

Regression Techniques and an Applied Laboratory in Logistic Classification

Part I: The Theory of Regression Analysis

Section 1: Foundational Concepts in Predictive Modeling

1.1 Defining Regression: Predicting Continuous Variables

Regression analysis stands as a foundation of supervised machine learning, a statistical technique engineered to model and study the relationship between variables. The primary objective of regression is to predict a continuous, numerical outcome, referred to as the dependent variable, based on the values of one or more input features, known as independent variables.² This predictive capability allows data scientists and analysts to move beyond historical analysis and into the realm of forecasting. The **core function of a regression model is to establish a mathematical relationship, often denoted as $Y=f(X)$** , where Y represents the continuous target variable and X signifies the set of predictor variables. By learning this function from a labeled dataset, the model can then predict the outcome for new, unseen data.

The applications of regression are vast and span numerous domains. For instance, in finance, **regression models can predict stock prices or assess risk**. In real estate, they are used to estimate the **price of a house based on its attributes, such as size, location, and the number of bedrooms**. In marketing, regression can *forecast future sales based on parameters like advertising expenditure* and market trends. Other common examples include predicting a person's age, estimating the impact of **academic scores on university admissions**, or even **forecasting weather conditions**. In all these scenarios, the goal is to predict a value on a continuous spectrum, which is the defining characteristic of a regression task.

1.2 Dependent vs. Independent Variables

To formalize the study of regression, it is essential to establish a clear and consistent terminology for the variables involved. The structure of a regression problem is built around two primary types of variables.

The **dependent variable** is the main variable of interest that we aim to understand, explain, or predict. It is the effect or the outcome whose variation is being studied. In literature, it is referred to by many names, including the response variable, target variable, or outcome

variable. In a graphical representation, **the dependent variable is conventionally plotted on the vertical y-axis**. For example, in a model predicting student exam scores based on hours studied, the "exam score" is the dependent variable because it is hypothesized to depend on the amount of study time.

The **independent variable(s)** are the factors that are believed to influence the dependent variable. These are the **inputs to the model, which are manipulated or observed to see their effect on the outcome**. They are also known by several synonyms, such as *explanatory variables*, *predictor variables*, *features*, or *covariates*. In a two-dimensional plot, the independent variable is placed on the horizontal x-axis. In the study time example, "hours studied" is the independent variable. A regression analysis can involve a single independent variable (simple regression) or multiple independent variables (multiple regression). The fundamental regression equation, in its simplest form, is often written as

$Y = \beta_0 + \beta_1 X + \epsilon$, where Y is the dependent variable, X is the independent variable, β_0 and β_1 are the model parameters (coefficients), and ϵ is the error term representing the variability in Y not explained by X .

1.3 Regression vs. Classification

Within the domain of supervised machine learning, problems are categorized into two types: *regression and classification*. The distinction between them is fundamental and hinges on the nature of the output variable the model is designed to predict.

Regression models are tasked with predicting a continuous quantity. The output is a numerical value that can exist anywhere within a given range. For example, a regression model might predict the exact temperature tomorrow in degrees Celsius, the price of a stock in dollars, or the amount of rainfall in millimeters. **The objective of the regression algorithm is to learn a mapping function that best fits a line or curve to the data points, representing the underlying trend.**

Conversely, **classification models are designed to predict a discrete class label or category.** The output is not a continuous value but rather a membership to a specific group. For instance, a classification model might predict whether an email is "spam" or "not spam," or classify a medical image as showing a "malignant" or "benign" tumor. **The objective of a classification algorithm is to find a decision boundary—a line, plane, or hyperplane—that effectively separates the data points into their respective classes.**

While distinct, these two tasks share similarities. Both are forms of supervised learning, meaning they require a labeled dataset for training. Furthermore, the lines can sometimes blur; a regression algorithm can predict an integer quantity (a discrete value), and a classification algorithm can output a continuous probability of class membership. The key

differentiator remains the ultimate nature of the target variable: continuous for regression and categorical for classification.

Table 1: Regression vs. Classification - A Comparative Overview.

Feature	Regression	Classification
Output Type	Continuous (e.g., price, temperature, age)	Discrete/Categorical (e.g., spam/not spam, cat/dog, yes/no)
Primary Goal	To predict a numerical value and model the relationship between variables.	To predict a class label and categorize data into predefined groups.
Model Output	A continuous value (e.g., 25.4, 150000.75)	A class label (e.g., 1, 0, 'spam') or a probability of class membership.
Algorithmic Aim	To find a best-fit line or curve that minimizes the error between predictions and actual values.	To find a decision boundary that separates the different classes in the feature space.
Example Question	"How much will this house sell for?"	"Will this customer churn or not?"
Evaluation Metrics	Mean Squared Error (MSE), R-Squared, Mean Absolute Error (MAE)	Accuracy, Precision, Recall, F1-Score, AUC-ROC

1.4 The Bias-Variance Tradeoff: A Conceptual Framework for Model Complexity

The central challenge in all predictive modeling is navigating the bias-variance tradeoff. This concept provides a framework for understanding model error and guides the process of model selection and tuning. Total error in a model can be decomposed into three components: bias, variance, and irreducible error. While irreducible error is due to inherent noise in the data that cannot be modeled, bias and variance are sources of error that can be managed by the data scientist.

Bias refers to the error introduced by approximating a real-world problem, which may be complex, with a much simpler model. It represents the systematic error of the model, or the

*difference between the average prediction of our model and the correct value which we are trying to predict. **A model with high bias pays little attention to the training data and oversimplifies the underlying patterns.*** This leads to **underfitting**, where the model performs poorly on both the training data and unseen test data because it fails to capture the true relationship between the variables.

