# Machine Learning Pipeline for MAGIC Gamma Telescope Dataset

## (a) Data & Pre-processing

### Dataset Selection

For this machine learning pipeline, I selected the MAGIC Gamma Telescope dataset from the UCI Machine Learning Repository. This dataset was chosen for several compelling reasons:

1. **Appropriate complexity and size**: With 19,020 samples and 10 features, the dataset is substantial enough to demonstrate meaningful patterns while remaining manageable computationally.
2. **Clear classification objective**: The dataset presents a well-defined binary classification problem - distinguishing gamma particle signals (signal) from hadron particles (background) based on telescope image data.
3. **Real-world scientific application**: The dataset represents measurements from the MAGIC (Major Atmospheric Gamma Imaging Cherenkov) telescope, which detects gamma rays in astrophysical observations. This gives the analysis practical significance beyond academic exercise.
4. **Balance of features**: The dataset contains continuous numerical features representing various aspects of the detected Cherenkov radiation, allowing for comprehensive modeling approaches.
5. **Class imbalance consideration**: With approximately 65% gamma (class 1) and 35% hadron (class 0) samples, the dataset presents a moderate class imbalance that requires careful handling but isn't extreme.

### Data Description

The MAGIC Gamma Telescope dataset contains 10 continuous features that describe the parameters of detected Cherenkov radiation in a telescope, with the target variable indicating whether the radiation came from a gamma particle (signal) or hadron particle (background). The features include:

- **fLength**: Major axis length of the elliptical image (mm)
- **fWidth**: Minor axis length of the elliptical image (mm)
- **fSize**: 10-log of sum of content of all pixels (# photons)
- **fConc**: Ratio of sum of two highest pixels over fSize
- **fConc1**: Ratio of highest pixel over fSize
- **fAsym**: Distance from highest pixel to center, projected onto major axis
- **fM3Long**: 3rd root of third moment along major axis (mm)
- **fM3Trans**: 3rd root of third moment along minor axis (mm)
- **fAlpha**: Angle of major axis with vector to origin (degrees)
- **fDist**: Distance from origin to center of ellipse (mm)

### Pre-processing Steps

I implemented the following pre-processing steps to prepare the data for modeling:

1. **Data loading and inspection**: I loaded the dataset and verified there were no missing values, which simplified preprocessing.
2. **Class encoding**: Converted class labels from categorical ('g' for gamma, 'h' for hadron) to binary (1 for gamma, 0 for hadron).
3. **Train-test split**: Split the data into 70% training and 30% testing, using stratification to maintain the same class distribution in both sets.
4. **Feature standardization**: Applied StandardScaler to normalize features to zero mean and unit variance. This step was essential because:

- Features have vastly different scales (e.g., fDist has values up to 495, while fConc has values between 0.013 and 0.893)
- Many algorithms (especially SVM, logistic regression, and neural networks) perform better with standardized features
- Distance-based metrics used in certain algorithms would be dominated by features with larger scales without standardization

The preprocessing approach ensures that the models are trained on properly scaled data with representative proportions of each class, reducing the risk of bias and improving convergence during training.

# (b) Exploratory Data Analysis & Visualization

## Distribution Analysis

The initial exploration of the dataset revealed several interesting patterns:

1. **Class distribution**: The dataset contains approximately 65% gamma signals (class 1) and 35% hadron background (class 0), showing a moderate imbalance that needs to be considered during model evaluation.
2. **Feature distributions**:
   - Many features show distinct distributions between the classes, particularly fAlpha, which shows a clear separation (hadrons have much higher values).
   - fLength and fWidth show that gamma signals tend to have smaller elliptical images compared to hadrons.
   - fConc and fConc1 have bimodal distributions, with gamma signals having higher concentration values.
3. **Outliers**: Several features contain outliers, particularly fLength, fWidth, and fAsym, which could potentially impact model performance if not properly handled through robust algorithms or standardization.

