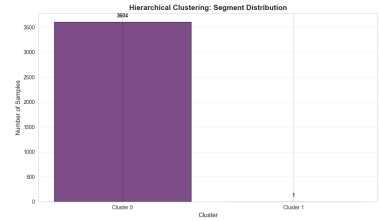


Figure: By cluster

Best DBSCAN

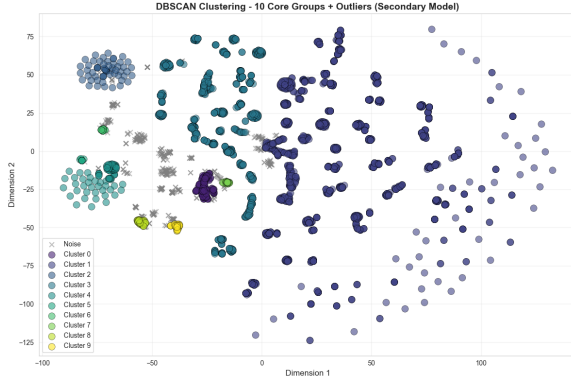


Figure: Segmentation

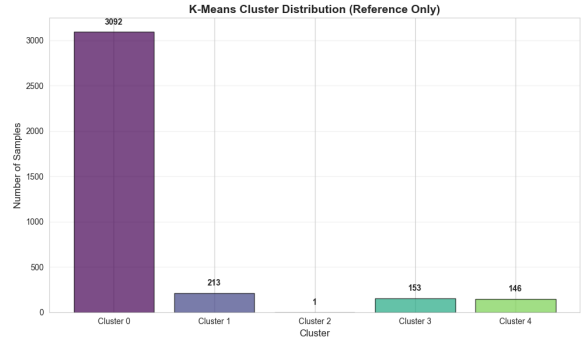


Figure: K-Means as baseline

Clusters Interpretation I

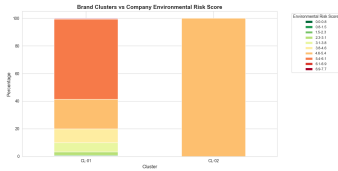


Figure: Brands x Risk score



Figure: Brands x Risk bucket

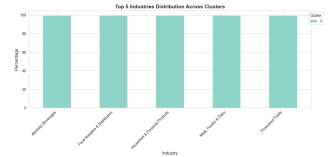


Figure: Industries Distribution

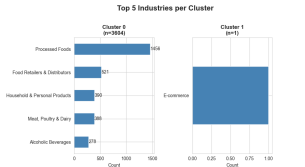


Figure: Top 5 Industries per Cluster



Figure: Top 5 Companies per Cluster

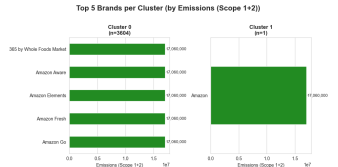


Figure: Top 5 Brands per Cluster - By scope

Clustering Results: My Honest Assessment

Algorithm	Clusters	Silhouette	Distribution	Key Observation
K-Means	5	0.658	3,092 / 213 / 153 / 146 / 1	One dominant cluster (86%) + one singleton
Hierarchical	2	0.943	3,604 / 1	Effectively a single cluster + 1 outlier
DBSCAN	10	0.683	Various (20-1,777)	13.3% noise (481 brands)

The results diverge from my initial expectations. Rather than discovering distinct, well-separated market segments, the algorithms suggest that **the brand space is largely continuous rather than discretely segmented.**

Conclusion - Failure or Not Failure?

- **The Finding IS the Insight:** Brands don't segment into distinct groups; differentiation operates on continuous spectra, not discrete categories.
- **Rigorous Methodology:** 25 features, 255 HP combinations tested, PCA with interpretable latent dimensions.
- **Outlier Detection:** DBSCAN identified 481 outlier brands (13.3%) with higher revenue and emissions; worth individual investigation.
- **Future Work:** Clean dataset ready for Deep Learning analysis.

