**ML_U2**

Regression is a **statistical and machine learning technique** used to model and analyze the relationship between a **dependent variable** (target) and one or more **independent variables** (predictors).

Its goal is to predict or estimate the dependent variable based on known values of the independent variables.

- **Dependent Variable (Y):** The value we want to predict (e.g., house price, sales revenue).
- **Independent Variables (X):** The input features used for prediction (e.g., area, location, number of bedrooms).



General Equation of Regression:

$$Y = f(X) + \varepsilon$$

Where:

- $f(X)$ = functional relationship between variables
- $\varepsilon$ = error term (difference between actual and predicted values)

Example:

Predicting a student's final score (Y) based on hours studied (X).

Simple equation:

$$\text{Score} = a + b \times \text{Hours studied}$$

**Need for Regression**

Regression is important in many real-life scenarios because it provides both **prediction** and **relationship analysis**.

1. **Prediction of future outcomes**
   o Example: Predicting rainfall based on humidity and temperature.
2. **Relationship understanding**
   o Shows how the dependent variable changes when independent variables change.
3. **Forecasting trends**
   o Example: Estimating future sales from past sales data.
4. **Decision-making & planning**
   o Helps businesses choose strategies (e.g., which features increase product sales).

5. **Quantifying impact of factors**
   o Example: Finding which factor most influences customer satisfaction.

**Regression and Correlation**

| Aspect | Regression | Correlation |
|---|---|---|
| Purpose | Predicts dependent variable from independent variables and shows the nature of relationship. | Measures strength and direction of relationship between two variables. |
| Direction | One-way (predict Y from X). | No direction; relationship is mutual. |
| Output | Regression equation ($Y = a + bX$). | Correlation coefficient (r) between -1 and +1. |
| Variable Types | One dependent (Y) and one/more independent (X). | Two variables treated equally; no dependent/independent distinction. |
| Application | Used for prediction and forecasting. | Used for measuring degree of association. |
| Example | Predicting salary from years of experience. | Measuring relationship between salary and years of experience. |

**Types of Regression**

**1. Univariate vs. Multivariate Regression**

**Univariate Regression**

A regression model with only one independent variable (X) used to predict a single dependent variable (Y).

**Purpose:** Helps to understand and quantify the relationship between two variables.

**Equation:** $Y = a + bX$, where $a =$ intercept, $b =$ slope (change in Y per unit change in X).

Example: Predicting student marks based on hours studied.

**Graph:** Straight line in a 2D plot (X vs Y).

**Multivariate Regression**

A regression model with two or more independent variables (X1, X2, …) used to predict one dependent variable (Y).

**Purpose:** Captures more complex relationships and accounts for multiple influencing factors.

**Equation:** $Y = a + b1X1 + b2X2 + … + bnXn$, where b1, b2 … = coefficients for each predictor.

Example: Predicting house price from size, location, and age.
**Graph**: Can be visualized as a plane (for 2 predictors) or hyperplane (for more predictors).

## 2. Linear vs. Nonlinear Regression

### Linear Regression
A regression model where the relationship between dependent and independent variables is linear in the parameters.
**Key Point**: Variables can be transformed (e.g., polynomial), but coefficients appear in a linear way.
**Equation**: Y = a + bX.
Example: Predicting salary from years of experience.
Advantages: Simple, fast to compute, easy to interpret.
**Graph**: Straight line in 2D, flat plane in higher dimensions.

### Nonlinear Regression
**Definition**: A regression model where the relationship is nonlinear in parameters. Parameters may appear in exponents, products, or other nonlinear forms.
Equation Example: Y = a * e^(bX).
**Example**: Modeling population growth using an exponential curve.
Advantages: Captures complex, curved patterns.
Disadvantages: Harder to compute, may need iterative methods.
**Graph**: Curved fit line (e.g., exponential, logarithmic, S-shape).

## 3. Simple Linear vs. Multiple Linear Regression

### Simple Linear Regression (SLR)
Linear regression with one independent variable.
**Equation**: Y = a + bX.
Example: Predicting marks from hours studied.
Interpretation: b tells how much Y changes for a one-unit increase in X.
**Graph**: Straight line on 2D plot.

### Multiple Linear Regression (MLR)
Linear regression with two or more independent variables.
**Equation:** $Y = a + b_1X_1 + b_2X_2 + \ldots + b_nX_n$.
**Example:** Predicting salary based on experience, education level, and skill score.
Interpretation: Each $b_i$ represents the effect of $X_i$ on Y, keeping other variables constant.

