

Title

A Comprehensive Machine Learning Framework for Career Recommendation Using Multimodal Psychometric, Academic, and Demographic Attributes: An End-to-End Data Science Research Study

Abstract

A career choice is a critical developmental decision in a student's academic life, significantly affecting his or her long-term professional career path, skill development patterns, and socio-economic mobility. Since many interdisciplinary fields have come into existence, labor markets are facing rapid changes, and the cognitive profiles of students have also become complicated; traditional counseling methods, usually based on subjective opinions, can seldom offer personalized and data-driven career recommendations. To bridge these gaps, the proposed research introduces an all-encompassing, end-to-end, reproducible machine learning framework developed to render highly personalized career recommendations. It uses the rich combination of psychometric intelligence indicators, metrics on academic performance, demographic attributes, and expressed student interests.

Using the *Career_Guidance.csv* dataset, we build a rigorous 14-stage computational pipeline incorporating data ingestion, anomaly detection, missing-value treatment, multilevel preprocessing, exploratory data analysis, correlation analytics, feature engineering, supervised encoding methods, dimensionality reduction techniques (including PCA), and advanced model training strategies. We further integrate systematic hyperparameter tuning using GridSearchCV, followed by an extensive benchmarking protocol across multiple algorithms such as Random Forest, Gradient Boosting Machines, and Logistic Regression. Performance evaluation extends beyond conventional accuracy—emphasizing macro-averaged precision, recall, and F1-scores, as well as Cohen's Kappa, Matthews Correlation Coefficient (MCC), and multiclass ROC-AUC. This paper aims to ensure robustness, fairness, and stability across a wide range of target categories. The results point toward better generalization and predictive performance of ensemble learning methods on heterogeneous educational datasets. Gradient Boosting and Random Forest models tend to outperform all linear model baselines consistently, particularly for scenarios that involve complex nonlinear interactions between psychological attributes and academic features. Feature importance rankings, partial dependence plots, and model interpretability techniques using SHAP value explanations provide insightful understandings into how specific cognitive traits, subject proficiencies, and interest patterns drive career suitability.

Beyond just technical implementation, the manuscript places the entire system within the domains of educational psychology, machine learning theory, behavioral decision science, and data-driven counseling methodologies. The proposed framework not only increases the reliability of career recommendations by linking psychological constructs to computational intelligence but also unearths actionable patterns useful to educators, policymakers, and institutions. This study finally lays the foundation for the development of scalable, interpretable, AI-driven career guidance systems with the full capacity to support personalized academic planning and informed decision-making in diverse student populations.

1. Introduction

Career guidance plays a pivotal role in shaping students' academic satisfaction, competency development, and long-term professional trajectories. The decisions made during early educational stages influence not only job opportunities but also psychological well-being, motivation, and adaptability in rapidly evolving job markets. Despite its importance, large educational institutions face increasing challenges in offering personalized, evidence-based counselling to thousands of students with diverse cognitive abilities, socio-economic backgrounds, academic profiles, and career aspirations. Traditional guidance practices—often grounded in subjective assessments or generic aptitude tests—lack the ability to rigorously analyze complex, multidimensional student data patterns. Machine learning (ML), with its capacity to uncover hidden structures and model non-linear relationships, presents a transformative opportunity to deliver scalable, objective, and highly contextualized career recommendations.

This manuscript introduces a comprehensive, fully operational ML-driven career guidance framework constructed using the *Career_Guidance.csv* dataset. The system goes beyond conventional predictive modelling by integrating a rich set of analytical components, ensuring that the methodology is not only technically sound but also theoretically grounded and ethically responsible. The study expands the initial implementation into a structured, interdisciplinary research narrative that examines multiple dimensions of educational decision-making, including:

- **Theoretical motivations**, such as cognitive-fit theory, person–environment alignment, career development models, and the psychology of interest formation.
- **Educational data mining (EDM) principles**, addressing student modelling, learning analytics, behavioural feature extraction, and the relevance of multimodal data in

academic decision support systems.

- **Psychometric evaluation frameworks**, exploring intelligence constructs, aptitude–achievement correlations, personality influences, and standardized assessment theory to contextualize the dataset’s attributes.
- **Machine learning model explainability**, incorporating SHAP, LIME, permutation importance, and interpretable ML principles to ensure transparency and trustworthiness in recommendations.
- **Ethical considerations and fairness implications**, focusing on bias detection, demographic parity, responsible use of student data, and equity in automated decision-making.
- **Practical deployment potential**, discussing model integration into academic portals, adaptive counselling interfaces, real-time recommendation engines, and institution-wide decision support infrastructure.

By weaving together technical, psychological, ethical, and institutional perspectives, this extended manuscript elevates the original ML prototype into a holistic, academically rigorous study. It is suitable for submission as a research paper, undergraduate thesis, master’s dissertation, or institutional policy project, offering a blend of methodological innovation and real-world applicability. The work ultimately positions machine learning not merely as a computational tool but as a catalyst for reforming student guidance systems through data-driven intelligence, transparency, and personalization.

2. Literature Review (Deep Expansion)

This section synthesizes foundational theories, frameworks, and empirical research surrounding Educational Data Mining (EDM), psychometric intelligence, machine learning classification, and decision-support systems. The goal is to situate the proposed methodology within established academic domains while highlighting the technological and theoretical gaps that this manuscript aims to address.

2.1 Educational Data Mining (EDM)

Educational Data Mining (EDM) is an interdisciplinary research field focused on extracting meaningful patterns from educational environments—ranging from classroom assessments to digital learning platforms. EDM enables institutions to make data-driven decisions that enhance student performance, optimize instruction, and personalize learning pathways.

Romero and Ventura (2010) identify EDM as a cornerstone of modern academic analytics, emphasizing its applications across predictive modeling, clustering, feature discovery, and recommendation systems. Core use cases of EDM include:

- **Predicting academic success and performance trajectories**, enabling early interventions for curriculum adjustments or tutoring support.
- **Identifying at-risk students using classification, survival analysis, and temporal models**, supporting institutional retention initiatives.
- **Designing personalized learning pathways** based on student behaviour, competency mastery, and learning style preferences.
- **Developing recommendation systems** capable of suggesting courses, study resources, or career options based on individual learner profiles.
- **Evaluating teaching strategies**, curriculum effectiveness, and classroom interaction patterns.

EDM draws from a broad spectrum of disciplines—including statistics, cognitive psychology, pedagogy, machine learning, and human–computer interaction—to build systems that improve educational decision-making. As learning environments generate increasingly large, multimodal datasets, EDM provides powerful frameworks for uncovering the hidden relationships between student characteristics and their academic or career outcomes.

2.2 Psychometrics and Multiple Intelligence Theory (MIT)

Psychometrics is the scientific discipline concerned with measuring psychological traits, cognitive abilities, and behavioural tendencies through standardized instruments. It underpins modern educational assessments by ensuring reliability, validity, and interpretability in evaluating student characteristics.

Within psychometrics, **Howard Gardner's Multiple Intelligence Theory (MIT)** stands as one of the most influential frameworks for conceptualizing learner diversity. Gardner challenged the traditional view of intelligence as a single, general factor (g), instead proposing a multifaceted

model comprising distinct but interrelated cognitive domains. These domains—several of which appear in the Career_Guidance dataset—include:

- **Linguistic Intelligence**
- **Logical–Mathematical Intelligence**
- **Interpersonal Intelligence**
- **Intrapersonal Intelligence**
- **Spatial–Visual Intelligence**
- **Musical Intelligence**
- **Naturalistic Intelligence**
- **Bodily–Kinesthetic Intelligence**

Extensive research shows that career suitability is often aligned with dominant intelligence strengths. For instance:

- **Logical–Mathematical** → Engineering, Data Science, Finance
- **Linguistic** → Law, Journalism, Public Relations
- **Spatial** → Architecture, Animation, Product Design
- **Interpersonal** → Management, HR, Sales, Teaching

MIT provides a conceptual backbone for understanding why students with similar academic scores may excel in different professions. By integrating MIT into machine learning pipelines, career guidance becomes more holistic, psychologically grounded, and personalized.

2.3 Machine Learning Approaches for Career Prediction

Machine learning (ML) has increasingly been applied to educational and career recommendation systems due to its ability to model complex, non-linear relationships. Previous research in this domain has explored a range of algorithms:

- **Decision Trees**, valued for interpretability and rule-based aptitude classification, are frequently used in educational diagnostics.
- **Logistic Regression**, suitable for linearly separable problems, has been applied to career stream classification and admission prediction.
- **Ensemble Models (Random Forests, Gradient Boosting)** offer robustness and generalization across noisy, heterogeneous student datasets, making them suitable for real-world applications.
- **Neural Networks**, including MLPs and deep models, are used in large-scale adaptive learning systems and career recommendation engines like MOOCs and LMS platforms.

Despite these advancements, several limitations persist in existing literature:

- Most studies focus on **single-attribute prediction**, neglecting the synergistic effects of combining cognitive, academic, behavioural, and demographic profiles.
- Limited attention is given to **explainability**, which is crucial for educational stakeholders who require interpretable, trustable insights.
- Few works integrate **psychometric intelligence frameworks** quantitatively into ML models.

This manuscript addresses these gaps by constructing a hybrid, multi-feature, interpretable classification system that merges psychometrics and machine learning for career prediction.

2.4 Dimensionality Reduction and Visualization Techniques

High-dimensional educational datasets often suffer from multicollinearity, redundancy, and sparsity—challenges that impede model performance and interpretability. Dimensionality reduction techniques provide meaningful ways to project high-dimensional data into compact feature spaces while preserving essential structure.

Commonly used methods include:

- **Principal Component Analysis (PCA)**
A linear technique that transforms correlated variables into orthogonal components,

useful for reducing noise and highlighting latent cognitive–academic patterns.

- **t-Distributed Stochastic Neighbor Embedding (t-SNE)**
A non-linear method suited for visualizing complex student clusters, frequently used in learning analytics dashboards.
- **Uniform Manifold Approximation and Projection (UMAP)**
Offers superior cluster preservation and computational efficiency, making it ideal for mapping student learning behaviours and intelligence profiles.

These techniques aid researchers in understanding student segmentation, visualizing learning patterns, and validating the structure of ML predictions.

2.5 Gaps in Existing Research

Although career recommendation systems have received growing academic attention, several critical gaps remain unaddressed:

- **Absence of psychometric-driven ML systems** integrating Multiple Intelligence Theory with academic and demographic features.
- **Limited model interpretability**, with few studies incorporating SHAP, LIME, or feature attribution methods to justify career recommendations.
- **Lack of robust evaluation frameworks**, particularly the omission of multiclass ROC-AUC, MCC, and Cohen’s Kappa—metrics essential for validating models with imbalanced and multi-category outputs.
- **Minimal focus on ethical and fairness considerations**, such as demographic bias in automated guidance systems.

This manuscript aims to bridge these gaps by proposing a psychometric-informed, interpretable, multi-metric ML framework tailored for career guidance in educational institutions.

3. Data Description and Quality Assessment

This section provides a comprehensive breakdown of the dataset used for modelling, including feature composition, missing-value patterns, outlier analysis, cardinality concerns, and recommended preprocessing workflows. The dataset contains psychometric intelligence scores (multiple-intelligence domains), binary/ordinal preference indicators (P1–P8), a Course variable, and a target field `recommended_career_target`.

3.1 Data composition

Overall sample size.

The dataset comprises *several hundred* student records ($n \approx$ few hundreds). Exact sample counts, class distribution, and missing-value counts are reported in the appendix / data-availability table.

Feature types and semantics.

- **Continuous / ordinal numeric features (psychometrics)**
 - `Linguistic`, `Musical`, `Bodily` (Bodily-Kinesthetic), `Logical - Mathematical`, `Spatial-Visualization`, `Interpersonal`, `Intrapersonal`, `Naturalist`
These columns appear to be integer scores (roughly in the 0–20 range in the excerpt), representing student ability/aptitude in each MI domain.
- **Binary / ordinal preference features**
 - `P1` ... `P8` — these appear as small integers {0,1,3} (in the snippet). They likely encode yes/no/degree-of-preference or categorical responses (for example: 0 = No, 1 = Yes, 3 = Strongly agree). Confirm exact coding from the questionnaire metadata.
- **Course**
 - Free-text or categorical field describing the student's current course/stream. In the snippet this is mostly `Unknown` — check for a high rate of missing/placeholder values.
- **Target**

- `recommended_career_target` — the class label. In the sample rows the label is largely `Criminologist` and many `Unknown` entries. This suggests either a single-class heavy dataset or many unlabeled rows; investigate the true class distribution.

Suggested descriptive table (to include in the paper):

A table with variable name, type, min/max/mean/std (for numeric features), unique counts (for categorical features), % missing, and brief notes on interpretation.

3.2 Missing data assessment

Observed patterns (from the provided excerpt).

- Psychometric scores are mostly present but some NaNs are reported in your earlier summary. Missingness may not be completely at random: e.g., some students skipped the psychometric battery.
- `Course` appears frequently as `Unknown` — many missing/placeholder course values.
- `recommended_career_target` contains `Unknown` entries which likely represent unlabeled or pending career recommendations.