Variance refers to the error introduced by the model's sensitivity to small fluctuations in the training data. It measures how much the model's predictions would change if it were trained on a different training dataset. *A model with high variance pays too much attention to the training data, modeling not only the underlying patterns but also the noise.* This leads to **overfitting**, where the model performs exceptionally well on the training data but poorly on unseen test data because it has effectively "memorized" the training set instead of learning a generalizable pattern.

The relationship between bias and variance is typically inverse. As a model's complexity increases, its bias tends to decrease, but its variance tends to increase. A very simple model (e.g., simple linear regression on a complex, non-linear dataset) will have high bias and low variance. A very complex model (e.g., a high-degree polynomial regression) will have low bias but high variance. *The ultimate goal is to find a model that strikes an optimal balance, minimizing the total error by achieving both low bias and low variance.* Every decision in the modeling process—from choosing an algorithm to tuning its hyperparameters—is an act of managing this fundamental tradeoff. The entire process is not about finding a "perfect" model that eliminates all error, but rather about finding the optimal balance of these two unavoidable sources of error for a given dataset and problem.

Section 2: Linear Regression

Linear regression is one of the simplest and most widely used regression algorithms, serving as a foundational technique in both statistics and machine learning. It is often the first model to be considered for a regression problem due to its simplicity, interpretability, and computational efficiency. The core assumption of linear regression is that a linear relationship exists between the independent and dependent variables.

2.1 Simple Linear Regression: Modeling the Relationship Between Two Variables

Simple linear regression is the most basic form of linear regression, used when there is a single independent variable (x) to predict a single dependent variable (y). The goal is to find the best-fitting straight line through the data points on a scatter plot, which mathematically represents the relationship between the two variables.

The equation for the simple linear regression line is analogous to the standard equation of a line from algebra:

$$\hat{y} = \beta_0 + \beta_1 x$$

Here, \hat{y} (read "y-hat") is the predicted value of the dependent variable. The two parameters, β_0 and β_1 , define the line ⁶:

- **Intercept (β_0):** This is the predicted value of y when the independent variable x is equal to zero. It is the point where the regression line crosses the y -axis.⁹ The practical interpretation of the intercept depends heavily on the context of the data. For instance, in predicting a person's weight based on their height, the intercept would represent the predicted weight of a person with zero height, which is a nonsensical extrapolation. In such cases, the intercept serves primarily as a mathematical anchor for the line rather than a meaningful standalone prediction.¹⁹
- **Slope (β_1):** This is the most critical component for interpretation. The slope coefficient represents the average change in the dependent variable y for a one-unit increase in the independent variable x .⁹ A positive slope indicates a positive relationship (as x increases, y tends to increase), while a negative slope indicates a negative relationship (as x increases, y tends to decrease).²³ The magnitude of the slope quantifies the strength of this effect.

To determine the "best-fit line," the most common method is **Ordinary Least Squares (OLS)**. The principle of OLS is to find the values of β_0 and β_1 that minimize the total prediction error.⁹ This error is quantified by calculating the

residuals, which are the vertical distances between each observed data point (y_i) and the value predicted by the regression line (\hat{y}_i).⁶ OLS specifically minimizes the

sum of the squared residuals (SSR), also known as the residual sum of squares (RSS).⁴

$$SSR = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

The residuals are squared for two primary reasons: first, it ensures that both positive and negative errors (points above and below the line) contribute positively to the total error, preventing them from canceling each other out. Second, squaring the residuals penalizes larger errors more heavily than smaller ones, effectively forcing the model to pay more attention to significant deviations.⁴ The optimal values for

β_0 and β_1 are found using calculus, by taking the partial derivatives of the SSR with respect to each parameter, setting them to zero, and solving the resulting system of equations.²²

2.2 Multiple Linear Regression: Expanding to Multiple Predictors

Simple linear regression is often insufficient for real-world problems, where an outcome is typically influenced by multiple factors. **Multiple linear regression** extends the model to accommodate two or more independent variables.² The equation becomes a linear

combination of all predictor variables:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

In this model, p is the number of predictors, and each predictor x_j has its own corresponding slope coefficient β_j .⁴

The interpretation of the coefficients in a multiple regression model is more nuanced and powerful than in a simple model. The coefficient β_j represents the expected change in the dependent variable y for a one-unit increase in the predictor variable x_j , **while holding all other independent variables in the model constant**.²³ This allows analysts to isolate the unique contribution of each predictor to the outcome, controlling for the effects of the other variables. For example, in a model predicting house prices based on size and age, the coefficient for size would estimate the change in price for an additional square foot, assuming the age of the house remains the same.

A significant challenge in multiple linear regression is **multicollinearity**. This issue arises when two or more independent variables in the model are highly correlated with each other.³ When multicollinearity is present, it becomes difficult for the model to disentangle the individual effects of the correlated predictors. This can lead to unstable and unreliable coefficient estimates, where small changes in the data can cause large fluctuations in the coefficient values. While multicollinearity does not necessarily reduce the overall predictive accuracy of the model, it severely undermines the interpretability of the individual coefficients.³ This problem is a key motivation for using more advanced techniques like regularized regression.

2.3 Critical Assumptions of Linear Models

Linear regression is a parametric model, meaning it makes several assumptions about the data. For the model's results to be valid, reliable, and unbiased, these assumptions must be reasonably met. Violations of these assumptions can lead to misleading conclusions.⁹

Assumption	Description	How to Check
Linearity	The relationship between the independent variables (X) and the mean of the dependent variable (Y) is linear. The model can only capture straight-line relationships. ⁹	Scatter plots of each predictor vs. the target variable. Residual plots (residuals vs. fitted values) should show no discernible pattern.

