

Fitting Generalized Additive Models for very large datasets with Apache Spark

An exposé for a bachelor thesis

Kai Thomas Brusch

05.11.2015

Summary This exposé describes content, goals and motivation and gives a general introduction to my bachelor thesis titled "Fitting General Additive Models for very large datasets with Apache Spark".

Contact `kai.brusch@gmail.com`

Location Hochschule für Angewandte Wissenschaften Hamburg

Department Dept. Informatik

Examiner Prof. Dr. Michael Köhler-Bußmeier

Second examiner TBA

Contents

1	Introduction	3
1.1	Context	3
2	Content	3
2.1	General Linear Models	3
2.2	Generalized Linear Models	4
2.3	Generalized Additive Models	5
2.4	Smoothing parameter Estimation with generalized cross validation	6
2.5	Generalized Additive Models for large data sets	7
2.6	Apache Spark	8
2.7	General additive models in Apache Spark	8
3	Guiding Questions	8
3.1	What is novel in this approach?	9
3.2	When is the thesis regarded as completed?	9
3.3	Which tasks are realistic to be accomplished within the scope of this thesis?	9
4	Time table	9
5	References	10

1 Introduction

1.1 Context

Regression analysis stands at the heart of modern statistical analysis. Discovering the relationship between dependent and independent variables offers great insight into the underlying system and can even provide far reaching predictive capabilities. The traditional linear models explain the dependent variable as the sum of the independent variables and the linear interaction with an estimated coefficient. [Wood, 2006] describes the generalized additive model (GAM) as a powerful extension of the generalized linear model, explaining the relationship between dependent and independent variables as the sum of smoothing functions. This modeling approach offers more flexibility than linear models and has been successfully established in academia and industry. However, the additional flexibility of GAMs comes at the cost of additional computational and memory requirements, placing strict constraints on the model and data size. [Wood et al., 2015] suggest two optimizations to alleviate these constraints: parallelization and numerical optimization. These optimizations have been implemented and used in high-level statistical languages such as R, SPSS and Stata. While these implementations already allow for fitting of larger models with more data they have not been implemented in a truly distributed environment. [Zaharia et al., 2010] introduced Apache Spark which has recently established itself as one of the most promising general purpose cluster computing engines. It's general and distributed nature offers the ideal environment to implement the suggested optimizations by [Wood et al., 2015]. Ideally proving better than previous implementations.

Many relevant concepts of my bachelor thesis have not been part of my bachelor curriculum and are mostly self-taught. I therefore would like to use this exposé to also introduce some of the relevant theory.

2 Content

2.1 General Linear Models

Linear models are statistical models in which an univariate response is modeled as the sum of a 'linear predictor' and a zero mean random error term. The linear predictor depends on some predictor variables, measured with the response variable, and some unknown parameters, which must be estimated. A key feature of linear models is that the linear predictor depends linearly on these parameters. Statistical inference with such models is usually based on the assumption that the response variable has a normal distribution. Linear models are used widely in most branches of science.

$$y_i = \beta x_i + \varepsilon_i \quad (1)$$

Equation 1 introduces formal notation on how to describe the dependent variable y as linear combination of x and the unknown parameter β plus an error term ε . We want β to fit all data as closely as possible and hence introduce S as the squared difference between an estimated $x_i\beta$ and y_i

$$S = \sum_{i=1}^n (y_i - x_i\beta)^2 \quad (2)$$

A good choice of β should bring S close to 0. The Markov-Gauss Theorem states that the minimization of S yields $\hat{\beta}$ which is the best possible estimation for b . 1 represents the univariate case where y is explained with only one variable. The process of minimizing S is commonly referred

to as ordinary least squares (OLS). This model can be extended to multiple independent variables yielding in a scalar matrix form.

$$y = \mathbf{X}\beta \quad (3)$$

The dependent variable vector y is the linear combination of model matrix X and the vector of unknown parameter β . Finding the best possible estimator $\hat{\beta}$ for 3 is formalized stated as the minimal length of the distance between y and $\mathbf{X}\beta$:

$$\|y - \mathbf{X}\beta\|^2 \quad (4)$$

Minimizing 4 will yield in the best possible estimator $\hat{\beta}$. OLS stands at the heart of many statistical models and will be relevant later in this exposé, it is important to note however that OLS only works to estimate $\hat{\beta}$ for normally distributed data.