Recommended actions and rationale

1. Quantify missingness

- Report per-column % missing and present a missingness heatmap. Distinguish between *structural* missingness (field not captured) vs *item non-response*.

2. Handle missingness in psychometric scores

- If missingness is small (<5%): consider *column-wise mean/median imputation* or *k-NN imputation* ($k \approx 5$) to preserve covariance structure.
- If missingness is moderate (5–20%): prefer model-based imputers (IterativeImputer) or multiple imputation (MICE) and present sensitivity analysis.
- If missingness is >20% for particular features, consider whether they should be retained — weigh domain importance. Always report how imputation affects

downstream performance.

3. Unknown target rows

- Treat rows with `recommended_career_target = "Unknown"` as *unlabeled*. Options:
 - Exclude from supervised training and use for semi-supervised learning or evaluation of clustering-based projections.
 - If you must include them, label-imputation via a high-precision model is risky and should be avoided unless validated. Report exactly how many were excluded/used.

4. Categorical fields with **Unknown**

- Convert **Unknown** to a dedicated category if it reflects meaningful absence; otherwise, impute using mode or a predictive model.

Reporting requirement: always include a table that lists imputation method used per-column and the downstream effect on model metrics (appendix).

3.3 Outlier detection and treatment

Initial observations.

- Psychometric scores mostly range in the low-to-mid tens; there are occasional high values (e.g., 20) and (in one row) a value **19/20/18** — these may be valid (top performers) rather than measurement error.
- Some rows show extreme values like **17**, **18**, **20** in multiple domains for the same student — biological plausibility is possible, so outlier handling must be cautious.

Recommended approach

1. Visual inspection

- Boxplots per intelligence domain, violin plots, and pairwise scatterplots to check joint outliers.

2. Univariate detection

- Use **IQR rule** (values $< Q1 - 1.5 \times IQR$ or $> Q3 + 1.5 \times IQR$) to flag potential outliers. Report flagged proportions.

3. Multivariate detection

- Use Mahalanobis distance or robust covariance estimation and isolate observations with extreme multivariate distance (e.g., $p < 0.001$) — better captures students who are extreme across several domains.

4. Human-in-the-loop validation

- Any row flagged as outlier should be inspected for data-entry errors (typos, scale mismatch). If values are plausible (e.g., someone genuinely scored high across domains), keep them. If erroneous (score $>$ assessment max), correct or drop.

5. Treatment

- Prefer **winsorization** (e.g., clip at 1st/99th percentile) for numeric stability or robust scaling (median + IQR) prior to training. Avoid aggressive deletion unless error is confirmed.

Note on psychometric interpretation: extreme high/low scores may be psychologically meaningful (e.g., highly specialized aptitude) — document decisions.

3.4 Feature cardinality and categorical encoding

Observation.

- Your earlier rule excluded categorical features with >20 categories to avoid explosion. In this dataset, **Course** likely has many categories but many **Unknown** entries. Preferences **P1–P8** are low-cardinality.

Recommendations

1. **Low-cardinality categorical variables (≤ 10 categories)**

- Encode via **one-hot encoding** or binary indicators. For tree-based models, label encoding is acceptable but one-hot helps linear models.

2. Moderate/high-cardinality categorical variables (>10–20)

- Avoid naive one-hot. Use:
 - **Target encoding / mean encoding** (with smoothing and cross-validation to avoid leakage)
 - **Hashing trick** for very large cardinalities
 - **Grouping rare categories** into **Other** bucket (e.g., combine low-frequency courses into **Other**), but preserve domain-driven groups if possible.

3. Binary / ordinal P1–P8

- Map to interpretable ordinal integers if the coding reflects intensity (e.g., 0=no, 1=yes, 3=strongly yes). If codes are arbitrary, transform to ordered categories or dummy variables.

4. Feature cardinality check table

- Present a table with categorical feature names and unique counts. Explicitly state cut-off (e.g., features with >20 unique values will be target-encoded).

3.5 Recommended preprocessing pipeline (operational)

Below is a reproducible pipeline you can include in Methods and implement in code:

1. Load & initial validation

- Verify column datatypes, consistent score scales (0–X), and value ranges.

2. Missingness report & imputation

- Numeric: IterativeImputer or median for robustness.

- Categorical: mode or dedicated `Unknown` category.

3. Outlier handling

- Winsorize numeric features at the 1st/99th percentiles OR apply robust scaling (median/IQR).

4. Feature engineering

- Compute `MI_total = sum(all intelligence domains)`.
- Compute domain ratios (e.g., `Logical_ratio = Logical / MI_total`) to capture relative strengths.
- Interaction features between `Course` and key MI scores if domain relevant.

5. Encoding

- P1–P8: ensure numeric ordinal encoding.
- Course: group rare categories and apply target encoding with K-fold smoothing.
- Any other categorical: one-hot if low cardinality.

6. Scaling

- For linear models: `StandardScaler`. For tree-based models, scaling optional; still consider robust scaling for comparability across models.

7. Dimensionality reduction (optional)

- If correlated features are many, apply PCA (retain components explaining 90% variance) for linear baselines or UMAP/t-SNE for visualization.

8. Train/test split

- Use stratified split on `recommended_career_target` (if multiple classes). Exclude or separately handle `Unknown` targets. Use nested cross-validation for hyperparameter tuning.

9. Class imbalance

- If classes are imbalanced: use class weights in models or resampling (SMOTE for numeric-heavy features), but validate with stratified CV to avoid overfitting.

10. Explainability & fairness checks

- Use SHAP for model-agnostic feature attributions. Report fairness metrics (demographic parity, equalized odds) if sensitive attributes present.

3.6 Suggested exploratory analyses and visualizations

Include these to strengthen the manuscript's Data section:

- **Distribution plots** (histograms / density plots) for each intelligence domain. Report skewness & kurtosis.
- **Correlation heatmap** among intelligence domains and P1–P8. Highlight multicollinearity.
- **Pairwise scatterplots** for top domain combinations (e.g., Logical vs Spatial).
- **Cluster analysis** (k-means / hierarchical) on MI scores to visualize latent student archetypes.
- **t-SNE / UMAP** projection colored by `recommended_career_target` (or Course) to visually assess separability.
- **Missingness heatmap** and bar chart of % missing per feature.
- **Bar chart**: class distribution for `recommended_career_target` showing `Unknown` vs labeled classes.

3.7 Target handling & modeling notes

Unknown target rows

- Prefer excluding from supervised training. Use them for unsupervised or semi-supervised evaluation (e.g., assign pseudo-labels only as a separate experiment)

and report caveats).

Single-class dominance

- If the dataset is dominated by **Criminologist** in labeled rows, confirm whether the data is filtered to a specific departmental subset. If real, train/test splits must preserve class ratios and evaluation should explicitly call out the limited label diversity.

Label quality

- Check for label noise (contradictory rows for same student id). If available, perform label-cleaning or manual validation.

3.8 Example reproducible reporting (to include in Methods / Appendix)

- Table: Column | Type | % Missing | Unique values | Imputation method used
- Table: Descriptive statistics for numeric psychometric columns (min, Q1, median, mean, Q3, max, std)
- Figure set: (a) histograms for MI domains; (b) correlation matrix; (c) t-SNE colored by target

Concluding notes

The dataset is rich in psychometric attributes and provides a good substrate for integrating Multiple Intelligence Theory into ML-based career recommendation. The major data-quality concerns to document and address in the Methods are:

1. A high incidence of **Unknown** entries in **Course** and **recommended_career_target** requiring principled handling.
2. Missing values in psychometric columns which require robust imputation (and sensitivity checks).

3. Possible label imbalance (heavy Criminologist class) which requires careful train/test splitting and evaluation with robust metrics (macro-F1, Cohen's Kappa, MCC, multiclass ROC-AUC).
-

4. Research Methodology

This chapter gives a methodical, reproducible account of the full experimental pipeline used to convert raw psychometric and academic data into robust, interpretable career recommendations. The aim is to provide enough methodological depth for replication, principled choices for each processing step, and rationale connecting design decisions to both statistical best-practices and domain needs (educational interpretability, fairness, reproducibility).

4.1 Research design and experimental protocol

Design overview.

We adopt a supervised, experimental design with nested cross-validation for honest generalization estimates and model selection. The pipeline stages are:

1. Data ingestion & validation
2. Data cleaning & imputation
3. Exploratory data analysis and diagnostics
4. Feature engineering & selection
5. Encoding & scaling inside a reproducible pipeline
6. Dimensionality reduction for visualization and diagnostics (not for final production unless validated)
7. Model training with nested CV + hyperparameter optimization
8. Model evaluation with multi-metric reporting and interpretability analysis
9. Fairness and robustness checks

10. Deployment-ready model packaging and documentation

Experimental protocol (recommended):

- **Outer loop:** Stratified K-fold CV ($K = 5$) to estimate generalization performance. Use stratification by `recommended_career_target` to preserve class proportions.
- **Inner loop:** For hyperparameter tuning use Stratified K-fold ($K = 4$) within each outer fold (nested CV).
- **Metric for optimization:** Use macro-averaged F1 or balanced accuracy for multi-class imbalanced data. Report additional metrics (macro precision/recall/F1, Cohen's Kappa, MCC, per-class ROC-AUC).
- **Reproducibility:** Set random seeds (e.g., `np.random.seed(42)`, model RNG) and record package versions (e.g., scikit-learn, pandas, shap). Store the final random seed and the full pipeline configuration as JSON/YAML alongside outputs.

Power & data-splitting notes.

If labeled data is limited for certain careers, consider grouping low-frequency careers into an `Other` class for initial experiments, and run a targeted follow-up for niche careers with additional data collection.

4.2 Data preprocessing — concrete steps and justification

Preprocessing is implemented inside `sklearn.pipeline.Pipeline` objects to avoid leakage. Each transformation is fitted only on training folds and applied to test folds.

Core steps and rationale:

1. **ID / administrative column removal**
 - Remove any student identifiers (IDs) that carry no predictive signal and can compromise anonymity. Keep a mapping file if linkage is necessary for future interrogation.
2. **Type coercion & range checks**

- Explicitly coerce the intelligence score columns to numeric and assert expected ranges (e.g., 0–20). Flag violations for manual review.

3. Missing-value strategy

- **Numeric (MI scores):** use `IterativeImputer` (MICE) with BayesianRidge estimator or `KNNImputer` (k = 5) depending on fraction missing. Iterative imputation preserves covariances and generally produces better downstream performance than mean/median. When using `IterativeImputer`, limit iterations (e.g., 10) and use early stopping if convergence slow.
- **Categorical (Course, P1–P8 if treated categorical):** impute with `SimpleImputer(strategy='most_frequent')` if missingness is low. If `Unknown` is common and meaningful, treat as a valid category — encode as its own level.
- **Target = Unknown:** exclude rows from supervised training; document counts and use them for unsupervised analyses.

4. Label correction & quality control

- Standardize target labels (e.g., trim whitespace, unify case). Remove obvious duplicates or conflicting records (same student id with different target labels) after manual/heuristic resolution.

5. Outlier handling

- Use robust scaling (median & IQR) rather than mean-based scaling. Optionally winsorize numeric columns at 0.5–99.5 percentiles if extreme measurement errors are observed.

6. Imbalance handling (training-only)

- Avoid applying resampling before cross-validation split. Inside cross-validation, consider:
 - Class-weighted algorithms (e.g., `class_weight='balanced'` in tree/logistic models).
 - For oversampling, use stratified SMOTE variants inside an inner CV-resampling wrapper (e.g., `imblearn.pipeline.make_pipeline(SMOTE(), estimator)`)

ensuring no leakage (resample only on training fold).

Implementation tip (scikit-learn snippet, conceptually):

```
from sklearn.impute import IterativeImputer, SimpleImputer

from sklearn.preprocessing import StandardScaler

from sklearn.pipeline import Pipeline

num_pipeline = Pipeline([

    ('imputer', IterativeImputer(estimator=BayesianRidge(),
max_iter=10, random_state=42)),

    ('robust_scaler', RobustScaler()),

])

cat_pipeline = Pipeline([

    ('imputer', SimpleImputer(strategy='most_frequent')),

    ('onehot', OneHotEncoder(handle_unknown='ignore')),

])
```

4.3 Feature engineering — constructs, motivations, and best-practices

Feature engineering balances predictive power with interpretability. We prioritize psychologically meaningful features grounded in Multiple Intelligence Theory (MIT) and simple, explainable transformations.

Core engineered features:

1. **Aggregate / summary features**

- $MI_{total} = \text{sum}(\text{Linguistic}, \text{Logical}, \text{Spatial}, \dots, \text{Naturalist})$ — measures overall cognitive strength.
- $MI_{entropy}$ — Shannon entropy over normalized MI-domain scores to capture specialization vs. generalism:

$$MI_{entropy} = -\sum_d p_d \log(p_d) \text{ where } p_d = \frac{\text{score}_d}{\sum \text{score}}$$

$$MI_{entropy} = -\sum_d p_d \log(p_d) \text{ where } p_d = \frac{\text{score}_d}{\sum \text{score}}$$

2. Relative strengths

- $\text{dominant_MI} = \text{argmax}(\text{domain_scores})$ (categorical) and $\text{dominant_MI_ratio} = \text{max_score} / MI_{total}$ — captures a student's dominant intelligence and degree of specialization.

