

Synthetic Brand Generation V2 - Enhanced with Ensemble Methods

University of Colorado Boulder - CSCA 5642: Introduction to Deep Learning

Generation

University of Colorado Boulder

December 6, 2025

Outline

- 1 Synthetic Brand Generation V2 - Enhanced with Ensemble Methods
- 2 Problem Statement
- 3 ML Approach Architecture
- 4 Phase 1: Setup Installation
- 5 Configuration
- 6 Phase 0: Hyperparameter Tuning (Optional)
- 7 Phase 1: Data Preparation
- 8 Phase 2: Tabular Ensemble Training
- 9 Phase 3: LLM Ensemble Training
- 10 Phase 4: Synthetic Data Generation
- 11 Phase 5: Synthetic Data Quality Evaluation
- 12 Clean Up
- 13 Conclusion
- 14 References Bibliography
- 15 Appendix - Mermeid Code
- 16 Tools Libraries Used

Dyego Fernandes de Sousa

To showcase and demonstrate the application of **Generative Deep Learning** techniques for synthesizing realistic tabular data with associated text fields. This project is a continuation of my previous works on **Supervised Learning** and **Unsupervised Learning**, you can find more information in the appendix.

Dataset Overview:

Detailed ESG (Environmental Social and Governance) dataset to the brand-level Observations:
3605 Features: 77

1. **Generative Adversarial Networks (GANs)**: Implement and train CTGAN for mixed-type tabular data
2. **Variational Autoencoders (VAEs)**: Apply TVAE as an alternative generative approach
3. **Statistical Methods**: Integrate Gaussian Copula for correlation structure preservation
4. **Large Language Models (LLMs)**: Fine-tune GPT-2 and Flan-T5 for conditional text generation
5. **Ensemble Methods**: Combine multiple generators with optimized weighting
6. **Evaluation Metrics**: Apply statistical tests (KS, correlation) to assess synthetic data quality

The following diagram illustrates the complete pipeline for synthetic brand data generation:

Deep Learning Models Used

Model	Type	Purpose	Key Features	——- ——- ——- ——-	CTGAN
Conditional GAN	Tabular synthesis	Mode-specific normalization, conditional vector			
TVAE	Variational Autoencoder	Distribution learning	KL divergence, latent space regularization	Gaussian Copula	Statistical Correlation preservation Multivariate dependencies
GPT-2 Medium	Transformer LLM	Brand name generation	355M params, fine-tuned	Flan-T5 Small	Encoder-Decoder Conditional text gen Instruction-following

!mermeid1.png

!mermeid2.png

!mermeid3.png

```
1 # Clone repository and install dependencies
2 !git clone https://github.com/dyegofern/csca5642-deep-learning.git
3 !pip install -q sdv transformers torch pandas numpy scikit-learn matplotlib
   seaborn plotly scipy
4 !pip install -q peft bitsandbytes accelerate sentencepiece
5 !pip install -q optuna
6
7 import sys
8 import os
9 from google.colab import drive
10
11 MAPPED_DIR = '/content/csca5642-deep-learning'
12
13 # Mount Google Drive
14 print("Mounting Google Drive...")
15 if not os.path.exists('/content/drive'):
16     drive.mount('/content/drive')
17 else:
18     print("Google Drive already mounted")
19
20 DATA_PATH = MAPPED_DIR + '/data/raw/brand_information.csv'
```

```

1  # Import V2 modules
2  from data_processor import BrandDataProcessor
3  from tabular_gan_v2 import (
4      EnsembleSynthesizer,
5      CTGANSynthesizerWrapper,
6      TVAESynthesizerWrapper,
7      GaussianCopulaSynthesizerWrapper,
8      calculate_generation_targets
9  )
10 from brand_name_generator_v2 import BrandNameGeneratorV2
11 from evaluator import BrandDataEvaluator
12
13 # Standard libraries
14 import pandas as pd
15 import numpy as np
16 import matplotlib.pyplot as plt
17 import seaborn as sns
18 import torch
19 import gc
20
21 # Set style
22 sns.set_style('whitegrid')

```

```
1 # Configuration for V2
2 FROM_PRETRAINED = False # Set to True to load pre-trained models
3
4 # Tabular Ensemble Config
5 CTGAN_EPOCHS = 300
6 TVAE_EPOCHS = 300
7 BATCH_SIZE = 500
8 ENSEMBLE_WEIGHTS = {
9     'ctgan': 0.40,
10    'tvae': 0.35,
11    'gaussian_copula': 0.25
12 }
13
14 # LLM Ensemble Config
15 LLM_MODELS = ['gpt2-medium', 'flan-t5-base'] # Can add 'phi-2', 'tinylama' if
        memory allows
16 LLM_EPOCHS = 3
17
18 # Generation Config
19 MIN_BRANDS_PER_COMPANY = 10
20 DIVERSITY_TEMPERATURE = 0.7
21 ADD_DIVERSITY_NOISE = True
```

This phase uses Optuna to find optimal hyperparameters for the tabular synthesizers. Set `RUN_HYPERPARAMETER_TUNING = True` to run optimization, or use previously saved best parameters. **Tuned Parameters:**

- **CTGAN/TVAE:** epochs, batch_size, embedding_dim, generator/discriminator dimensions
- **Ensemble:** weights for each model
- **Generation:** noise_level, diversity_temperature

```
1 # Hyperparameter Tuning Configuration
2 RUN_HYPERPARAMETER_TUNING = True # Set to True to run Optuna optimization
3 N_TUNING_TRIALS = 20 # Number of Optuna trials (more = better but slower)
4 TUNING_TIMEOUT = 300 # Maximum time for tuning in seconds (60 minutes)
5
6 # Path to save/load best hyperparameters
7 HYPERPARAMS_PATH = os.path.join(MODEL_DIR, 'best_hyperparameters.json')
8
9 print(f"Hyperparameter tuning: {'ENABLED' if RUN_HYPERPARAMETER_TUNING else 'DISABLED'}")
10 print(f"Trials: {N_TUNING_TRIALS}, Timeout: {TUNING_TIMEOUT}s")
11 print(f"Hyperparameters path: {HYPERPARAMS_PATH}")
```

Figure: Code from Cell 12


```
1 # Import hyperparameter tuner from external module
2 from hyperparameter_tuner_v2 import HyperparameterTunerV2
3
4 # Initialize tuner (will be used after data preparation)
5 tuner = None
6 best_hyperparams = None
7
8 print("HyperparameterTunerV2 imported successfully.")
9 print("Tuner will be initialized after data preparation (Phase 1).")
```

Figure: Code from Cell 13

```

1  # Check for saved hyperparameters
2  if os.path.exists(HYPERPARAMS_PATH):
3      print("Found saved hyperparameters!")
4      best_hyperparams = HyperparameterTunerV2.load(HYPERPARAMS_PATH)
5
6      print(f"Previous best score: {best_hyperparams.get('best_score', 'N/A')}")
7      print(f"Trials completed: {best_hyperparams.get('n_trials', 'N/A')}")
8      print(f"Timestamp: {best_hyperparams.get('tuning_timestamp', 'N/A')}")
9
10     # Update configuration with loaded hyperparameters
11     CTGAN_EPOCHS = best_hyperparams.get('ctgan_epochs', CTGAN_EPOCHS)
12     TVAE_EPOCHS = best_hyperparams.get('tvae_epochs', TVAE_EPOCHS)
13     BATCH_SIZE = best_hyperparams.get('batch_size', BATCH_SIZE)
14     if 'ensemble_weights' in best_hyperparams:
15         ENSEMBLE_WEIGHTS = best_hyperparams['ensemble_weights']
16
17     print(f"\nUsing optimized configuration:")
18     print(f"    CTGAN_EPOCHS: {CTGAN_EPOCHS}")
19     print(f"    TVAE_EPOCHS: {TVAE_EPOCHS}")
20     print(f"    BATCH_SIZE: {BATCH_SIZE}")
21     print(f"    ENSEMBLE_WEIGHTS: {ENSEMBLE_WEIGHTS}")
22 else:

