

Adaptive Basis Function Construction: An Approach for Adaptive Building of Sparse Polynomial Regression Models

Gints Jekabsons
Riga Technical University
Latvia

1. Introduction

The task of learning useful models from available data is common in virtually all fields of science, engineering, and finance. The goal of the learning task is to estimate unknown (input, output) dependency (or model) from training data (consisting of a finite number of samples) with good prediction (generalization) capabilities for future (test) data (Cherkassky & Mulier, 2007; Hastie et al., 2003). One of the specific learning tasks is regression – estimating an unknown real-valued function. The process of regression model learning is also called regression modelling or regression model building.

Many practical regression modelling methods use basis function representation – these are also called dictionary methods (Friedman, 1994; Cherkassky & Mulier, 2007; Hastie et al., 2003), where a particular type of chosen basis functions constitutes a “dictionary”. Further distinction is then made between non-adaptive methods and adaptive (also called flexible) methods.

The most widely used form of basis function expansions is polynomial of a fixed degree. If a model always includes a fixed (predetermined) set of basis functions (i.e. they are not adapted to training data), the modelling method is considered non-adaptive (Cherkassky & Mulier, 2007; Hastie et al., 2003). Using adaptive modelling methods however the basis functions themselves are adapted to data (by employing some kind of search mechanism). This includes methods where the restriction of fixed polynomial degree is removed and the model’s degree now becomes another parameter to fit. Adaptive methods use a very wide dictionary of candidate basis functions and can, in principle, approximate any continuous function with a pre-specified accuracy. This is also known as the universal approximation property (Kolmogorov & Fomin, 1975, Cherkassky & Mulier, 2007).

However, in polynomial regression the increase in the model’s degree leads to exponential growth of the number of basis functions in the model (Cherkassky & Mulier, 2007; Hastie et al., 2003). With finite training data, the number of basis functions along with the number of model’s parameters (coefficients) quickly exceeds the number of data samples, making model’s parameter estimation impossible. Additionally the model should not be overly

complex even if the number of its basis functions is lower than the number of data samples, as too complex models will overfit the data and produce large prediction errors.

To obtain a polynomial regression model that does not overfit (nor underfit) and describes the relations in data sufficiently well, typically the subset selection approach (Hastie et al., 2003; Reunanen, 2006) is used where the goal is from a fixed full predetermined dictionary of basis functions to find a subset which corresponds to a model (a sparse polynomial) with the best predictive performance. This is done via combinatorial optimization. However, for the subset selection approach still the two issues remain – deficiency of adaptation as well as computational inefficiency.

Searching through all the possible combinations of basis functions takes double-exponential runtime as the number of combinations grows exponentially in the number of basis functions of the predetermined dictionary while the number of the basis functions in the dictionary grows exponentially in the number of input variables and “full” model’s degree (Hastie et al., 2003). This makes the exhaustive search through all the combinations impractical. The heuristic greedy search algorithms, such as forward selection (Hastie et al., 2003; Reunanen, 2006), substantially reduce the time and make it practical for not too large number of input variables and not too high degree. Nevertheless, the search time actually is still exponential, hindering their use in problems of larger dimensionality and hindering the removal of the restriction of a fixed degree.

The approach of subset selection assumes that the chosen fixed finite dictionary of the predefined basis functions contains a subset that is sufficient to describe the target relation sufficiently well. However, in most practical situations the required dictionary (and “full” model’s degree) is not known beforehand and needs to be either guessed or found by an additional search loop over the whole model building process, since it will differ from one regression task to another. In many cases, especially when the studied data dependencies are complex and not well studied, this means either a non-trivial and long trial-and-error process or acceptance of a possibly inadequate model.

This chapter presents a sparse polynomial regression model building approach which enables adaptive model building without restrictions on model’s degree and does it in polynomial time instead of exponential time (in the number of input variables, required degree, and target model’s complexity) as well as without the requirement to repeat the model building process. The required basis functions are automatically iteratively constructed using heuristic search specifically for the particular data at hand instead of choosing a subset from a very restricted finite user-defined dictionary (hence the approach is called Adaptive Basis Function Construction, ABFC). The basis function dictionary now becomes infinite and polynomials of arbitrary complexity can be generated bringing the desired flexibility to the model building process.

The remainder of this chapter is organized as follows. The next two sections give brief overview of polynomial regression and the subset selection approach. In Section 4 the ABFC approach is described. Section 5 outlines the related work. The results of the empirical evaluations of the proposed methods and their comparison to other well-known regression modelling methods are presented in Section 6. Section 7 concludes this chapter.

2. Polynomial regression

In standard regression formulation (Vapnik, 1995; Cherkassky & Mulier, 2007; Hastie et al., 2003) the goal is to estimate unknown real-valued function in the relationship

$$y = G(x) + \varepsilon, \tag{1}$$

where ε is independent and identically distributed random noise with zero mean, $x = (x_1, x_2, \dots, x_d)$ is d -dimensional input, and y is scalar output. The estimation is made based on a finite number of samples (training data) provided in form of matrix \mathbf{x} of input values for each sample and vector \mathbf{y} of output values for each corresponding sample. Using the finite number n of training samples (\mathbf{x}_j, y_j) , $j = 1, 2, \dots, n$ one wants to build a model F that allows predicting the output values for yet unseen input values as closely as possible.