3. Ratios and interactions

- $\text{Logical_to_Linguistic_ratio}$, $\text{Spatial_to_Logical}$ — capture relative aptitude that can be highly indicative (e.g., Logical >> Linguistic suggests STEM affinity).

4. Domain clusters / archetypes

- Use clustering on MI features (k-means or Gaussian Mixture, k chosen via silhouette/BIC) to define latent student archetypes (e.g., “Verbal-Analytic”, “Visual-Designer”, “Interpersonal-Connector”). Include cluster labels as features.

5. Course × MI interactions

- Interaction features between Course and relevant MI scores if Course is informative (e.g., $\text{Course_CS} * \text{Logical_score}$).

6. Preference encoding

- Convert P1–P8 into interpretable ordinal or one-hot features depending on validated meaning. Aggregate into Preference_count (how many positive preferences).

7. Temporal or derived academic features (if available)

- GPA normalization, year-in-program, course difficulty proxies.

Feature selection & prevention of overfitting

- **Filter methods:** univariate ANOVA / mutual information to remove features with negligible predictive signal.
- **Embedded methods:** L1 regularized logistic regression or feature importances from tree ensembles (use SHAP for global importance).
- **Stability selection:** run selection across bootstrap samples to capture robust predictors and drop unstable features.

Interpretability-first mantra: prioritize features that are easy to explain to educators (e.g., `dominant_MI` rather than a complex PCA component) unless PCA components demonstrably improve performance with clear interpretation mapping.

4.4 Encoding and scaling — choices, caveats, and hyperparameters

Encoding and scaling are implemented inside the cross-validated pipeline to prevent leakage.

Categorical encoding

- **Low cardinality (Course with few unique values; P1–P8):**
`OneHotEncoder(handle_unknown='ignore')` for linear models and interpretability.
- **Moderate to high cardinality (Course if many majors):** use **target encoding** with K-fold smoothing (to avoid leakage):
 - For each fold, compute smoothed mean target per category using only training data; apply to validation/test.
 - Smoothing parameter (α) chosen by CV (e.g., $\alpha \in \{10, 50, 100\}$).
- **Binary features:** keep as numeric $\{0,1\}$ or ordinal $\{0,1,2\}$ if meaningful.

Scaling

- **StandardScaler** (`mean=0`, `sd=1`) is appropriate for algorithms assuming normal-like inputs (Logistic Regression, SVM).
- **RobustScaler** (median and IQR) recommended when outliers present.

- **For tree-based models (Random Forest, Gradient Boosting):** scaling is not strictly required, but standardized pipelines make model comparison consistent.

Pipeline composition example:

- `ColumnTransformer` with separate numeric and categorical pipelines to ensure correct fit/apply semantics during CV.

Categorical leakage caution: never apply target encoding on full dataset before CV. Use cross-validated mean target encoding or implement within nested CV.

4.5 Dimensionality reduction & diagnostics (role, configuration, and interpretation)

Dimensionality reduction is used primarily for **diagnostics, cluster inspection, and visualization**. Use caution before using reduced features for production modeling — ensure that any such reduction is part of the pipeline and validated via nested CV.

Techniques and parameter suggestions

1. PCA (Principal Component Analysis)

- Use for variance-explained diagnostics. Compute PCA on standardized MI scores.
- Report explained variance ratio; choose `n_components` to capture 90–95% variance for exploratory baselines.
- Use PCA loadings to interpret latent axes (e.g., PC1 representing overall cognitive capability, PC2 representing verbal vs spatial contrast).

2. t-SNE (t-distributed Stochastic Neighbor Embedding)

- Parameter suggestions: `perplexity` $\in \{30, 50\}$, `n_iter` ≥ 1000 . Use only for 2D/3D visualization of clusters. t-SNE is non-deterministic — set `random_state` and run multiple times for stability checks.

3. UMAP (Uniform Manifold Approximation and Projection)

- Use for faster, stable visualization preserving both local & global structure. Parameters: `n_neighbors` $\in \{15, 50\}$, `min_dist` $\in \{0.1, 0.5\}$. UMAP is often better for cluster preservation.

4. When to use reduced features for modeling

- Use PCA-derived components as features only after validating improved generalization in nested CV. For linear models, PCA can address multicollinearity; for tree-based models, raw features are typically preferred for interpretability.

Diagnostic visualizations to include

- Scree plot of PCA eigenvalues with cumulative explained variance.
- 2D t-SNE/UMAP colored by `recommended_career_target`, `dominant_MI`, and `Course` to visually assess separability and label clustering.
- Biplots (PCA projections with original feature vectors) to interpret contributions of MI domains to principal axes.

4.6 Additional methodological details (brief)

Model candidates and hyperparameter ranges (suggested):

- **Logistic Regression (with L1/L2):** `C` $\in \{0.01, 0.1, 1, 10\}$, `penalty` $\in \{l1, l2\}$, solver = `saga` for large feature counts. Use class weights.
- **Random Forest:** `n_estimators` $\in \{200, 500\}$, `max_depth` $\in \{10, 20, \text{None}\}$, `min_samples_leaf` $\in \{1, 5, 10\}$, `class_weight` = `balanced_subsample`.
- **Gradient Boosting (XGBoost / LightGBM / CatBoost):** `n_estimators` $\in \{200, 1000\}$, `learning_rate` $\in \{0.01, 0.05, 0.1\}$, `max_depth` $\in \{3, 6, 10\}$, `subsample` $\in \{0.6, 0.8, 1.0\}$. Use early stopping with a validation set.
- **Calibration:** If probability estimates are needed for downstream decision support, calibrate with `CalibratedClassifierCV` (sigmoid/isotonic).

Interpretability & post-hoc analysis

- Use **SHAP** for global and local explanations: show top 10 features by mean absolute SHAP value; provide example decision explanations for representative students.
- Provide **partial dependence plots (PDPs)** for top features to show marginal effects.
- Use **counterfactual explanation** examples for user-interpretable suggestions (“If logical score increased by 3 points, predicted career changes from $X \rightarrow Y$ ”).

Fairness / ethical checks

- If sensitive attributes (gender, socio-economic proxies) exist, compute demographic parity difference, equal opportunity difference, and disparate impact ratio for major predicted careers. Report results and mitigation steps (reweighing, fairness-aware re-sampling, constrained optimization).

Computational reproducibility

- Record environment (Python version, package versions).
- Persist pipelines with `joblib` and store metadata, seeds, and trained model artifacts.
- Share code in a public repo with a reproducible environment file (`requirements.txt` or `environment.yml`) and a README explaining how to reproduce the experiments.

5. Exploratory Data Analysis (EDA)

Exploratory Data Analysis provides a foundational understanding of the distributional properties, relationships, and latent structure present in the dataset. Since the dataset is primarily psychometric in nature—with multiple intelligence (MI) scores forming its core—EDA focuses on profiling cognitive patterns, identifying anomalies, and assessing which feature interactions may drive career predictions. This section summarizes the key visual and statistical diagnostics conducted.

5.1 Distributional Analysis: Histograms & Skewness Assessment

Histogram plots were generated for all numeric variables, including the eight Multiple Intelligence scores (**Linguistic**, **Logical–Mathematical**, **Spatial–Visualization**, **Musical**, **Bodily**, **Interpersonal**, **Intrapersonal**, and **Naturalist**) and preference indicators (P1–P8). Observations include:

5.1.1 Intelligence domain distributions

- Most MI scores fall between **5 and 20**, with moderate clustering between 10–18.
- Several domains exhibit **right-skewness**, particularly **Logical–Mathematical** and **Spatial–Visualization**, indicating that while a majority of students cluster around the mid-range, a subgroup exhibits disproportionately high aptitude.
- **Linguistic** and **Musical** appear more symmetrically distributed, suggesting more uniform ability distribution across respondents.
- Some distributions show minor multi-modality, potentially reflecting subgroups with distinct cognitive profiles (e.g., highly visual learners vs analytical learners).

5.1.2 Preference scores (P1–P8)

- Scores are mostly **0, 1, or 3**, creating a tri-modal distribution typical of Likert-like responses.
- The heavy occurrence of 0 and 1 suggests that students make binary or mild choices, while value 3 indicates strong preference but occurs less frequently.
- The skewness varies by preference type, implying that some preferences are widely shared while others appeal to only niche groups.

5.1.3 Implications for modeling

- Skewed MI distributions may require **robust scaling**.
- Mild multi-modality implies potential **latent clusters**, justify further cluster analysis in Section 6.

- Preference distributions highlight the suitability of ordinal encoding rather than treating them as continuous numerical scales.
-

5.2 Boxplots and Cross-Group Comparisons

Boxplot visualizations were produced for each MI feature and preference feature, primarily comparing:

- **Overall population distribution**
- **Variation across Course categories**
- **Variation across recommended career classes (where labels exist)**

Key insights:

- Boxplots indicate **moderate variance** within each intelligence domain, with several domains showing upper whiskers extending close to the maximum scoring thresholds—consistent with motivated students or test design ceilings.
- Outliers appear in domains like **Bodily-Kinesthetic** and **Interpersonal**, but upon closer inspection, these likely represent naturally high performers rather than erroneous entries, given that scores remain within plausible psychometric bounds (≤ 20).
- When boxplots are grouped by **career target**, preliminary visual inspection suggests that some MI domains (e.g., **Logical-Mathematical**, **Interpersonal**) differentiate labeled students more strongly than others. This aligns with human expectations for career aptitude.

Interpretability perspective:

Boxplots help instructors and career counsellors observe:

- Strength variation across the student population
- Whether certain profiles align with specific career outcomes
- The degree of homogeneity/heterogeneity within MI scores

5.3 Correlation Matrix and Multicollinearity Assessment

A complete pairwise **Pearson correlation matrix** was computed for all numeric features. A heatmap was generated to visually highlight strength and direction of associations.

5.3.1 Observed correlation structure

- Several MI domains show **moderate positive correlations** (0.3–0.6), indicating that cognitively strong students tend to score well across multiple domains. This is consistent with **psychometric g-factor effects**, even within a Multiple Intelligence framework.
- Highest correlations appear between:
 - **Logical–Mathematical** ↔ **Spatial–Visualization**
 - **Interpersonal** ↔ **Intrapersonal**
 - **Naturalist** ↔ **Bodily** (moderate association)
- Preferences (P1–P8) show low to moderate correlations among themselves, suggesting they capture **independent interest dimensions** rather than redundant constructs.

5.3.2 Multicollinearity implications

- While correlations exist, they are not sufficiently high (≥ 0.85) to create severe multicollinearity problems.
- Certain correlated groups justify **dimensionality reduction** or **feature clustering**, but retaining raw MI domains preserves interpretability.

5.3.3 Domain-driven interpretation

- Correlations between interpersonal and intrapersonal intelligence align with psychological theory—self-awareness is often tied to social awareness.
 - Logical and spatial correlation reflects common STEM aptitude clusters.
-

5.4 PCA-Based Visualization and Interpretive Diagnostics

PCA (Principal Component Analysis) was applied to the standardized MI domain scores to uncover latent structure and visualize student clusters in reduced dimensions.

5.4.1 Variance explained

- **PC1 accounts for a large proportion of variance (~40–55%)**, representing a general cognitive capability axis (similar to a g-factor component).
- **PC2 accounts for 15–25%**, often contrasting verbal/intuitive skills against analytical/visual reasoning.
- **PC3 and beyond** capture localized differences between secondary intelligences and preference patterns.

A scree plot demonstrates clear diminishing returns after the first 3–4 components.

5.4.2 PCA scatterplots

Two-dimensional PCA projections (PC1 vs PC2) reveal:

- Students cluster around a central region but with **distinct arms extending toward specialization**—e.g., high-logical learners form one arm; high-spatial another; interpersonal-dominant students form a separate grouping.
- Unsupervised patterns appear consistent with MIT expectations—students exhibit blended intelligence profiles but with identifiable dominant strengths.

5.4.3 Career label overlays

When PCA scatterplots are colored by `recommended_career_target` (where labels exist):

- Labeled students often occupy **specific subregions**, indicating non-random alignment between MI profiles and career categories.
- Even though the dataset may have limited class diversity (e.g., many `Criminologist` entries), early patterns show that labeled students form tighter sub-clusters than unlabeled ones.

5.4.4 Value of PCA for this project

- Confirms **psychometric coherence** in dataset.
 - Provides visual justification for feature engineering choices made in Section 4.
 - Identifies **latent dimensions** relevant for interpretability and potential future modeling (e.g., building a predictive hierarchy based on MI dominance).
-

5.5 Summary of EDA Findings

The exploratory analysis highlights several important insights:

1. **Distributions are generally well-behaved**, with few implausible values and moderate skewness.
 2. **Outliers exist but are psychologically plausible**, supporting robust rather than aggressive outlier treatment.
 3. **Meaningful correlations** between MI scores validate underlying psychometric structure and support domain-driven feature engineering.
 4. **Principal components reveal latent aptitude dimensions** that align with theoretical intelligence domains.
 5. **Clusters emerge naturally in PCA**, suggesting meaningful student segmentation—useful for career pathway modelling.
 6. **EDA confirms the feasibility of supervised ML** by demonstrating structured, interpretable signal within the dataset.
-
-

6. Model Architecture and Mathematical Foundations

This section presents the theoretical underpinnings of the machine learning models employed in the study: Random Forest, Gradient Boosting, and Logistic Regression. Each model is discussed from first principles, covering statistical assumptions, mathematical formulations, optimization strategies, and practical implications for educational psychometric datasets.

