

Okay, here's a concept for a complex dataset designed to cover the machine learning topics you've listed. This dataset will focus on Employee Performance and Attrition Prediction.

Dataset Name: EmployeeAnalyticsHR.csv

Dataset Goal: To provide a rich environment for exploring multiple linear regression, logistic regression, decision trees, preprocessing techniques, model evaluation, and fundamental machine learning concepts.

Number of Rows (Suggested): 750 - 1500 records (to allow for robust training/testing splits and demonstrate various phenomena like the curse of dimensionality with OHE, or benefits of regularization).

Number of Columns (Features + Targets): ~25-30 (allowing for multicollinearity, feature selection, and a good mix of data types).

I. Core Features:

Here's a breakdown of the features, their types, example values, and their relevance to your topics:

Feature Name	Type	Description & Example Values	ML Topic Relevance
EmployeeID	Identifier	Unique ID for each employee (e.g., EMP001, EMP002)	Basic identification
Age	Numerical	Employee's age in years (e.g., 22-60)	Regression, Transformations, Ethics (potential bias)
Gender	Categorical	Male, Female, Other, PreferNotToSay	Categorical Encoding, Ethics (bias detection,

			fairness), Decision Tree Splits
EducationLevel	Ordinal	High School, Bachelor's, Master's, PhD	Categorical Encoding (One-Hot, Ordinal), Decision Tree Splits
Department	Categorical	Sales, Engineering, HR, Marketing, Finance, Operations	One-Hot Encoding, Interpreting OHE Coefficients, Decision Tree Splits
JobRole	Categorical	Analyst, Engineer, Senior Engineer, Manager, Specialist, Director	One-Hot Encoding (many categories to show benefits of regularization), Decision Tree Splits
OfficeLocation	Categorical	HQ-CityA, Branch-CityB, Branch-CityC, Remote	One-Hot Encoding, Decision Tree Splits
YearsAtCompany	Numerical	Total years employee has been with the company (e.g., 0.5 - 30)	Regression, Potential multicollinearity with TotalWorkingYears, Transformations
TotalWorkingYears	Numerical	Total years of professional	Regression, Potential multicollinearity with

		experience (e.g., 1 - 40)	YearsAtCompany, Transformations
PreviousSalaryUSD	Numerical	Employee's salary at their previous job (if applicable, else 0 or NaN)	Regression, Transformations (likely skewed, good for Log), Missing data handling
MonthlyIncomeUSD	Numerical	Current gross monthly income	Primary Target for Regression (transformed to Annual Salary for some exercises), Skewed (good for Log)
TrainingHoursLastYear	Numerical	Hours spent in company- provided training last year (e.g., 0- 200)	Regression, Feature Importance
ProjectsCompletedLastCycle	Numerical	Number of projects completed in the last review cycle (e.g., 0-10)	Regression, Feature Importance
LastPerformanceRating	Ordinal	1-Poor, 2-Needs Improvement, 3-Meets Expectations, 4-	Categorical Encoding, Potential target or strong predictor for other targets

		Exceeds, 5- Outstanding	
SatisfactionScore	Ordinal	Employee- reported job satisfaction (1- Low, 2-Medium, 3-High, 4-Very High)	Categorical Encoding, Predictor for churn, Decision Tree Splits
WorkLifeBalanceScore	Ordinal	Employee- reported WLB (1-Poor, 2-Fair, 3-Good, 4- Excellent)	Categorical Encoding, Predictor for churn, Decision Tree Splits
OverTimeHoursLastMonth	Numerical	Hours of overtime worked last month (e.g., 0- 80)	Regression, Skewed data potential, Predictor for churn
PromotionLast2Years	Binary	0 (No), 1 (Yes)	Categorical Encoding, Predictor for churn/salary
AvgDailyCommuteTimeMin	Numerical	Average daily commute time in minutes (0 for Remote)	Regression, Transformations
TechnicalSkillScore	Numerical	Score from a standardized technical	Regression, Feature Selection

