# Examples of Linear Regression models and Non-linear Regression models

Linear regression models and nonlinear regression models are both used in machine learning to model the relationship between input variables and output variables. Here are some examples of each:

<u>Linear regression models:</u>

Simple linear regression: A linear model that uses a single input variable to predict a continuous output variable. For example, predicting housing prices based on square footage.

Multiple linear regression: A linear model that uses multiple input variables to predict a continuous output variable. For example, predicting customer churn based on age, income, and purchase history.

Polynomial regression: A linear model that fits a polynomial function to the input data. For example, fitting a quadratic function to data that has a curved relationship between the input and output variables.

<u>Nonlinear regression models:</u>

Logistic regression: A nonlinear model used for classification problems that predict the probability of a

binary output variable. For example, predicting whether a customer will make a purchase based on their demographic information.

Support vector regression (SVR): A nonlinear model used for regression problems that use support vector machines to fit a curve to the input data. For example, predicting the stock prices of a company based on historical data.

Decision trees: A nonlinear model that makes predictions by recursively partitioning the input space into smaller regions based on the input variables. For example, predicting customer churn based on their demographic information and purchase history.

In summary, linear regression models are used to model linear relationships between input and output variables, while nonlinear regression models are used to model more complex relationships that are not linear.

## Examples of Polynomial Regression

Polynomial regression is a type of linear regression that fits a polynomial function to the input data, allowing for a curve to be fitted to the data instead of a straight line. Here are some examples of polynomial regression:

Fitting a quadratic function to data that has a curved relationship between the input and output variables. For example, predicting the sales of a product based on the price and advertising spend.

Fitting a cubic function to data that has an S-shaped relationship between the input and output variables. For example, predicting the amount of rainfall based on the temperature and humidity.

Fitting a higher-degree polynomial function to data that has a more complex relationship between the input and output variables. For example, predicting the lifespan of a machine based on the operating temperature and pressure.

In general, polynomial regression can be used to fit any degree of polynomial to the data, depending on the complexity of the relationship between the input and output variables. However, it's important to keep in mind that overfitting can occur if the degree of the polynomial is too high, so it's important to use techniques like regularization and cross-validation to ensure that the model is not overfitting to the training data.

In the context of Linear Regression, which situations does require OLS and which doesn't

OLS (ordinary least squares) is a method used in linear regression to estimate the parameters of a linear model by minimizing the sum of the squared residuals. OLS is a widely used method in linear regression, but there are situations where it may not be appropriate to use it. Here are some situations where OLS may or may not be suitable:

## When to use OLS:

The relationship between the independent and dependent variables is linear or approximately linear.

The residuals are normally distributed.

The residuals have constant variance (i.e., homoscedasticity).

The independent variables are not highly correlated with each other (i.e., there is no multicollinearity).

## When not to use OLS:

The relationship between the independent and dependent variables is not linear.

The residuals are not normally distributed (i.e., there is non-normality).

The residuals do not have constant variance (i.e., there is heteroscedasticity).

The independent variables are highly correlated with each other (i.e., there is multicollinearity).

In summary, OLS is a commonly used method in linear regression, but it is important to consider whether the assumptions of the method are met in a particular situation. If the assumptions are not met, other regression methods, such as nonlinear regression or robust regression, may be more appropriate.

When the assumptions of OLS are violated, there are several techniques that can be used in linear regression to address the issues. Here are some examples:

Weighted least squares regression: A variation of OLS that assigns different weights to the data points based on their variance or importance. Weighted least squares regression can be used to address heteroscedasticity, where the variance of the residuals

is not constant across the range of the independent variables.

Nonlinear regression: A regression technique that models the relationship between the independent and dependent variables using a nonlinear function. Nonlinear regression can be used when the relationship between the variables is not linear.

Generalized linear models: A family of regression models that can handle non-normal response variables and non-constant variance. Generalized linear models include regression models such as logistic regression, Poisson regression, and gamma regression.

Robust regression: A regression technique that is less sensitive to outliers than OLS. Robust regression methods, such as the Huber and Tukey bisquare methods, can be used to address the issue of non-normality and outliers.

Ridge regression and Lasso regression: Regularization techniques that can be used to address the issue of multicollinearity. Ridge regression and Lasso regression add a penalty term to the OLS cost function that encourages smaller parameter estimates, which can help to reduce the impact of multicollinearity.

In summary, when the assumptions of OLS are violated, there are several techniques that can be used in linear regression to address the issues. These techniques include weighted least squares regression, nonlinear regression, generalized linear models, robust regression, and regularization techniques such as ridge regression and Lasso regression.

## What are GLMs and how are they different from Linear Models

GLM (Generalized Linear Models) is an extension of linear models that allows for non-normally distributed response variables and non-linear relationships between the independent and dependent variables. Linear models assume that the response variable has a normal distribution, while GLMs allow for other types of response distributions, such as binomial, Poisson, or gamma distributions.

GLMs are similar to linear models in that they involve a linear combination of the independent variables, but they differ in the following ways:

Link function: In a GLM, a link function is used to relate the linear combination of the independent variables to

the response variable. The link function transforms the response variable to a linear scale, allowing for non-linear relationships between the independent and dependent variables.