## Correlation Analysis

The correlation heatmap revealed important relationships between features:

1. **High correlations**: fLength and fWidth are highly correlated (0.92), suggesting potential redundancy.
2. **Target correlations**: fAlpha has the strongest negative correlation with the target class (-0.45), indicating it may be one of the most discriminative features. This aligns with physics, as gamma rays tend to come from the center of the field of view (low alpha).
3. **Feature groups**: The features can be grouped into related clusters based on correlations, which aligns with their physical interpretation:
   - Size-related features (fLength, fWidth, fSize)
   - Concentration-related features (fConc, fConc1)
   - Position-related features (fAlpha, fDist)

## Dimensionality Reduction

PCA and t-SNE were applied to visualize the data in lower dimensions:

1. **PCA**: The first two principal components explained 58% of the variance (42.2% and 15.8% respectively). The PCA projection showed some separation between classes, with a significant overlap in the center region. This suggests that linear separation might be challenging.
2. **t-SNE**: The t-SNE visualization showed better class separation than PCA, revealing clustered regions of each class. This suggests that non-linear transformations may be beneficial in modeling this data.
3. **KMeans clustering**: Unsupervised KMeans clustering with k=2 was applied to see if natural clusters aligned with the classes. The clustering achieved very low alignment with actual classes (Adjusted Rand Index of 0.0062), indicating that simple clustering is insufficient and supervised methods are necessary.

## Feature Importance by Variance

Analysis of feature variance showed:

1. **High variance features**: fDist (5584.8), fAsym (3505.4), and fM3Long (2601.0) have the highest variance.
2. **Low variance features**: fSize (0.223), fConc (0.033), and fConc1 (0.012) have much lower variance.

## EDA Insights for Modeling

The exploratory analysis provided several insights that influenced the modeling approach:

1. **Model selection**: The moderate class imbalance and complex decision boundaries visible in the PCA and t-SNE plots suggested that ensemble methods like Random Forest and Gradient Boosting might perform well on this dataset.
2. **Feature engineering**: The high correlation between some features suggested that dimensionality reduction might be beneficial, but the strong relationship between individual features and the target suggested that preserving the original features was important.
3. **Evaluation metrics**: The class imbalance indicated that accuracy alone would be insufficient as an evaluation metric; ROC AUC and F1 score would provide more balanced assessments of model performance.
4. **Hyperparameter focus**: The complex relationships between features suggested that models with greater flexibility (deeper trees, more complex kernels in SVM) might be necessary to capture the classification boundaries.
5. **Standardization impact**: The large differences in feature scales confirmed that standardization was a necessary preprocessing step, especially for distance-based algorithms and SVM.

The EDA process clearly showed that this dataset requires sophisticated modeling approaches to effectively separate gamma signals from hadron background, with particular attention to feature scaling and non-linear relationships.

# (c) Modeling & Model Validation

## Model Selection

Based on the insights from EDA, I selected several classification algorithms to evaluate:

1. **Logistic Regression**: As a baseline linear model to assess if linear separation is possible and to serve as a benchmark for more complex models.
2. **K-Nearest Neighbors**: Chosen because the t-SNE visualization suggested local neighborhood structures might be informative for classification.
3. **Support Vector Machine**: Selected for its ability to find complex decision boundaries using kernel transformations, which seemed necessary given the PCA and t-SNE visualizations.
4. **Random Forest**: Chosen for its robustness to outliers, ability to capture non-linear relationships, and capacity to handle the moderate class imbalance.
5. **Gradient Boosting**: Selected for its typically strong performance on structured data and ability to sequentially improve on difficult-to-classify instances.

This selection provides a diverse set of algorithms with different strengths, allowing for a comprehensive evaluation of approaches for this dataset.