		assessment (0-100)	
<b>SoftSkillScore</b>	<b>Numerical</b>	Score from a standardized soft skills assessment (0-100)	Regression, Feature Selection
<b>EngagementSurveyScore</b>	<b>Numerical</b>	Score from last employee engagement survey (1-5)	Regression, Predictor for churn
<b>NumDirectReports</b>	<b>Numerical</b>	Number of direct reports (0 for individual contributors)	Regression, Feature Engineering (e.g., IsManager flag)
<b>UsesCompanyLaptop</b>	<b>Binary</b>	0 (No), 1 (Yes)	Can be an example of a less impactful feature for feature selection exercises
<b>LeftCompanyWithinYear</b>	<b>Binary</b>	TARGET: 0 (No), 1 (Yes) - for Churn Prediction	Primary Target for Classification, Intentionally create class imbalance (e.g., 15-25% 'Yes') for Class Imbalance Problems, Logistic Regression, Decision Trees, Classification Metrics.

PerformanceScoreOutOf100	Numerical	TARGET: Composite performance score (0-100)	Primary Target for Regression, can be engineered from other features or given directly. For Multiple Linear Regression, CART Regression.
--------------------------	-----------	--	---

---

## II. How this Dataset Addresses Your Topics:

### 1. Multiple Linear Regression (MLR) - Introduction & StatsModels:

- Predict PerformanceScoreOutOf100 or MonthlyIncomeUSD using Age, YearsAtCompany, TrainingHoursLastYear, TechnicalSkillScore, SoftSkillScore.
- StatsModels can be used for detailed statistical output (R-squared, p-values, coefficients).

### 2. MLR Model Evaluation:

- Calculate R-squared, Adjusted R-squared, MSE, MAE for the regression models.

### 3. Dealing with Categorical Variables & One-Hot Encoding (OHE):

- Department, Gender, EducationLevel, JobRole, OfficeLocation will require OHE for linear models.
- JobRole having many categories will demonstrate the increase in dimensionality.

### 4. Interpreting One-Hot Encoded Coefficients:

- Analyze how belonging to a specific Department (e.g., Engineering vs. a baseline) affects salary or performance score, after OHE.

### 5. Regression Diagnostics - Introduction:

- Examine residuals for patterns (heteroscedasticity).
- Check for multicollinearity (e.g., between YearsAtCompany and TotalWorkingYears).

- Identify outliers or influential points (e.g., an employee with exceptionally high `OverTimeHoursLastMonth` or very low `SatisfactionScore`).

#### 6. Linear & Log Transformations:

- `PreviousSalaryUSD` and `MonthlyIncomeUSD` are likely right-skewed; apply log transformations and observe improvements in model assumptions (e.g., normality of residuals).
- Create a new feature like  $\text{AvgSkillScore} = (\text{TechnicalSkillScore} + \text{SoftSkillScore}) / 2$  (linear transformation).

#### 7. Object-Oriented Programming (OOP) - Introduction, Classes, Instances, Methods, Initialization, Inheritance:

- This is a programming concept, but the scikit-learn section below relies heavily on it. You can create custom transformer classes.

#### 8. OOP with Scikit-Learn & Preprocessing with scikit-learn:

- Use scikit-learn classes like `StandardScaler`, `MinMaxScaler`, `OneHotEncoder`, `SimpleImputer`. These are all objects instantiated from classes with methods like `.fit()`, `.transform()`.
- Build Pipeline objects.

#### 9. Machine Learning Fundamentals-Introduction, Data Science Processes, Statistical Learning Theory:

- The entire workflow of using this dataset (from data loading, cleaning, preprocessing, modeling, to evaluation) covers these.

#### 10. Data Ethics-Training On Bad Data:

- The Gender feature can be used to discuss fairness. If the historical data (synthetically generated) shows Gender correlating with `MonthlyIncomeUSD` even after accounting for other factors, it indicates bias. Training a model on this could perpetuate it.
- Similarly for Age.

#### 11. Regression Model Validation, Introduction to Cross-Validation, Bias-Variance Tradeoff:

- Split data into train/test.

- Implement k-fold cross-validation for both regression (PerformanceScoreOutOf100) and classification (LeftCompanyWithinYear).
- Discuss how model complexity (e.g., number of features, polynomial degree in regression, tree depth) affects bias and variance.