2.2 Generalized Linear Models

Generalized linear models (GLMs) relax the strict linearity assumption of linear models, by allowing the expected value of the response to depend on a smooth monotonic function of the linear predictor. Similarly the assumption that the response is normally distributed is relaxed by allowing it to follow any distribution from the exponential family (normal, Poisson, Binomial, Gamma etc.). While OLS was sufficient for estimating β for normal distributed data we have to generalize this notion to account for an arbitrary amount of distribution parameters. Instead of looking for a β that will minimize the squared error we are now looking for a set of parameters given an exponential distribution that fits our data the best. This process of finding the best parameters for a given distribution and data points is called maximum likelihood estimation (MLE). A maximum likelihood estimate, estimates some hidden parameter λ for identical, independent distributed data X_1, \dots, X_i for a given probability density distribution of the exponential family. Let a given distribution be governed by a probability density function $f(x; \lambda_1, \dots, \lambda_i)$, we can then formally describe MLE as:

$$L(X_1, \dots, X_i) = \prod_{i=1}^n f(X_i; \lambda_1, \dots, \lambda_i) \quad (5)$$

MLE is the essential method in estimation GLMs and finding the best model parameter λ stands at the core of fitting any statistical model beyond simple linear ones. Fitting models via MLE involves iterating over possible λ values to maximize the likelihood. With some work we can state MLE as a general case of OLS. Modern applications use a form of the iterative reweighted least squares method in conjunction with a Gauss-Newton method to maximize the likelihood. Details of these methods shall be discussed in my bachelor thesis.

$$g(\mu_i) = \mathbf{X}\beta \quad (6)$$

Equation 6 introduces formal notation for GLMs and illustrates the similarities to the general linear model. The change of parameter estimation method does not become obvious in the usual GLM notation but is essential. It is important to recognize the smooth monotonic function $g(\cdot)$. The smooth monotonic function, also known as link function, is the key extension which enables GLMs to model members of the exponential family.

2.3 Generalized Additive Models

Generalized Additive Models (GAMs) extends the GLM by specifying the linear prediction in terms of the summation of smooth functions. This allows for a more flexible modeling of the influence for each explanatory variable. The gained flexibility comes at the cost of additional questions concerning the smooth function:

- The smooth functions must be represented somehow.
- The degree of smoothness of the functions must be made controllable.
- Some means for estimating the most appropriate degree of smoothness from data is required.

GAMs are formally described by the following equation:

$$g(\mu_i) = \mathbf{X}_i\Theta + f_1(x_{1i}) + f_2(x_{2i}) + f_3(x_{3i}, x_{4i})\dots \quad (7)$$

7 explains y_i as the model matrix for this row and the smooth functions $f_j(x_{1j})$ of the x values for this row. X_i is a row of the model matrix with parametric component θ . Unlike the linear model we can now specify a smooth function for each explanatory variable. This proves to be way more flexible than only allowing for a constant influence per explanatory variable. The natural question that arises now are: How do I find proper smoothing functions? Finding the right smooth function stands at the heart of GAM fitting and can be best illustrated in the univariate case 8.

$$y_i = f(x_i) + \epsilon_i \quad (8)$$

Smooth functions form a vector space, which can be approximated using a linear basis. Only allowing linear basis allow us to heavily leverage the theory already developed for linear models and S as the optimal model fit. For the sake of illustration we assume that 8 can be rewritten as the following equation if $b_i(x)$ is the i th basis function:

$$f(x) = \sum_{i=1}^q b_i(x)\beta_i \quad (9)$$

In 9 we already know $f()$ is linear in regard to 9. We now have to specify a basis function to represent b_i . We can choose from many basis functions for b_i , each with advantages and disadvantages. A common choice however is a fourth order polynomial basis function. 9 represented by a fourth order polynomial yields the following model:

$$f(x) = \beta_1 + x\beta_2 + x^2\beta_3 + x^3\beta_4 + x^4\beta_5 \quad (10)$$

Applying 10 to 8 we get the modeling of y_i as the sum of smoothing functions.