Graph: Plane for two predictors, hyperplane for more.

**Bias–Variance Trade-Off**

The **Bias–Variance Trade-Off** describes the balance between a model's ability to **fit the training data** and its ability to **generalize to unseen data**. The goal is to minimize the **Total Error** (Generalization Error), which comes from three sources:

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Err}$$

| Model Complexity | Bias | Variance | Total Error |
|---|---|---|---|
| Too Simple | High | Low | High |
| Too Complex | Low | High | High |
| Optimal | Balanced | Balanced | Low |

**High Bias** → Underfitting → Simple model.

**High Variance** → Overfitting → Complex model.

### Bias

- Error caused by **overly simplistic assumptions** in the model.
- High bias → **Underfitting** (model misses important patterns).
- Example: Using a straight line to fit non-linear data.

**Characteristics of High Bias:**

- Poor training accuracy.
- Poor test accuracy.
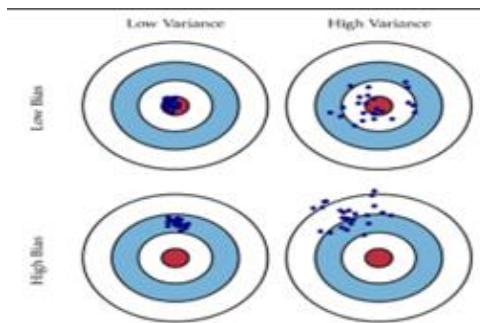- Model is too simple.

### Variance

- Error caused by **too much sensitivity to training data fluctuations**.
- High variance → **Overfitting** (model memorizes training data but fails to generalize).
- Example: Using a high-degree polynomial to fit simple data.

**Characteristics of High Variance:**

- High training accuracy.

- Low test accuracy.
- Model is too complex.



## Overfitting

When a model learns the training data too well, including noise, it performs well on training data but poorly on new (test) data.

- Low training error, high test error
- Predictions follow overly complex curves or patterns

**Causes:**

- Too many features or parameters
- Too little training data
- Long training without regularization
- Data leakage

**Detection:**

- Large gap between training and test performance
- High variation in cross-validation scores
- Learning curve shows low training error but high test error

**Remedies:**

- Add regularization
- Use simpler model
- Prune decision trees or limit depth
- Increase k in k-NN
- Early stopping or dropout (for neural networks)
- Data augmentation or cleaning
- Use cross-validation
- Apply ensembling (bagging)

## Underfitting

When a model is too simple to capture the actual pattern in data, it performs poorly on both training and test sets.

- High training error and high test error
- Misses obvious patterns or trends

**Causes:**

- Model too simple
- Not enough features
- Too much regularization
- Training stopped too early

**Detection:**

- Learning curve shows both training and test errors are high and close
- Residual plots show clear patterns

**Remedies:**

- Use a more complex model
- Add more useful features
- Reduce regularization
- Train for longer
- Add interaction terms or nonlinear transformations

## Polynomial Regression

Polynomial Regression is an **extension of Linear Regression** that allows us to model **non-linear relationships** between the independent variable(s) $XXX$ and dependent variable $YYY$.

While **Linear Regression** fits a straight line to the data, Polynomial Regression fits a **curved line** by adding higher-degree terms of the predictor variable.

**Why Polynomial Regression?**

- In real-world problems, relationships are rarely perfectly linear.
- Example: Growth rate of bacteria, housing prices vs. area, learning curves, etc. often follow a curved trend.
- Linear regression would **underfit** such data because it assumes a constant slope, but polynomial regression can capture bends and curves.

**How It Works**

- We **transform** the original features by adding polynomial terms.
- Then we perform **Linear Regression** on these transformed features.
- Even though the curve looks non-linear in $xxx$, the regression is **linear in coefficients**,

so it can be solved using the same least squares method.

Choosing the polynomial degree (n) is critical:

- **Too low degree** → Underfitting.
- **Too high degree** → Overfitting.
- Use **cross-validation** to choose optimal n.