Acknowledgement of Limitations

- **Feature Space:** The 25 features may not capture the dimensions that truly differentiate brands in meaningful ways
- **Data Homogeneity:** Brands within this ESG database proved to be inherently similar; The reason is that they were pre-selected for sustainability reporting purposes
- **Scale Imbalance:** My finding of one dominant cluster across all algorithms confirms at least one of the hypotheses: the data is either homogeneous or suboptimal

The project is not a waste of effort. While the clustering results fell short of discovering dramatically distinct brand segments, the work accomplished legitimate goals:

- Demonstrated mastery of unsupervised learning techniques and proper application
- Revealed an honest truth about the data
- Produced a clean, engineered dataset ready for advanced analysis
- Identified the class imbalance problem that GANs can directly address

The journey continues in the Deep Learning project.

- ① **Jain, A. K. (2010).** *Data clustering: 50 years beyond K-means*. Pattern Recognition Letters, 31(8), 651-666.
- ② **Hartigan, J. A., & Wong, M. A. (1979).** *Algorithm AS 136: A K-Means Clustering Algorithm*. Journal of the Royal Statistical Society. Series C, 28(1), 100-108.
- ③ **Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996).** *A density-based algorithm for discovering clusters in large spatial databases with noise*. KDD-96 Proceedings, 226-231.
- ④ **van der Maaten, L., & Hinton, G. (2008).** *Visualizing Data using t-SNE*. JMLR, 9, 2579-2605.
- ⑤ **Jolliffe, I. T., & Cadima, J. (2016).** *Principal component analysis: a review and recent developments*. Phil. Trans. R. Soc. A, 374(2065).
- ⑥ **Rousseeuw, P. J. (1987).** *Silhouettes: A graphical aid to the interpretation and validation of cluster analysis*. J. Comput. Appl. Math., 20, 53-65.

- 7 **Pedregosa, F., et al. (2011).** *Scikit-learn: Machine Learning in Python*. JMLR, 12, 2825-2830.
- 8 **scikit-learn Documentation.** *Clustering*.
<https://scikit-learn.org/stable/modules/clustering.html>
- 9 **scikit-learn Documentation.** *Decomposition (PCA)*.
<https://scikit-learn.org/stable/modules/decomposition.html>
- 10 **Caliński, T., & Harabasz, J. (1974).** *A dendrite method for cluster analysis*. Comm. in Statistics, 3(1), 1-27.
- 11 **Davies, D. L., & Bouldin, D. W. (1979).** *A Cluster Separation Measure*. IEEE Trans. PAMI, 1(2), 224-227.
- 12 **Google Developers.** *Clustering in Machine Learning*.
<https://developers.google.com/machine-learning/clustering>
- 13 **Towards Data Science.** *The 5 Clustering Algorithms Data Scientists Need to Know*.

- 14 StatQuest with Josh Starmer. *K-means clustering*. YouTube.
- 15 The GHG Shopper application. <https://ghgshopper.org>

The implementation uses modular Python files:

- `clustering_models.py` - K-Means, Hierarchical, DBSCAN + HyperparameterTuner
- `dimensionality_reduction.py` - PCA and t-SNE with component analysis
- `visualization_utils.py` - Plotting and visualization functions