$$y_i = \beta_1 + x_i\beta_2 + x_i^2\beta_3 + x_i^3\beta_4 + x_i^4\beta_5 + \epsilon_i \quad (11)$$

Equation 11 gives a formal full description of y_i as the sum of a fourth degree polynomial. Fourth order polynomial basis does a fine job of illustrating the used concepts but enforce rather strict limitations. Instead of trying to fit a single fourth order polynomial to all our data points we will now try to represent f with cubic splines. This approach is called regression splines and follows this general pattern:

- Divide the curve in to a fixed number of sections, the point between two sections is called a knot.
- Find a cubic polynomial for each section between two knots

- Each cubic polynomial must match it's neighbors in value, first and second derivative at the location of the knot.

Using cubic splines as a basis for f means that 8 becomes a linear model identical to 3. We can hence rely well established linear model theory to fit a good polynomial. The number and location of knots will influence the fit of our splines and are carefully chosen. The right number and location of knots can be estimated but are usually a problem specific design decision. After establishing how to represent the smooth functions and the use of knots we will now turn to the smoothing parameter λ . λ is a penalty factor that controls the smoothness by scaling the second derivative of the smooth function, this can be formalized by:

$$\|y - X\beta\|^2 + \lambda \int_0^1 [f''(x)]^2 dx \quad (12)$$

The integrated square root of second derivative penalizes models that have a too 'rough' function. As λ approaches 0 f will become a straight regression line while λ of 0 becomes an unpenalized regression spline. The trade off between model fit and model smoothness is controlled by the smoothing parameter λ . The linear nature of f allows us to rewrite 12 in a more concise form.

$$\|y - X\beta\|^2 + \lambda \beta^T S \beta \quad (13)$$

13 shows that describing an unknown variable as the sum of smooth functions can be restated as the minimization of penalized regression splines and hence becomes problem of minimizing equation 13. While the the problem of finding a smooth function for a univariate function has almost been solved by 13 we are still left with an unknown parameter of λ . Finding the optimal λ is subject of the following section. It is important to remember that this λ is not the distribution parameter from MLE but the smoothing parameter λ for GAMs.

2.4 Smoothing parameter Estimation with generalized cross validation

Finding a good smoothing parameter through minimizing 13 the essential question for GAM research. Finding an optimal choice for λ means to trade goodness of fit to smoothness. A high λ will overfit the data while a low λ will create undersmoothed splines. The common approach is to define a measure for the predictive capabilities of our estimated spline. For our purposes we will use the generalized cross validation (GCV) which formalizes the notion that an ideal splines minimizes the average distance between our splines and the function f . Another interpretation is to see the GCV as the prediction error and that finding the best smoothing parameter is equivalent to find the smoothing parameter with the minimal prediction error. A formal definition of the GCV score yields:

$$V_g(\lambda) = \frac{n \|y - X\hat{\beta}_\lambda\|^2}{\{n - \text{tr}(F_\lambda)\}^2} \quad (14)$$

14 formalizes that the generalized cross validation score for a given λ , formally $V_g(\lambda)$ is n times the squared distance for a currently estimated parameter $\hat{\beta}$ given a λ . This distance get put in perspective to degrees of freedom for the current λ : $\{n - \text{tr}(F_\lambda)\}^2$. Hence the GCV score expresses goodness of fit as the tradeoff between fitting the points and degrees of freedom for a given λ . Finding λ is done via penalized likelihood iteratively reweighted least squares (PIRLS) which proceeds as follows, where V is the function such that $\text{var}(y_i) = \phi V(\mu_i)$ and $\mu_i = E(y_i)$. First initialize $\hat{\mu}_i = y_i + \epsilon_i$ and $\eta_i = g(\mu_i)$ where ϵ_i is a small quantity (often 0) added to ensure that $g(\mu_i)$ exists.