```

```
1 # Load and process data
2 processor = BrandDataProcessor(DATA_PATH)
3 raw_data = processor.load_data()
4 print(f"Loaded {len(raw_data)} brands with {len(raw_data.columns)} features")
```

Figure: Code from Cell 16

```
1 # Clean data
2 cleaned_data = processor.clean_data()
3 print(f"\nCleaned data: {len(cleaned_data)} rows, {len(cleaned_data.columns)}
      columns")
```

Figure: Code from Cell 17

```
1 # Prepare for GAN training
2 train_df, val_df = processor.prepare_for_gan(test_size=0.2)
3
4 print(f"\nTraining set: {len(train_df)} brands")
5 print(f"Validation set: {len(val_df)} brands")
6
7 # Get column types
8 discrete_cols = processor.categorical_features
9 binary_cols = [col for col in train_df.columns if train_df[col].nunique() == 2
10                and set(train_df[col].unique()).issubset({0, 1})]
11 numerical_cols = [col for col in train_df.columns if col not in discrete_cols
12                  and col not in binary_cols]
13
14 print(f"\nColumn types:")
15 print(f"    Numerical: {len(numerical_cols)}")
16 print(f"    Categorical: {len(discrete_cols)}")
17 print(f"    Binary: {len(binary_cols)}")
```

Figure: Code from Cell 18

```
1 !pip install optuna
```

Figure: Code from Cell 19

```
1 # Execute hyperparameter tuning if enabled and no saved params exist
2 if RUN_HYPERPARAMETER_TUNING and best_hyperparams is None:
3     print("Initializing hyperparameter tuner...")
4
5     # Create tuner instance
6     tuner = HyperparameterTunerV2(
7         train_data=train_df,
8         discrete_cols=discrete_cols,
9         binary_cols=binary_cols,
10        eval_sample_size=min(1000, len(train_df)),
11        gen_sample_size=500,
12        verbose=True
13    )
14
15    # Run optimization
16    best_hyperparams = tuner.tune(
17        n_trials=N_TUNING_TRIALS,
18        timeout=TUNING_TIMEOUT,
19        seed=42,
20        show_progress_bar=True
21    )
22
```

Analysis Result



Figure: Output from Cell 20

Training CTGAN, TVAE, and Gaussian Copula models

```
1  # Initialize Ensemble Synthesizer
2  tabular_ensemble = EnsembleSynthesizer(
3      ctgan_epochs=CTGAN_EPOCHS,
4      ctgan_batch_size=BATCH_SIZE,
5      tvae_epochs=TVAE_EPOCHS,
6      tvae_batch_size=BATCH_SIZE,
7      gc_default_distribution='beta',
8      weights=ENSEMBLE_WEIGHTS,
9      verbose=True,
10     cuda=True
11 )
12
13 print("Tabular Ensemble initialized with:")
14 print(f"    - CTGAN (epochs={CTGAN_EPOCHS})")
15 print(f"    - TVAE (epochs={TVAE_EPOCHS})")
16 print(f"    - Gaussian Copula (distribution=beta)")
```

Figure: Code from Cell 22

```
1 if FROM_PRETRAINED:
2     # Load pre-trained models
3     print("Loading pre-trained tabular ensemble...")
4     tabular_ensemble.load_models(os.path.join(MODEL_DIR, 'tabular_ensemble'))
5 else:
6     # Train all models
7     print("Training tabular ensemble (this will take ~30-60 minutes)...")
8     training_times = tabular_ensemble.train(
9         data=train_df,
10        discrete_columns=discrete_cols,
11        binary_columns=binary_cols
12    )
13
14    # Save models
15    tabular_ensemble.save_models(os.path.join(MODEL_DIR, 'tabular_ensemble'))
16
17    print(f"\nTraining times:")
18    for model, time in training_times.items():
19        print(f"    {model}: {time:.1f} seconds")
```

Figure: Code from Cell 23

```
1 # Compare individual model quality
2 print("Evaluating individual model quality...")
3 comparison_df = tabular_ensemble.compare_all_models(train_df, n_samples=1000)
4 print("\nModel Comparison:")
5 display(comparison_df)
```

Figure: Code from Cell 24

```
1 # Optionally optimize weights based on quality  
2 optimized_weights = tabular_ensemble.optimize_weights(train_df, n_eval_samples  
    =1000)  
3 print(f"Optimized weights: {optimized_weights}")
```

Figure: Code from Cell 25

Training GPT-2 Medium and Flan-T5 for brand name generation

```
1 # Prepare brand name training data
2 brands_df = processor.df[['brand_name', 'company_name', 'industry_name']].
    dropna()
3 print(f"Brand name training data: {len(brands_df)} examples")
4 brands_df.head()
```

Figure: Code from Cell 27

```
1 # Initialize LLM Ensemble Generator
2 llm_generator = BrandNameGeneratorV2(
3     models=LLM_MODELS,
4     memory_efficient=True,
5     verbose=True
6 )
7
8 print(f"LLM Ensemble initialized with models: {LLM_MODELS}")
```

Figure: Code from Cell 28


```
1 if FROM_PRETRAINED:
2     # Load pre-trained models
3     print("Loading pre-trained LLM ensemble...")
4     llm_generator.load_model(os.path.join(MODEL_DIR, 'llm_ensemble'))
5 else:
6     # Fine-tune all models
7     print(f"Fine-tuning LLM ensemble (epochs={LLM_EPOCHS})...")
8     print("This will train each model sequentially to save memory.")
9
10    llm_generator.fine_tune(
11        brands_df=brands_df,
12        epochs=LLM_EPOCHS,
13        output_dir=os.path.join(MODEL_DIR, 'llm_ensemble')
14    )
15
16    # Save ensemble config
17    llm_generator.save_model(os.path.join(MODEL_DIR, 'llm_ensemble'))
```

Figure: Code from Cell 29

```
1 # Test LLM generation
2 print("Testing LLM ensemble generation...")
3 llm_generator.prepare_model()
4
5 test_companies = [
6     ("PepsiCo", "Non-Alcoholic Beverages"),
7     ("Nestle", "Processed Foods"),
8     ("Mars, Incorporated", "Processed Foods")
9 ]
10
11 for company, industry in test_companies:
12     names = llm_generator.generate_brand_names(company, industry, n_names=3)
13     print(f"\n{company} ({industry}): {names}")
```

Figure: Code from Cell 30

```
1 # Calculate generation targets
2 generation_targets = calculate_generation_targets(
3     data=train_df,
4     company_column='company_name',
5     min_brands_per_company=MIN_BRANDS_PER_COMPANY
6 )
7 generation_targets = dict(list(generation_targets.items())[:150])
```

Figure: Code from Cell 32

```
1 # Generate synthetic tabular features using ensemble
2 print("Generating synthetic features with ensemble...")
3 synthetic_features, failed_companies = tabular_ensemble.generate_stratified(
4     company_distribution=generation_targets,
5     verbose=True
6 )
7
8 print(f"\nGenerated {len(synthetic_features)} synthetic brand features")
9 if failed_companies:
10     print(f"Failed companies: {len(failed_companies)}")
```

Figure: Code from Cell 33

```
1 # Add diversity noise if enabled
2 if ADD_DIVERSITY_NOISE:
3     print("Adding diversity noise to numerical features...")
4     synthetic_features = tabular_ensemble.add_diversity_noise(
5         synthetic_features,
6         noise_level=0.07
7     )
```

Figure: Code from Cell 34

```
1 # Decode categorical features back to original values  
2 print("Decoding categorical features...")  
3 synthetic_decoded = processor.decode_categorical(synthetic_features)  
4 print(f"Decoded {len(synthetic_decoded)} synthetic brands")
```

Figure: Code from Cell 35

```
1 # Generate brand names using LLM ensemble
2 print("\nGenerating brand names with LLM ensemble...")
3 llm_generator.reset_uniqueness_tracker()
4
5 synthetic_with_names = llm_generator.generate_for_dataframe(
6     synthetic_df=synthetic_decoded,
7     temperature=DIVERSITY_TEMPERATURE,
8     verbose=True
9 )
10
11 print(f"\nFinal synthetic dataset: {len(synthetic_with_names)} brands")
```