## Hyperparameter Tuning Approach

For each model, I implemented a systematic hyperparameter tuning process:

1. **Grid search with cross-validation**: Used GridSearchCV with 5-fold stratified cross-validation to ensure consistent class distributions across folds.
2. **Model-specific hyperparameters**: Tuned the following key parameters for each model:
   - Logistic Regression: Regularization strength (C) and solver algorithm
   - KNN: Number of neighbors, weighting scheme, and distance metric
   - SVM: Kernel type, regularization parameter (C), and gamma parameter
   - Random Forest: Number of trees, maximum depth, minimum samples for splits and leaves
   - Gradient Boosting: Number of trees, learning rate, maximum depth, and subsampling ratio

3. **Optimization metric**: ROC AUC was chosen as the primary optimization metric to account for class imbalance and prioritize ranking performance.
4. **Parallelization**: Utilized n_jobs=-1 to parallelize the tuning process and reduce computational time.

This approach ensures a thorough exploration of the hyperparameter space while maintaining computational efficiency.

## Validation Strategy

The validation strategy consisted of:

1. **Stratified k-fold cross-validation**: Used 5-fold stratified cross-validation to ensure representative class distributions in each fold.
2. **Train/validation/test separation**: Maintained a strict separation between:
   - Training data (used for model fitting)
   - Validation data (used for hyperparameter tuning via cross-validation)
   - Test data (used only for final evaluation)
3. **Multiple evaluation metrics**: Assessed performance using multiple metrics to gain a comprehensive understanding:
   - Accuracy: Overall correctness of predictions
   - Precision: Proportion of true positives among positive predictions
   - Recall: Proportion of true positives identified correctly
   - F1 Score: Harmonic mean of precision and recall
   - ROC AUC: Area under the ROC curve, measuring ranking performance

This validation approach ensures robust model selection and reliable performance estimates.

## Model Performance

The performance of each model on the test set:

| Model | Accuracy | Precision | Recall | F1 Score | ROC AUC |
|---|---|---|---|---|---|
| Logistic Regression | 0.7808 | 0.7923 | 0.8970 | 0.8414 | 0.8276 |
| K-Nearest Neighbors | 0.8412 | 0.8214 | 0.9649 | 0.8874 | 0.8990 |
| Support Vector Machine | 0.8747 | 0.8613 | 0.9616 | 0.9087 | 0.9284 |
| Random Forest | 0.8836 | 0.8843 | 0.9441 | 0.9132 | 0.9374 |
| Gradient Boosting | 0.8873 | 0.8898 | 0.9430 | 0.9156 | 0.9377 |

## Best Model Analysis

The Gradient Boosting model emerged as the best performer with:

- Highest ROC AUC: 0.9377
- Highest Accuracy: 0.8873
- Highest F1 Score: 0.9156

The optimal hyperparameters for this model were:

- learning_rate: 0.1
- max_depth: 7
- n_estimators: 200
- subsample: 0.8

## Model Comparison and Justification

Several observations about the model performance:

1. **Progression of complexity**: There is a clear improvement in performance as model complexity increases, from Logistic Regression (simplest) to Gradient Boosting (most complex).
2. **Non-linear advantage**: The substantial performance gap between Logistic Regression and the other models confirms that non-linear decision boundaries are critical for this dataset, as suggested by the EDA.
3. **Ensemble superiority**: Both ensemble methods (Random Forest and Gradient Boosting) outperformed single models, demonstrating the benefit of combining multiple weak learners for this problem.
4. **Gradient Boosting edge**: Gradient Boosting likely performed best because:
   - It sequentially focuses on harder-to-classify examples
   - The moderate depth (7) allows for complex patterns while avoiding overfitting
   - The subsample parameter (0.8) introduces randomness that improves generalization
   - The learning rate (0.1) provides a good balance between convergence speed and precision
5. **KNN limitations**: Despite the promising t-SNE visualization, KNN didn't perform as well as tree-based models, likely due to the curse of dimensionality in the original 10-dimensional space.

The validation approach confirms that Gradient Boosting is robustly the best model, with its superior performance consistent across multiple evaluation metrics.