Then iterate the following steps to convergence.

1. form $z_i = g'(\hat{\mu}_i)(y_i - \hat{\mu}_i) + \eta_i$ and $w_i = V(\hat{\mu}_i)^{1/2} g'(\hat{\mu}_i)^{-1}$.
2. putting the w_i in a diagonal matrix W , minimize the weighted version of expression 13:

$$\|Wz - WX\beta\|^2 + \lambda\beta^T S\beta \quad (15)$$

with respect to β to obtain $\hat{\beta}_i$ and the updates $\hat{\eta}_i = X\hat{\beta}_i$, and $\hat{\mu}_i = g^{-1}(\hat{\eta}_i)$.

For moderate size data PIRLS will converge for each trial λ . Iterating 15 to obtain the new estimates can place a heavy burden on the memory requirements. Each step of the PIRLS requires to form the whole model matrix X . The difficulty is simply that the model matrix for the model can become too big: if n and p are respectively the number of rows and columns of the model matrix, and M is the number of smoothing parameters, then the memory requirements of GAM fitting methods are typically $O(Mnp^2)$. The aim of my bachelor thesis is to discover and implement the in [Wood et al., 2015] suggested optimizations for PIRLS. Though I am able to reproduce the wording of this method I currently do not understand this method in detail, understanding this method will be one of the core tasks in the next weeks, table 1 for reference.

2.5 Generalized Additive Models for large data sets

[Wood et al., 2015] introduces two extensions of the development methods for fitting GAMs. Woods methods aim at numerical optimization and parallelization to eliviate the rescriptions on data size.

The numerical optimization propossed by [Wood et al., 2015] rewrites the general cross validation step to not require the full model matrix X in each step. Woods describes a strategy on how model matrix factorization can be used to avoid formation of the whole model matrix in each iteration.

Suppose that the model matrix is first QR decomposed into a column orthogonal $n * p$ factor Q and an upper triangular $p * p$ factor R so that $X = QR$: If we also form $f = Q^T y$ and then expression 13 becomes

$$\|f - R\beta\|^2 + \|r\|^2 + \sum_j \lambda_j \beta^T S_j \beta \quad (16)$$

16 can be rewritten in a form that does not require the full model matrix X .

$$V_g(\lambda) = \frac{n \|f - R\hat{\beta}_\lambda\|^2 + \|r\|^2}{\{n - tr(F_\lambda)\}^2} \quad (17)$$

17 brings a tramendous advantage over 14 in terms of memory and computational usage. Receiving all relevant information to find λ for a given X only forming subsections $R f$ and $\|r\|^2$ gives a significant advantage over the proposed method in 14.

The parallelization optimimiztion of 14 envolves the computations of $R f$ and $\|r\|^2$ in a way that only requires small, disjunct subsections of X . This allows us to distribute the computation by grouping X into equally sized non-overlapping sets and allocate them to different processors. I do not understand the details of this optimization at this moment, understanding them will be part of my studies in the upcoming weeks.

2.6 Apache Spark

Distributed computing is a special interest of mine. I've had the pleasure of setting up a SparkR cluster and have been following its development closely. [Zaharia et al., 2010] advertises Spark as a fast and general cluster computing engine for large-scale data processing. Here a overview of it's key features:

- fast, Apache Spark currently holds the world record in sorting 1 TB on data, beating the previous record held by Hadoop. Spark reaches this high speed by removing I/O to disc and maintaining all relevant in memory.
- general, Apache Spark offers a high level API called Resilient Data Set. RDDs are designed to facilitate a resilient, distributed data set as a optimal level of interaction with data. Each of the key features of RDDs are outlined below:
 - Resilient: Every RDD is evaluated lazy and constructed from it's lineage, point to different previous RDD.
 - Distributed: RDDs are partiioned accross many machines
 - Data set: RDDs are ment to think about in terms of matrixes, this should feel very natural for people familar with data.
- cluster computing, Apache Spark is build to be executed in a cluster environment like Mesos or Yarn. The management and design of these systems are not subject of my bachelor thesis.
- large scale data processing, the data managable by Apache Spark is limited by the amount of memory available in a cluster.