Independence	The errors (residuals) are independent of each other. The error of one observation does not influence the error of another. This is often a concern with time-series data. ⁹	Durbin-Watson test. Residual plots against time or observation order should show no patterns (e.g., autocorrelation).
Homoscedasticity	The variance of the errors is constant across all levels of the independent variables. This is also known as constant variance. ⁹	Residual plots (residuals vs. fitted values) should show a random scatter of points with constant spread. Funnel shapes indicate heteroscedasticity.
Normality	The errors (residuals) of the model are normally distributed with a mean of zero. This is important for hypothesis testing and constructing confidence intervals. ²⁶	Q-Q (Quantile-Quantile) plots of the residuals should form a straight line. Histograms of residuals should appear bell-shaped. Statistical tests like Shapiro-Wilk can also be used.

Table 2: Key Assumptions of Linear Regression Models.

Section 3: Modeling Curvature and Complexity

While linear models are powerful, they are constrained by the assumption of a linear relationship. Many real-world phenomena exhibit non-linear trends. Polynomial regression provides a straightforward way to extend linear regression to model these more complex, curved relationships.³

3.1 Polynomial Regression: Capturing Non-Linear Relationships

Polynomial regression models the relationship between the independent variable x and the dependent variable y as an n -th degree polynomial in x .³³ Instead of fitting a straight line, it fits a curve to the data. This is achieved by adding new features to the model that are polynomial transformations of the original feature(s), such as the square (

x^2), cube (x^3), and so on.⁵

For example, a second-order (quadratic) polynomial regression model with one independent

variable takes the form:

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2$$

This equation can capture a simple parabolic curve.⁵ By increasing the degree of the polynomial, the model can fit increasingly complex curves to the data. For example, a third-order (cubic) model would be

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3.$$

A crucial conceptual point is that despite its ability to model non-linear relationships, polynomial regression is still considered a special case of **multiple linear regression**.³³ The term "linear" in this context refers to the model being linear in its

parameters (β_j), not necessarily in its variables (x). The model $\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2$ is non-linear with respect to the variable x . However, if we create new features, let's say $x_1 = x$ and $x_2 = x^2$, the equation becomes $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. This is the standard form of a multiple linear regression equation. Because the model is a linear combination of its coefficients, the same estimation techniques, such as Ordinary Least Squares (OLS), can be used to find the optimal parameter values.³³

3.2 The Risk of Overfitting

The flexibility of polynomial regression comes at a cost: a high risk of **overfitting**. As the degree of the polynomial increases, the model's complexity and flexibility grow, allowing it to fit the training data more closely.³⁷ A very high-degree polynomial can create a curve that passes through or very near every single point in the training set. While this results in a very low error on the training data, the model is likely capturing not just the underlying trend but also the random noise inherent in the data.⁵

When this overfitted model is exposed to new, unseen data, its performance is typically very poor. It has failed to generalize the true relationship and instead has "memorized" the idiosyncrasies of the training set.¹⁸ This is a classic manifestation of the high-variance problem discussed in the bias-variance tradeoff.

Visually, this can be demonstrated by plotting models of different degrees on the same data. A first-degree polynomial (a straight line) might underfit the data if the true relationship is curved. A second or third-degree polynomial might capture the trend well. However, a tenth-degree polynomial would likely produce a highly convoluted curve that wiggles erratically to hit every training point, leading to nonsensical predictions for new data points that fall between the training examples.³⁷ The challenge, therefore, is to select a polynomial degree that is complex enough to capture the underlying trend without being so complex that it models the noise. This dilemma directly motivates the need for regularization techniques, which provide a more systematic way to control model complexity.

Section 4: Advanced Models - Regularization for Robustness and Feature Selection

When a model becomes too complex, such as a multiple regression model with many predictors or a high-degree polynomial regression, it is prone to overfitting. Regularization is a class of techniques designed to combat this problem by adding a penalty for model complexity to the optimization process, thereby promoting simpler, more generalizable models.¹⁸

4.1 The Philosophy of Regularization

The core idea of regularization is to modify the model's cost function. In standard linear regression, the goal is to minimize the Residual Sum of Squares (RSS). Regularization introduces an additional component, a **penalty term**, that constrains the magnitude of the model's coefficients.³⁹ The new objective function becomes:

Minimize: $(RSS + \lambda \times \text{Penalty Term})$

In this formulation, the penalty term is a function of the regression coefficients (β_j). By penalizing large coefficient values, regularization discourages the model from assigning excessive importance to any single feature, thus reducing its complexity and variance.³⁹

The **regularization parameter**, denoted by λ (lambda) or α (alpha), is a hyperparameter that controls the strength of this penalty.³⁹

- If $\lambda=0$, the penalty term has no effect, and the model is equivalent to a standard OLS regression.
- As λ increases, the penalty becomes more influential, forcing the coefficients to shrink towards zero. This increases the model's bias but decreases its variance.
- If λ is very large, the penalty dominates, and all coefficients will be pushed very close to (or exactly to) zero, resulting in a highly biased, underfit model.⁴⁰

The optimal value of λ is typically determined through cross-validation, finding the value that yields the best performance on unseen data.³⁹

4.2 Ridge Regression (L2 Regularization)

Ridge regression is a regularization technique that uses the **L2 norm** as its penalty term. The L2 norm is the sum of the squared magnitudes of the coefficients.³ The cost function for Ridge regression is:

\$\$ \text{Minimize: } \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right) \$\$

The effect of the L2 penalty is to shrink the coefficients of the model towards zero. However, due to the nature of the squared term, the penalty for a large coefficient is substantial, while the penalty for a small coefficient is minimal. As a result, Ridge regression will make large coefficients smaller but will **never shrink any coefficient exactly to zero**.⁴¹ All features are retained in the final model, although their influence may be significantly reduced.