## 3. Mathematical Representation

Linear Regression Equation:

$$y = \beta_0 + \beta_1 x + \epsilon$$

Polynomial Regression Equation (degree n):

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \cdots + \beta_n x^n + \epsilon$$

where:
- $y$ = dependent variable
- $x, x^2, x^3, \ldots, x^n$ = independent variable and its powers
- $\beta_0, \beta_1, \ldots, \beta_n$ = coefficients to be estimated
- $\epsilon$ = error term

### Advantages

✅ Captures non-linear trends.
✅ Easy to implement with existing linear regression algorithms.
✅ Works well when the true relationship is polynomial-like.

## 8. Limitations

✖ High-degree polynomials are sensitive to noise → Overfitting.
✖ Can oscillate wildly between points (Runge's phenomenon).
✖ Poor extrapolation — predictions outside the data range can be unrealistic.

**Example:** Predicting exam scores based on study hours where improvement rate slows after a certain point.

### Stepwise Regression

Stepwise regression is a **variable selection method** for multiple linear regression that iteratively adds or removes predictors based on statistical criteria such as **p-values, AIC, BIC, or adjusted R²**. It aims to find a parsimonious model that explains the data without including unnecessary variables.

**Mathematical Expression:**
Starts with:

y=β0+β1X1+⋯+βpXp+ϵ

The algorithm chooses the subset of XXX that minimizes an error metric.

**Working:**

1. **Forward Selection** – Start with no variables, add the one with the highest significance until no improvement.
2. **Backward Elimination** – Start with all variables, remove the least significant one at each step.
3. **Stepwise** – Combination of both methods.

**Advantages:**

- Reduces overfitting by eliminating irrelevant variables.
- Automates model building.
- Improves interpretability.

**Limitations:**

- Can ignore variable interactions.
- May overfit if selection is based solely on training data.
- Unstable when predictors are highly correlated.

**Example:** Selecting important economic indicators for predicting GDP growth.

### Decision Tree Regression

Decision Tree Regression is a **non-parametric supervised learning** method used for predicting continuous values.
It works by **splitting the dataset into regions** based on decision rules, where each region has similar target values.
The final prediction for a region is the **mean** (or sometimes median) of its target values.

Mathematical Representation:
For a region $R_m$:

$$\hat{y}_{R_m} = \frac{1}{|R_m|} \sum_{x_i \in R_m} y_i$$

The best split minimizes:

$$\text{MSE} = \sum_{m=1}^{M} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2$$

**Working Process:**

1. Start with the full dataset as the root node.
2. Find the feature and split point that **maximally reduces Mean Squared Error (MSE)**.

3. Recursively split each subset into child nodes until stopping criteria are met (max depth, min samples per leaf, or no further improvement).
4. Prediction = **mean target value** of the leaf node where the sample falls.

**Advantages:**

- Simple to **interpret and visualize**.
- Can capture **non-linear relationships**.
- Works with both numerical and categorical features.
- No need for feature scaling.

**Limitations:**

- **Prone to overfitting** if grown too deep.
- **High variance** – small changes in data can alter the tree structure.
- Predictions are **piecewise constant**, which may be less smooth than other models.

**Example:**
Predicting **house prices** based on size, location, and number of bedrooms by splitting data into ranges (e.g., "if size > 1500 sqft and location = city center → higher price").

**Random Forest Regression**

Random Forest Regression is an **ensemble learning** method that combines predictions from multiple **decision trees** to improve accuracy and reduce overfitting.
It uses **bagging** (Bootstrap Aggregating) and **random feature selection** to create diverse trees and averages their outputs for final prediction.

**Mathematical Representation:**
If there are $T$ trees $h_t(X)$:

$$\hat{y} = \frac{1}{T}\sum_{t=1}^{T} h_t(X)$$

**Working Process:**

1. **Bootstrap Sampling** – Create random samples (with replacement) from the training data.
2. **Random Feature Selection** – At each tree split, consider only a random subset of features.
3. Train each decision tree **without pruning**.
4. For prediction, take the **average** of predictions from all trees.

**Advantages:**

- **Reduces overfitting** compared to a single decision tree.
- Works well with **non-linear data** and complex relationships.
- Robust to **outliers and noise**.
- Can handle **missing values** effectively.

**Limitations:**

- Less **interpretable** than a single tree.
- **Computationally expensive** for large datasets.
- Requires careful tuning of parameters (number of trees, max depth, etc.).

**Example:**
Predicting **car prices** using multiple features (mileage, age, brand, condition) where each tree learns different patterns and the forest combines them for a stable result.

