# Week 08 Assignment

**Name:** Bala Sai Phani Krishna Yadamreddy

**Introduction:**

In modern data science, generalized linear models (GLMs) are widely used for various regression tasks, offering flexibility in modeling different types of relationships between variables. The optimization methods used to estimate the parameters of GLMs are crucial for ensuring the accuracy and efficiency of these models. Different frameworks and packages, such as **Base R**, **Big Data Version of R**, **Dask-ML**, **Spark R**, and **Scikit-learn**, offer unique approaches to fitting GLMs. The efficiency of these implementations can vary depending on the size of the data, the complexity of the model, and the computational resources available. This assignment explores the optimization methods used in these frameworks, highlighting their strengths and limitations. By examining the details of the algorithms and providing practical examples, we can better understand how to choose the most suitable GLM implementation for different use cases.

| Module/Framework/Package | Name and Description of the Algorithm | An example of a situation where using the provided GLM implementation provides superior performance compared to that of base R or its equivalent in Python (identify the equivalent in Python) |
|---|---|---|
| Base R (stats library) | Iteratively Reweighted Least Squares (IRLS): An optimization technique that iteratively applies weighted least squares to estimate parameters, adjusting weights based on the residuals. | Base R (lm function): Suitable for small to moderate datasets where the computational load is not significant. In Python, the equivalent would be statsmodels for GLMs, which also relies on IRLS but is more memory-intensive in large datasets. |
| Big Data Version of R | Distributed IRLS: Extends IRLS by distributing computations across multiple nodes, handling large-scale data efficiently. | Bigstatsr: Outperforms base R in memory-intensive tasks, particularly for large datasets that do not fit into memory. This is beneficial when compared to R's default GLM implementation or statsmodels in Python when handling very large datasets. |
| Dask-ML | Stochastic Gradient Descent (SGD): An optimization method that approximates the gradient of the loss function using a subset of data, suitable for large-scale data processing. | Dask-ML Logistic Regression: Superior for handling very large datasets in a distributed manner, offering faster performance than base R or Python's scikit-learn when scaling to large datasets. |
| Spark R | Stochastic Gradient Descent (SGD): Similar to Dask-ML, Spark R employs SGD for optimizing GLMs, benefiting | Spark's GLM Implementation: Outperforms base R when dealing with very large-scale data on a cluster, offering more efficient scaling compared to R's base implementation or Python's |

| Module/Framework/Package | Name and Description of the Algorithm | An example of a situation where using the provided GLM implementation provides superior performance compared to that of base R or its equivalent in Python (identify the equivalent in Python) |
|---|---|---|
| | from Spark's distributed computing capabilities. | scikit-learn for datasets that require distributed computing. |
| Scikit-learn | Coordinate Descent and SGD: Utilizes coordinate descent for Lasso regression and SGD for various GLMs, optimizing parameters through iterative updates. | Scikit-learn Logistic Regression: Performs better in large-scale machine learning tasks than R's base GLM implementation, particularly in high-dimensional datasets, as scikit-learn can efficiently handle parallelization and sparse datasets. |

**Conclusion:**

In conclusion, each GLM implementation across different frameworks and packages offers distinct advantages depending on the scale and computational requirements of the task at hand. **Base R** is a solid choice for smaller datasets, using **Iteratively Reweighted Least Squares (IRLS)** for efficient parameter estimation. However, when dealing with large-scale data, frameworks such as **Bigstatsr**, **Dask-ML**, and **Spark R** provide optimized solutions that distribute the computation across multiple nodes, significantly reducing processing time. **Scikit-learn** stands out for its efficient handling of high-dimensional datasets and scalable machine learning tasks. By comparing these methods, we can see that distributed computing and specialized algorithms like **Stochastic Gradient Descent (SGD)** play a crucial role in improving the performance of GLMs in large-scale and resource-intensive environments. The choice of the appropriate framework should be guided by factors like data size, computational resources, and specific performance requirements.

**References:**

- Base R (in the stats library) https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/glmLinks to an external site.
- Big data version of R (look here: https://cran.r-project.org/web/views/HighPerformanceComputing.htmlLinks to an external site.
- Dask ML: https://ml.dask.org/glm.htmlLinks to an external site.
- SparkR: https://spark.apache.org/docs/3.5.0/api/R/reference/spark.glm.htmlLinks to an external site.
- Spark optimization: https://github.com/apache/spark/blob/master/docs/mllib-optimization.mdLinks to an external site.
- Scikit-learn: https://scikit-learn.org/stable/modules/linear_model.htmlLinks to an external site.