The primary use case for Ridge regression is to handle **multicollinearity**.⁴³ When independent variables are highly correlated, OLS estimates can be highly variable. Ridge regression stabilizes these estimates by shrinking the coefficients of correlated predictors together, distributing their effect more evenly.⁴⁵

4.3 Lasso Regression (L1 Regularization)

Lasso regression, which stands for Least Absolute Shrinkage and Selection Operator, employs the **L1 norm** as its penalty term. The L1 norm is the sum of the absolute values of the coefficients.⁴⁰ The cost function for Lasso regression is:

$$\text{Minimize: } \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j| \right)$$

The key difference from Ridge lies in the absolute value penalty. This form of penalty allows Lasso to shrink some coefficients **exactly to zero** when λ is sufficiently large.³⁹ This property is a direct consequence of the geometry of the L1 constraint, which is diamond-shaped and has "corners" at the axes, making it possible for the optimal solution to lie on an axis where one or more coefficients are zero.⁴⁷

This ability to zero out coefficients makes Lasso regression an extremely powerful tool for **automatic feature selection**.⁴² By eliminating irrelevant or redundant features from the model, Lasso produces a

sparse model—one with fewer parameters. Sparse models are often easier to interpret and can be more computationally efficient.⁴⁰

4.4 Elastic Net Regression: A Hybrid Approach

Elastic Net regression was developed to combine the strengths of both Ridge and Lasso regression.¹⁸ It includes both the L1 and L2 penalty terms in its cost function, controlled by a mixing parameter,

α (often referred to as `l1_ratio`), in addition to the overall regularization strength, λ .

$$\text{Minimize: } \left(\text{RSS} + \lambda \left[(1-\alpha) \sum_{j=1}^p \beta_j^2 + \alpha \sum_{j=1}^p |\beta_j| \right] \right)$$

- When $\alpha=0$, the model is equivalent to Ridge regression.

- When $\alpha=1$, the model is equivalent to Lasso regression.
- For $0<\alpha<1$, the model is a hybrid of the two.⁴⁹

Elastic Net is particularly useful in scenarios where there are multiple highly correlated features. In such cases, Lasso tends to arbitrarily select one feature from the group and eliminate the others. Ridge, on the other hand, would shrink the coefficients of the correlated group together. Elastic Net bridges this gap; it can perform feature selection like Lasso but also exhibits a "grouping effect," where it tends to select or discard a group of correlated variables together.⁴⁵ This often leads to better performance and stability in high-dimensional datasets with correlated predictors.

The choice between these regularization methods often reflects a prior assumption about the nature of the data. Opting for Ridge suggests a belief that most features are relevant and contribute to the outcome. Choosing Lasso implies a belief that the underlying signal is sparse, meaning only a subset of features are truly important. Elastic Net is chosen when there is uncertainty or a belief that features are correlated in meaningful groups. This elevates the decision from a purely technical one to a strategic choice informed by domain knowledge.

Feature	Ridge Regression (L2)	Lasso Regression (L1)	Elastic Regression Net
Penalty Term	Sum of squared coefficients ($\lambda \sum \beta_j^2$)	Sum of absolute coefficients ($\lambda \sum \beta_j $)	$\lambda \sum \beta_j^2$
Effect on Coefficients	Shrinks coefficients towards zero.	Can shrink coefficients exactly to zero.	Can shrink coefficients exactly to zero.
Feature Selection	No. All features are retained in the model.	Yes. Performs automatic feature selection, creating sparse models.	Yes. Performs feature selection.
Handling Multicollinearity	Very effective. Shrinks coefficients of correlated variables together.	Can be unstable. May arbitrarily select one variable from a correlated group.	Effective. Exhibits a grouping effect, selecting or dropping correlated variables together.

Primary Use Case	When most predictors are useful and multicollinearity is present.	When you suspect many predictors are irrelevant and want a simpler, more interpretable model.	When there are many correlated predictors and you want the benefits of both Ridge and Lasso.
-------------------------	---	---	--

Table 3: Comparison of L1, L2, and Elastic Net Regularization.

Section 5: Evaluating Regression Model Performance

After training a regression model, it is crucial to evaluate its performance to understand how well it predicts outcomes on unseen data. Several metrics are used for this purpose, each providing a different perspective on the model's accuracy and goodness of fit.⁵²

5.1 Quantifying Error: Metrics on the Scale of the Target Variable

These metrics measure the average error between the predicted values (\hat{y}_i) and the actual values (y_i). They are expressed in the same units as the target variable (or its square), making them directly interpretable in the context of the problem.

- **Mean Absolute Error (MAE):** MAE calculates the average of the absolute differences between the predictions and the actual values. It provides a straightforward measure of the average magnitude of the errors.⁵³

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Because it does not square the errors, MAE is less sensitive to outliers than other metrics like MSE. It treats all errors equally in the average.¹⁶

- **Mean Squared Error (MSE):** MSE is one of the most common metrics, calculating the average of the squared differences between predictions and actual values.⁵⁶

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

The squaring of the error term has two important consequences: it ensures all errors are positive, and it gives significantly more weight to larger errors (outliers).¹⁶ This property makes MSE a good choice when large errors are particularly undesirable. Additionally, its mathematical properties (being differentiable) make it a preferred loss function for optimization algorithms like gradient descent.⁵⁶

- **Root Mean Squared Error (RMSE):** RMSE is simply the square root of the MSE.