## 12. Ridge and Lasso Regression, Feature Selection Methods:

- With many features (especially after OHE of JobRole and Department), use Ridge (L2) and Lasso (L1) to predict PerformanceScoreOutOf100 or MonthlyIncomeUSD.
- Observe how Lasso performs feature selection by shrinking some coefficients to zero.
- Compare with other feature selection methods (e.g., recursive feature elimination, select K best).

## 13. Logistic Regression-Introduction, Linear to Logistic, Fitting, scikit-learn:

- Predict LeftCompanyWithinYear (binary target) using a mix of numerical and OHE categorical features.

## 14. MLE Review, MLE and Logistic Regression, Gradient Descent Review, Applying Gradient Descent - Lab, Coding Logistic Regression From Scratch- Lab:

- The dataset provides the features and target for students to understand the mathematical underpinnings of logistic regression, and to implement it from scratch or observe its iterative optimization.

## 15. Classification Metrics-Introduction, Confusion Matrices, Visualizing-Lab, Evaluation Metrics, Evaluating Logistic Regression Models Lab:

- For the LeftCompanyWithinYear prediction:
  - Generate confusion matrices.
  - Calculate Accuracy, Precision, Recall, F1-score, Specificity.

## 16. ROC Curves and AUC, ROC Curves and AUC-Lab:

- Plot ROC curves and calculate AUC for the churn prediction model.

## 17. Class Imbalance Problems, Class Imbalance Problems-Lab:

- Ensure LeftCompanyWithinYear has an imbalanced distribution (e.g., 80% "No", 20% "Yes").



- Demonstrate how accuracy can be misleading.
- Apply techniques like oversampling (SMOTE), undersampling, or using class weights in models.

#### **18. Logistic Regression Model Comparisons -Lab:**

- Compare logistic regression models with different sets of features, regularization strengths, or after applying class imbalance techniques.

#### **19. Decision Trees -Introduction, Entropy and Information Gain, ID3 Classification Trees: Perfect Split-Lab:**

- Use LeftCompanyWithinYear as the target.
- Features like SatisfactionScore, OverTimeHoursLastMonth, Department, LastPerformanceRating would be good candidates for demonstrating information gain and splits.

#### **20. Building Trees using scikit-learn, Building Trees using scikit-learn Lab:**

- Implement DecisionTreeClassifier for churn and DecisionTreeRegressor for performance/salary.

#### **21. Hyperparameter Tuning and Pruning In Decision Trees, Hyperparameter Tuning and Pruning in Decision Trees-Lab:**

- Tune max\_depth, min\_samples\_split, min\_samples\_leaf, ccp\_alpha for both classification and regression trees.

#### **22. Regression with CART Trees, Regression with CART Trees-Lab, Regression Trees and Model Optimization-Lab:**

- Predict PerformanceScoreOutOf100 or MonthlyIncomeUSD using decision trees.
- Optimize these regression trees.

---

### **III. Generating Synthetic Data - Considerations:**

- **Correlations:** Intentionally introduce correlations:
  - **Positive:** YearsAtCompany and MonthlyIncomeUSD, TechnicalSkillScore and PerformanceScoreOutOf100.

- Negative: SatisfactionScore and LeftCompanyWithinYear, OverTimeHoursLastMonth and WorkLifeBalanceScore.
- Skewness: Make features like MonthlyIncomeUSD, PreviousSalaryUSD, OverTimeHoursLastMonth skewed.
- Outliers: Introduce a few realistic outliers (e.g., an employee with very high training hours but average performance).
- Missing Values (Optional): Introduce a small percentage of missing values in features like PreviousSalaryUSD or AvgDailyCommuteTimeMin to cover imputation techniques.
- Interaction Terms: Design it so that interactions might be present (e.g., the effect of TrainingHoursLastYear on PerformanceScoreOutOf100 might differ across Departments).
- Class Imbalance for LeftCompanyWithinYear: Ensure roughly 15-25% of employees are marked as 1 (Yes).

This dataset schema provides a comprehensive playground. You would then need to synthetically generate data points that adhere to these characteristics, perhaps using Python libraries like NumPy and Pandas, and introducing relationships using statistical functions or rules.