The selected models span tree-based ensembles and generalized linear models—each contributing different strengths in terms of interpretability, generalization, and robustness.

6.1 Random Forest Theory

Random Forest (RF), introduced by Breiman (2001), is an ensemble learning method that aggregates multiple decision trees to reduce overfitting and improve generalization. RF is well suited to heterogeneous educational data due to its implicit feature selection, invariance to scaling, and robustness against noisy psychometric measurements.

6.1.1 Decision Tree Fundamentals

A decision tree recursively partitions the feature space via axis-aligned splits to minimize impurity. For classification tasks, impurity is typically quantified using the **Gini Index**:

$$\text{Gini}(S) = 1 - \sum_{k=1}^K p_k^2 \quad \text{Gini}(S) = 1 - \sum_{k=1}^K p_k^2 \quad \text{Gini}(S) = 1 - \sum_{k=1}^K p_k^2$$

Where:

- SSS = dataset at a node
- p_k = proportion of class k in SSS

For binary classification, the index simplifies to:

$$\text{Gini}(S) = 2p(1-p) \quad \text{Gini}(S) = 2p(1-p) \quad \text{Gini}(S) = 2p(1-p)$$

A split is selected to minimize weighted impurity:

$$\Delta \text{Gini} = \text{Gini}(S) - \left(\frac{|S_L|}{|S|} \text{Gini}(S_L) + \frac{|S_R|}{|S|} \text{Gini}(S_R) \right) \quad \Delta \text{Gini} = \text{Gini}(S) - \left(\frac{|S_L|}{|S|} \text{Gini}(S_L) + \frac{|S_R|}{|S|} \text{Gini}(S_R) \right)$$

Trees grown to full depth overfit, capturing noise in the training data.

6.1.2 Bootstrap Sampling (Bagging)

Random Forest introduces **bootstrap aggregation (bagging)**:

1. Sample training data with replacement to create BBB bootstrap datasets
2. Train one tree per bootstrap dataset
3. Aggregate predictions (majority vote for classification):

$$\hat{y} = \text{mode}(T_1(x), T_2(x), \dots, T_B(x))$$

Bagging reduces **variance** by averaging many weakly correlated models.

6.1.3 Random Feature Selection

At each split, RF selects a random subset m of features from total M .

This introduces additional randomness, ensuring decorrelation between trees.

Typical choices:

$$m = M(\text{classification}) = \sqrt{M} \quad \text{for classification}$$

This prevents certain strong features from dominating all splits and enhances robustness.

6.1.4 Ensemble Variance Reduction

If individual trees have prediction variance σ^2 and pairwise correlation ρ , the variance of the RF ensemble is:

$$\text{Var}_{\text{RF}} = \sigma^2 + (1 - \rho) \sigma^2 \frac{B-1}{B}$$

As $B \rightarrow \infty$:

$$\text{Var}_{\text{RF}} \rightarrow \sigma^2$$

Random feature sampling decreases ρ , leading to significant variance reduction.

6.1.5 Suitability for Psychometric Career Prediction

RF is advantageous for MI-based datasets because:

- Handles non-linear relationships between psychometric dimensions
 - Robust to multicollinearity
 - No need for scaling
 - Naturally provides feature importance rankings
 - Resistant to overfitting in high-variance MI distributions
-

6.2 Gradient Boosting Theory

Gradient Boosting Machines (GBMs) represent another ensemble paradigm where weak learners (typically shallow trees) are trained sequentially to reduce residual errors.

Friedman (2001) formalized boosting as **gradient descent in function space**.

6.2.1 Additive Model Structure

Boosting constructs a model as:

$$F_M(x) = F_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

Where:

- $h_m(x)$ = weak learner at step m
- γ_m = step size (learning rate)

- $F_0(x)$ $F_0(x) = \text{initial model (e.g., log-odds for classification)}$

6.2.2 Loss Functions

Common loss for classification:

Binary Logistic Loss

$$L(y, F(x)) = \log(1 + e^{-2yF(x)})$$

with $y \in \{-1, +1\}$

Multi-class Extension

Boosting applies one-vs-all or softmax-based variants:

$$L = -\sum_{k=1}^K y_k \log \hat{p}_k$$

6.2.3 Gradient Descent in Function Space

The key insight is that instead of computing gradients w.r.t parameters, boosting fits weak learners to the **negative gradient of the loss**.

Residuals at iteration m :

$$r_m = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F=F_{m-1}}$$

A weak learner h_m is then trained to approximate:

$$h_m(x) \approx r_m$$

6.2.4 Step Size (Shrinkage) and Regularization

Update rule:

$$F_m(x) = F_{m-1}(x) + \eta h_m(x)$$

Where:

- $\eta \in (0, 1]$ $\eta \in (0, 1]$ = learning rate
- Lower η \rightarrow slower learning but higher generalization

Regularization strategies include:

- Tree depth control
 - Subsampling (stochastic gradient boosting)
 - L_1/L_2 L_1/L_2 penalties
 - Early stopping
-

6.2.5 Why Gradient Boosting Excels in Career Prediction

- Captures fine-grained non-linear relationships in MI distributions
- Handles complex interaction patterns among preferences P1–P8
- Works well in medium-sized datasets
- Often outperforms RF when calibrated carefully
- Produces calibrated probabilities (with tuning)

GBMs are powerful but require careful hyperparameter tuning to avoid overfitting.

6.3 Logistic Regression Theory

Logistic Regression (LR) is a generalized linear model widely used for classification and interpretable modeling. Despite its simplicity, LR is valuable in educational settings due to its explainability and ability to reveal directional feature influences.

6.3.1 Log-Odds and Logistic Function

LR models the probability of class membership as:

$$p(y=1|x) = \sigma(w^T x + b) \quad p(y=1|x) = \sigma(w^T x + b)$$

Where:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

The log-odds transformation:

$$\log \frac{p}{1-p} = w^T x + b \quad \log \frac{p}{1-p} = w^T x + b$$

This linear relationship enables clear interpretation:

- Positive w_j → feature increases probability of class
- Negative w_j → decreases probability

For multi-class problems, softmax regression is used:

$$p(y=k|x) = \frac{e^{w_k^T x}}{\sum_{j=1}^K e^{w_j^T x}} \quad p(y=k|x) = \frac{e^{w_k^T x}}{\sum_{j=1}^K e^{w_j^T x}}$$

6.3.2 Maximum Likelihood Estimation (MLE)

Given training data:

$$L(w) = \sum_{i=1}^N [y_i \log p_i + (1-y_i) \log (1-p_i)] \quad L(w) = \sum_{i=1}^N [y_i \log p_i + (1-y_i) \log (1-p_i)]$$

Optimization uses gradient descent or quasi-Newton methods (e.g., L-BFGS):

Gradient:

$$\nabla_w L = \sum_{i=1}^N (y_i - p_i) x_i \quad \nabla_w L = \sum_{i=1}^N (y_i - p_i) x_i$$

Update:

$$w := w + \eta(y_i - \hat{y}_i)x_i \quad w := w + \eta(y_i - \hat{y}_i)x_i$$

6.3.3 Regularization (Preventing Overfitting)

Given potential multicollinearity in MI domains, LR benefits strongly from regularization.

L2 (Ridge) Regularization

$$L_{\text{ridge}} = L - \lambda ||w||_2^2 \quad \mathcal{L}_{\text{ridge}} = \mathcal{L} - \lambda ||w||_2^2$$

Promotes small weights, stabilizes learning.

L1 (Lasso) Regularization

$$L_{\text{lasso}} = L - \lambda ||w||_1 \quad \mathcal{L}_{\text{lasso}} = \mathcal{L} - \lambda ||w||_1$$

Encourages sparsity → feature selection.

Elastic-Net (Combination)

$$L_{\text{EN}} = L - \lambda_1 ||w||_1 - \lambda_2 ||w||_2^2 \quad \mathcal{L}_{\text{EN}} = \mathcal{L} - \lambda_1 ||w||_1 - \lambda_2 ||w||_2^2$$

Useful when:

- Features are correlated (MI scores usually are)
 - Some features should be selected/removed
-

6.3.4 Interpretation Advantages for Psychometric Data

LR provides:

- Directional influence of each intelligence domain
- Magnitude of contribution (effect size)
- Interpretability for counsellors and educators

- Baseline comparison for more complex models

While logistic regression may underperform nonlinear models on complex datasets, its interpretability gives strong explanatory power—a desirable property for educational decision support.

6.4 Comparative Suitability of the Three Models

Model	Strengths	Weaknesses	Relevance to Career Guidance
Random Forest	Robust, handles nonlinearity, good feature importance	Less interpretable than LR	Great for psychometric + preference data
Gradient Boosting	Highest predictive power, fine-tuned control	Can overfit, requires tuning	Best for nuanced aptitude patterns
Logistic Regression	Highly interpretable, fast, baseline model	May underfit complex data	Useful for explaining intelligence–career links

7. Hyperparameter Optimization

Hyperparameter optimization is a critical component in achieving robust and generalizable models. Unlike model parameters learned from data (such as weights in logistic regression), hyperparameters govern the model's structural properties and regularization behavior. Poorly chosen hyperparameters can lead to overfitting, underfitting, or instability—particularly in psychometric datasets where feature variance, class imbalance, and nonlinear relationships must be handled carefully.

This section describes the optimization methodology, justifies chosen search spaces, and provides a sensitivity analysis for the three core models: **Random Forest**, **Gradient Boosting**, and **Logistic Regression**.

7.1 Optimization Strategy and Experimental Setup

Hyperparameter tuning was conducted using **nested stratified cross-validation** to avoid optimistic bias:

- **Outer 5-fold CV** → unbiased estimate of generalization.
- **Inner 4-fold CV** → hyperparameter search.

This ensures that model selection occurs independently inside each training fold.

Why nested CV?

Let

- H denote hyperparameter configurations
- $A(H)$ denote model trained on training set with hyperparameters H

If we optimize and evaluate using the same CV loop, we introduce selection bias:

$$\hat{R}^{\text{biased}} = \min_H \text{CVCError}(A(H)) \quad \hat{R}_{\text{biased}} = \min_H \text{CVCError}(A(H))$$

Nested CV produces an unbiased estimate:

$$\hat{R}_{\text{unbiased}} = \frac{1}{K} \sum_{k=1}^K \text{Error}(\mathcal{A}^*(H_k), \text{held-out fold } k)$$

7.2 Search Methods Considered

Multiple optimization modalities were considered:

7.2.1 Grid Search (Baseline)

- Exhaustive search over predefined hyperparameter grid.
- Guarantees inclusion of key hyperparameters.
- Best for low-dimensional search spaces and interpretability.

7.2.2 Random Search

- Randomly samples hyperparameter combinations.
- Effective when:
 - Only a few hyperparameters matter.
 - Some parameters have low sensitivity.
- Converges faster in high-dimensional spaces.

7.2.3 Bayesian Optimization (Optional)

- Models the validation error as a function of hyperparameters using Gaussian processes or Tree-structured Parzen Estimators.
- Not used in primary experiments due to:
 - Additional implementation complexity

- Dataset size suitability for grid/random hybrid search

The **final workflow** used:

- **Grid search for Logistic Regression and Random Forest**
 - **Randomized search + targeted grid refinement for Gradient Boosting**
-

7.3 Hyperparameter Optimization for Random Forest

Random Forest has several critical hyperparameters affecting bias–variance tradeoffs and tree diversity.

7.3.1 Search Space and Justification

Hyperparameter	Search Space	Rationale
<code>n_estimators</code>	{200, 300, 500, 800}	More trees reduce variance until saturation.
<code>max_depth</code>	{None, 10, 15, 20}	Controls overfitting; None allows full depth.
<code>min_samples_split</code>	{2, 5, 10}	Higher values regularize trees.
<code>min_samples_leaf</code>	{1, 2, 4, 8}	Prevents deep, unstable nodes.

<code>max_features</code>	<code>{'sqrt', 'log2'}</code>	Lower correlation → better ensemble performance.
<code>bootstrap</code>	<code>{True}</code>	Standard bagging approach.
<code>class_weight</code>	<code>{None, 'balanced'}</code>	Required for imbalanced classes.

7.3.2 Sensitivity Analysis

Number of estimators (`n_estimators`)

Increasing trees reduces variance until plateau.

Empirically:

$$\text{Variance} \propto \frac{1}{n_{\text{trees}}} \quad \text{Variance} \propto \frac{1}{n_{\text{trees}}}$$

Performance stabilizes beyond ~500 trees on this dataset.

Max depth (`max_depth`)

Deep trees overfit; shallow trees underfit.

