# Smarter Models, Fewer Features: Why PCA Matters in Feature Selection
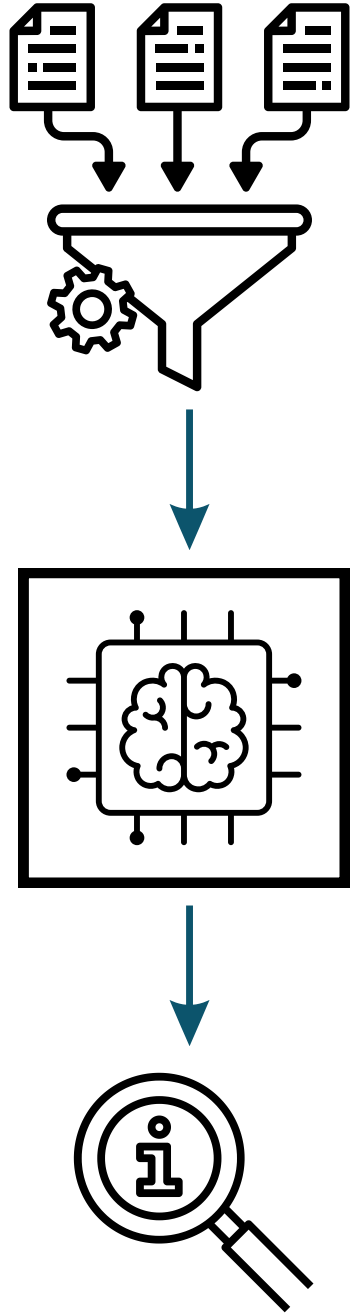
Correlation Heatmap - ABEV3
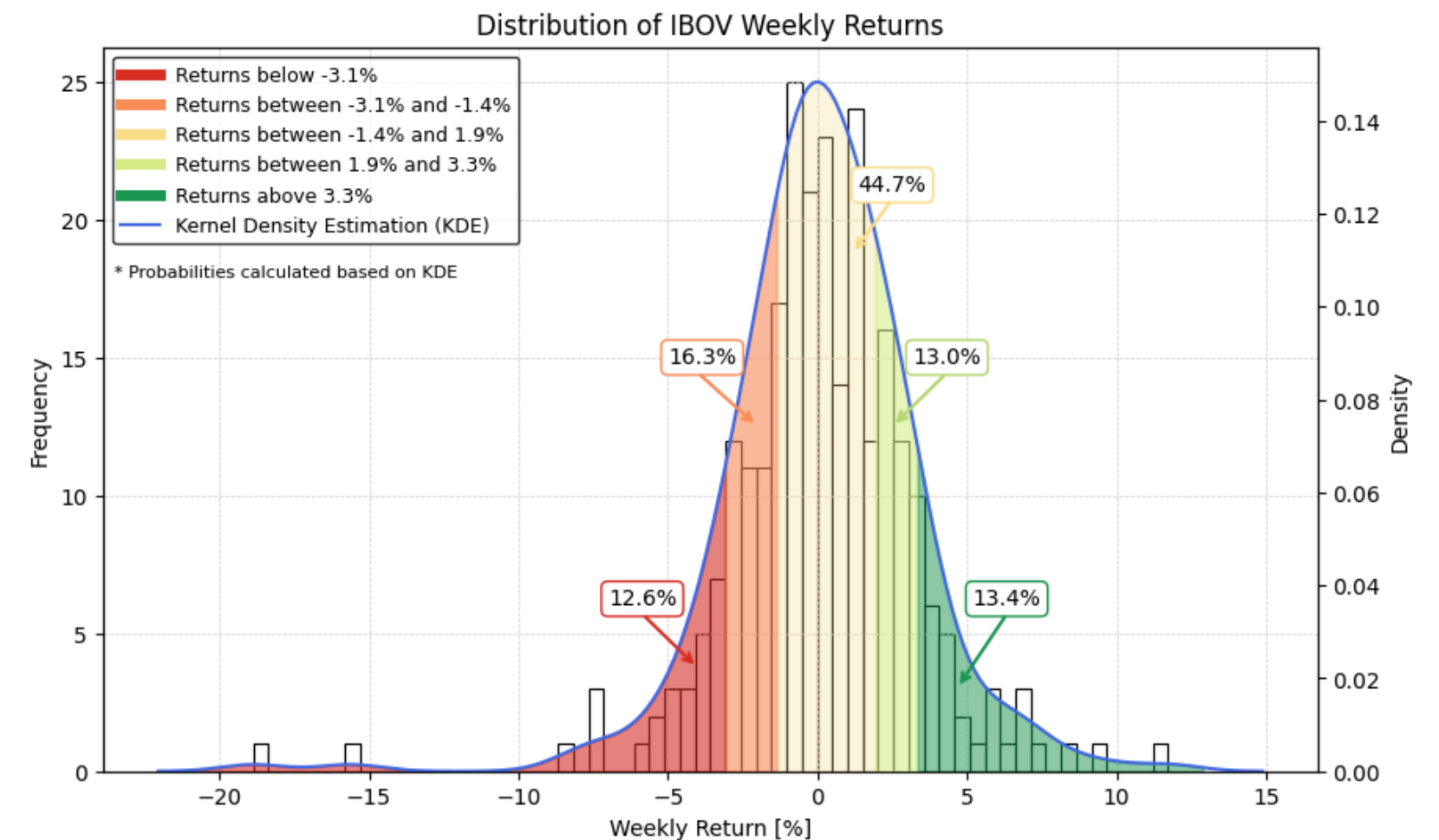
**'Raw' data as input:**
- Total of 19 features.
- Dataset combines technical, fundamental, and sentiment analysis data from the ABEV3 stock.