2.7 General additive models in Apache Spark

I strongly believe that Apache Spark offers an ideal environment for gitting GAMs for large data sets. Spark itself offers high-level APIs for R, Pyhton and Scala. Implementing the suggested methods in Scala would allow them to be used in all three languages. The optimizations suggested by [Wood et al., 2015] are optimal for an general purpose in-cluster computing engine for the following reasons:

- The iterative nature of PIRLS is optimal for Spark's in memory computing.
- Spark's high level RDD allows for a flexible data structure to implement the suggested methods.
- Spark is designed to facilitate the parallel nature of finding Rf and $\|r\|^2$ in subblocks of X .
- Spark supports [Hall, 2009], a linear algebra library for Scala with far reaching capabilities.

3 Guiding Questions

The leading question for my bachelor thesis can be stated as:

Can the distributed, numerically optimized, in-memory version of GAM fitting out-perform previous benchmark implementation in model size, data size and fitting speed?

Subquestions arising are:

- Can the suggested methods be implemented in Apache Spark?
- How much faster is a parallel variant against the sequential one?
- How much faster is the Spark in-memory variant against the R one?
- How do the limit for data and model size change?
- How much does horizontal scaling effect performance?

3.1 What is novel in this approach?

The suggested optimizations have been implemented in high-level statistical languages, proving its effectiveness. The suggested methods have not been implemented in an environment designed for in-memory cluster computing. As I lined out earlier the iterative nature of the proposed methods is well suited for a distributed in-memory environment and hence should be discovered

3.2 When is the thesis regarded as completed?

The thesis is considered finished, once the optimized version of GAM fitting has been implemented in Spark made available in SparkR and compared to other implementation of GAMs.

Additional tasks

- Develop a deep theoretical understanding of GAMs with [Wood, 2006]
- Develop a better understanding of used matrix decompositions with [Strang, 2009]
- Set up an Apache Spark environment.
- Familiarize with Breeze in Apache Spark.
- Coordinate with FSB and my employer to request my employer as second examiner

3.3 Which tasks are realistic to be accomplished within the scope of this thesis?

I strongly believe that the suggested methods can be implemented in Apache Spark within this thesis. In an ideal situation I would like to make my implementation available to the Apache project.

4 Time table

My time plan is visible at table 1.

Table 1: Time table

CW	thesis work
48	Develop a deep theoretical understanding relevant linear algebra
49	Develop a deep theoretical understanding of relevant linear algebra
50	Develop a deep theoretical understanding of GAMs and relevant linear algebra methods
51	Develop a deep theoretical understanding of GAMs and relevant linear algebra methods
52	Winter break
1	Implementation of matrix decomposition and Exams
2	Implementation of spline fitting and Exams
3	Implementation and Exams
4	Implementation and Exams
5	Post Exams Break
6	Start implementing GAM fitting
7	Implementation of GAM fitting
8	Implementation of GAM fitting and run benchmark tests
9	Run benchmarks tests
10	Buffer
11	writing bachelor thesis
12	writing bachelor thesis
13	writing bachelor thesis
14	writing bachelor thesis
15	writing bachelor thesis
16	Done

5 References

- [Hall, 2009] Hall, D. (2009). Breeze.
- [Strang, 2009] Strang, G. (2009). *Linear Algebra and Its Applications*. Wellesley-Cambridge Press, Wellesley, MA.
- [Wood, 2006] Wood, S. (2006). *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC.
- [Wood et al., 2015] Wood, S. N., Goude, Y., and Shaw, S. (2015). Generalized additive models for large data sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 64(1):139–155.
- [Zaharia et al., 2010] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., and Stoica, I. (2010). Spark: Cluster computing with working sets. In *Proceedings of the 2Nd USENIX Conference on Hot Topics in Cloud Computing*, HotCloud’10, pages 10–10, Berkeley, CA, USA. USENIX Association.