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

The primary advantage of RMSE over MSE is that its units are the same as the target variable's units, making it more intuitive to interpret.¹⁶ For example, if predicting house prices in dollars, an RMSE of 20,000 means the model's predictions are, on average, about \$20,000 away from the actual prices.

5.2 Explaining Variance: The Coefficient of Determination (R-Squared)

R-Squared, or the coefficient of determination, is a statistical measure that represents the proportion of the variance in the dependent variable that is explained by the independent variables in the model.⁵ It provides a measure of the "goodness of fit" of the model.

The formula for R-Squared is:

$$R^2 = 1 - \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

where SSR is the Sum of Squared Residuals (unexplained variance) and SST is the Total Sum of Squares (total variance in the data).⁵³

R-Squared values range from 0 to 1.

- An R² of 1 indicates that the model perfectly explains all the variability of the response data around its mean.
- An R² of 0 indicates that the model explains none of the variability.
- An R² of 0.82, for example, means that 82% of the variance in the target variable can be explained by the model's inputs.⁵

However, R-Squared has a significant flaw: its value will always increase or stay the same whenever a new predictor is added to the model, even if that predictor is completely irrelevant.⁵³ This can be misleading and may encourage the creation of overly complex models (overfitting).

5.3 A More Honest Metric: Adjusted R-Squared

To address the limitation of R-Squared, the Adjusted R-Squared metric was developed. It modifies the R-Squared value to account for the number of predictors in the model.⁶⁰

The formula for Adjusted R-Squared is:

$$R_{adj}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$$

where n is the number of observations and p is the number of predictors.⁵³

Adjusted R-Squared effectively penalizes the model for adding predictors that do not improve the model's explanatory power more than would be expected by chance.⁵⁸ Unlike R-Squared, the Adjusted R-Squared value can decrease if a new, non-significant predictor is added to the model. This makes it a much more reliable metric for comparing the goodness of fit of models

with different numbers of independent variables, as it favors more parsimonious models that achieve good performance with fewer features.⁵⁵ A higher Adjusted R-Squared generally indicates a better model.

Part II: Applied Laboratory - Logistic Regression for Binary Classification

Section 6: Introduction to Logistic Regression

While the first part of this lecture focused on predicting continuous values, many data mining problems require predicting categorical outcomes. Logistic Regression is a fundamental and widely used algorithm for such classification tasks, particularly for binary classification where the outcome has only two possible classes (e.g., Yes/No, True/False, 1/0).⁶²

6.1 The Conceptual Bridge: Why a "Regression" Technique is Used for Classification

A common point of confusion is the name "Logistic Regression," as it is fundamentally a classification algorithm.⁶⁴ The name arises because the underlying technique is an extension of linear regression.⁶⁵ Logistic regression does not directly predict a discrete class label. Instead, it builds a regression model to predict a continuous outcome: the

probability that a given instance belongs to the positive class (usually denoted as class '1').⁶⁶

This predicted probability, which is a value between 0 and 1, is then transformed into a discrete class prediction by applying a **decision threshold**. The most common threshold is 0.5: if the predicted probability is greater than or equal to 0.5, the instance is classified as '1'; otherwise, it is classified as '0'.⁶³ So, while its application is classification, its mechanism is rooted in regressing for a probability, hence the name.⁶⁴

6.2 From Linear Output to Probability: The Role of the Sigmoid Function

A standard linear regression model, $z = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$, produces an output that can range from negative infinity to positive infinity. This is unsuitable for modeling a probability, which must be constrained between 0 and 1.⁷¹

To solve this, logistic regression passes the output of the linear equation through a special function called the **Sigmoid function**, also known as the logistic function.⁶² The Sigmoid

function takes any real-valued number as input and maps it to a value between 0 and 1.⁷¹

The formula for the Sigmoid function is:

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

where z is the output of the linear model. The resulting probability is thus:

$$P(Y=1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)}}$$

The Sigmoid function has a characteristic "S-shaped" curve. As the linear output z approaches positive infinity, e^{-z} approaches 0, and the probability $\sigma(z)$ approaches 1. As z approaches negative infinity, e^{-z} approaches infinity, and the probability $\sigma(z)$ approaches 0. When $z=0$, the probability is exactly 0.5.⁷⁰ This elegant transformation allows the linear model to produce valid probability estimates.

6.3 Understanding the Logit: The Mathematics of Log-Odds

To interpret the coefficients (β_j) of a logistic regression model, one must understand the transformation that links the linear equation to the probability. This involves the concepts of odds and log-odds.

- **Probability:** The likelihood of an event occurring, $p = P(Y=1)$.
- **Odds:** The ratio of the probability of an event occurring to the probability of it not occurring.

$$\text{Odds} = \frac{p}{1-p}$$

While probability ranges from 0 to 1, odds can range from 0 to infinity.⁷⁴

- **Log-Odds (or Logit):** The natural logarithm of the odds.

$$\text{Logit}(p) = \ln\left(\frac{p}{1-p}\right)$$

The logit transformation is crucial because it maps a probability value from the range to the entire real number line $[-\infty, +\infty]$, which matches the output range of a linear equation.⁷⁵

By algebraically rearranging the sigmoid function, we can see the fundamental relationship in logistic regression:

$$\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

This equation reveals that logistic regression is a linear model for the log-odds of the

outcome.⁶⁶ This provides a direct way to interpret the coefficients. A one-unit increase in a predictor variable x_j is associated with an increase of β_j in the log-odds of the outcome, holding all other predictors constant.⁷⁶