Appendix - Feature Engineering and Scaling

```
1 full_df['emissions_intensity'] = full_df['scope12_total'] / (  
    full_df['revenues'] + 1)  
2 full_df['esg_investment_ratio'] = full_df['  
    r_and_d_spend_percent_revenue'] / (full_df['  
    market_cap_billion_usd'] + 1)  
3 current_year = 2025  
4 full_df['brand_age'] = current_year - full_df['year_of_foundation'  
    ']  
5 full_df['brand_age_log'] = np.log1p(full_df['brand_age'])  
6 full_df['revenue_per_employee'] = (full_df['revenues'] / (full_df  
    ['employees'] + 1)) / 1000  
7 age_cols = ['age_prenatal', 'age_0_5', 'age_6_12', 'age_teens', '  
    age_young_adults', 'age_seniors']  
8 full_df['demographic_breadth'] = full_df[age_cols].sum(axis=1)  
9 income_cols = ['income_low', 'income_middle', 'income_high', '  
    income_premium']  
10 full_df['income_diversity'] = full_df[income_cols].sum(axis=1)  
11 lifestyle_cols = ['lifestyle_family', 'lifestyle_youth', '  
    lifestyle_seniors',  
12     'lifestyle_health_focused', '  
    lifestyle_convenience',  
13     'lifestyle_tech_savvy', '  
    lifestyle_sustainability_conscious']  
14 full_df['lifestyle_complexity'] = full_df[lifestyle_cols].sum(  
    axis=1)  
15 full_df['operational_complexity'] = (  
16     full_df['has_franchises'].astype(float) +  
17     full_df['owns_fleet'].astype(float) +  
18     full_df['has_drive_through'].astype(float) +
```

```
1     full_df['online_sales'].astype(float) +  
2     full_df['product_portfolio_diversity']  
3 )  
4 full_df['market_positioning'] = (  
5     np.log1p(full_df['market_cap_billion_usd']) * 0.5 +  
6     full_df['customer_loyalty_index'] * 0.3 +  
7     full_df['branding_innovation_level'] * 0.2  
8 )  
9 env_risk_cols = ['deforestation_risk', 'chemical_pollution_risk',  
    'fossil_fuel_reliance_encoded']  
10 full_df['env_risk_concentration'] = full_df[env_risk_cols].mean(  
    axis=1)  
11  
12 interaction_cols = ['emissions_intensity', 'esg_investment_ratio',  
    'brand_age_log',  
13     'revenue_per_employee', 'demographic_breadth',  
    'income_diversity',  
14     'lifestyle_complexity', '  
    operational_complexity', '  
    market_positioning',  
15     'env_risk_concentration']  
16 print(full_df[interaction_cols].describe().round(3))  
17  
18 print(full_df.describe())
```

Appendix - Interpreting Brand-Company Divergence Scores

These divergence scores measure the alignment (or misalignment) between a brand's positioning and its parent company's actual environmental performance. **Sustainability Divergence:**

- **Positive values (> 0):** Brand positions strongly on sustainability while parent company has higher environmental risk → **High divergence** (brand is more environmentally conscious than company)
- **Negative values (< 0):** Brand under-emphasizes sustainability despite parent company's better environmental performance → **Reverse divergence** (company is more environmentally conscious than brand positioning suggests)
- **Near zero (0):** Brand positioning aligns with company environmental performance → **Aligned**

ESG-Premium Divergence:

- **Positive values:** Premium brand from company with stronger ESG performance → **Aligned premium positioning**
- **Negative values:** Premium brand from company with weaker ESG performance → **Misaligned premium positioning**
- **Near zero:** Premium positioning matches ESG performance level → **Consistent**

Key Insight: These scores help identify patterns in how brands strategically position themselves relative to their parent companies, which is central to understanding brand-company environmental alignment.

Visualizations

```
1 divergence_cols = ['brand_name', '
    sustainability_positioning', '
    company_env_risk', '
    sustainability_divergence']
2 top_divergence = full_df.nlargest
    (10, 'sustainability_divergence'
    )[divergence_cols]
3 print(top_divergence.to_string(index
    =False))
4
5 print("\n" + "="*60)
6 print("BRANDS WITH LOWEST
    SUSTAINABILITY DIVERGENCE")
7 print("(Parent company more
    environmentally conscious than
    brand positioning)")
8 print("="*60)
9 bottom_divergence = full_df.
    nsmallest(10, '
    sustainability_divergence')[
    divergence_cols]
```

```
1 print("\n" + "="*60)
2 print("PREMIUM-ESG DIVERGENCE
    ANALYSIS")
3 print("="*60)
4 prem_div_cols = ['brand_name', '
    premium_positioning', '
    esg_premium_divergence']
5 print("\nHighest Premium-ESG
    Divergence (Premium brands from
    strong ESG companies):")
6 print(full_df.nlargest(10, '
    esg_premium_divergence')[
    prem_div_cols].to_string(index=
    False))
7
8 print("\nLowest Premium-ESG
    Divergence (Premium brands from
    weak ESG companies):")
9 print(full_df.nsmallest(10, '
    esg_premium_divergence')[
    prem_div_cols].to_string(index=
```