Figure: Code from Cell 36

```
1 # Preview synthetic data
2 print("\nSample of generated synthetic brands:")
3 display(synthetic_with_names[['company_name', 'industry_name', 'brand_name']].
    head(20))
```

Figure: Code from Cell 37


```
1 # Save synthetic data
2 synthetic_path = os.path.join(OUTPUT_DIR, 'synthetic_brands_v2.csv')
3 synthetic_with_names.to_csv(synthetic_path, index=False)
4 print(f"Synthetic data saved to {synthetic_path}")
5
6 # Create augmented dataset
7 original_decoded = processor.decode_categorical(train_df)
8 augmented_df = pd.concat([original_decoded, synthetic_with_names], ignore_index=True)
9
10 augmented_path = os.path.join(OUTPUT_DIR, 'augmented_brands_v2.csv')
11 augmented_df.to_csv(augmented_path, index=False)
12 print(f"Augmented data saved to {augmented_path}")
13 print(f"Total augmented size: {len(augmented_df)} brands")
```

Figure: Code from Cell 38

This phase comprehensively evaluates the quality of generated synthetic data through: 1. **Statistical Fidelity:** Distribution matching (KS tests), correlation preservation 2. **Visual Analysis:** Distribution overlays, QQ plots, feature comparisons 3. **Dimensionality Reduction:** PCA and t-SNE projections 4. **Summary Metrics:** Quality scorecards and radar charts

```
1 # Initialize evaluator and prepare data
2 evaluator = BrandDataEvaluator()
3
4 # Get numerical columns for evaluation
5 eval_numerical_cols = [col for col in numerical_cols if col in
    synthetic_features.columns and col in train_df.columns]
6
7 # Prepare augmented numerical data for later use
8 augmented_numerical = pd.concat([train_df[eval_numerical_cols],
    synthetic_features[eval_numerical_cols]], ignore_index=True)
9
10 print(f"Evaluation columns: {len(eval_numerical_cols)} numerical features")
11 print(f"Real data samples: {len(train_df)}")
12 print(f"Synthetic data samples: {len(synthetic_features)}")
13 print(f"Augmented data samples: {len(augmented_numerical)}")
```

Figure: Code from Cell 40

```

1  ### 5.1 Statistical Distribution Comparison (KS Test)
2
3  from scipy import stats
4
5  # Compute KS statistics for all features
6  print("="*70)
7  print("KOLMOGOROV-SMIRNOV TEST RESULTS")
8  print("="*70)
9  print(f"{'Feature':<35} {'KS Stat':>10} {'P-Value':>12} {'Result':>10}")
10 print("-"*70)
11
12 ks_results = {}
13 for col in eval_numerical_cols:
14     if col in train_df.columns and col in synthetic_features.columns:
15         stat, pvalue = stats.ks_2samp(
16             train_df[col].dropna(),
17             synthetic_features[col].dropna()
18         )
19         ks_results[col] = {'statistic': stat, 'pvalue': pvalue}
20         result = "PASS" if pvalue > 0.05 else "FAIL"
21         print(f"{'col':<35} {'stat':>10.4f} {'pvalue':>12.4f} {'result':>10}")
22

```

```

1  ### 5.2 KS Statistics Visualization
2
3  # Create a bar chart of KS statistics
4  fig, axes = plt.subplots(1, 2, figsize=(16, 6))
5
6  # Sort features by KS statistic
7  sorted_features = sorted(ks_results.items(), key=lambda x: x[1]['statistic'],
8                           reverse=True)
9  features = [f[0] for f in sorted_features]
10 ks_stats = [f[1]['statistic'] for f in sorted_features]
11 colors = ['#e74c3c' if f[1]['pvalue'] <= 0.05 else '#27ae60' for f in
12           sorted_features]
13
14 # Bar chart of KS statistics
15 ax1 = axes[0]
16 bars = ax1.barh(range(len(features)), ks_stats, color=colors, edgecolor='black',
17                 , alpha=0.8)
18 ax1.set_yticks(range(len(features)))
19 ax1.set_yticklabels(features, fontsize=9)
20 ax1.set_xlabel('KS Statistic', fontsize=12)
21 ax1.set_title('Distribution Similarity (KS Test)\nLower is Better', fontsize
22               =14, fontweight='bold')

```

Analysis Result

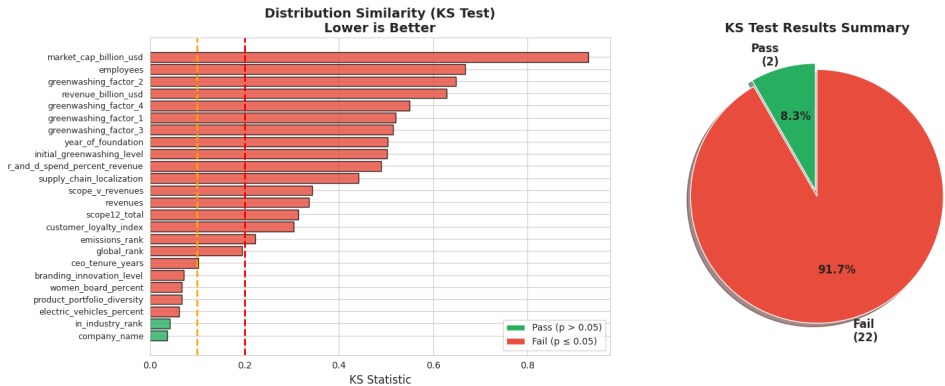


Figure: Output from Cell 42

```

1  ### 5.3 Distribution Comparison - Histograms with KDE
2
3  # Select top features for visualization (mix of good and bad performers)
4  best_features = [f[0] for f in sorted_features[-4:]] # Best 4 (lowest KS)
5  worst_features = [f[0] for f in sorted_features[:4]] # Worst 4 (highest KS)
6  viz_features = worst_features + best_features
7
8  n_features = len(viz_features)
9  fig, axes = plt.subplots(2, 4, figsize=(20, 10))
10 axes = axes.flatten()
11
12 for idx, feature in enumerate(viz_features):
13     ax = axes[idx]
14
15     # Get data
16     real_data = train_df[feature].dropna()
17     synth_data = synthetic_features[feature].dropna()
18
19     # Plot histograms with KDE
20     ax.hist(real_data, bins=30, alpha=0.5, label='Real', color='steelblue',
21             density=True, edgecolor='black')
21     ax.hist(synth_data, bins=30, alpha=0.5, label='Synthetic', color='coral',