For a more intuitive interpretation, we can exponentiate the coefficient, e^{β_j} . This value is the **odds ratio**. It represents the multiplicative factor by which the odds of the outcome change for a one-unit increase in x_j . For example, if $\beta_j=0.693$, the odds ratio is $e^{0.693}\approx 2$. This means that a one-unit increase in x_j doubles the odds of the outcome occurring.⁷⁵

6.4 The Decision Boundary

The **decision boundary** is the surface (a line in 2D, a plane in 3D, or a hyperplane in higher dimensions) that separates the feature space into regions where the model predicts different classes.¹³

Using a standard threshold of 0.5, the model predicts class '1' when $P(Y=1|X)\geq 0.5$ and class '0' when $P(Y=1|X)< 0.5$. The decision boundary is the point where the probability is exactly 0.5. Looking at the sigmoid function, this occurs precisely when its input, z , is equal to 0.⁷¹

Therefore, the equation of the decision boundary is given by setting the linear part of the model to zero:

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p = 0$$

This is the equation of a hyperplane.⁷⁸ For a model with two features, x_1 and x_2 , the decision boundary is the line $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$. All points on one side of this line will be classified as '1', and all points on the other side will be classified as '0'. It is important to note that logistic regression inherently produces a linear decision boundary. To model non-linear separation, one would need to add polynomial or interaction terms as features, similar to polynomial regression.⁷⁹

Section 7: Practical Implementation: Predicting Survival on the Titanic

For this laboratory exercise, we will implement a logistic regression model to tackle a classic binary classification problem: predicting passenger survival on the RMS Titanic.

7.1 Dataset Introduction: The Titanic Dataset

The Titanic dataset, widely available on platforms like Kaggle, provides data on the passengers of the Titanic, which sank in 1912.⁸² The objective is to build a predictive model

that answers the question: "what sorts of people were more likely to survive?".⁸² This is a supervised learning task, as the training data includes the ground truth for survival.

The dataset contains several features about each passenger. For our model, we will focus on a few key predictors ⁸²:

- **Survived:** The target variable (0 = No, 1 = Yes).
- **Pclass:** Ticket class (1 = 1st, 2 = 2nd, 3 = 3rd), a proxy for socio-economic status.
- **Sex:** The gender of the passenger.
- **Age:** The age of the passenger in years.
- **SibSp:** The number of siblings or spouses aboard the Titanic.
- **Parch:** The number of parents or children aboard the Titanic.
- **Fare:** The fare paid by the passenger.

This problem is an ideal application for logistic regression because the outcome variable is binary (survived or not), and we are interested in modeling the probability of survival based on a set of predictor variables.⁸⁶

7.2 Model Training: Implementing Logistic Regression

The following section provides a step-by-step guide to building the logistic regression model using R.

Step 1: Load Data and Libraries

First, we import the necessary libraries and load the training data

```
# Install packages if you haven't already
# install.packages(c("dplyr", "caTools", "caret", "pROC", "ggplot2"))

# Load necessary libraries
library(dplyr)
library(caTools)
library(caret)
library(pROC)
library(ggplot2)

# Load the dataset (assuming 'train.csv' is in your working directory)
titanic_df <- read.csv('Titanic-Dataset.csv')
cat("--- 1. Data Loading Complete ---\n")
print(head(titanic_df, 3))
```

Step 2: Data Preparation

Real-world datasets often require cleaning and preprocessing. While this lecture assumes prior knowledge of these steps, some minimal preparation is essential for the model to function. The 'Age' column has missing values, which we will fill with the median age. The 'Sex' column is categorical and must be converted to a numerical format.⁸⁸ We will also select our features and define the target variable.

```
# Step 2: Data Preprocessing and Exploration
cat("\n--- 2. Starting Data Preprocessing ---\n")

# Impute missing 'Age' values with the median. This is a robust measure against outliers.
median_age <- median(titanic_df$Age, na.rm = TRUE)
titanic_df$Age[is.na(titanic_df$Age)] <- median_age
cat("Missing 'Age' values imputed with median:", median_age, "\n")

# Impute missing 'Embarked' values with the mode (the most common port).
# First, find the mode. 'S' is overwhelmingly the most common.
embarked_mode <- names(sort(table(titanic_df$Embarked), decreasing = TRUE))[1]
titanic_df$Embarked[titanic_df$Embarked == ""] <- NA
titanic_df$Embarked[is.na(titanic_df$Embarked)] <- embarked_mode
cat("Missing 'Embarked' values imputed with mode:", embarked_mode, "\n")

# Convert variables to the correct factor data type for the model
titanic_df$Survived <- as.factor(titanic_df$Survived)
titanic_df$Pclass <- as.factor(titanic_df$Pclass)
titanic_df$Sex <- as.factor(titanic_df$Sex)
titanic_df$Embarked <- as.factor(titanic_df$Embarked)

# Select the final, cleaned features for our model. We exclude identifiers like Name and Ticket.
titanic_clean <- titanic_df %>%
  select(Survived, Pclass, Sex, Age, SibSp, Parch, Fare, Embarked)
cat("Data cleaning and feature selection complete.\n")
```

Step 3: Splitting the Data

To evaluate our model's performance on unseen data, we split the dataset into a training set and a testing set. A common split is 80% for training and 20% for testing.⁹⁰

```
# Step 3: Split Data into Training and Testing Sets
cat("\n--- 3. Splitting Data (80% Train, 20% Test) ---\n")
set.seed(123) # for reproducible results
split <- sample.split(titanic_clean$Survived, SplitRatio = 0.8)
train_set <- subset(titanic_clean, split == TRUE)
test_set <- subset(titanic_clean, split == FALSE)
```

```
cat("Training set size:", nrow(train_set), "\n")
cat("Testing set size:", nrow(test_set), "\n")
```