# (d) Communication of Results & Interpretation

## Model Performance Summary

The Gradient Boosting model achieved excellent performance on this gamma/hadron classification task:

- **Accuracy: 88.7%** - Almost 9 out of 10 particles are correctly classified
- **Precision: 89.0%** - When the model predicts gamma, it's right 89% of the time
- **Recall: 94.3%** - The model captures 94.3% of all actual gamma particles
- **F1 Score: 91.6%** - Strong balanced performance between precision and recall
- **ROC AUC: 93.8%** - Excellent discrimination ability between classes

## Confusion Matrix Analysis

The confusion matrix reveals:

- **True Negatives (Hadrons correctly identified)**: 1574 instances
- **False Positives (Hadrons misclassified as gamma)**: 432 instances
- **False Negatives (Gamma misclassified as hadrons)**: 211 instances
- **True Positives (Gamma correctly identified)**: 3489 instances

Key insights:

1. **Higher hadron misclassification**: The model has more difficulty correctly identifying hadrons (78.5% accuracy) compared to gamma particles (94.3% accuracy).
2. **Real-world impact**: In astronomical observations, the higher recall for gamma particles means we're less likely to miss actual gamma ray signals, which is often prioritized in detection scenarios.

## Feature Importance Analysis

The feature importance from the Gradient Boosting model provides valuable insights:

1. **Top features**:
   - fAlpha (28.1%): Angle of major axis with vector to origin
   - fLength (22.5%): Major axis length of elliptical image
   - fSize (12.1%): 10-log of sum of content of all pixels

- ○ fWidth (10.6%): Minor axis length of elliptical image
2. **Physics interpretation**:
    - ○ The dominant importance of fAlpha aligns with astrophysical theory - gamma rays tend to come directly from the observed source (low alpha angle), while hadrons can come from any direction.
    - ○ The importance of size parameters (fLength, fWidth, fSize) confirms that the shower shape is highly discriminative between gamma and hadron events.
3. **Least important features**:
    - ○ fM3Trans (2.6%)
    - ○ fAsym (3.2%)
    - ○ fConc1 (4.2%)

This feature importance analysis provides valuable feedback to physicists about which measurements are most crucial for discrimination, potentially informing future telescope designs or measurement priorities.

## ROC Curve Interpretation

The ROC curve (with AUC of 0.938) illustrates:

1. **Excellent discrimination ability**: The curve bows strongly toward the upper-left corner, indicating the model can achieve high true positive rates while maintaining low false positive rates.
2. **Operational flexibility**: The curve offers multiple potential operating points depending on scientific priorities:
    - ○ A threshold yielding ~80% true positive rate with only ~10% false positive rate
    - ○ A high-sensitivity setting with ~95% true positive rate at the cost of ~30% false positives
    - ○ A high-specificity setting with <5% false positive rate while still capturing ~60% of gamma signals

This allows astronomers to adjust classification thresholds based on their specific research needs - whether prioritizing pure samples or maximizing detection of rare events.

## Model Limitations

Despite the strong performance, several limitations should be acknowledged:

1. **Misclassified instances**: About 11% of particles are still misclassified, with a higher error rate for hadron particles.
2. **Feature correlations**: Some features show high correlation (like fLength and fWidth), indicating potential redundancy that could be addressed in future modeling.
3. **Complex black-box model**: While Gradient Boosting offers excellent performance, its complex ensemble nature makes detailed interpretation challenging beyond feature importance.

## Practical Applications

The developed model can be applied to:

1. **Gamma-ray astronomy**: Efficiently filtering telescope data to identify gamma ray signals from various cosmic sources
2. **Automated analysis pipelines**: Integration into real-time data processing systems for telescopes
3. **Educational demonstrations**: Teaching the application of machine learning to particle physics problems

## Conclusion

The machine learning pipeline successfully developed a high-performing Gradient Boosting model for gamma/hadron discrimination with 88.7% accuracy and 93.8% ROC AUC. The analysis confirmed that tree-based ensemble methods are particularly well-suited to this classification task, and identified the most discriminative features aligned with astrophysical understanding of the phenomena.

The angle of the elliptical image (fAlpha) and the size parameters (fLength, fWidth, fSize) are the most important features for discrimination, which provides valuable scientific validation of the model. The resulting classifier offers flexible operational points via threshold adjustment, making it adaptable to different research priorities in gamma-ray astronomy.