```

Analysis Result

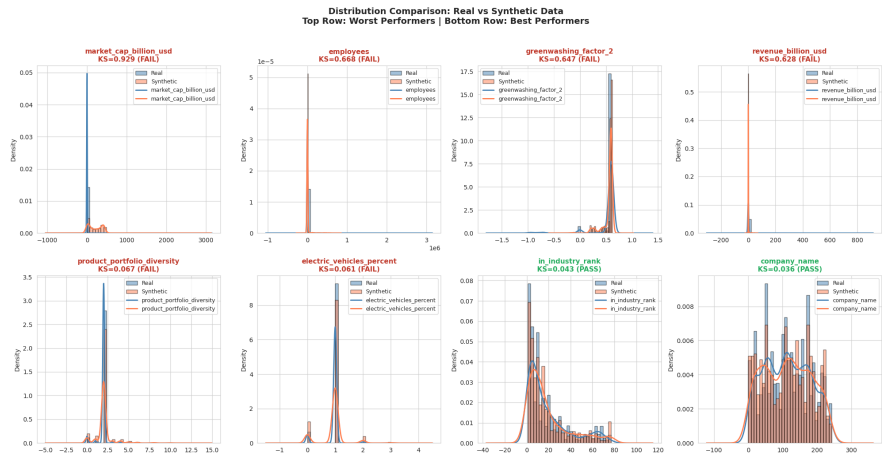


Figure: Output from Cell 43


```

1  ### 5.4 Correlation Structure Preservation
2
3  # Compute correlation matrices
4  real_corr = train_df[eval_numerical_cols].corr()
5  synth_corr = synthetic_features[eval_numerical_cols].corr()
6  corr_diff = np.abs(real_corr - synth_corr)
7
8  # Create figure with 3 subplots
9  fig, axes = plt.subplots(1, 3, figsize=(20, 6))
10
11 # Real data correlations
12 sns.heatmap(real_corr, ax=axes[0], cmap='RdBu_r', center=0, vmin=-1, vmax=1,
13             square=True, cbar_kws={'label': 'Correlation', 'shrink': 0.8},
14             xticklabels=True, yticklabels=True)
15 axes[0].set_title('Real Data Correlations', fontsize=14, fontweight='bold')
16 axes[0].tick_params(axis='both', labelsize=8, rotation=45)
17
18 # Synthetic data correlations
19 sns.heatmap(synth_corr, ax=axes[1], cmap='RdBu_r', center=0, vmin=-1, vmax=1,
20             square=True, cbar_kws={'label': 'Correlation', 'shrink': 0.8},
21             xticklabels=True, yticklabels=True)
22 axes[1].set_title('Synthetic Data Correlations', fontsize=14, fontweight='bold')

```

Analysis Result

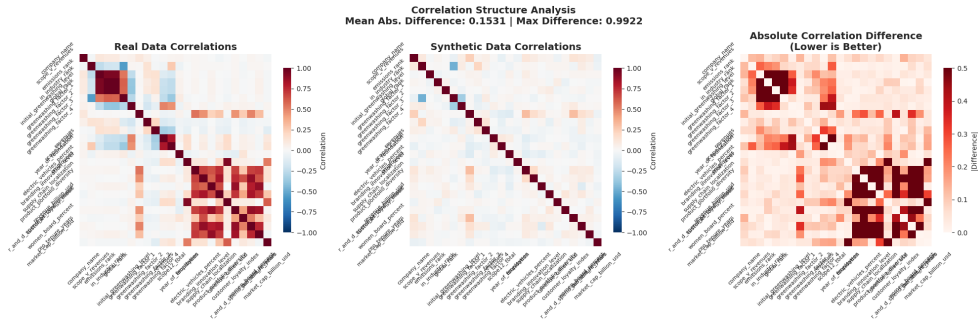


Figure: Output from Cell 44

```

1  ### 5.5 QQ Plots - Quantile Comparison
2
3  # QQ plots for selected features
4  from scipy import stats as scipy_stats
5
6  fig, axes = plt.subplots(2, 4, figsize=(20, 10))
7  axes = axes.flatten()
8
9  for idx, feature in enumerate(viz_features):
10     ax = axes[idx]
11
12     real_data = np.sort(train_df[feature].dropna().values)
13     synth_data = np.sort(synthetic_features[feature].dropna().values)
14
15     # Interpolate to same length for QQ plot
16     n_points = min(len(real_data), len(synth_data), 100)
17     real_quantiles = np.percentile(real_data, np.linspace(0, 100, n_points))
18     synth_quantiles = np.percentile(synth_data, np.linspace(0, 100, n_points))
19
20     # Plot QQ
21     ax.scatter(real_quantiles, synth_quantiles, alpha=0.6, s=30, c='steelblue',
                edgecolor='black')

```

Analysis Result

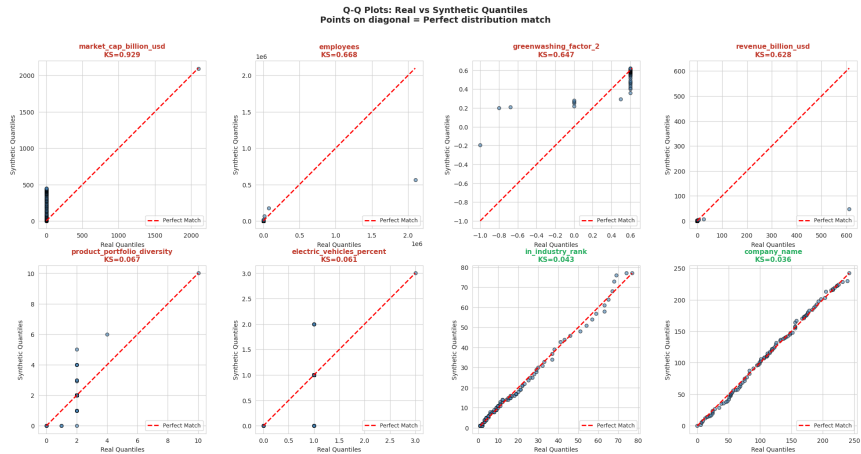


Figure: Output from Cell 45

```
1  ### 5.6 PCA & t-SNE Dimensionality Reduction
2
3  from sklearn.decomposition import PCA
4  from sklearn.manifold import TSNE
5  from sklearn.preprocessing import StandardScaler
6
7  # Prepare data
8  X_real = train_df[eval_numerical_cols].fillna(0).values
9  X_synth = synthetic_features[eval_numerical_cols].fillna(0).values
10
11 # Standardize
12 scaler = StandardScaler()
13 X_real_scaled = scaler.fit_transform(X_real)
14 X_synth_scaled = scaler.transform(X_synth)
15
16 # PCA
17 pca = PCA(n_components=2)
18 X_real_pca = pca.fit_transform(X_real_scaled)
19 X_synth_pca = pca.transform(X_synth_scaled)
```

```

20
21 # t-SNE (on combined data for fair comparison)
22 print("Computing t-SNE projection (this may take a moment)...")
23 X_combined = np.vstack([X_real_scaled, X_synth_scaled])
24 labels = np.array(['Real'] * len(X_real) + ['Synthetic'] * len(X_synth))
25
26 tsne = TSNE(n_components=2, perplexity=30, random_state=42, n_iter=1000)
27 X_combined_tsne = tsne.fit_transform(X_combined)
28 X_real_tsne = X_combined_tsne[:len(X_real)]
29 X_synth_tsne = X_combined_tsne[len(X_real):]
30
31 # Create visualization
32 fig, axes = plt.subplots(1, 2, figsize=(16, 7))
33
34 # PCA plot
35 ax1 = axes[0]
36 ax1.scatter(X_real_pca[:, 0], X_real_pca[:, 1], alpha=0.5, s=30, c='steelblue',
37             label='Real', edgecolor='white', linewidth=0.5)
38 ax1.scatter(X_synth_pca[:, 0], X_synth_pca[:, 1], alpha=0.5, s=30, c='coral',
39             label='Synthetic', edgecolor='white', linewidth=0.5)

```

```

38 ax1.set_xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.1%} variance)',
    fontsize=12)
39 ax1.set_ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.1%} variance)',
    fontsize=12)
40 ax1.set_title('PCA Projection\nReal vs Synthetic Data', fontsize=14, fontweight
    ='bold')
41 ax1.legend(fontsize=11, loc='upper right')
42 ax1.grid(True, alpha=0.3)
43
44 # t-SNE plot
45 ax2 = axes[1]
46 ax2.scatter(X_real_tsne[:, 0], X_real_tsne[:, 1], alpha=0.5, s=30, c='steelblue',
    label='Real', edgecolor='white', linewidth=0.5)
47 ax2.scatter(X_synth_tsne[:, 0], X_synth_tsne[:, 1], alpha=0.5, s=30, c='coral',
    label='Synthetic', edgecolor='white', linewidth=0.5)
48 ax2.set_xlabel('t-SNE Dimension 1', fontsize=12)
49 ax2.set_ylabel('t-SNE Dimension 2', fontsize=12)
50 ax2.set_title('t-SNE Projection\nReal vs Synthetic Data', fontsize=14,
    fontweight='bold')
51 ax2.legend(fontsize=11, loc='upper right')