**Support Vector Regression (SVR)**

Support Vector Regression is the regression counterpart of **Support Vector Machines (SVM)**. Instead of finding a classification boundary, SVR tries to find a function that fits the data **within an ε-margin of tolerance**, ignoring small deviations and penalizing only large errors. It can model **linear and non-linear** relationships using **kernels**.

**Mathematical Representation:**
Optimization problem:

$$\min \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{n}(\xi_i + \xi_i^*)$$

subject to:

$$y_i - (w^T x_i + b) \le \epsilon + \xi_i$$
$$(w^T x_i + b) - y_i \le \epsilon + \xi_i^*$$

where:

- $\epsilon$ = tolerance margin
- $C$ = penalty parameter
- $\xi_i, \xi_i^*$ = slack variables for points outside the margin

**Working Process:**

1. **Choose a kernel** (linear, polynomial, RBF) to transform data if non-linear.
2. Fit a function such that most data points lie **within ±ε** from it.
3. Penalize only the points lying **outside the ε-tube**.
4. Use **support vectors** (critical boundary points) for prediction.

**Advantages:**

- Effective for **high-dimensional** and **non-linear** data.
- Robust to overfitting due to margin control.
- Uses only support vectors, making it memory efficient for small datasets.

**Limitations:**

- Requires careful tuning of **C, ε, and kernel parameters**.
- Computationally expensive for very large datasets.
- Performance sensitive to choice of kernel.

**Example:**
Predicting **electricity demand** from temperature and time-of-day where the relationship is complex and non-linear.

**Ridge Regression**

**Concept:**
Ridge Regression is a **regularized form of linear regression** that adds an **L2 penalty** to the loss function.
This penalty shrinks large coefficient values, which helps reduce **overfitting** and handle **multicollinearity** (high correlation between predictors).
It keeps all predictors in the model but with reduced magnitudes.

Cost function:

$$\min_{\beta} \left[ \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \right]$$

where:
- $\lambda$ = regularization strength (larger $\lambda \rightarrow$ more shrinkage)
- $\beta_j$ = model coefficients

**Working Process:**

1. Start with the **Ordinary Least Squares (OLS)** regression formulation.
2. Add an **L2 penalty term** to the cost function.
3. Optimize to find coefficients that minimize both **prediction error** and **coefficient magnitude**.
4. Larger $\lambda \rightarrow$ more shrinkage; $\lambda = 0$ gives standard linear regression.

**Advantages:**

- Handles **multicollinearity** effectively.

- Reduces model complexity and **overfitting**.
- Works well when most predictors have small to moderate effects.

**Limitations:**

- Does **not** perform feature selection (keeps all predictors).
- Choice of $\lambda$ requires **cross-validation**.
- Coefficients become biased due to shrinkage.

**Example:**
Predicting **sales** using multiple highly correlated marketing variables (e.g., TV ads, online ads, print ads).

**Lasso Regression**

Lasso Regression (**Least Absolute Shrinkage and Selection Operator**) is a **regularized linear regression** method that uses an **L1 penalty** on the coefficients.
Unlike Ridge Regression, Lasso can shrink some coefficients **exactly to zero**, effectively performing **feature selection** along with regularization.

Mathematical Representation:
Cost function:

$$\min_{\beta} \left[ \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right]$$

where:
- $\lambda$ = regularization strength
- $|\beta_j|$ = absolute value of coefficient $j$

**Working Process:**

1. Start with **Ordinary Least Squares (OLS)** regression.
2. Add an **L1 penalty term** to the cost function.
3. Optimization via **coordinate descent** or **LARS algorithm** forces some coefficients to be exactly **zero**.
4. Larger $\lambda \rightarrow$ more coefficients shrink to zero.

**Advantages:**

- Performs **automatic feature selection**.
- Produces simpler, more interpretable models.
- Useful when only a subset of predictors are important.

**Limitations:**

- In presence of highly correlated predictors, it tends to select only one and ignore others.
- Choice of $\lambda$ is critical and requires cross-validation.
- May be unstable for small datasets.