Visualizations

```
1 fig, axes = plt.subplots(1, 2,  
    figsize=(14, 5))  
2  
3 axes[0].hist(full_df['  
    sustainability_divergence'],  
    bins=50, edgecolor='black',  
    alpha=0.7, color='green')  
4 axes[0].axvline(x=0, color='red',  
    linestyle='--', linewidth=2,  
    label='Zero divergence')  
5 axes[0].set_xlabel('Sustainability  
    Divergence Score', fontsize=11)  
6 axes[0].set_ylabel('Number of Brands  
    ', fontsize=11)  
7 axes[0].set_title('Sustainability  
    Divergence Distribution\n(  
    Positive = Potential  
    Greenwashing)',  
8  
    fontsize=12,  
    fontweight='  
    bold'))
```

```
1 axes[1].hist(full_df['  
    esg_premium_divergence'], bins  
    =50, edgecolor='black', alpha  
    =0.7, color='purple')  
2 axes[1].axvline(x=0, color='red',  
    linestyle='--', linewidth=2,  
    label='Zero divergence')  
3 axes[1].set_xlabel('ESG-Premium  
    Divergence Score', fontsize=11)  
4 axes[1].set_ylabel('Number of Brands  
    ', fontsize=11)  
5 axes[1].set_title('ESG-Premium  
    Alignment Distribution\n(  
    Negative = Premium Brand, Poor  
    ESG)',  
6  
    fontsize=12,  
    fontweight='  
    bold')  
7 axes[1].legend()  
8 axes[1].grid(True, alpha=0.3)  
9
```

Visualizations

```
1 n_components_to_interpret = min(10,
    X_pca.shape[1])
2 component_interpretations = {}
3 for i in range(
    n_components_to_interpret):
4     pc_name = f'PC{i+1}'
5
6     top_features_pc = top_features[
7         (top_features['Component']
8          == pc_name) &
9         (top_features['Rank'] <= 5)
10    ].sort_values('Abs_Loading',
11                  ascending=False)
12
13     print(f"\n{pc_name} (explains {
14           variance_df.iloc[i]['
15             Variance_Explained']:.1f}%
16           variance):")
17
18     print("-" * 80)
```

```
1     for _, row in top_features_pc.
2       iterrows():
3         direction = 'positively' if
4           row['Loading'] > 0 else
5           'negatively'
6         print(f"    {row['Feature']}:50
7              s} {row['Loading']:.3f}
8              ({direction})")
9
10    component_names = {
11        'PC1': 'To be determined by
12              loadings', # Placeholder
13        'PC2': 'To be determined by
14              loadings',
15        'PC3': 'To be determined by
16              loadings',
17        'PC4': 'To be determined by
18              loadings',
19        'PC5': 'To be determined by
20              loadings',
```


Appendix - Latent Features Interpretation - Code

```
1 n_top_features = 20
2 top_feature_names = (top_features.
    groupby('Feature')['Abs_Loading']
    ).max().sort_values(ascending=
    False).head(n_top_features).
    index)
3
4 loading_matrix = pd.DataFrame(
5     pca_model.components_[
6         n_components_to_interpret,
7         :],
8     columns=feature_columns,
9     index=[f'PC{i+1}' for i in range
10            (n_components_to_interpret)]
11 ) [top_feature_names].T
```

```
1 fig, ax = plt.subplots(figsize=(14,
2     10))
3 sns.heatmap(loading_matrix, cmap='
4     RdBu_r', center=0, annot=True,
5     fmt='.2f', cbar_kws={'label': '
6     Loading'}, ax=ax)
7 ax.set_title(f'Feature Loadings
8     Heatmap: Top {n_top_features}
9     Features vs Principal Components
10     ', fontweight='bold', fontsize
11     =14, pad=20)
12 ax.set_xlabel('Principal Components'
13     , fontweight='bold')
14 ax.set_ylabel('Features', fontweight
15     ='bold')
16 plt.tight_layout()
17 plt.show()
```

Visualizations

```
1 pca_scores_df = pd.DataFrame(  
2     X_pca[:, :  
3     n_components_to_interpret],  
4     columns=[f'PC{i+1}' for i in  
5         range(  
6             n_components_to_interpret)],  
7     index=full_df.index  
8 )  
9  
10 pca_scores_df['brand_name'] =  
11     full_df['brand_name']  
12  
13 pca_scores_df['company_name'] =  
14     full_df['company_name']  
15  
16 pca_scores_df['industry_name'] =  
17     full_df['industry_name']
```

```
1 for pc in ['PC1', 'PC2', 'PC3']:  
2     print(f"\n{pc}:")  
3     print("\nHighest scoring brands:  
4         ")  
5     print(pca_scores_df.nlargest(5,  
6         pc)[['brand_name', '  
7             company_name', pc]])  
8  
9     print("\nLowest scoring brands:"  
10         )  
11     print(pca_scores_df.nsmallest(5,  
12         pc)[['brand_name', '  
13             company_name', pc]])  
14     print("-" * 80)  
15  
16 full_df[[f'PC{i+1}' for i in range(  
17     n_components_to_interpret)]] =  
18     X_pca[:, :  
19     n_components_to_interpret]
```