```

```
52 ax2.grid(True, alpha=0.3)
53
54 plt.suptitle('Dimensionality Reduction: Synthetic Data Should Overlap with Real
    Data', fontsize=14, fontweight='bold', y=1.02)
55 plt.tight_layout()
56 plt.savefig(os.path.join(OUTPUT_DIR, 'pca_tsne_comparison.png'), dpi=150,
    bbox_inches='tight')
57 plt.show()
58
59 print(f"\nPCA Explained Variance: PC1={pca.explained_variance_ratio_[0]:.1%},
    PC2={pca.explained_variance_ratio_[1]:.1%}")
```

Figure: Code from Cell 46

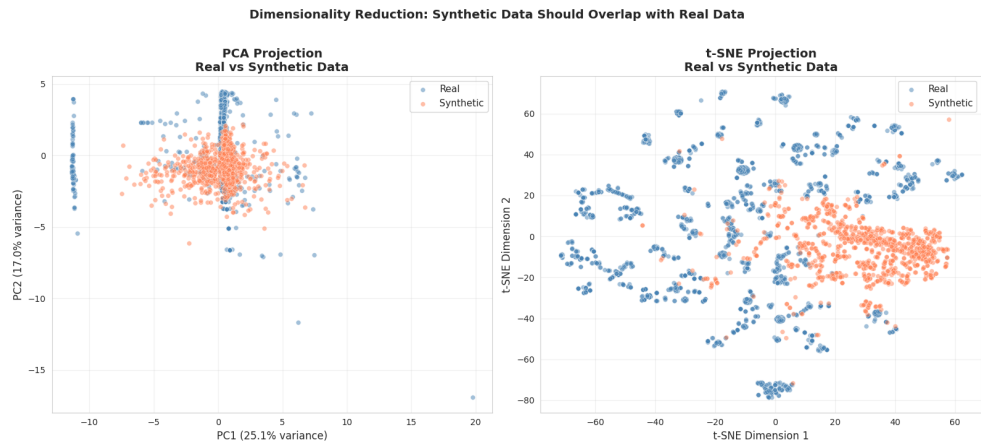


Figure: Output from Cell 46

```
1 ### 5.7 Feature-wise Statistics Comparison
2
3 # Compute statistics for real and synthetic data
4 stats_comparison = []
5
6 for col in eval_numerical_cols:
7     real_col = train_df[col].dropna()
8     synth_col = synthetic_features[col].dropna()
9
10    stats_comparison.append({
11        'Feature': col,
12        'Real Mean': real_col.mean(),
13        'Synth Mean': synth_col.mean(),
14        'Mean Diff %': abs(real_col.mean() - synth_col.mean()) / (abs(real_col.
15            mean()) + 1e-10) * 100,
16        'Real Std': real_col.std(),
17        'Synth Std': synth_col.std(),
18        'Std Diff %': abs(real_col.std() - synth_col.std()) / (abs(real_col.std
19            ()) + 1e-10) * 100,
```

```
18         'Real Min': real_col.min(),
19         'Synth Min': synth_col.min(),
20         'Real Max': real_col.max(),
21         'Synth Max': synth_col.max(),
22         'KS Stat': ks_results[col]['statistic']
23     })
24
25 stats_df = pd.DataFrame(stats_comparison)
26
27 # Display summary table
28 print("="*80)
29 print("FEATURE-WISE STATISTICS COMPARISON")
30 print("="*80)
31 display(stats_df[['Feature', 'Real Mean', 'Synth Mean', 'Mean Diff %', 'Real
    Std', 'Synth Std', 'Std Diff %', 'KS Stat']].round(3))
32
33 # Visualize mean and std differences
34 fig, axes = plt.subplots(1, 2, figsize=(16, 6))
35
36 # Mean difference
```

```

37 ax1 = axes[0]
38 sorted_by_mean = stats_df.sort_values('Mean Diff %', ascending=False)
39 colors_mean = ['#e74c3c' if x > 50 else '#f39c12' if x > 20 else '#27ae60' for
    x in sorted_by_mean['Mean Diff %']]
40 ax1.barh(range(len(sorted_by_mean)), sorted_by_mean['Mean Diff %'], color=
    colors_mean, edgecolor='black', alpha=0.8)
41 ax1.set_yticks(range(len(sorted_by_mean)))
42 ax1.set_yticklabels(sorted_by_mean['Feature'], fontsize=9)
43 ax1.set_xlabel('Mean Difference (%)', fontsize=12)
44 ax1.set_title('Mean Value Difference\n(Lower is Better)', fontsize=14,
    fontweight='bold')
45 ax1.axvline(x=20, color='orange', linestyle='--', linewidth=2, alpha=0.7)
46 ax1.axvline(x=50, color='red', linestyle='--', linewidth=2, alpha=0.7)
47 ax1.invert_yaxis()
48
49 # Std difference
50 ax2 = axes[1]
51 sorted_by_std = stats_df.sort_values('Std Diff %', ascending=False)
52 colors_std = ['#e74c3c' if x > 50 else '#f39c12' if x > 20 else '#27ae60' for x
    in sorted_by_std['Std Diff %']]

```

```
53 ax2.barh(range(len(sorted_by_std)), sorted_by_std['Std Diff %'], color=
    colors_std, edgecolor='black', alpha=0.8)
54 ax2.set_yticks(range(len(sorted_by_std)))
55 ax2.set_yticklabels(sorted_by_std['Feature'], fontsize=9)
56 ax2.set_xlabel('Std Deviation Difference (%)', fontsize=12)
57 ax2.set_title('Standard Deviation Difference\n(Lower is Better)', fontsize=14,
    fontweight='bold')
58 ax2.axvline(x=20, color='orange', linestyle='--', linewidth=2, alpha=0.7)
59 ax2.axvline(x=50, color='red', linestyle='--', linewidth=2, alpha=0.7)
60 ax2.invert_yaxis()
61
62 plt.tight_layout()
63 plt.savefig(os.path.join(OUTPUT_DIR, 'statistics_comparison.png'), dpi=150,
    bbox_inches='tight')
64 plt.show()
```

Figure: Code from Cell 47

Analysis Result

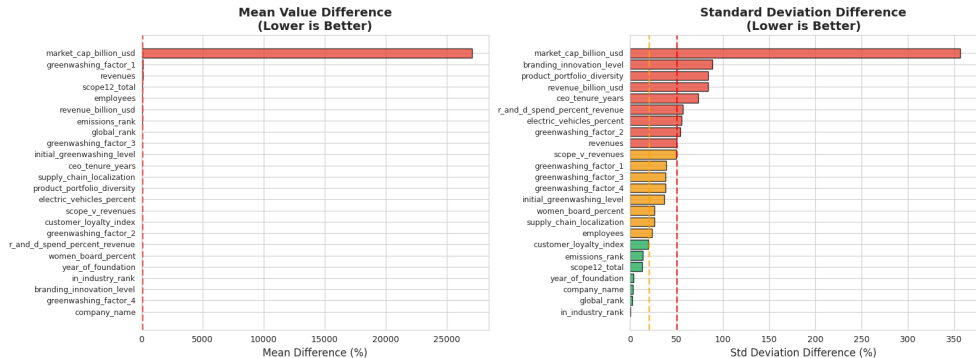


Figure: Output from Cell 47

```

1  ### 5.8 Synthetic Data Quality Scorecard & Radar Chart
2
3  # Calculate quality metrics
4  ks_pass_rate = passes / len(ks_results) * 100
5  mean_ks_stat = np.mean([v['statistic'] for v in ks_results.values()])
6  mean_mean_diff = stats_df['Mean Diff %'].mean()
7  mean_std_diff = stats_df['Std Diff %'].mean()
8
9  # Normalize metrics to 0-100 scale (higher is better)
10 distribution_score = max(0, 100 - mean_ks_stat * 200)  # KS of 0 = 100, KS of
    0.5 = 0
11 correlation_score = max(0, 100 - mean_corr_diff * 200)  # Diff of 0 = 100, Diff
    of 0.5 = 0
12 mean_preservation = max(0, 100 - mean_mean_diff)  # Lower diff = higher score
13 variance_preservation = max(0, 100 - mean_std_diff)  # Lower diff = higher
    score
14 coverage_score = ks_pass_rate  # Percentage of features passing KS test
15
16 # Overall score (weighted average)