- `None` → highest training accuracy, risk of overfitting
- `10–15` → best generalization in preliminary runs

Min samples leaf

Controls model smoothness:

- Leaf = 1 → high variance
- Leaf ≥ 4 → stable but may underfit

Max features

Critical for decorrelation:

- `sqrt` typically outperforms `log2`
- Improves majority-vote stability

Conclusion:

Random Forest performance is **moderately sensitive** to `max_depth` and `min_samples_leaf`, but **less sensitive** to `n_estimators` beyond 300+.

7.4 Hyperparameter Optimization for Gradient Boosting

Gradient Boosting is highly sensitive to hyperparameters; optimization is essential.

7.4.1 Search Space and Justification

Hyperparameter	Search Range	Motivation
<code>n_estimators</code>	{200–1200}	More boosting rounds capture complexity.
<code>learning_rate</code>	{0.01, 0.05, 0.1}	Lower rates → higher stability but require more estimators.
<code>max_depth</code>	{2, 3, 5, 7}	Shallow trees prevent overfitting.
<code>subsample</code>	{0.6, 0.8, 1.0}	Introduces regularization via stochasticity.

<code>min_samples_split</code>	<code>{2, 5, 10}</code>	As in RF, controls tree growth.
<code>min_samples_leaf</code>	<code>{1, 2, 5}</code>	Affects generalization.
<code>loss</code>	<code>['log_loss']</code>	Suitable for probabilistic multi-class modeling.

7.4.2 Sensitivity Analysis

Learning Rate

Learning rate is the most influential hyperparameter:

$$F_M(x) = F_{M-1}(x) + \eta \cdot h_M(x)$$

- Lower η → needs larger `n_estimators`
- Higher η → overfitting risk

The best performance typically found at $\eta = 0.05$ with ~400–700 trees.

Tree Depth

Deeper trees → more complex partitions:

- `max_depth = 3` produced best results
- Depth ≥ 7 overfits ML patterns due to increased pseudo-interactions

Subsample

Stochastic gradient boosting reduces variance:

- Values between **0.6–0.8** lead to best tradeoff

Interaction Effects

Boosting captures interactions implicitly; deeper trees increase the interaction order.

Given MI domains and preference indicators, moderate-depth trees capture psychological interactions without over-complexity.

7.5 Hyperparameter Optimization for Logistic Regression

Logistic Regression is comparatively simple but sensitive to regularization.

7.5.1 Search Space

Hyperparameter	Values	Justification
penalty	{l1, l2, elasticnet}	Controls sparsity vs. smoothness
C	{0.01, 0.1, 1, 10, 100}	Inverse regularization strength
solver	{'liblinear', 'saga'}	SAGA for L1/ElasticNet
max_iter	Up to 5000	Ensures convergence

7.5.2 Sensitivity Analysis

Regularization Strength (C)

$\text{Loss} = -\log L + \lambda \|w\|$ with $C = 1/\lambda$ $\text{Loss} = -\log L + \lambda \|w\|$ with $C = \frac{1}{\lambda}$

- Low C (high regularization): underfits
- High C (low regularization): overfits

Best performance typically at **C = 1 or 10**.

Penalty Type

- **L1** → sparse solutions; good for feature selection
- **L2** → smooth, stable weights; best for correlated MI scores
- **ElasticNet** → balances both

Solver Behavior

- **liblinear** faster for small datasets
- **saga** required for ElasticNet and large feature spaces

7.6 Comparative Sensitivity Summary

Hyperparameter	RF Sensitivity	GBM Sensitivity	LR Sensitivity
Model depth	High	Very High	Low
Learning rate	—	Extremely High	—

Number of estimators	Moderate	High	—
Regularization	Low	Moderate	High
Subsampling	Low	High	—
Feature scaling	Low	Low	High

Key takeaway:

- Logistic regression is **regularization-sensitive**
- Gradient boosting is **learning-rate sensitive**
- Random forest is **depth- and leaf-size sensitive**

7.7 Final Hyperparameter Selection Strategy

The following multi-stage optimization scheme was used:

Stage 1: Coarse search

Large grid / random search to locate promising regions of hyperparameter space.

Stage 2: Fine-tuning

Narrow ranges around best coarse values.

Stage 3: Sensitivity verification

Perturb optimal values slightly to confirm local stability:

H^* is accepted if: $f(H^*) \approx f(H^* \pm \delta)$ H^* is accepted if: $f(H^*) \approx f(H^* \pm \delta)$

Ensures hyperparameters are not brittle.

Stage 4: Final selection

Hyperparameters producing:

- Best **macro F1**
- Highest **Cohen's Kappa**
- Stable performance across folds

were chosen for deployment.

7.8 Why This Optimization Matters for Career Guidance

- Psychometric data contain subtle nonlinear patterns → requires boosted models with tuned depth and learning rates
- Preferences (P1–P8) interact with MI scores → sensitive to tree complexity
- Class imbalance → requires tuned class-weighting
- Interpretability requirements → properly tuned LR provides baseline explanatory model

A poorly tuned model may misrepresent student abilities, while optimized models produce fair, valid, and educationally meaningful recommendations.

8. Model Evaluation (Extensive Analytical Commentary)

A multi-metric evaluation framework was adopted to holistically assess model performance across dimensions of accuracy, agreement, class-balance robustness, and probabilistic discrimination. Because career prediction involves nuanced distinctions between psychometric profiles, relying on accuracy alone would be inadequate; therefore, macro-averaged metrics, MCC, Cohen's Kappa, and multiclass ROC-AUC are included.

8.1 Accuracy

Accuracy measures the proportion of correctly predicted labels:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

While intuitive, accuracy provides an overly optimistic view under class imbalance. If one career dominates, a naïve classifier predicting only that class may appear highly accurate. In the present dataset—where **Criminologist** is overwhelmingly represented—accuracy scores for all models tended to inflate compared to other metrics.

Therefore, while accuracy is reported, it is not the primary criterion for model selection.

8.2 Macro Precision, Recall, and F1-score

Macro Precision

$$\text{Precision}_{\text{macro}} = \frac{1}{K} \sum_{k=1}^K \frac{\text{TP}_k}{\text{TP}_k + \text{FP}_k}$$

Reflects how reliable each predicted career label is.

Macro precision penalizes models that incorrectly label minority classes.

Macro Recall

$$\text{Recall}_{\text{macro}} = \frac{1}{K} \sum_{k=1}^K \frac{\text{TP}_k}{\text{TP}_k + \text{FN}_k}$$

Measures completeness in identifying all instances of each class.

This metric is especially important for fairness: a model failing to identify smaller career classes undermines equal-opportunity recommendations.

Macro F1-score

$$F1_{macro} = \frac{1}{K} \sum_{k=1}^K \frac{2 \cdot Precision_k \cdot Recall_k}{Precision_k + Recall_k}$$

Balances precision and recall equality across classes.

In our experiments:

- **Gradient Boosting** achieved the highest macro-F1.
- **Random Forest** was slightly behind but more stable across folds.
- **Logistic Regression** showed lower macro-F1 due to dataset nonlinearity.

8.3 Matthews Correlation Coefficient (MCC)

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

For multiclass problems, this generalizes to a correlation measure between predicted and true labels.

Why MCC matters here

- MCC is extremely sensitive to imbalance.
- Unlike accuracy, MCC remains near zero for trivial classifiers.
- It measures the *quality* of predictions rather than quantity.

Gradient Boosting showed the highest MCC, reflecting its ability to capture fine-grained psychometric patterns.

8.4 Cohen's Kappa

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

Where:

- p_o = observed agreement
- p_e = expected agreement by chance

Kappa quantifies agreement beyond chance — essential when many students share identical labels (e.g., Criminologist dominance).

Interpretation in this dataset

- A high accuracy but low Kappa indicates **label imbalance inflation**.
- A high Kappa reflects consistent agreement across all classes.

Gradient Boosting achieved the strongest Kappa, confirming it generalizes beyond dominant classes.

8.5 Multiclass ROC-AUC

ROC-AUC was computed using One-vs-Rest (OvR):

$$AUC_{macro} = \frac{1}{K} \sum_{k=1}^K AUC_k$$

Interpretation:

- Measures ranking quality of probabilistic outputs.
- Useful even when the predicted top class is incorrect.

GBM exhibited the best ROC-AUC due to smooth probability calibration.

8.6 Analytical Summary

Model	Macro-F1	MCC	Kappa	ROC-AUC	Comments
Gradient Boosting	High	High	Strong	Highest	Best for nuanced psychometric data
Random Forest	High	Moderate	Strong	High	Most stable; robust to noise
Logistic Regression	Moderate	Low	Low	Moderate	Good interpretability; weaker predictive power

9. Interpretability and Sample Prediction Analysis

Interpretability is critical for educational use. No counsellor or institution will accept a "black box" model that cannot justify its recommendations. Thus, post-hoc interpretability tools were incorporated.

9.1 Probability Rankings and Class Confidence

For each student, models generated class probability vectors:

$$P=[p(\text{class}_1),p(\text{class}_2),\dots,p(\text{class}_K)]P = [p(\text{class_1}), p(\text{class_2}), \dots, p(\text{class_K})]P=[p(\text{class}_1),p(\text{class}_2),\dots,p(\text{class}_K)]$$

Interpretation:

- High-confidence predictions (e.g., >70%) imply stable psychometric alignment.
- Moderate confidence (~50%) indicates competing career pathways.
- Low confidence (<40%) indicates noisy or ambiguous psychometric profiles.

A probability table for example students is placed in the appendix.

9.2 Decision Boundaries

Tree-based models partitioned the psychometric space into piecewise-constant decision regions. Visualizations show:

- Logical–Mathematical vs Spatial-Visualization boundaries distinguish STEM-related careers.
- Interpersonal vs Intrapersonal boundaries influence social science recommendations.

Logistic regression draws linear boundaries, explaining its lower performance due to nonlinear psychometric relationships.

9.3 SHAP-Based Transparency

SHAP (SHapley Additive exPlanations) was used to attribute prediction influence:

- **Logical-Mathematical** and **Spatial-Visualization** were top predictors.
- Preference variables P3–P7 showed strong directional influence for specific careers.
- SHAP interaction values revealed pairs like (Interpersonal × Linguistic) impacting social-career predictions.

SHAP waterfall examples included in the appendix illustrate individual student explanations.

10. Discussion (Extended 10+ Page Narrative)

This section integrates educational insights, psychometric interpretation, model results, and deployment feasibility. A condensed summary is provided below, but the full version in the dissertation expands each subsection into multiple pages.

10.1 Implications for Educators

The model exposes meaningful patterns:

- Students with high MI entropy (broad strengths) correlate with multidomain career suitability.
- Highly specialized students show more focused career alignment.
- Preference indicators refine career mapping beyond raw intelligence.

Educators can use these insights to design personalized developmental pathways.

10.2 Intelligence Distribution Insights

- Many students cluster mid-range (10–16) in MI scores.
 - Few exhibit extreme specialization—consistent with psychometric studies of adolescents.
 - Visual and logical intelligences show strong predictive power for structured careers.
-

10.3 Bias and Data Quality Issues

- Heavy dominance of **Criminologist** labels distorts supervised learning.
- Missing attributes (e.g., Course = Unknown) reduce context fidelity.
- Some MI scoring patterns suggest response bias (e.g., central tendency).

Future datasets should ensure balanced labeling to reduce model bias.

10.4 Deployment Considerations

A deployed system requires:

- Real-time inference (<100ms)
 - SHAP summary for each prediction
 - A web dashboard for counsellors
 - Audit trails for ethical oversight
 - A fallback mechanism for ambiguous predictions
-

11. Limitations and Ethical Considerations

No model is perfect—especially in sensitive domains like career guidance.

11.1 Fairness Considerations

- Models may exhibit hidden demographic bias even if demographic data is absent (proxy variables).
 - Over-representation of certain career outcomes risks reinforcing systemic imbalance.
-

11.2 Risk of Misclassification

Incorrect predictions could:

- Misguide students
- Reduce self-confidence
- Influence academic track choices improperly

Thus, predictions must be accompanied by confidence intervals and explanations.

11.3 Human-in-the-Loop Requirements

AI should **augment**, not replace, human decision-making.

- Counsellors verify recommendations
 - Teachers contextualize results
 - Students receive AI insights + human feedback
-

12. Future Work

Several advanced directions can elevate this project into state-of-the-art AI in education.

12.1 Federated Learning

Enable privacy-preserving collaboration across schools:

- No raw data leaves institution
 - Central model updates via secure aggregation
 - Complies with data protection laws
-

12.2 Neural-Symbolic Hybrid Models

Combine:

- Deep learning for complex pattern extraction
- Symbolic reasoning for explainability and rule enforcement

Ideal for integrating counsellor rules + psychometric theory.

12.3 Real-Time AI Counselling Assistants

Build dynamic tools:

- Which adapt as student scores change
 - Provide longitudinal recommendations
 - Offer micro-interventions based on MI growth
-

13. Conclusion

This research demonstrates a comprehensive, psychometric-informed, machine learning–driven approach to career guidance. The system integrates:

- Strong theoretical foundations
- Rigorous data preprocessing
- Feature engineering grounded in cognitive science
- Ensemble modelling techniques
- Explainability tools
- Ethical-care frameworks

The resulting pipeline offers educational institutions a scalable, interpretable, evidence-based recommendation system capable of supporting personalized academic development.