Response distribution: In a GLM, the response variable can have any distribution from the exponential family of distributions, such as the binomial, Poisson, or gamma distributions, which are not normally distributed. This allows for modeling of count data, binary outcomes, and continuous outcomes that have a skewed distribution.

Non-constant variance: In a GLM, the variance of the response variable can vary as a function of the mean, allowing for non-constant variance across the range of the independent variables.

In summary, GLMs are an extension of linear models that allow for non-normally distributed response variables and non-linear relationships between the independent and dependent variables. GLMs use a link function to relate the linear combination of the independent variables to the response variable, and they allow for a wider range of response distributions and non-constant variance.

# Can we use Least squares for GLMs?

Least squares is a method used in linear regression to estimate the parameters of a linear model by minimizing the sum of the squared residuals. However, it cannot be used for GLMs (Generalized Linear Models) because GLMs involve non-normally distributed response variables and non-linear relationships between the independent and dependent variables.

Instead of least squares, GLMs use a maximum likelihood estimation (MLE) approach to estimate the model parameters. MLE is a method used to find the parameters that maximize the likelihood of observing the data given the model. The likelihood function is determined by the chosen response distribution and link function of the GLM.

The maximum likelihood estimation approach in GLMs is based on iteratively reweighted least squares (IRLS), where at each iteration, the linear predictor is updated based on the current estimates of the model parameters, and the weights are updated based on the response distribution and the current estimate of the

mean. This iterative process continues until convergence is reached and the estimates of the model parameters are stable.

In summary, least squares cannot be used for GLMs due to their non-normal response variables and non-linear relationships between the independent and dependent variables. Instead, GLMs use maximum likelihood estimation (MLE) through iteratively reweighted least squares (IRLS) to estimate the model parameters.

Examples of ML models under GLM

- GLM (Generalized Linear Models) is a framework for building a variety of models that can handle non-normally distributed response variables and non-linear relationships between the independent and dependent variables. Here are some examples of models that fall under GLM:

- Logistic regression: A GLM that models the probability of a binary outcome, such as whether a customer will buy a product or not

- Poisson regression: A GLM that models count data, such as the number of insurance claims filed by a policyholder.

- Gamma regression: A GLM that models continuous positive data with a skewed distribution, such as healthcare costs.

- Inverse Gaussian regression: A GLM that models continuous positive data with a highly skewed distribution, such as the time to failure of a component.

- Negative binomial regression: A GLM that models over dispersed count data, where the variance is higher than the mean, such as the number of accidents at an intersection.

- In summary, GLM is a flexible framework that includes a variety of models, such as logistic regression, Poisson regression, gamma regression, inverse Gaussian regression, and negative binomial regression, among others. Each of these models can handle specific types of response variables and relationships between the independent and dependent variables.

# Does Decision Trees, Random forests come under GLMs?

<u>Decision Trees and Random Forests are not examples of GLMs</u> (Generalized Linear Models). GLMs are a class of models that can handle non-normally distributed response variables and non-linear relationships between the independent and dependent variables, whereas Decision Trees and Random Forests are examples of non-linear models that can model complex relationships between the variables without explicitly assuming a functional form.

<u>Decision Trees are a non-linear model</u> that recursively splits the data based on the most informative feature at each node to create a tree-like structure of decision rules. Random Forests are an ensemble method that combines multiple Decision Trees, each trained on a random subset of the features and/or data, to reduce overfitting and improve generalization performance.

<u>Both Decision Trees and Random Forests are flexible</u> and can handle a wide range of data types and relationships between variables. They are often used in problems where the relationships between variables

are complex and cannot be easily captured by linear or GLM models.

In summary, Decision Trees and Random Forests are examples of non-linear models that are different from GLMs. While GLMs are suitable for handling non-normally distributed response variables and non-linear relationships between the independent and dependent variables, Decision Trees and Random Forests are flexible and can capture complex relationships between variables without assuming a functional form.

<span style="color:red">Among GLMs and Decision Trees, which models are more suitable for non-linear relationships</span>

<u>Decision Trees are generally more suitable for non-linear relationships</u> between variables compared to GLMs. Decision Trees are non-parametric models that can capture non-linear and complex relationships between the independent and dependent variables without assuming a specific functional form.

<u>GLMs</u>, on the other hand, assume a linear relationship between the independent and dependent variables, which may not be appropriate for data with non-linear

relationships. While GLMs can use link functions to capture non-linearity in the relationship, they are generally more suited for data with linear relationships and may not be as effective in capturing complex non-linear relationships.

In summary, while both GLMs and Decision Trees can handle non-linear relationships between variables, Decision Trees are generally considered to be more suitable for capturing complex non-linear relationships, while GLMs are more appropriate for linear relationships.

