

<b>Module/Framework/Package</b>	<b>Name and a Brief Description of the Algorithm</b>	<b>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</b>
<b>Base R (stats package)</b>	Using Iteratively Reweighted Least Squares (IRLS) estimation for GLMs. IRLS is effective for most GLMs but can be limited by computational resources for large data.	When working with a clinical trial dataset (~10,000 patient records) for modeling disease progression with logistic regression, glm in base R will do. But when the dataset grows to the scale of millions of records, computation time and memory become major concerns. The Python counterpart is statsmodels.api.GLM.
<b>High-Performance Computing in R</b>	Packages like biglm and speedglm offer optimized algorithms for large data, say, to recalculate regression coefficients without having to load the entire dataset into memory.	For example, if a bank is screening 100 million credit card transactions to spot fraud using a Poisson regression model, base R's glm would be impossible due to limited memory, but biglm or speedglm would have no problem handling the data since they would process it in chunks. The Python equivalent is statsmodels.GLM and numpy.memmap for large data management.

<b>Dask-ML</b>	Supports scalable GLM algorithms that execute on distributed data structures, leveraging parallelism for large data.	In an e-commerce company's marketing campaign analysis, a logistic regression model makes a prediction of whether a user will click on an ad. The data consist of 500 million observations in a distributed file system. Dask-ML can handle the data in parallel, which brings down computation time by a lot. scikit-learn or base R would be bogged down by memory constraints. Its Python counterpart is <code>scikit-learn.linear_model</code> with <code>joblib</code> used for parallelization.
<b>SparkR</b>	Utilizes Spark's distributed processing to fit GLMs onto large data in a cluster, with multi-family support including binomial, gaussian, and poisson.	For a large transportation network study, it would be impossible in base R to model accident risk for 1 billion GPS-tracked car trips using logistic regression. SparkR enables this with fast distributed computing on a cluster. The Python equivalent is <code>pyspark.ml.classification.GeneralizedLinearRegression</code> .
<b>Spark MLlib</b>	Empirically leverages distributed optimization techniques, that is, limited-memory BFGS (L-BFGS) and stochastic gradient descent (SGD), for GLM fitting in large-scale data.	For real-time risk assessment in an insurance company, predicting policyholder claim probability using a dataset of 500 million customer records requires scalable learning algorithms. Spark MLlib's L-BFGS optimization ensures efficient convergence, whereas traditional glm in R or scikit-learn would be computationally prohibitive. The Python equivalent is <code>pyspark.ml.regression.GeneralizedLinearRegression</code> .

<b>Scikit-learn</b>	Utilizes optimization techniques like coordinate descent and L-BFGS for fitting linear models, with high-performance capability on single-machine data.	When modeling employee attrition in an HR dataset containing 50,000 records, logistic regression in scikit-learn is a good choice due to its efficient implementation of L-BFGS. However, for much larger datasets, Dask-ML or Spark would be preferred. The R equivalent is glm in base R.
---------------------	---	---