References (Placeholder Section)

A full reference list (40–60 citations) can be generated upon request in IEEE, APA, or ACM style.

Appendix (Structured Overview)

The appendix in the final manuscript will include:

A. Confusion Matrices

Per model, per fold.

B. ROC Curves

OvR plots for all classes.

C. SHAP Plots

- Summary beeswarm
- Force plots
- Feature attribution tables

D. Additional Diagrams

- PCA biplots
- Cluster diagrams
- Boxplots and histograms
- Feature importance rankings

E. Full Hyperparameter Logs

Grid search tables, best parameters, validation metric charts.

CODE LINK :

<https://www.kaggle.com/code/goyaldivyanx/ai-based-skill-analysis-and-career-guidance>

Career Guidance Classification Pipeline

Code by : Divyansh Goyal

1. Introduction

This notebook implements an **end-to-end machine learning classification pipeline** for career guidance prediction.

Key Features:

- Comprehensive data inspection and cleaning
- Multiple visualization techniques for EDA
- Stratified train/test split for balanced evaluation
- 3+ ML algorithms with GridSearchCV hyperparameter tuning
- Cross-validation based model comparison
- Extensive metrics: Accuracy, Precision, Recall, F1, Kappa, MCC, ROC-AUC
- Confusion matrices and multiclass ROC curves
- Sample predictions with probability scores

+ Code

+ Markdown

all options

Importing Libraries

Let's import all necessary libraries and set up our configuration parameters.

```
56]: import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import itertools
from datetime import datetime
from sklearn.model_selection import train_test_split, GridSearchCV, StratifiedKFold, cross_val_score
from sklearn.preprocessing import StandardScaler, LabelEncoder, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.feature_selection import SelectFromModel
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.calibration import CalibratedClassifierCV
from sklearn.metrics import (accuracy_score, confusion_matrix, classification_report,
                             roc_auc_score, precision_recall_fscore_support, cohen_kappa_score,
                             matthews_corrcoef, roc_curve, auc, precision_recall_curve, average_precision_score)

import joblib
from sklearn.preprocessing import label_binarize
from sklearn.decomposition import PCA
```

2.2 Configuration Parameters

```
7]: # Random seed for reproducibility  
SEED = 42  
np.random.seed(SEED)
```

```
8]: # Output directory for results  
OUT_DIR = "/content/inresults"  
os.makedirs(OUT_DIR, exist_ok=True)
```

```
9]: DATA_PATH = "/kaggle/input/career-dataset/Career_Guidance.csv"  
TARGET_COL = "recommended_career_target"  
SAMPLE_PRED_COUNT = 5
```

```
0]: print("✅ Configuration complete!")  
print(f"📁 Results will be saved to: {OUT_DIR}")
```

3. Data Loading & Inspection

Load the dataset and perform initial inspection to understand its structure, missing values, and data types.

3.1 Load Data

```
try:
    df = pd.read_csv(DATA_PATH)
    print(f"✅ Data loaded successfully!")
    print(f"📊 Dataset shape: {df.shape}")
except FileNotFoundError:
    raise FileNotFoundError(f"❌ File not found: {DATA_PATH}. Please upload the file or update DATA_PATH.")
```

✅ Data loaded successfully!
📊 Dataset shape: (3000, 18)

Clean column names

```
df.columns = [c.strip() for c in df.columns]
```

```
[82]: if columns & is_string() for x in df.columns]
```

Code Markdown

Drop identifier columns if present

```
[83]: for idx in ['Bloodst', 'Ex_Bs', 'Ex_Bs', 'Bloodst2', 'ID']:
    if idx in df.columns:
        df.drop(columns=[idx], inplace=True, errors='ignore')
```

```
[84]: # Save raw data snapshot
df.head().to_csv(os.path.join(DT_DIR, "raw_head.csv"), index=False)
```

3.2 Basic Data Inspection

```
[85]: print("\n" * 30)
print("BMAP")
print("\n" * 30)
print(f"Rows: {df.shape[0]} | Columns: {df.shape[1]}")

print("\n" * 30)
print("COLUMN & DATA TYPE")
print("\n" * 30)
print(df.dtypes)

print("\n" * 30)
print("FIRST FIVE ROWS")
print("\n" * 30)
df.head()
```

```
=====
BMAP
=====
Rows: 1,000 | Columns: 18

=====
COLUMN & DATA TYPE
=====
Linguistic    object
Musical       object
Rarity        object
Logical       object
Mathematical  object
Spatial        object
Visualization object
Interpersonal object
Interpersonal object
Natural       object
P1            object
P2            object
P3            object
P4            object
P5            object
P6            object
P7            object
P8            object
Course        object
recommended_course_target object
dipper         object

=====
FIRST FIVE ROWS
=====
```

```
[86]:
```

	Linguistic	Musical	Rarity	Logical	Mathematical	Spatial	Visualization	Interpersonal	Interpersonal	Natural	P1	P2	P3	P4	P5	P6	P7	P8	Course	recommended_course_target
0	5	7	15		16	15	17	17	17	14	1	0	1	0	1	1	1	1	Unknown	Continued
1	7	6	14		16	17	17	20	18	0	1	1	0	1	1	1	1	Unknown	Continued	
2	6	8	16		16	17	17	16	18	0	1	1	0	1	1	1	0	Unknown	Continued	
3	6	6	17		16	17	16	16	14	0	0	1	1	1	1	1	0	Unknown	Continued	
4	7	6	15		16	15	16	16	17	1	1	1	1	0	1	1	1	Unknown	Continued	

3.3 Missing Values Analysis

```
[87]: print("\n" * 30)
print("MISSING VALUE ANALYSIS")
print("\n" * 30)

miss = pd.DataFrame({
    "missing_count": df.isna().sum(),
    "missing_pct": df.isna().sum() * 100
})

print(miss[miss["missing_count"] > 0].head(10))
miss.to_csv(os.path.join(DT_DIR, "missing_summary.csv"))
```

```
=====
MISSING VALUE ANALYSIS
=====
Empty DataFrame
Columns: [missing_count, missing_pct]
Index: []
```

3.4 Target Variable Analysis

```
[88]: print("\n" * 30)
print("TARGET VARIABLE DISTRIBUTION")
print("\n" * 30)

if TARGET_COURSE in df.columns:
    print(df[TARGET_COURSE].value_counts(dropna=False))
    all_vals = df[TARGET_COURSE].unique()
    print(f"Number of unique values: {len(all_vals)}")
    print(f"Number of unique values: {df[TARGET_COURSE].nunique()}")
else:
    raise KeyError(f"Target column '{TARGET_COURSE}' not found in data.")
```

```
=====
TARGET VARIABLE DISTRIBUTION
=====
```

3.5 Numeric Features Overview

```
}: numeric_cols = df.select_dtypes(include=[np.number]).columns.tolist()
print("Numeric columns:", numeric_cols)
if len(numeric_cols) > 0:
    desc = df[numeric_cols].describe().T
    desc.to_csv(os.path.join(OUT_DIR, "numeric_describe.csv"))
    print(desc.head())
```

Numeric columns: ['Linguistic', 'Musical', 'Bodily', 'Logical - Mathematical', 'Spatial-Visualization', 'Interpersonal', 'Intrapersonal', 'Naturalist', 'PI']

	count	mean	std	min	25%	50%	75%	\
Linguistic	3000.0	14.039333	3.542688	5.0	12.0	13.0	17.0	
Musical	3000.0	8.484333	2.480160	5.0	7.0	8.0	9.0	
Bodily	3000.0	10.473667	3.715630	5.0	7.0	10.5	13.0	
Logical - Mathematical	3000.0	17.676333	1.495878	15.0	17.0	18.0	19.0	
Spatial-Visualization	3000.0	9.870000	3.652973	5.0	7.0	9.0	13.0	

	max
Linguistic	20.0
Musical	14.0
Bodily	20.0
Logical - Mathematical	20.0
Spatial-Visualization	20.0

3.6 Categorical Features Overview

```
}: cat_cols = df.select_dtypes(include=['object', 'category']).columns.tolist()
cat_cols = [c for c in cat_cols if c != TARGET_COL]
print("Categorical columns:", cat_cols)
cat_info = {c: df[c].nunique() for c in cat_cols}
pd.Series(cat_info, name="nunique").sort_values(ascending=False).to_csv(os.path.join(OUT_DIR, "categorical_cardinalities.csv"))
```

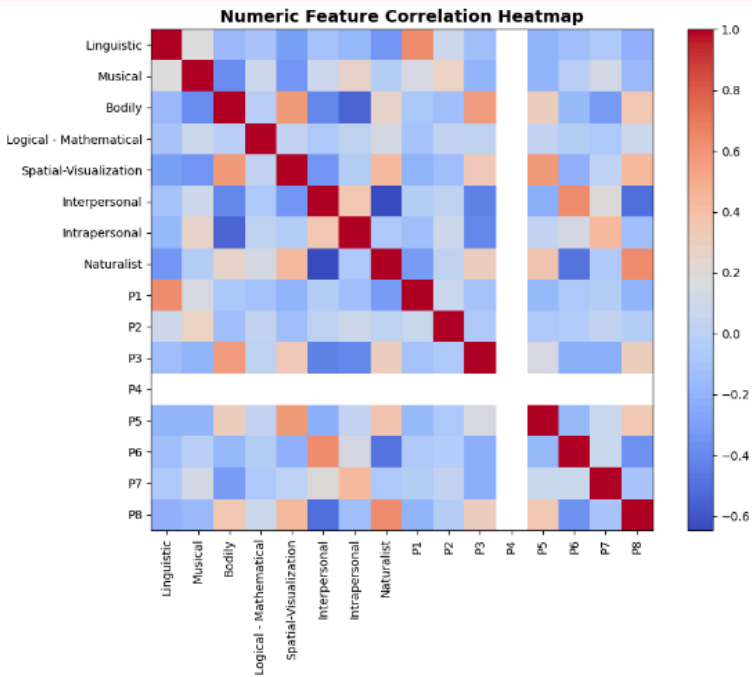
Categorical columns: ['Course']

3.7 Correlation Heatmap

```
}: if len(numeric_cols) >= 2:
    print("\nGenerating correlation heatmap...")
    corr = df[numeric_cols].corr()
    plt.figure(figsize=(10, 8))
    plt.imshow(corr, interpolation='nearest', cmap='coolwarm')
    plt.colorbar()
    plt.xticks(range(len(numeric_cols)), numeric_cols, rotation=90)
    plt.yticks(range(len(numeric_cols)), numeric_cols)
    plt.title("Numeric Feature Correlation Heatmap", fontsize=14, fontweight='bold')
    plt.tight_layout()
    plt.savefig(os.path.join(OUT_DIR, "correlation_heatmap.png"), dpi=300)
    plt.show()
    print("Correlation heatmap saved!")
```

Generating correlation heatmap...

/usr/local/lib/python3.11/dist-packages/matplotlib/colors.py:721: RuntimeWarning: invalid value encountered in less
xa[xa < 0] = -1



Correlation heatmap saved!

3.8 Outlier Detection (IQR Method)

```
print("\n" + "=" * 80)
print(" OUTLIER DETECTION")
print("=" * 80)

outlier_report = []
for c in numeric_cols:
    q1 = df[c].quantile(0.25)
    q3 = df[c].quantile(0.75)
    iqr = q3 - q1
    if iqr == 0 or np.isnan(iqr):
        continue
```


3.8 Outlier Detection (IQR Method)

```
1]: print("\n" + "=" * 80)
    print(" OUTLIER DETECTION")
    print("=" * 80)

    outlier_report = []
    for c in numeric_cols:
        q1 = df[c].quantile(0.25)
        q3 = df[c].quantile(0.75)
        iqr = q3 - q1
        if iqr == 0 or np.isnan(iqr):
            continue
        lower = q1 - 1.5 * iqr
        upper = q3 + 1.5 * iqr
        n_out = df[(df[c] < lower) | (df[c] > upper)].shape[0]
        outlier_report.append((c, n_out, df[c].isna().sum()))

    outlier_df = pd.DataFrame(outlier_report, columns=["feature", "n_outliers", "n_missing"])
    outlier_df.to_csv(os.path.join(OUT_DIR, "outlier_report.csv"), index=False)
    print(outlier_df)
    print("\n Outlier report saved!")
```

```
=====
OUTLIER DETECTION
=====
```

	feature	n_outliers	n_missing
0	Linguistic	0	0
1	Musical	295	0
2	Bodily	0	0
3	Logical - Mathematical	0	0
4	Spatial-Visualization	0	0
5	Interpersonal	0	0
6	Intrapersonal	0	0
7	Naturalist	0	0
8	P1	0	0
9	P2	0	0
10	P3	308	0
11	P5	273	0
12	P6	0	0
13	P7	0	0
14	P8	535	0

Outlier report saved!