Step 4: Model Instantiation and Training

We now create an instance of the LogisticRegression model and train it.

```
# Step 4: Train the Logistic Regression Model
cat("\n--- 4. Training the Logistic Regression Model ---\n")
# We use the glm() function with family = "binomial"
# The formula `Survived ~ .` predicts 'Survived' using all other variables in the dataset.
titanic_model <- glm(Survived ~ ., data = train_set, family = "binomial")
print(summary(titanic_model))
```

7.3 Making Predictions

With the model trained, we can use it to make predictions:

```
# Step 5: Make Predictions on the Test Set
cat("\n--- 5. Making Predictions on the Test Set ---\n")
predicted_probabilities <- predict(titanic_model, newdata = test_set, type = "response")
predicted_classes <- ifelse(predicted_probabilities > 0.5, "1", "0")
predicted_classes <- factor(predicted_classes, levels = levels(test_set$Survived))
cat("Predictions generated.\n")
```

The model is now ready for a thorough performance evaluation.

Section 8: Evaluating Classification Model Performance

Evaluating a classification model requires a more nuanced approach than simply calculating a single error score. We need metrics that describe the different types of correct and incorrect predictions the model makes.

8.1 The Confusion Matrix: A Deeper Look at Prediction Accuracy

A simple accuracy score can be highly misleading, especially in datasets with imbalanced classes. For example, if 95% of passengers did not survive, a model that always predicts 'did not survive' would be 95% accurate but completely useless. The

confusion matrix provides a much more detailed and insightful summary of model

performance by breaking down the predictions into four categories.

For a binary classification problem, the confusion matrix is a 2x2 table that tabulates the number of correct and incorrect predictions for each class.

	Predicted: No (0)	Predicted: Yes (1)
Actual: No (0)	True Negative (TN)	False Positive (FP)
Actual: Yes (1)	False Negative (FN)	True Positive (TP)

Table 4: The Confusion Matrix Deconstructed.

The four components are:

- **True Positives (TP):** The number of passengers who actually survived and were correctly predicted as having survived.
- **True Negatives (TN):** The number of passengers who did not survive and were correctly predicted as not having survived.
- **False Positives (FP) (Type I Error):** The number of passengers who did not survive but were incorrectly predicted as having survived.
- **False Negatives (FN) (Type II Error):** The number of passengers who did survive but were incorrectly predicted as not having survived.

8.2 Core Metrics Derived from the Confusion Matrix

From these four values, we can calculate several key performance metrics:

- **Accuracy:** The proportion of all predictions that were correct. It is a good general measure but can be misleading for imbalanced datasets.
 $Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$
- **Precision:** Measures the accuracy of the positive predictions. It answers the question: "Of all the passengers the model predicted would survive, what percentage actually did?" High precision is important when the cost of a False Positive is high (e.g., wrongly flagging a legitimate email as spam).
 $Precision = \frac{TP}{TP + FP}$
- **Recall (Sensitivity or True Positive Rate):** Measures the model's ability to find all the actual positive instances. It answers the question: "Of all the passengers who actually survived, what percentage did the model correctly identify?" High recall is crucial when the cost of a False Negative is high (e.g., failing to detect a cancerous tumor).
 $Recall = \frac{TP}{TP + FN}$
- **F1-Score:** The harmonic mean of precision and recall. It provides a single score that

balances both metrics. The F1-score is particularly useful when you need to find a balance between precision and recall, especially in the case of imbalanced classes.

$$F1\text{-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

8.3 The ROC Curve and Area Under the Curve (AUC)

The **Receiver Operating Characteristic (ROC) curve** is a powerful tool for visualizing and evaluating the performance of a binary classifier across all possible decision thresholds.⁹⁸ Instead of picking a single threshold (like 0.5), the ROC curve shows the model's performance at every threshold.

It plots the **True Positive Rate (TPR)**, which is the same as Recall, on the y-axis against the **False Positive Rate (FPR)** on the x-axis. The FPR is the proportion of actual negatives that were incorrectly classified as positive:

$$FPR = FP + TNFP$$

- A perfect classifier would have a curve that goes straight up the y-axis to (0, 1) and then across to (1, 1), indicating a TPR of 1 and an FPR of 0.
- A model that performs no better than random guessing will have a diagonal line from (0, 0) to (1, 1).

The **Area Under the Curve (AUC)** summarizes the ROC curve into a single value. The AUC represents the probability that the model will rank a randomly chosen positive instance higher than a randomly chosen negative instance.

- An AUC of 1.0 represents a perfect model.
- An AUC of 0.5 represents a model that is no better than random chance.

AUC is a very useful metric because it is threshold-independent and provides a comprehensive measure of the model's discriminative power.