```

```

17 overall_score = (distribution_score * 0.3 + correlation_score * 0.2 +
18                 mean_preservation * 0.2 + variance_preservation * 0.2 +
19                 coverage_score * 0.1)
20
21 # Create figure with scorecard and radar chart
22
23 # Left side: Scorecard
24 ax1 = fig.add_subplot(121)
25 ax1.axis('off')
26
27 # Create scorecard table
28 scorecard_data = [
29     ['Metric', 'Value', 'Score', 'Grade'],
30     ['KS Test Pass Rate', f'{ks_pass_rate:.1f}%', f'{coverage_score:.1f}', 'A'
31         if coverage_score >= 70 else 'B' if coverage_score >= 50 else 'C' if
32         coverage_score >= 30 else 'D'],
33     ['Mean KS Statistic', f'{mean_ks_stat:.3f}', f'{distribution_score:.1f}', '
34         A' if distribution_score >= 70 else 'B' if distribution_score >= 50
35         else 'C' if distribution_score >= 30 else 'D'],

```



```

32 ['Correlation Preservation', f'{mean_corr_diff:.3f}', f'{correlation_score
    :.1f}', 'A' if correlation_score >= 70 else 'B' if correlation_score >=
    50 else 'C' if correlation_score >= 30 else 'D'],
33 ['Mean Value Preservation', f'{mean_mean_diff:.1f}%', f'{mean_preservation
    :.1f}', 'A' if mean_preservation >= 70 else 'B' if mean_preservation >=
    50 else 'C' if mean_preservation >= 30 else 'D'],
34 ['Variance Preservation', f'{mean_std_diff:.1f}%', f'{variance_preservation
    :.1f}', 'A' if variance_preservation >= 70 else 'B' if
    variance_preservation >= 50 else 'C' if variance_preservation >= 30
    else 'D'],
35 ['    ' * 20, '    ' * 10, '    ' * 10, '    ' * 5],
36 ['OVERALL QUALITY', '', f'{overall_score:.1f}', 'A' if overall_score >= 70
    else 'B' if overall_score >= 50 else 'C' if overall_score >= 30 else 'D
    '],
37 ]
38
39 table = ax1.table(cellText=scorecard_data, loc='center', cellLoc='center',
40                  colWidths=[0.35, 0.2, 0.15, 0.1])
41 table.auto_set_font_size(False)
42 table.set_fontsize(11)

```

```
43 table.scale(1.2, 2)
44
45 # Style the table
46 for i in range(len(scorecard_data)):
47     for j in range(4):
48         cell = table[(i, j)]
49         if i == 0: # Header
50             cell.set_facecolor('#34495e')
51             cell.set_text_props(color='white', fontweight='bold')
52         elif i == len(scorecard_data) - 1: # Overall row
53             cell.set_facecolor('#2ecc71' if overall_score >= 70 else '#f39c12'
54                                     if overall_score >= 50 else '#e74c3c')
55             cell.set_text_props(fontweight='bold')
56         elif i == len(scorecard_data) - 2: # Separator
57             cell.set_facecolor('#ecf0f1')
58         elif j == 3: # Grade column
59             grade = scorecard_data[i][3]
60             if grade == 'A':
61                 cell.set_facecolor('#27ae60')
62             elif grade == 'B':
```

```

62         cell.set_facecolor('#f39c12')
63     elif grade == 'C':
64         cell.set_facecolor('#e67e22')
65     else:
66         cell.set_facecolor('#e74c3c')
67         cell.set_text_props(color='white', fontweight='bold')
68
69 ax1.set_title('Synthetic Data Quality Scorecard', fontsize=16, fontweight='bold',
70              ', pad=20)
71
72 # Right side: Radar chart
73
74 # Radar chart data
75 categories = ['Distribution\nMatching', 'Correlation\nPreservation', 'Mean\nPreservation',
76              'Variance\nPreservation', 'KS Test\nPass Rate']
77 values = [distribution_score, correlation_score, mean_preservation,
78           variance_preservation, coverage_score]
79 values += values[:1] # Close the polygon

```

```

79
80 angles = np.linspace(0, 2 * np.pi, len(categories), endpoint=False).tolist()
81 angles += angles[:1]
82
83 ax2.plot(angles, values, 'o-', linewidth=2, color='#3498db', markersize=8)
84 ax2.fill(angles, values, alpha=0.25, color='#3498db')
85 ax2.set_xticks(angles[:-1])
86 ax2.set_xticklabels(categories, fontsize=10)
87 ax2.set_ylim(0, 100)
88 ax2.set_yticks([20, 40, 60, 80, 100])
89 ax2.set_yticklabels(['20', '40', '60', '80', '100'], fontsize=9)
90 ax2.set_title('Quality Metrics Radar Chart\n(Higher is Better)', fontsize=14,
               fontweight='bold', pad=20)
91
92 # Add reference circles
93 for val in [30, 50, 70]:
94     circle = plt.Circle((0, 0), val, transform=ax2.transData._b, fill=False,
95                          linestyle='--', alpha=0.3, color='gray')
96

```

```
97 plt.suptitle('SYNTHETIC DATA QUALITY ASSESSMENT', fontsize=18, fontweight='bold', y=1.02)
98 plt.tight_layout()
99 plt.savefig(os.path.join(OUTPUT_DIR, 'quality_scorecard.png'), dpi=150,
    bbox_inches='tight')
100 plt.show()
101
102 # Print summary
103 print("\n" + "="*70)
104 print("QUALITY ASSESSMENT SUMMARY")
105 print("="*70)
106 print(f"Overall Quality Score: {overall_score:.1f}/100 ({'Excellent' if
    overall_score >= 70 else 'Good' if overall_score >= 50 else 'Fair' if
    overall_score >= 30 else 'Poor'})")
107 print(f"Distribution Matching: {distribution_score:.1f}/100")
108 print(f"Correlation Preservation: {correlation_score:.1f}/100")
109 print(f"Statistical Fidelity: {(mean_preservation + variance_preservation)/2:.1
    f}/100")
110 print("="*70)
```

Figure: Code from Cell 48

SYNTHETIC DATA QUALITY ASSESSMENT

Synthetic Data Quality Scorecard

Metric	Value	Score	Grade
KS Test Pass Rate	8.3%	8.3	D
Mean KS Statistic	0.357	28.7	D
Correlation Preservation	0.153	69.4	B
Mean Value Preservation	1145.0%	0.0	D
Variance Preservation	51.6%	48.4	C
OVERALL QUALITY		33.0	C