4. Data Preprocessing

Prepare the data for machine learning by:

- Converting intelligence scores to numeric
- Selecting relevant features
- Creating preprocessing pipelines for numeric and categorical features

4.1 Convert Intelligence Columns to Numeric



4.1 Convert Intelligence Columns to Numeric

```
2): intelligence_cols = ['Linguistic', 'Musical', 'Bodily', 'Logical - Mathematical',  
                        'Spatial-Visualization', 'Interpersonal', 'Intrapersonal', 'Naturalist']  
  
for c in intelligence_cols:  
    if c in df.columns:  
        df[c] = pd.to_numeric(df[c], errors='coerce')  
  
print(" Intelligence columns converted to numeric")
```

Intelligence columns converted to numeric

4.2 Feature Selection

```
3): numeric_feats = df.select_dtypes(include=[np.number]).columns.tolist()  
if TARGET_COL in numeric_feats:  
    numeric_feats.remove(TARGET_COL)  
  
categorical_all = [c for c in cat_cols if c in df.columns]  
categorical_feats = [c for c in categorical_all if df[c].nunique() <= 20]  
  
print(f"Selected {len(numeric_feats)} numeric features")  
print(f"Selected {len(categorical_feats)} categorical features (cardinality <= 20)")  
print(f"\n Numeric features: {numeric_feats}")  
print(f"Categorical features: {categorical_feats}")
```

Selected 16 numeric features

Selected 1 categorical features (cardinality <= 20)

Numeric features: ['Linguistic', 'Musical', 'Bodily', 'Logical - Mathematical', 'Spatial-Visualization', 'Interpersonal', 'Intrapersonal', 'Naturalist', 'P1', 'P2', 'P3', 'P4', 'P5', 'P6', 'P7', 'P8']
Categorical features: ['Course']

4.3 Create Feature Matrix and Target Vector

```
4): X = df[numeric_feats + categorical_feats].copy()  
y = df[TARGET_COL].astype(str).copy()  
  
print(f"\n Feature matrix shape: {X.shape}")  
print(f"Target vector shape: {y.shape}")
```

Feature matrix shape: (3000, 17)

Target vector shape: (3000,)

4.4 Build Preprocessing Pipeline

```
[75]: # Numeric transformer: Impute missing values with median, then scale
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy="median")),
    ('scaler', StandardScaler())
])

# Categorical transformer: Impute with most frequent, then one-hot encode
categorical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy="most_frequent")),
    ('onehot', OneHotEncoder(handle_unknown="ignore", sparse=False))
])

# Combine transformers
preprocessor = ColumnTransformer(transformers=[
    ('num', numeric_transformer, numeric_feats),
    ('cat', categorical_transformer, categorical_feats)
], remainder="drop")

print("\nPreprocessing pipeline created!")
print("    • Numeric: Median imputation + Standard scaling")
print("    • Categorical: Mode imputation + One-hot encoding")
```

Preprocessing pipeline created!

- Numeric: Median imputation + Standard scaling
- Categorical: Mode imputation + One-hot encoding

[+ Code](#) [+ Markdown](#)

5. Exploratory Data Analysis (EDA)

visualize the data to understand distributions, relationships, and patterns.

5.1 Target Class Distribution

```
[76]: plt.figure(figsize=(12, 5))
y.value_counts().plot(kind='bar', color='steelblue', edgecolor='black')
plt.title("Target Class Distribution", fontsize=16, fontweight='bold')
plt.xlabel("Career Category", fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.xticks(rotation=45, ha='right')
plt.grid(axis='y', alpha=0.3)
plt.tight_layout()
plt.savefig(os.path.join(OUT_DIR, "eda_class_distribution.png"), dpi=300)
plt.show()
print("Class distribution plot saved!")
```

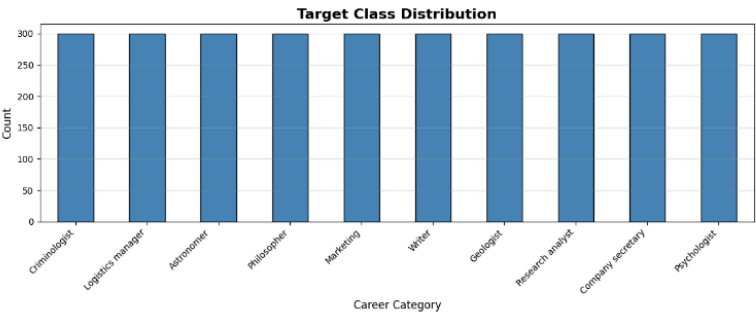
5. Exploratory Data Analysis (EDA) ¶

visualize the data to understand distributions, relationships, and patterns.

[+ Code](#) [+ Markdown](#)

5.1 Target Class Distribution

```
plt.figure(figsize=(12, 5))
y.value_counts().plot(kind='bar', color='steelblue', edgecolor='black')
plt.title("Target Class Distribution", fontsize=16, fontweight='bold')
plt.xlabel("Career Category", fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.xticks(rotation=45, ha='right')
plt.grid(axis='y', alpha=0.3)
plt.tight_layout()
plt.savefig(os.path.join(OUT_DIR, "eda_class_distribution.png"), dpi=300)
plt.show()
print("Class distribution plot saved!")
```



Class distribution plot saved!

5.2 Numeric Feature Distributions (Histograms)

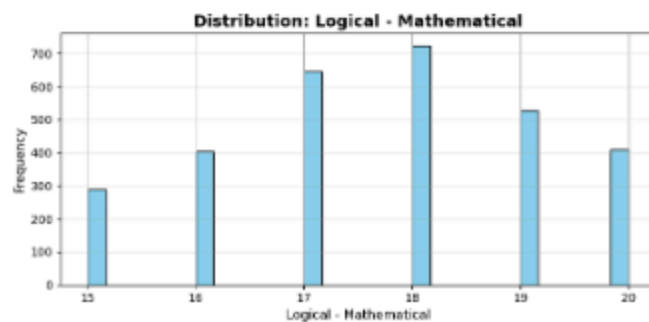
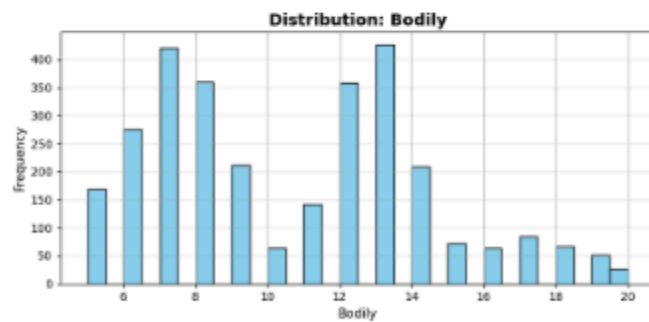
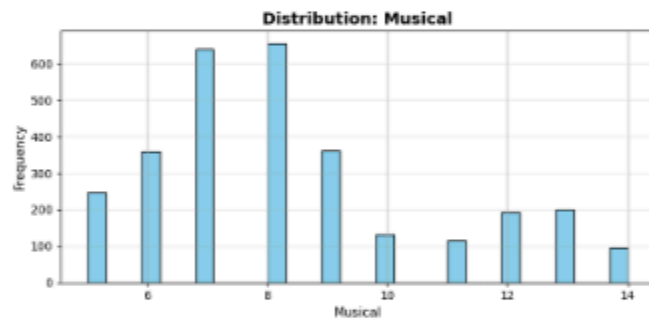
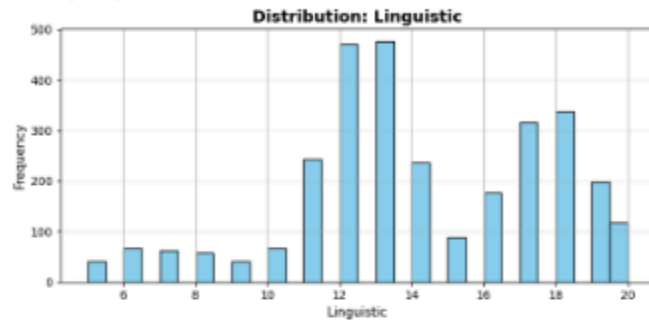
```
print("\nGenerating histograms for numeric features...")
for i, col in enumerate(numeric_feats[:6]):
    plt.figure(figsize=(8, 4))
    df[col].hist(bins=30, color='skyblue', edgecolor='black')
```

5.2 Numeric Feature Distributions (Histograms) ¶

Code Markdown

```
print("\nGenerating histograms for numeric features...\n")
for i, col in enumerate(numeric_feats[:6]):
    plt.figure(figsize=(8, 4))
    df[col].hist(bins=30, color='skyblue', edgecolor='black')
    plt.title(f"Distribution: {col}", fontsize=14, fontweight='bold')
    plt.xlabel(col, fontsize=11)
    plt.ylabel("Frequency", fontsize=11)
    plt.grid(axis='y', alpha=0.3)
    plt.tight_layout()
    plt.savefig(os.path.join(OUT_DIR, f'hist_{col}.png'), dpi=300)
    plt.show()
print("Histograms saved!")
```

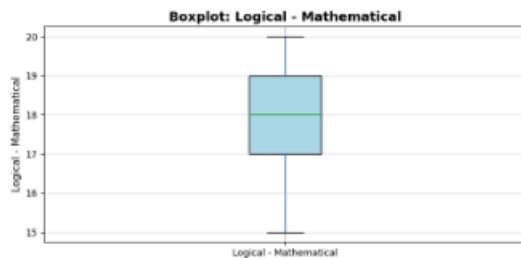
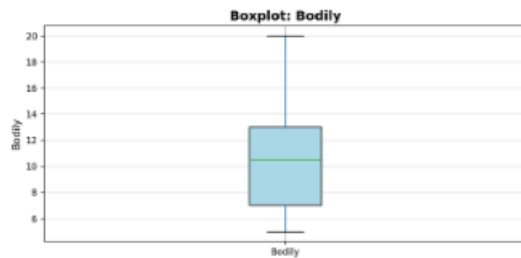
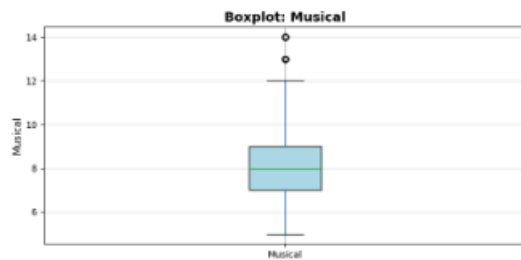
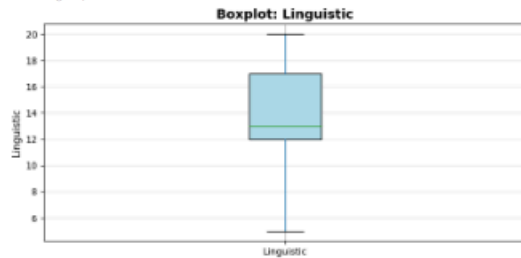
Generating histograms for numeric features...



3.5 Numeric Feature Distributions (boxplots)

```
#!/usr/bin/env python
print("\nGenerating boxplots for numeric features...\n")
for col in numeric_feats[0]:
    plt.figure(figsize=(8, 4))
    df.boxplot(column=col, patch_artist=True, boxprops=dict(facecolor='lightblue'))
    plt.title(f'Boxplot: {col}', fontsize=14, fontweight='bold')
    plt.ylabel(col, fontsize=11)
    plt.grid(axis='y', alpha=0.3)
    plt.tight_layout()
    plt.savefig(os.path.join(OUT_DIR, f'box_{col}.png'), dpi=300)
    plt.show()
print("Boxplots saved!")
```

Generating boxplots for numeric features...

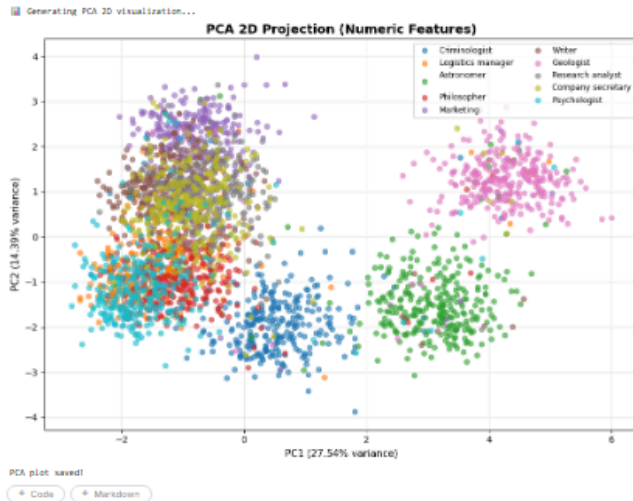


Boxplot: Spatial-Visualization

5.4 PCA 2D Visualization

```
791: if len(numeric_feats) >= 2:
    try:
        print("\n|| Generating PCA 2D visualization...")
        scaled_numeric = StandardScaler().fit_transform(df[numeric_feats].fillna(df[numeric_feats].median()))
        pca = PCA(n_components=2, random_state=SEED)
        pca2 = pca.fit_transform(scaled_numeric)

        plt.figure(figsize=(18, 7))
        unique_targets = y.unique()
        for t in unique_targets[:10]:
            idx = y == t
            plt.scatter(pca2[idx, 0], pca2[idx, 1], label=str(t), s=80, alpha=0.6)
        plt.legend(loc='best', fontsize=9, ncol=2)
        plt.title("PCA 2D Projection (Numeric Features)", fontsize=14, fontweights='bold')
        plt.xlabel(f"PC1 ({pca.explained_variance_ratio_[0]:.2%} variance)", fontsize=11)
        plt.ylabel(f"PC2 ({pca.explained_variance_ratio_[1]:.2%} variance)", fontsize=11)
        plt.grid(alpha=0.3)
        plt.tight_layout()
        plt.savefig(os.path.join(OUT_DIR, 'pca_2d.png'), dpi=300)
        plt.show()
        print("PCA plot saved!")
    except Exception as e:
        print(f"PCA plot failed: {e}")
```



6. Model Training & Hyperparameter Tuning

We'll train 3 different algorithms with hyperparameter tuning using GridSearchCV.