Visualizations

```
1 fig, axes = plt.subplots(1, 2,  
2     figsize=(16, 6))  
3 scatter1 = axes[0].scatter(  
4     pca_scores_df['PC1'],  
5     pca_scores_df['PC2'], c=full_df['  
6     'environmental_risk_score'],  
7     cmap='RdYlGn_r', alpha=0.6, s  
8     =50, edgecolors='black',  
9     linewidth=0.5)  
10 axes[0].set_xlabel('PC1', fontweight  
11     = 'bold', fontsize=12)  
12 axes[0].set_ylabel('PC2', fontweight  
13     = 'bold', fontsize=12)  
14 axes[0].set_title('Brands in Latent  
15     Space (PC1 vs PC2)\nColored by  
16     Environmental Risk', fontweight=  
17     'bold', fontsize=13)  
18 axes[0].grid(True, alpha=0.3)  
19 plt.colorbar(scatter1, ax=axes[0],  
20     label='Environmental Risk Score')
```

```
1 scatter2 = axes[1].scatter(  
2     pca_scores_df['PC1'],  
3     pca_scores_df['PC3'], c=full_df['  
4     'revenues'], cmap='viridis',  
5     alpha=0.6, s=50, edgecolors='  
6     black', linewidth=0.5, norm=plt.  
7     cm.colors.LogNorm())  
8 axes[1].set_xlabel('PC1', fontweight  
9     = 'bold', fontsize=12)  
10 axes[1].set_ylabel('PC3', fontweight  
11     = 'bold', fontsize=12)  
12 axes[1].set_title('Brands in Latent  
13     Space (PC1 vs PC3)\nColored by  
14     Revenue',  
15     fontweight='bold',  
16     fontsize=13)  
17 axes[1].grid(True, alpha=0.3)  
18 plt.colorbar(scatter2, ax=axes[1],  
19     label='Revenue (log scale)')  
20 plt.tight_layout()
```

Appendix - Business Interpretation of Latent Features

Would this ever be used in the real business world, the Clusters should be named to represent what they mean. For the scope of this project I simply named them sequentially with placeholders, as I don't claim to be an ESG expert.

Using the following names for Principal Components:

- **PC1:** PC-1
- **PC2:** PC-2
- **PC3:** PC-2

Using the following names for Clusters:

- **CL_01:** CL_01
- **CL_02:** CL_02
- **CL_n:** CL_n

```
1 # Apply t-SNE on PCA-reduced data (recommended workflow)  
2 # This reduces computational cost and noise  
3 X_tsne, tsne_model = reducer.apply_tsne(X_pca, n_components=2)
```

Figure: Applying t-SNE on PCA-reduced data

Hyperparameter Tuning variables for ranges

```
1 # K-MEANS CONFIGURATION
2 KMEANS_MIN_CLUSTERS = 2
3 KMEANS_MAX_CLUSTERS = 21
4 KMEANS_N_INIT = [10, 20]
5 KMEANS_MAX_ITER = [100, 300, 500]
6 KMEANS_BALANCED_MIN = 5
7 KMEANS_BALANCED_MAX = 15
8 KMEANS_BALANCED_MIN_SILHOUETTE = 0.2
9 # HIERARCHICAL CLUSTERING CONFIGURATION
10 HIER_MIN_CLUSTERS = 2
11 HIER_MAX_CLUSTERS = 21
12 HIER_LINKAGE_METHODS = ['ward', 'complete', 'average']
13 HIER_METRIC = ['euclidean']
14 HIER_BALANCED_MIN = 2
15 HIER_BALANCED_MAX = 15
16 HIER_BALANCED_MIN_SILHOUETTE = 0.2
17 # DBSCAN CONFIGURATION
18 DBSCAN_EPS_MIN_PERCENTILE = 60
19 DBSCAN_EPS_MAX_PERCENTILE = 95
20 DBSCAN_EPS_N_VALUES = 12
21 DBSCAN_MIN_SAMPLES = [3, 5, 10, 15, 20, 25, 35]
22 DBSCAN_MIN_CLUSTERS = 5
```