8.4 Interpreting the Performance of our Titanic Survival Model

Now, we apply these evaluation techniques to our trained logistic regression model using the test data.

```
# Step 6: Evaluate Model Performance
```

```
cat("\n--- 6. Evaluating Model Performance ---\n")
```

```
# Generate the confusion matrix and key statistics
```

```
# We set 'positive = "1"' to ensure metrics are calculated for the "Survived" class
```

```
conf_matrix <- confusionMatrix(data = predicted_classes, reference = test_set$Survived,  
positive = "1")
```

```
cat("\n--- Confusion Matrix and Statistics ---\n")
```

```

print(conf_matrix)

# Calculate and display the AUC score
roc_curve <- roc(response = test_set$Survived, predictor = predicted_probabilities)
auc_score <- auc(roc_curve)
cat("\n--- AUC Score ---\n")
cat("AUC:", round(auc_score, 4), "\n")

# Generate and save the ROC curve plot
roc_plot <- ggroc(roc_curve, colour = 'steelblue', size = 1.5) +
  geom_segment(aes(x = 1, xend = 0, y = 0, yend = 1), color="darkgrey", linetype="dashed") +
  ggtitle(paste0("ROC Curve (AUC = ", round(auc_score, 4), ")")) +
  xlab("False Positive Rate (1 - Specificity)") +
  ylab("True Positive Rate (Sensitivity)") +
  theme_minimal()
ggsave("roc_curve.png", roc_plot, width = 7, height = 5)
cat("\nROC curve plot has been saved as 'roc_curve.png'.\n")

```

Interpretation of Results (Example):

Assuming the model yields an accuracy of around 80%, an AUC of 0.85, and a classification report, the interpretation would be as follows:

- **Accuracy:** The model correctly predicts survival status for approximately 80% of the passengers in the test set.
- **Confusion Matrix:** The heatmap would visually show the exact counts of TP, TN, FP, and FN. For instance, we might see that the model is better at correctly identifying passengers who did not survive (high TN) than those who did (lower TP), or vice versa.
- **Classification Report:**
 - **Precision for 'Survived' (1):** This tells us that when our model predicts a passenger survived, it is correct X% of the time. A high value here means the model's "survived" predictions are reliable.
 - **Recall for 'Survived' (1):** This tells us that our model successfully identified Y% of all the people who truly survived. A high value means the model is good at finding the survivors.
 - The **F1-Score** provides the balance between this precision and recall.
- **ROC/AUC:** An AUC score of 0.85 indicates that there is an 85% chance that the model will correctly rank a randomly chosen survivor higher than a randomly chosen non-survivor. This is significantly better than random chance (0.5) and suggests good discriminative ability across all thresholds.

This comprehensive evaluation provides a much richer understanding of the model's strengths and weaknesses than a single accuracy score could.

Laboratory Exercises

1 Practical Coding Challenge: California Housing Price Prediction

In this exercise, you will apply various regression techniques to predict median house values using the California Housing dataset.

1. Load and Prepare the Data:

- Load the California Housing dataset.
- Separate the data into features (X) and the target variable (y, which is MedHouseVal).
- Split the data into an 80% training set and a 20% testing set. For reproducibility, use a `random_state`.

2. Baseline Model (Linear Regression):

- Train a standard `LinearRegression` model on the training data.

3. Evaluate the Baseline Model:

- Use the trained model to make predictions on the test set.
- Calculate the Root Mean Squared Error (RMSE) and the Adjusted R-Squared value for the model's predictions on the test set.
- Interpret these metrics. What does the RMSE value tell you about the typical prediction error in terms of dollars (remember the target is in units of \$100,000)? What does the Adjusted R-Squared value say about how well your model explains the variance in housing prices?

4. Train Regularized Models:

- Train three additional models on the same training data: Ridge, Lasso, and ElasticNet.
- For each of these models, you need to choose a value for the regularization hyperparameter `alpha`. Experiment with a few different values, such as `alpha` in `[0.01, 0.1, 1.0, 10.0]`. For ElasticNet, also experiment with the `l1_ratio` (e.g., `[0.25, 0.5, 0.75]`). Select the best `alpha` (and `l1_ratio`) for each model based on its performance on the test set (in a real-world scenario, you would use cross-validation on the training set for this).

5. Compare Models and Conclude:

- For your best-performing Ridge, Lasso, and ElasticNet models, calculate the RMSE and Adjusted R-Squared on the test set.
- Create a summary table that compares all four models (Linear, Ridge, Lasso, ElasticNet) based on their Test RMSE, Test Adjusted R-Squared, and the number of non-zero coefficients.
- Analyze the coefficients of the Lasso model. Did it set any feature coefficients to zero? If so, what does this suggest about the importance of those features in predicting California housing prices?
- Based on your table and analysis, which model would you conclude is the best for this task? Justify your choice by considering both predictive performance (RMSE) and model complexity/interpretability (number of non-zero coefficients).

2 Theoretical Questions

1. **Ordinary Least Squares (OLS):** Explain the core principle of the OLS method used in linear regression. Why are the residuals squared in the cost function? What are the potential consequences of using the sum of absolute residuals (Least Absolute Deviations) as the cost function instead?
2. **Bias-Variance Tradeoff:** Describe the bias-variance tradeoff in the context of model complexity. Specifically, how do the bias and variance of a polynomial regression model change as you increase the degree of the polynomial? Use diagrams to illustrate the concepts of underfitting, optimal fit, and overfitting.
3. **Regularization Comparison:** Contrast L1 (Lasso) and L2 (Ridge) regularization. Provide the mathematical form of their respective penalty terms. Explain the geometric intuition (using the concept of constraint regions) that underlies why Lasso can perform feature selection by setting coefficients to exactly zero, while Ridge cannot.
4. **Model Evaluation Metrics:** A colleague has built a multiple linear regression model with 20 features and proudly reports a high R-Squared value of 0.95. Why might this value be misleading? What specific metric would you request to get a more reliable assessment of the model's goodness of fit, and why is that metric superior in this context?
5. **Logistic Regression's Name:** Explain the seeming paradox of why Logistic Regression, an algorithm used for classification, has "regression" in its name. What continuous quantity is the model actually "regressing" on? Describe the mathematical relationship between the model's linear component and the final probability output.