**Example:**
Selecting important **genes** from thousands of genomic variables to predict disease risk.

**ElasticNet                                    Regression**
**ElasticNet Regression is a** hybrid regularization method **that combines** L1 penalty **(from Lasso) and** L2 penalty **(from Ridge). It performs** feature selection **like Lasso while maintaining** stability **in the presence of highly correlated predictors like Ridge.**



**Mathematical Representation:**
Cost function:

$$\min_{\beta} \left[ \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{j=1}^{p} |\beta_j| + \lambda_2 \sum_{j=1}^{p} \beta_j^2 \right]$$

Alternatively, using mixing parameter $\alpha$ ($0 \leq \alpha \leq 1$):

$$\text{Penalty} = \lambda \left[ \alpha \sum |\beta_j| + (1 - \alpha) \sum \beta_j^2 \right]$$

**Working Process:**

1. Transform problem into a mix of **L1 and L2 regularization**.
2. $\alpha$ controls the balance:
   - $\alpha = 1 \rightarrow$ pure Lasso
   - $\alpha = 0 \rightarrow$ pure Ridge
3. Solve using coordinate descent or similar optimization methods.
4. Select $\lambda$ and $\alpha$ using cross-validation.

**Advantages:**

- Handles **multicollinearity** better than Lasso.
- Performs feature selection while keeping correlated variables.
- More flexible than Ridge or Lasso alone.

**Limitations:**

- Requires tuning of **two parameters** ($\lambda$ and $\alpha$).
- More computationally intensive than Ridge/Lasso.
- Interpretability slightly reduced compared to Lasso.

**Example:**
Predicting **customer churn** using many correlated customer behavior features where both feature selection and stability are needed.

**Bayesian Linear Regression**

Bayesian Linear Regression is a **probabilistic approach** to linear regression where model parameters (coefficients) are treated as **random variables** with prior probability distributions. Using **Bayes' theorem**, the model updates these priors into **posterior distributions** after seeing the data, providing both predictions and **uncertainty estimates**.

**Mathematical                          Representation:**
Bayes' theorem for regression:

$p(\beta|X,y)=(p(y|X,\beta)\cdot p(\beta))/p(y|X)$

where:

- $p(\beta)$ = prior distribution of coefficients
- $p(y|X,\beta)$ = likelihood from data
- $p(\beta|X,y)$ = posterior distribution after observing data

Prediction distribution:

$p(y*|x*,X,y)=\int p(y*|x*,\beta) \, p(\beta|X,y) \, d\beta p$

**Working Process:**

1. **Define prior** for coefficients (commonly Gaussian with mean 0).
2. **Calculate likelihood** from observed data using the linear model.
3. Apply **Bayes' theorem** to obtain the posterior distribution.
4. Predict new values using the **posterior mean** (point estimate) or full distribution for uncertainty quantification.

**Advantages:**

- Produces **uncertainty intervals** for predictions.
- Can incorporate **prior knowledge** about parameters.
- Works well with small datasets.
- Naturally prevents overfitting through prior regularization.

**Limitations:**

- Requires **specifying priors** carefully.
- Computationally expensive for large datasets.
- Analytical solutions may not exist for complex priors (requires approximation methods like MCMC).

**Example:**
Forecasting **monthly sales** while accounting for uncertainty and prior knowledge about seasonal effects.

## 1. Mean Squared Error (MSE)

- **Definition**
  MSE is the average of the squared differences between actual and predicted values.
  Squaring ensures all errors are positive and penalizes larger errors more heavily.

Formula:

$$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}(y_i -$$

Where:
- $n$ = number of data points
- $y_i$ = actual value
- $\hat{y}_i$ = predicted value

- **Interpretation**
  - Lower MSE → better fit.
  - Value is **always ≥ 0**; 0 means perfect predictions.
- **Advantages**
  - Easy to compute mathematically.
  - Penalizes large errors strongly → useful when big mistakes are costly.
- **Disadvantages**
  - Unit is **squared** (not same as target variable).
  - Highly sensitive to outliers.

## 2. Mean Absolute Error (MAE)

- **Definition**
  MAE is the average of the absolute differences between actual and predicted values.
  Unlike MSE, it does not square errors, so it treats all errors equally.

Formula:

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

- **Interpretation**
  - Measures average magnitude of errors in the **same unit** as the target variable.
  - Gives an idea of how wrong predictions are, on average.
- **Advantages**
  - Easy to interpret.
  - Less sensitive to outliers than MSE.
- **Disadvantages**
  - Does not emphasize large errors → may be misleading if big mistakes are important.

## 3. Root Mean Squared Error (RMSE)

- **Definition**
  RMSE is the square root of MSE.
  It gives the error in the **same unit** as the target variable.