**Random Forest algorithm:**
- Since the dataset is small, parameters are set to **stress-test**, simulating computational intensity.
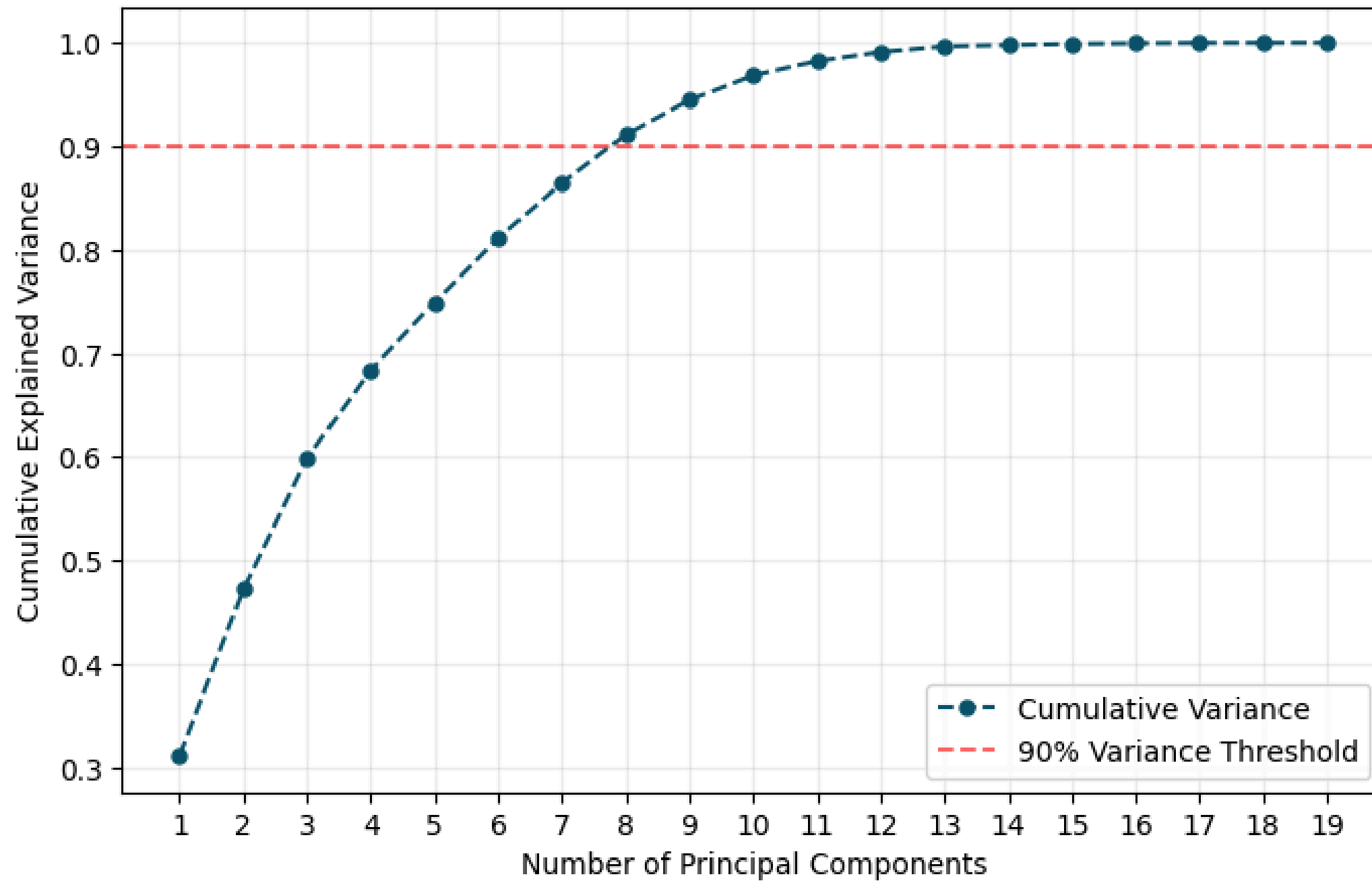
**Output**
- (Encoded) categorical variable representing the next week's return based on histogram-defined regions.