## Differences between Linear models and GLM in a table

|  | Linear Models | GLMs |
|---|---|---|
| variable distribution | Assumes normal distribution | Can handle a variety of distributions from the exponential family (e.g., binomial, Poisson, gamma) |
| Link function | Uses the identity link function | Uses a link function to relate the linear predictor to the expected value of the response variable |
| Non-constant variance | Assumes constant variance | Can model non-constant variance (i.e., heteroscedasticity) |

| Assumptions | Assumes linear relationship between independent and dependent variables | Can handle non-linear relationships between independent and dependent variables |
|---|---|---|
| Maximum likelihood estimation | Not used | Used to estimate the model parameters |
| Example | Simple linear regression | Logistic regression, Poisson regression, gamma regression |

In summary, GLMs extend the capabilities of linear models by allowing for non-normally distributed response variables, non-linear relationships between the independent and dependent variables, and non-constant variance. GLMs use a link function to relate the linear predictor to the expected value of the response variable and use maximum likelihood estimation to estimate the model parameters.

if GLMs allow non-linear relationships among variables, why is the word 'linear' still present in GLMs?

- The word 'linear' in GLMs refers to the linear relationship between the predictors (i.e.,

independent variables) and the linear predictor, which is a combination of the predictors weighted by their coefficients.

- In a GLM, the relationship between the predictors and the response (i.e., dependent variable) is not assumed to be linear, as GLMs can handle non-linear relationships through the use of link functions. However, the relationship between the linear predictor and the response is assumed to be linear. This linearity assumption is necessary for the estimation of the model parameters using maximum likelihood estimation (MLE), which is a key component of GLMs.

- The linearity assumption in GLMs refers to the linear relationship between the predictors and the linear predictor, rather than the relationship between the predictors and the response. By allowing for non-linear relationships between the predictors and the response, GLMs are able to capture a wide range of relationships in the data, making them a powerful tool in statistical modeling.

What do you mean by the word LINEAR PREDICTOR?

In a Generalized Linear Model (GLM), the linear predictor is a combination of the predictors (i.e., independent variables) weighted by their respective coefficients. The linear predictor is a linear function of the predictors and is used to estimate the expected value of the response variable.

The linear predictor can be expressed mathematically as follows:

$$\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_p x_p$$

where $\eta$ is the linear predictor, $\beta_0$ is the intercept, $\beta_1$, $\beta_2$, ..., $\beta_p$ are the coefficients corresponding to the predictors $x_1$, $x_2$, ..., $x_p$, respectively.

The linear predictor is then transformed using a link function to relate it to the expected value of the response variable. The link function is a non-linear function that maps the linear predictor to the expected value of the response variable, taking into account the distribution of the response variable.

In summary, the linear predictor in a GLM is a linear combination of the predictors, weighted by their coefficients, and is used to estimate the expected value of the response variable. The relationship between the

linear predictor and the response variable is then transformed using a non-linear link function.

| Distribution | Link Function |
|---|---|
| Normal | Identity |
| Binomial | Logit |
| Poisson | Log |
| Gamma | Inverse |
| Inverse Gaussian | Inverse squared |

Can you please provide a table of What link function should be used with respect to the distribution of the response variable?

Note that the choice of link function depends on the distribution of the response variable and can have a significant impact on the performance and interpretability of the model. It is important to choose

an appropriate link function that matches the distribution of the response variable and fits the specific needs of the analysis.

## Are GLM models interpretable?

- Yes, Generalized Linear Models (GLMs) are generally considered to be interpretable, as they provide a clear understanding of how the predictors (i.e., independent variables) influence the response variable.

- In GLMs, the coefficients for each predictor are typically interpreted as the change in the expected value of the response variable associated with a one-unit increase in the corresponding predictor, holding all other predictors constant. Additionally, the p-values associated with each coefficient can be used to assess the significance of each predictor in the model.

- GLMs allow for the examination of the residual deviance and goodness-of-fit statistics to evaluate how well the model fits the data. This can provide insights into whether the model adequately captures the relationships between the predictors and the response variable.

- However, it's important to note that the interpretability of GLMs can be affected by the choice of link function, as well as by interactions between predictors and non-linear relationships between predictors and the response variable. In some cases, the interpretation of the coefficients may be less straightforward. Additionally, GLMs can be limited by the assumptions made about the distribution of the response variable and the linear relationship between the predictors and the response variable.

## When to use OLS, GLS, and GLM in a table

| Situation | OLS | GLS | GLM |
| --- | --- | --- | --- |
| Constant variance errors | Yes | No | Yes |
| Non-normal response variable | No | No | Yes |
| Non-linear relationship between response and predictors | No | No | Yes |
| Heteroscedasticity | No | Yes | No |
| Wide range of response variables and relationships | No | No | Yes |

Is Logistic Regression a strictly linear classifier or a partially linear classifier?

Logistic Regression is a strictly linear classifier, as it models the relationship between the independent variables and the log odds of the dependent variable using a linear function.

But you said logistic Regression falls under GLMs which are not strictly linear.

You are correct that logistic regression is an example of a generalized linear model (GLM), which can model nonlinear relationships between the response and predictors through the use of a link function. However, the logistic regression model itself is a strictly linear classifier, meaning it separates the classes using a linear decision boundary. The nonlinearity comes from the use of a sigmoidal link function that maps the linear predictor to the probability of the positive class. So, in a sense, logistic regression can be considered a partially linear classifier.