Hyperparameter Tuning - Continued

```
1  if 'kmeans' in tuner.best_params:
2      kmeans_best = tuner.get_best_params('kmeans')
3      print("\nK-Means:")
4      print(f"    n_clusters: {int(kmeans_best['n_clusters'])}")
5      print(f"    n_init: {int(kmeans_best['n_init'])}")
6      print(f"    Silhouette Score: {kmeans_best['silhouette']:.4f}")
7
8  if 'hierarchical' in tuner.best_params:
9      hier_best = tuner.get_best_params('hierarchical')
10     print("\nHierarchical Clustering:")
11     print(f"    n_clusters: {int(hier_best['n_clusters'])}")
12     print(f"    linkage: {hier_best['linkage']}")
13     print(f"    Silhouette Score: {hier_best['silhouette']:.4f}")
14
15  if 'dbscan' in tuner.best_params:
16      dbscan_best = tuner.get_best_params('dbscan')
17      print("\nDBSCAN:")
18      print(f"    eps: {dbscan_best['eps']:.4f}")
19      print(f"    min_samples: {int(dbscan_best['min_samples'])}")
20      print(f"    n_clusters: {int(dbscan_best['n_clusters'])}")
21      print(f"    noise_pct: {dbscan_best['noise_pct']:.1f}%")
22      print(f"    Silhouette Score: {dbscan_best['silhouette']:.4f}")
```

Hyperparameter Tuning - Continued

```
1
2 best_params = tuner.get_best_params(
3     'kmeans')
4 n_clusters = int(best_params['
5     n_clusters'])
6 print(f"Using tuned K-Means
7     parameters: n_clusters={
8         n_clusters}")
9 kmeans_model, kmeans_labels =
10     clusterer.kmeans_clustering(
11         X_pca,
12         n_clusters=n_clusters,
13         find_optimal=False # Use
14             specified k from tuning
15     )
```

```
1 best_params = tuner.get_best_params(
2     'hierarchical')
3 n_clusters = int(best_params['
4     n_clusters'])
5 print(f"Using tuned Hierarchical
6     parameters: n_clusters={
7         n_clusters}, linkage={
8         best_params['linkage']}")
9
10 # Update config for hierarchical
11     clustering to use tuned linkage
12 clusterer.config['hierarchical']['
13     linkage'] = best_params['linkage']
14
15 clusterer.config['hierarchical']['
16     n_clusters'] = n_clusters
17
18 hierarchical_model,
19     hierarchical_labels = clusterer.
20     hierarchical_clustering(
21         X_pca,
```


Hierarchical

```
1 # Plot dendrogram
2 if X_pca.shape[0] <= 50:
3     fig = viz.plot_dendrogram(X_pca, method='ward')
4     plt.show()
5 else:
6     #plotting a random subset
7     subset_indices = np.random.choice(X_pca.shape[0], size=50, replace=False)
8     subset = X_pca[subset_indices]
9     fig = viz.plot_dendrogram(subset, method='ward')
10    plt.show()
```

```
1 # Visualize hierarchical clusters (PRIMARY MODEL)
2 # Focus on this visualization for business interpretation
3
4 fig = viz.plot_clusters_2d(
5     X_tsne,
6     hierarchical_labels,
7     title="Hierarchical Clustering - 5 Brand Segments (Primary Model)",
8     brand_names=full_df['brand_name'].tolist(),
9     show_labels=False
10 )
11 plt.show()
12
13 # Plot cluster sizes
14 fig = viz.plot_cluster_sizes(
15     hierarchical_labels,
16     title="Hierarchical Clustering: Segment Distribution"
17 )
18 plt.show()
19
20 print(f"\n{'='*80}")
21 print("HIERARCHICAL CLUSTERING RESULTS (PRIMARY MODEL)")
22 print(f"{'='*80}")
23 print(f"Number of clusters: {len(set(hierarchical_labels))}")
24 print(f"{'='*80}")
```