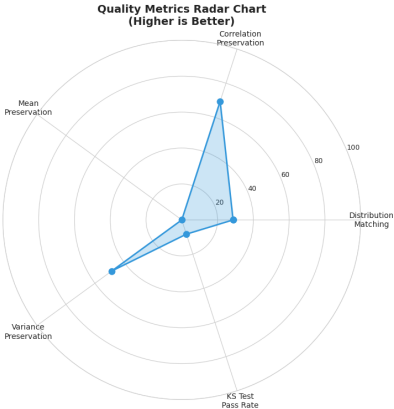


Figure: Output from Cell 48

```
1  ### 5.9 Augmented Data Clustering Analysis
2
3  from sklearn.cluster import AgglomerativeClustering
4  from sklearn.metrics import silhouette_score, davies_bouldin_score
5
6  # Cluster the augmented data (real + synthetic combined)
7  print("="*70)
8  print("AUGMENTED DATA CLUSTERING ANALYSIS")
9  print("="*70)
10
11 # Prepare augmented data with labels
12 augmented_with_labels = augmented_numerical.copy()
13 augmented_with_labels['source'] = ['Real'] * len(train_df) + ['Synthetic'] *
    len(synthetic_features)
14
15 # Standardize for clustering
16 X_augmented = scaler.fit_transform(augmented_numerical.fillna(0))
17
18 # Perform clustering
```



```
19 n_clusters = 5
20 clustering = AgglomerativeClustering(n_clusters=n_clusters)
21 cluster_labels = clustering.fit_predict(X_augmented)
22
23 # Calculate metrics
24 silhouette = silhouette_score(X_augmented, cluster_labels)
25 davies_bouldin = davies_bouldin_score(X_augmented, cluster_labels)
26
27 print(f"\nClustering Results (n_clusters={n_clusters}):")
28 print(f"    Silhouette Score: {silhouette:.4f} (higher is better, range: -1 to 1)
29       ")
30 print(f"    Davies-Bouldin Index: {davies_bouldin:.4f} (lower is better)")
31
32 # Analyze cluster composition
33 augmented_with_labels['cluster'] = cluster_labels
34
35 # Create visualization
36 fig, axes = plt.subplots(1, 2, figsize=(16, 6))
37
38 # Cluster composition (Real vs Synthetic distribution per cluster)
```

```

38 ax1 = axes[0]
39 cluster_composition = pd.crosstab(augmented_with_labels['cluster'],
    augmented_with_labels['source'], normalize='index') * 100
40 cluster_composition.plot(kind='bar', ax=ax1, color=['steelblue', 'coral'],
    edgecolor='black', alpha=0.8)
41 ax1.set_xlabel('Cluster', fontsize=12)
42 ax1.set_ylabel('Percentage (%)', fontsize=12)
43 ax1.set_title('Cluster Composition: Real vs Synthetic\n(Balanced = Good
    Integration)', fontsize=14, fontweight='bold')
44 ax1.legend(title='Source', fontsize=10)
45 ax1.set_xticklabels(ax1.get_xticklabels(), rotation=0)
46
47 # Add cluster sizes as text
48 cluster_sizes = augmented_with_labels['cluster'].value_counts().sort_index()
49 for i, (idx, size) in enumerate(cluster_sizes.items()):
50     ax1.annotate(f'n={size}', xy=(i, 105), ha='center', fontsize=9, fontweight=
        'bold')
51
52 # PCA visualization with clusters
53 ax2 = axes[1]

```

```

54 X_augmented_pca = pca.fit_transform(X_augmented)
55 scatter = ax2.scatter(X_augmented_pca[:, 0], X_augmented_pca[:, 1],
56                       c=cluster_labels, cmap='Set2', alpha=0.6, s=30, edgecolor
57                       ='white', linewidth=0.5)
58 # Mark synthetic points with different marker
59 synthetic_mask = np.array([False] * len(train_df) + [True] * len(
60     synthetic_features))
61 ax2.scatter(X_augmented_pca[synthetic_mask, 0], X_augmented_pca[synthetic_mask,
62     1],
63             c='none', edgecolor='red', s=50, linewidth=1.5, marker='o', label='
64     Synthetic', alpha=0.5)
65 ax2.set_xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.1%} variance)',
66               fontsize=12)
67 ax2.set_ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.1%} variance)',
68               fontsize=12)
69 ax2.set_title(f'Augmented Data Clusters (PCA View)\nSilhouette: {silhouette:.3f
70     }', fontsize=14, fontweight='bold')
71 ax2.legend(loc='upper right', fontsize=10)

```

```
67 |
68 | plt.colorbar(scatter, ax=ax2, label='Cluster')
69 | plt.suptitle('Augmented Data (Real + Synthetic) Clustering', fontsize=16,
70 |             fontweight='bold', y=1.02)
71 | plt.tight_layout()
72 | plt.savefig(os.path.join(OUTPUT_DIR, 'augmented_clustering.png'), dpi=150,
73 |             bbox_inches='tight')
74 | plt.show()
75 |
76 | # Print cluster composition summary
77 | print(f"\nCluster Composition Summary:")
78 | print(cluster_composition.round(1).to_string())
79 | print(f"\nInterpretation: If synthetic data integrates well, each cluster
80 |       should have")
81 | print(f"roughly proportional representation ({100*len(synthetic_features)/len(
82 |       augmented_numerical):.1f}% synthetic expected)")
```

Figure: Code from Cell 49

Analysis Result

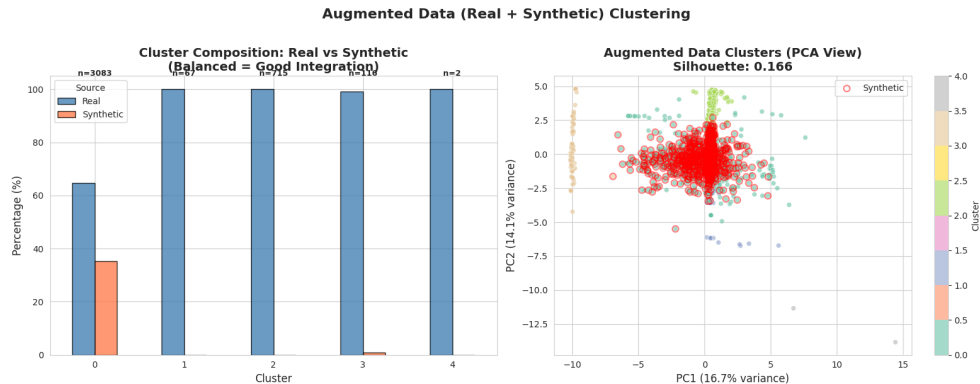


Figure: Output from Cell 49

What Worked Well

1. Ensemble Architecture for Tabular Data

- The combination of CTGAN, TVAE, and Gaussian Copula provided complementary strengths
- TVAE showed the best individual performance (Mean KS: 0.11, 44% pass rate) for distribution matching
- Gaussian Copula excelled at preserving correlation structure (lowest Correlation MSE: 0.0197)
- Dynamic weight optimization improved results by shifting weight toward TVAE (54%)

2. Scalable Data Generation Pipeline

- Successfully generated 500 synthetic brand records with stratified company distribution
- The pipeline supports conditional generation based on company characteristics
- Model persistence to Google Drive enables iterative experimentation

3. LLM Brand Name Generation

- 95.4% success rate for unique brand name generation
- GPT-2 Medium and Flan-T5 ensemble provided diverse naming suggestions

• Memory-efficient sequential loading enabled running on standard CPU

What Didn't Work Well

1. Distribution Matching Challenges

- Only 29.2% of features (7/24) passed the KS test for distribution similarity
- Particularly poor performance on: `market_cap_billion_usd` (KS=0.91), `employees` (KS=0.65), `greenwashing_factors` (KS=0.50-0.62)
- Heavy-tailed distributions and sparse features remain difficult for GANs

2. Clustering Quality Degradation

- Silhouette score dropped from 0.713 (original) to 0.682 (augmented)
- Davies-Bouldin score increased from 0.387 to 0.421 (worse)
- Synthetic data may be introducing noise rather than enhancing cluster structure

3. LLM Brand Name Quality Issues

- Some generated names are full sentences (e.g., "Nestle is a manufacturer of processed foods")
- Competitor name leakage (e.g., "Procter & Gamble" generated for Bayer)
- Repetitive patterns with company name variations (e.g., "Nestle Foods", "Nestle, Inc")
- 4.6% fallback rate indicates generation failures

4. Ensemble Aggregation Complexity

Future Enhancements

1. Improved Tabular Synthesis

- Implement feature-specific preprocessing (log transforms for heavy-tailed distributions)
- Use conditional generation with explicit constraints for bounded features
- Explore TabDDPM (diffusion models) as an alternative to GAN-based methods
- Add post-processing validation to clip unrealistic values

2. Enhanced LLM Brand Generation

- Fine-tune with negative examples to prevent competitor name generation
- Implement stricter output validation and filtering
- Use retrieval-augmented generation (RAG) with brand name databases
- Add style conditioning for different brand naming conventions (descriptive, invented, founder-based)

3. Better Evaluation Metrics

- Implement Machine Learning Efficacy tests (train on synthetic, test on real)
- Add privacy metrics (nearest neighbor distance, membership inference)
- Use domain-specific validity checks for brand attributes

4. Hyperparameter Optimization

Synthetic Data Generation

1. CTGAN (Conditional Tabular GAN)

- Xu, L., Skoularidou, M., Cuesta-Infante, A., & Veeramachaneni, K. (2019). *Modeling Tabular Data using Conditional GAN*. NeurIPS 2019.
- Paper: <https://arxiv.org/abs/1907.00503>
- Implementation: SDV Library