Formula:

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

- **Interpretation**
  - Shows how far predictions are from actual values, on average.
  - Lower RMSE = better model.
- **Advantages**
  - Same units as target variable → easy to compare.
  - Still penalizes large errors more than MAE.
- **Disadvantages**
  - Sensitive to outliers.

## 4. R-squared (R2R^2R2)

- **Definition**
  Measures the proportion of variance in the dependent variable explained by the model.
- **Formula**

Formula:

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

Where:
- $\bar{y}$ = mean of actual values

- **Interpretation**
  - $R^2 = 1 \rightarrow$ perfect fit
  - $R^2 = 0 \rightarrow$ model explains none of the variance
  - Can be **negative** if the model is worse than predicting the mean.
- **Advantages**
  - Intuitive measure of goodness-of-fit.
  - Easy to compare models on the same dataset.
- **Disadvantages**
  - Can increase artificially when adding more variables, even if they are irrelevant.
  - Does not indicate bias or overfitting.

## 5. Adjusted R-squared

- **Definition**
  Modified version of $R^2$ that adjusts for the number of predictors in the model. Penalizes adding predictors that don't improve the model.

  Formula:
  $$\text{Adjusted } R^2 = 1 - \left[ \frac{(1 - R^2)(n - 1)}{n - p - 1} \right]$$

  Where:
  - $n$ = number of observations
  - $p$ = number of predictors

- **Interpretation**
  - Increases only if the new variable improves the model significantly.
  - Useful for comparing models with different numbers of predictors.
- **Advantages**
  - Prevents misleading improvement from adding irrelevant predictors.
  - Better measure for model selection.
- **Disadvantages**
  - More complex formula.
  - Still doesn't guarantee the model is correct.

### Univariate Regression vs. Multivariate Regression

| Aspect | Univariate Regression | Multivariate Regression |
|---|---|---|
| 1. Number of Independent Variables | Involves **only one** independent (predictor) variable. | Involves **two or more** independent (predictor) variables. |
| 2. Purpose | Examines the relationship between a single predictor and the dependent variable. | Examines the combined effect of multiple predictors on the dependent variable. |
| 3. Complexity | Simple to model and interpret. | More complex due to multiple predictors and possible interactions. |
| 4. Equation Form | $Y = \beta_0 + \beta_1 X + \varepsilon$ | $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \varepsilon$ |
| 5. Use Case | Predicting sales based on **only** price. | Predicting sales based on **price, advertising spend, and store location** together. |

### 1. Ridge Regression vs. Lasso Regression (6 Points)

| Aspect | Ridge Regression | Lasso Regression |
|---|---|---|
| 1. Regularization Type | Uses **L2** regularization ($\lambda \sum \beta_j^2$). | Uses **L1** regularization (( \lambda \sum |
| 2. Effect on Coefficients | Shrinks coefficients but never exactly zero. | Can shrink some coefficients exactly to zero (feature selection). |
| 3. Feature Selection | Does **not** perform feature selection; retains all variables. | Performs automatic feature selection by eliminating irrelevant variables. |
| 4. Handling Multicollinearity | Good at handling multicollinearity while keeping all predictors. | Can remove redundant predictors entirely when they are highly correlated. |
| 5. Equation Bias | Adds small bias to reduce variance. | Adds bias but can also create sparse models. |
| 6. Best Use Case | Best when most features are useful but may have small effects. | Best when many features are irrelevant or dataset is high-dimensional. |

## 2. Regression vs. Correlation (6 Points)

| Aspect | Regression | Correlation |
| --- | --- | --- |
| 1. Definition | Predicts the value of a dependent variable from one or more independent variables. | Measures the degree and direction of relationship between two variables. |
| 2. Purpose | Establishes a predictive equation. | Quantifies strength of association. |
| 3. Direction of Relationship | One variable is dependent, others are independent. | No distinction between dependent and independent variables. |
| 4. Output | Gives a regression equation ($Y = \beta_0 + \beta_1 X + \varepsilon$). | Gives a correlation coefficient ($r$) between $-1$ and $+1$. |
| 5. Cause & Effect | Implies a cause-effect relationship (if model assumptions are valid). | Does not imply causation, only association. |
| 6. Use Case | Predicting sales from price and advertising spend. | Measuring how strongly price and sales are related. |