DBSCAN I

```
1 # DBSCAN (density-based) with tuned
  hyperparameters
2 # Secondary model - excellent for
  outlier detection
3
4 if 'tuner' in locals() and 'dbscan'
  in tuner.best_params:
5     # Use tuned parameters
6     dbscan_params = tuner.
        best_params['dbscan']
7     eps = dbscan_params['eps']
8     min_samples = dbscan_params['
        min_samples']
9 else:
10     # Default parameters
11     eps = 2.5
12     min_samples = 25
13
14 dbscan_labels = clusterer.fit_dbscan
    (X_pca, eps=eps, min_samples=int
    (min_samples))
```

```
1 # Compare clustering algorithms
2
3 comparison_df = clusterer.
    compare_algorithms()
4 print("\nClustering Algorithm
    Comparison:")
5 print("=" * 80)
6 print(comparison_df.to_string(index=
    False))
7 print("=" * 80)
```

Visualization

```
1 # Visualize DBSCAN clusters
2 fig = viz.plot_clusters_2d(
3     X_tsne,
4     dbscan_labels,
5     title="DBSCAN Clustering (t-SNE Visualization)",
6     show_legend=True # Hide legend if there are too many clusters
7 )
8 plt.show()
9
10 # Examine outliers (noise points)
11 outliers = full_df[full_df['dbscan_cluster'] == -1]
12 print(f"\nDBSCAN identified {len(outliers)} outlier brands:")
13 if len(outliers) > 0:
14     display_cols = ['brand_name', 'company_name', 'industry_name']
15     if 'initial_greenwashing_level' in outliers.columns:
16         display_cols.append('initial_greenwashing_level')
17     available_cols = [c for c in display_cols if c in outliers.columns]
18     print(outliers[available_cols].head(10))
```

- Successfully integrated **3,605 brands** across **243 unique companies** and **14 industries**
- Combined multiple data sources: industry ESG data, company emissions, brand demographics
- Engineered **environmental_risk_score** metric from multiple greenwashing and risk factors
- Handled missing values and normalized **25 features** across different scales

Dimensionality Reduction

- Applied **PCA** to reduce noise and preserve 95% of variance (25 \rightarrow 17 dimensions)
- Applied **t-SNE** for 2D visualization of high-dimensional brand space
- Identified key principal components driving brand differentiation:
- **PC1 (15.07%)**: Environmental Risk Dimension (deforestation, labor exploitation, greenwashing)
- **PC2 (12.21%)**: Scale & Sustainability Dimension (EV%, loyalty index, revenue, employees)
- **PC3 (10.71%)**: Corporate Size Dimension (revenue, employees, market cap)
- **PC4 (10.30%)**: Income Demographics Dimension (income targeting patterns)

- Implemented and compared **3 clustering algorithms** with systematic hyperparameter tuning
- Applied multi-metric evaluation (Silhouette, Calinski-Harabasz, Davies-Bouldin)
- Tested **114 K-Means**, **57 Hierarchical**, and **84 DBSCAN** parameter combinations

- Exportable CSV with cluster assignments for all brands
- Comprehensive visualizations (PCA variance, t-SNE plots, cluster profiles)

Future Enhancement: Deep Learning with GANs (Potentially Pretrained) I

This unsupervised learning project serves as a foundation for the upcoming **Deep Learning final project**, where the plan is to:

Use Clusters as Pseudo-Labels

The K-Means 5-cluster solution, despite its imbalance, provides initial class labels for supervised approaches. Address Class Imbalance with GANs

Generative Adversarial Networks can synthesize new samples for underrepresented clusters:

- The K-Means clusters 1-4 (213, 153, 146, 1 brands) are severely underrepresented vs Cluster 0 (3,092)
- GANs can generate realistic *synthetic* brand profiles to balance training data
- This will enable training to be more robust, and classifiers will work better for brand segment prediction

Future Enhancement: Deep Learning with GANs (Potentially Pretrained) II

Improved Segmentation:

With balanced, augmented data, deep learning models may discover finer-grained patterns that traditional clustering missed. Validation Pipeline:

The augmented dataset can train and evaluate classification models, with held-out original data serving as ground truth. Another round of Unsupervised Learning:

At the end, the generated augmented dataset can be re-analyzed using this pipeline to evaluate the Cluster distribution.