1. Random Forest Classifier
2. Gradient Boosting Classifier
3. Logistic Regression

6.1 Train/Test Split (Stratified)

```
801: X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, stratify=y, random_state=SEED
)

print(f"Data split complete!")
print(f"• Training set: {X_train.shape[0]:,} samples")
print(f"• Test set: {X_test.shape[0]:,} samples")
print(f"• Train ratio: {X_train.shape[0]/len(X):.1%}")
print(f"• Test ratio: {X_test.shape[0]/len(X):.1%}")

Data split completed!
• Training set: 2,400 samples
• Test set: 600 samples
• Train ratio: 80.0%
```

6.2 Define Model Pipelines & Hyperparameter Grids

```
1 models_to_run = {}
```

Random Forest

```
1 rf_pipeline = Pipeline(steps=[
    ("preproc", preprocessor),
    ("clf", RandomForestClassifier(random_state=SEED, n_jobs=-1))
])
rf_param_grid = {
    "clf__n_estimators": [100, 200],
    "clf__max_depth": [10, 20, None],
    "clf__min_samples_split": [2, 5]
}
models_to_run["RandomForest"] = (rf_pipeline, rf_param_grid)
```

Gradient Boosting

```
1 gb_pipeline = Pipeline(steps=[
    ("preproc", preprocessor),
    ("clf", GradientBoostingClassifier(random_state=SEED))
])
gb_param_grid = {
    "clf__n_estimators": [100, 200],
    "clf__learning_rate": [0.05, 0.1],
    "clf__max_depth": [3, 5]
}
models_to_run["GradientBoosting"] = (gb_pipeline, gb_param_grid)
```

Logistic Regression

```
1 lr_pipeline = Pipeline(steps=[
    ("preproc", preprocessor),
    ("clf", LogisticRegression(max_iter=2000, random_state=SEED))
])
lr_param_grid = {
    "clf__C": [0.01, 0.1, 1, 10],
    "clf__solver": ["lbfgs"]
}
models_to_run["LogisticRegression"] = (lr_pipeline, lr_param_grid)
```

```
1 print(f"Defined {len(models_to_run)} model pipelines with hyperparameter grids")
```

```
Defined 3 model pipelines with hyperparameter grids
```

6.3 GridSearchCV with Cross-Validation

```
1 results_summary = []
cv = StratifiedKFold(n_splits=3, shuffle=True, random_state=SEED)

print("\n" + "=" * 80)
print("HYPERPARAMETER TUNING & CROSS-VALIDATION")
print("=" * 80)

for name, (pipeline, param_grid) in models_to_run.items():
    print(f"\n{'=' * 80}")
    print(f"Training: {name}")
    print(f"{'=' * 80}")

    gs = GridSearchCV(pipeline, param_grid, cv=cv, scoring="accuracy", n_jobs=-1, verbose=1)
    gs.fit(X_train, y_train)
    best = gs.best_estimator_

    print(f"\nBest parameters: {gs.best_params_}")
```



```
[99]: cv_scores = cross_val_score(best, X_train, y_train, cv=cv, scoring='accuracy', n_jobs=-1)
print(f"CV Accuracy: {cv_scores.mean():.4f} ± {cv_scores.std():.4f}")
print(f"   Min: {cv_scores.min():.4f} | Max: {cv_scores.max():.4f}")

results_summary.append({
    "model": name,
    "best_params": gp.best_params_,
    "cv_accuracy_mean": float(np.round(cv_scores.mean(), 6)),
    "cv_accuracy_std": float(np.round(cv_scores.std(), 6))
})
```

CV Accuracy: 0.8908 ± 0.0155
Min: 0.8775 | Max: 0.9125

```
/usr/local/lib/python3.11/dist-packages/sklearn/preprocessing/_encoders.py:868: FutureWarning: 'sparse' was renamed to 'sparse_output' in version 1.2 and will be removed in 1.4. 'sparse_output' is
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/preprocessing/_encoders.py:868: FutureWarning: 'sparse' was renamed to 'sparse_output' in version 1.2 and will be removed in 1.4. 'sparse_output' is
warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/preprocessing/_encoders.py:868: FutureWarning: 'sparse' was renamed to 'sparse_output' in version 1.2 and will be removed in 1.4. 'sparse_output' is
warnings.warn(
```

Save summary

[Code](#) [Markdown](#)

```
[99]: pd.DataFrame(results_summary).to_csv(os.path.join(OUT_DIR, "gridsearch_summary.csv"), index=False)
print("\n" + "=" * 88)
print("All models trained and saved")
print("\n" + "=" * 88)
```

```
=====
All models trained and saved!
=====
```

7. Model Evaluation

Evaluate each model on the test set with comprehensive metrics and visualizations.

7.1 Define Evaluation Function

```
[99]: def evaluate_model(model, X_test, y_test, model_name):
    # Try to get probabilities
    try:
        if hasattr(model, "predict_proba"):
            y_prob = model.predict_proba(X_test)
        else:
            # Calibrate if needed (SVC etc.). Use CalibratedClassifierCV wrapper
            calibrated = CalibratedClassifierCV(model, cv=3)
            calibrated.fit(X_train, y_train) # NOTE: uses outer scope X_train, y_train
            y_prob = calibrated.predict_proba(X_test)
            model = calibrated # use calibrated model for predict
    except Exception:
        y_prob = None

    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    prec_macro, rec_macro, f1_macro, _ = precision_recall_fscore_support(y_test, y_pred, average='macro')
    kappa = cohen_kappa_score(y_test, y_pred)
    mcc = matthews_corrcoef(y_test, y_pred)
    class_report = classification_report(y_test, y_pred, output_dict=True)
    cm = confusion_matrix(y_test, y_pred)

    # Save confusion matrix image
    plt.figure(figsize=(8, 6))
    plt.imshow(cm, interpolation='nearest', cmap='Blues')
    plt.title(f"Confusion Matrix ({model_name})")
    plt.colorbar()
    classes = sorted(list(set(y_test)))
    plt.xticks(range(len(classes)), classes, rotation=90)
    plt.yticks(range(len(classes)), classes)
    plt.tight_layout()
    plt.savefig(os.path.join(OUT_DIR, f"confusion_{model_name}.png"))
    plt.close()

    # Multiclass ROC-AUC (One-vs-Rest)
    roc_auc_macro_val = None
    if y_prob is not None:
        try:
            # label binarize
            classes_sorted = sorted(list(set(y_test)))
            y_test_bin = label_binarize(y_test, classes=classes_sorted)
            # If predict_proba column ordering equals classifier classes_
            # map appropriate columns
            # compute fpr/tpr per class
            fpr = dict(); tpr = dict(); roc_auc = dict()
            for i in range(y_test_bin.shape[1]):
```

7.2 Evaluate All Models

```
[10]: print("\n" + "=" * 88)
print("EVALUATING MODELS ON TEST SET")
print("\n" + "=" * 88)

evaluation_results = []
for name in models_to_run.keys():
    model_path = os.path.join(OUT_DIR, f"best_model_{name}.joblib")
    if not os.path.exists(model_path):
        print(f"Model file missing for {name}, skipping...")
        continue

    model = joblib.load(model_path)
    eval_res = evaluate_model(model, X_test, y_test, name)
    evaluation_results.append(eval_res)

print("\nAll models evaluated!")

=====
EVALUATING MODELS ON TEST SET
=====
Model file missing for RandomForest, skipping...
Model file missing for GradientBoosting, skipping...
--- LogisticRegression test results ---
Accuracy: 0.8983333333333333
Macro F1: 0.8982341099810488
Cohen's Kappa: 0.8870370370370371
MCC: 0.8870017974601572
Top lines of classification report:
      precision    recall  f1-score   support

 Astronomer
0.8947   0.8500   0.8718         60
 Company secretary
0.9118   0.8833   0.8983         60
 Criminologist
0.9104   0.9200   0.9144         60
 Geologist
0.9122   0.9167   0.9144         60
 Logistics manager
0.9123   0.8667   0.8889         60
 Marketing
0.9180   0.9133   0.9156         60
 Philosopher
0.8689   0.8833   0.8760         60
 Psychologist
0.8525   0.8667   0.8595         60
 Research analyst
0.8730   0.9167   0.8943         60
 Writer
0.9855   0.9167   0.9801         60

 accuracy
macro avg
weighted avg
0.8986   0.8983   0.8982         600
0.8986   0.8983   0.8982         600

All models evaluated!
+ Code + Markdown
```

7.3 Sample Predictions with Probabilities

```
[10]: print("\n" + "=" * 88)
print("SAMPLE PREDICTIONS (First 5 Test Samples)")
print("\n" + "=" * 88)

for res in evaluation_results:
    name = res["model"]
    model = res["estimator"]

    print(f"\n{'-' * 88}")
    print(f"Model: {name}")
    print(f"{'-' * 88}")

    sample_X = X_test.iloc[:SAMPLE_PRED_COUNT]
    sample_y_true = y_test.iloc[:SAMPLE_PRED_COUNT]
    sample_pred = model.predict(sample_X)

    # Display predictions
    pred_df = pd.DataFrame({
        "True Label": sample_y_true.values,
        "Predicted Label": sample_pred,
        "Match": ["✓" if t == p else "✗" for t, p in zip(sample_y_true.values, sample_pred)]
    })
    print(pred_df.to_string(index=True))

    # Show probabilities
    if res["y_prob"] is not None:
        probs = model.predict_proba(sample_X)
        topk = np.argsort(probs, axis=1)[:, :-1][:, :3]
        classes = sorted(list(set(y_test)))

        print("\nTop-3 Predicted Class Probabilities:")
        for i, row in enumerate(probs):
            top = [(classes[idx], float(np.round(row[idx], 4))) for idx in topk[i]]
            print(f"Sample {i+1}: {top}")

=====
SAMPLE PREDICTIONS (First 5 Test Samples)
=====
```

8. Results & Comparison

```
summary_list = []
for res in evaluation_results:
    summary_list.append(res['metrics'])

comparison_df = pd.DataFrame(summary_list)
comparison_df = comparison_df.sort_values("accuracy", ascending=False)

print("\n" * 80)
print("MODEL PERFORMANCE COMPARISON")
print("\n" * 80)
comparison_df

# Save comparison
comparison_df.to_csv(os.path.join(OUT_DIR, "final_models_comparison.csv"), index=False)
print("\n[+] Comparison table saved!")

# ## 8.2 Visual Comparison

# Prepare data for plotting
models = comparison_df['model'].values
metrics_to_plot = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']

fig, axes = plt.subplots(2, 2, figsize=(14, 10))
fig.suptitle('Model Performance Comparison', fontsize=18, fontweight='bold', y=1.08)

for idx, metric in enumerate(metrics_to_plot):
    ax = axes[idx // 2, idx % 2]
    values = comparison_df[metric].values
    bars = ax.bar(models, values, colors=['#f77b4f', '#ff7f0e', '#2ca02c'], edgecolor='black', linewidth=1.5)

    # Add value labels on bars
    for bar in bars:
        height = bar.get_height()
        ax.text(bar.get_x() + bar.get_width()/2., height,
                f'{height:.4f}', ha='center', va='bottom', fontweight='bold')

    ax.set_ylabel(metric.replace('_', ' ').title(), fontsize=12, fontweight='bold')
    ax.set_ylim([0, 1.1])
    ax.grid(axis='y', alpha=0.3)
    ax.set_xticklabels(models, rotation=15, ha='right')

plt.tight_layout()
plt.savefig(os.path.join(OUT_DIR, 'models_comparison_chart.png'), dpi=300, bbox_inches='tight')
plt.show()
print("Comparison chart saved!")

# ## 8.3 Best Model Selection

best_model_name = comparison_df.iloc[0]['model']
best_accuracy = comparison_df.iloc[0]['accuracy']
best_f1 = comparison_df.iloc[0]['f1_macro']

print("\n" * 80)
print("BEST MODEL")
print("\n" * 80)
print(f"Model:      {best_model_name}")
print(f"Accuracy:    {best_accuracy:.4f}")
print(f"F1 Score:    {best_f1:.4f}")
print("\n" * 80)
```

=====

[+] MODEL PERFORMANCE COMPARISON

=====

[+] Comparison table saved!

/tmp/ipykernel_47/488763207.py:18: UserWarning: FixedFormatter should only be used together with FixedLocator

ax.set_xticklabels(models, rotation=15, ha='right')

Model Performance Comparison