2. TVAE (Tabular Variational Autoencoder)

- Xu, L., Skoularidou, M., Cuesta-Infante, A., & Veeramachaneni, K. (2019). *Modeling Tabular Data using Conditional GAN*. NeurIPS 2019.
- Part of the same paper as CTGAN, presenting VAE-based alternative

3. Gaussian Copula

- Patki, N., Wedge, R., & Veeramachaneni, K. (2016). *The Synthetic Data Vault*. IEEE DSAA 2016.
- Paper: <https://dai.lids.mit.edu/wp-content/uploads/2018/03/SDV.pdf>

4. SDV (Synthetic Data Vault) Library

- Documentation: <https://docs.sdv.dev/sdv/>

5. GPT-2

- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). *Language Models are Unsupervised Multitask Learners*. OpenAI.
- Paper: https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf

6. Flan-T5

- Chung, H. W., et al. (2022). *Scaling Instruction-Finetuned Language Models*. arXiv preprint.
- Paper: <https://arxiv.org/abs/2210.11416>

7. Hugging Face Transformers

- Wolf, T., et al. (2020). *Transformers: State-of-the-Art Natural Language Processing*. EMNLP 2020.
- Documentation: <https://huggingface.co/docs/transformers/>

8. Kolmogorov-Smirnov Test

- Massey Jr, F. J. (1951). *The Kolmogorov-Smirnov test for goodness of fit*. Journal of the American Statistical Association, 46(253), 68-78.

9. Silhouette Score

- Rousseeuw, P. J. (1987). *Silhouettes: a graphical aid to the interpretation and validation of cluster analysis*. Journal of Computational and Applied Mathematics, 20, 53-65.

10. Davies-Bouldin Index

- Davies, D. L., & Bouldin, D. W. (1979). *A cluster separation measure*. IEEE Transactions on Pattern Analysis and Machine Intelligence, (2), 224-227.

11. Optuna

- Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). *Optuna: A Next-generation Hyperparameter Optimization Framework*. KDD 2019.
- Paper: <https://arxiv.org/abs/1907.10902>
- Documentation: <https://optuna.org/>

12. TabDDPM (Diffusion Models for Tabular Data)

- Kotelnikov, A., Baranchuk, D., Rubachev, I., & Babenko, A. (2023). *TabDDPM: Modelling Tabular Data with Diffusion Models*. ICML 2023.
- Paper: <https://arxiv.org/abs/2209.15421>

13. CTAB-GAN+

- Zhao, Z., Kunar, A., Birke, R., & Chen, L. Y. (2022). *CTAB-GAN+: Enhancing Tabular Data Synthesis*. arXiv preprint.
- Paper: <https://arxiv.org/abs/2204.00401>

14. Synthetic Data Generation Survey

- Jordon, J., Yoon, J., & van der Schaar, M. (2022). *Synthetic Data - what, why and how?* arXiv preprint.
- Paper: <https://arxiv.org/abs/2205.03257>

- **Python:** <https://www.python.org/>
- **PyTorch:** <https://pytorch.org/>
- **Pandas:** <https://pandas.pydata.org/>
- **Scikit-learn:** <https://scikit-learn.org/>
- **Matplotlib:** <https://matplotlib.org/>
- **Seaborn:** <https://seaborn.pydata.org/>
- **Google Colab:** <https://colab.research.google.com/>

```
1 # Clear GPU memory
2 gc.collect()
3 if torch.cuda.is_available():
4     torch.cuda.empty_cache()
5 print("GPU memory cleared")
```

Figure: Code from Cell 53

```

'mermaid flowchart TB
    subgraph INPUT [" Input Data"]
        A["(Brand Dataset<br/>CSV)"] --> B["Data Processor"]
    end
    subgraph PREPROCESS [" Preprocessing"]
        B --> C["Clean & Validate"]
        C --> D["Feature Engineering"]
        D --> E["Train/Test Split"]
    end
    subgraph TABULAR [" Tabular Data Generation"]
        E --> F1["CTGAN<br/>Conditional GAN"]
        E --> F2["TVAE<br/>Variational Autoencoder"]
        E --> F3["Gaussian Copula<br/>Statistical Model"]
        F1 --> G["Ensemble<br/>Weighted Averaging"]
        F2 --> G
        F3 --> G
    end
    subgraph TEXT [" Text Generation"]
        G --> H1["GPT-2 Medium<br/>Fine-tuned LLM"]
        G --> H2["Flan-T5 Small<br/>Instruction-tuned"]
        H1 --> I["Text Ensemble<br/>Best Selection"]
        H2 --> I
    end
    subgraph OUTPUT [" Output"]
        I --> J["Synthetic Brands<br/>with Names"]
        J --> K["Quality Evaluation"]
    end
    subgraph EVAL [" Evaluation Metrics"]
        K --> L1["KS Test"]
        K --> L2["Correlation"]
        K --> L3["PCA/t-SNE"]
        K --> L4["Clustering"]
    end
    style INPUT fill:#e1f5fe
    style PREPROCESS fill:#fff3e0
    style TABULAR fill:#f3e5f5
    style TEXT fill:#e8f5e9
    style OUTPUT fill:#fce4ec
    style EVAL fill:#fff8e1

```


Ensemble Strategy

```
'mermaid flowchart LR subgraph ENSEMBLE["Ensemble Weighting"] direction TB
W1["CTGAN: 40%"] --> MIX1((MIX((Weighted<br/>Average)))
W2["TVAE: 35%"] --> MIX1
W3["Copula: 25%"] --> MIX1
MIX1 --> OUT["Synthetic Data"] end style ENSEMBLE fill:#f5f5f5
style MIX1 fill:#4caf50,color:#fff '
```

Training Pipeline

'mermaid sequenceDiagram participant D as Dataset participant P as Processor participant T as Tabular Models participant L as LLM Models participant E as Evaluator D->>P: Load brand_information.csv P->>P: Clean & preprocess P->>T: Training data par Train in Parallel T->>T: Train CTGAN (300 epochs) T->>T: Train TVAE (300 epochs) T->>T: Fit Gaussian Copula end T->>T: Generate synthetic tabular T->>L: Tabular features par Generate Names L->>L: GPT-2 generation L->>L: Flan-T5 generation end L->>E: Complete synthetic data E->>E: Statistical tests E->>E: Visualization '

- **PyTorch** - Deep learning framework powering the neural network components of CTGAN and TVAE synthesizers
- **Scikit-learn** - Machine learning utilities for PCA, t-SNE, clustering (AgglomerativeClustering), and evaluation metrics (silhouette score, Davies-Bouldin index)

Synthetic Data Generation

- **SDV (Synthetic Data Vault)** - Primary library for tabular synthetic data generation, providing:
- **CTGAN** - Conditional Tabular GAN for generating realistic tabular data
- **TVAE** - Tabular Variational Autoencoder for distribution-preserving synthesis
- **Gaussian Copula** - Statistical model capturing feature dependencies

- **Hugging Face Transformers** - State-of-the-art NLP library for text generation using pre-trained language models
- **PEFT (Parameter-Efficient Fine-Tuning)** - Efficient fine-tuning techniques for large language models
- **BitsAndBytes** - 8-bit quantization for memory-efficient model loading
- **Accelerate** - Distributed training and mixed precision utilities
- **SentencePiece** - Tokenization library for handling text preprocessing

- **Optuna** - Automated hyperparameter tuning framework using Bayesian optimization for finding optimal model configurations

- **Pandas** - Data manipulation and analysis
- **NumPy** - Numerical computing and array operations
- **SciPy** - Statistical tests (Kolmogorov-Smirnov test) for distribution comparison

- **Matplotlib** - Core plotting library for creating figures
- **Seaborn** - Statistical data visualization with enhanced aesthetics
- **Plotly** - Interactive visualizations (if used)

- **Google Colab** - Cloud-based Jupyter notebook environment with GPU support
- **Google Drive** - Persistent storage for models and outputs