Distribution of IBOV Weekly Returns

Returns below -3.1%
Returns between -3.1% and -1.4%
Returns between -1.4% and 1.9%
Returns between 1.9% and 3.3%
Returns above 3.3%
Kernel Density Estimation (KDE)

* Probabilities calculated based on KDE
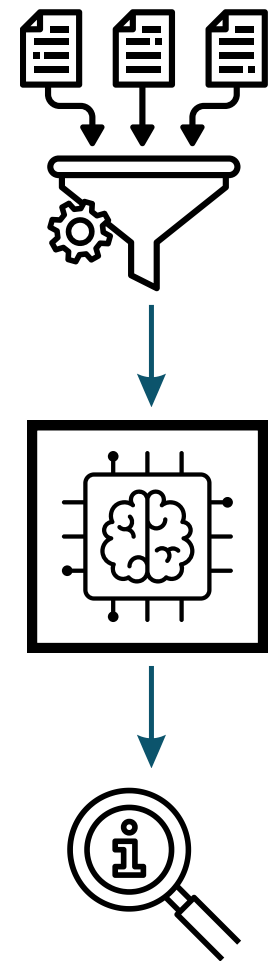
44.7%
16.3%
13.0%
12.6%
13.4%

# Principal Component Analysis (PCA)



Cumulative Explained Variance vs. Number of Principal Components

- The **cumulative explained variance** is the **sum** of the variance proportions explained by each principal component (PC).

- It indicates how much information is retained in **lower-dimensional space**.

- Helps decide the **optimal number** of principal components to retain.

# Example 2: Principal Components Classifier

**8 Principal Components as input:**
- ~90% of the variance.

**Random Forest algorithm**
- Same parameters as before.

| | Overall Results | |
|---|---|---|
| | **RAW inputs** | **PCA** |
| **Training time** | 6.38 s | 4.42 s |
| **Overall Accuracy** | 23.8 % | 33.3 % |