Generally, a linear regression model may be defined as a linear expansion of basis functions:

$$F(x) = \sum_{i=1}^k a_i f_i(x), \tag{2}$$

where $\mathbf{a} = (a_1, a_2, \dots, a_k)^T$ are model's parameters, k is the number of basis functions included in the model (equal to the number of model's parameters), and $f_i(x)$, $i = 1, 2, \dots, k$ are the included basis functions of the input x . As the model is linear in the parameters, the estimation of its parameters is typically done using the Ordinary Least-Squares (OLS) method (Hastie et al., 2003) minimizing the squared-error:

$$\mathbf{a} = \arg \min_{\mathbf{a}} \sum_{j=1}^n (y_j - F(\mathbf{x}_j))^2. \tag{3}$$

The basis function representation enables moving beyond pure linearity, by defining nonlinear transformations of x while still working with linear models (and employing OLS). For example, for $d = 1$ a polynomial model of fixed degree p can be defined as follows:

$$F(x) = \sum_{i=0}^p a_i x^i. \tag{4}$$

Generally for a given d and p the total number of basis functions in a "full" polynomial, i.e. the total number of basis functions in the dictionary, is

$$m = \prod_{i=1}^p (1 + d/i). \tag{5}$$

3. Subset selection

Models which are too complex (i.e. that fit the training data too well causing overfitting) or too simple (i.e. that fit the data poorly causing underfitting) provide poor predictive performance for the future data. The most popular approach of controlling model's complexity is subset selection. The goal of subset selection is from a fixed full predetermined dictionary of basis functions to find a subset that provides the best predictive performance of the model (Hastie et al., 2003; Reunanen, 2006). Now in addition to the estimation of model's parameters, the structure of the model itself needs to be found.

The total number of possible subsets from a dictionary of size m is 2^m . This means that searching through all the possible subsets is in most cases impractical. Hence in subset selection heuristic search algorithms are used. They efficiently traverse the space of subsets, by adding and deleting basis functions of the model, and use model evaluation measure to direct the search into areas of increased performance. The typical examples of heuristic search algorithms are the greedy hill-climbing algorithms – Forward Selection (also known as Sequential Forward Selection, SFS) and Backward Elimination (also known as Sequential Backward Selection, SBS) (Hastie et al., 2003; Reunanen, 2006). However, there exist also more recently developed search strategies, such as Beam Search, Floating Search, Simulated Annealing, Genetic Algorithms etc. (Reunanen, 2006; Pudil et al., 1994; Russel & Norvig, 2002).

Summarizing (Russel & Norvig, 2002; Molina et al., 2002; Kohavi & John, 1997), in order to characterize a heuristic search problem one must define the following: 1) initial state of the search; 2) available state-transition operators; 3) search strategy; 4) evaluation measure; 5) termination condition. Note that in the context of model building the “initial state” is also called “initial model” and the “state-transition operators” are also called “model refinement operators”.

In the subset selection approach for polynomial regression, typically the *initial state* is the model that corresponds to the empty subset, the subset with only the intercept term in it, full set of all the defined basis functions, or a randomly chosen subset. The typical *basic state-transition operators* are addition and deletion of a basis function. The typical *search strategy* is the hill-climbing (Russel & Norvig, 2002) which, in combination with the empty (or sufficiently small) subset as initial state and the addition operator, becomes SFS, but, in combination with the full subset as initial state and the deletion operator, becomes SBS. As the *evaluation measures* classically the statistical significance tests are used (Hastie et al., 2003; Dreyfus & Guyon, 2006). However, in model building currently two other strategies predominate (Cherkassky & Mulier, 2007; Dreyfus & Guyon, 2006): employment of complexity penalization criteria (also known as analytical criteria), e.g., the well-known Akaike's Information Criterion, AIC (Akaike, 1974; Burnham & Anderson, 2002), and the resampling techniques, e.g., Hold-Out, Cross-Validation (CV), and Bootstrap (Kohavi, 1995; Hastie et al., 2003; Dreyfus & Guyon, 2006). The *termination condition* typically corresponds to finding of a state in that none of the state-transition operators can lead to a better state (i.e. a local minimum).

In polynomial regression, increase in the model's degree leads to exponential growth of the number of basis functions in the dictionary, i.e. $O(m) = O(d^p)$ (Cherkassky & Mulier, 2007; Hastie et al., 2003) and to double-exponential growth of the number of all possible subsets (or the number of states in the state space): $O(2^m) = O(2^{d^p})$. When using one or both of the

two basic state-transition operators, the order of the branching factor of a state in the state space in the very first iteration of the search is already equal to the number of basis functions in the dictionary, i.e. it also increases exponentially.

Assuming that the “best” model found in the search process includes a total of k_* basis functions and that in each iteration the number of basis functions of the current model is increased by 1, the total number of evaluated models (subsets) is of order

$$O\left(\sum_{i=1}^{k_*} d^i\right) = O(d^p k_*). \quad (6)$$

Hence for larger values of d and p (e.g., when m reaches thousands) subset selection is rendered impractical. Additionally, because of the branching factor’s direct dependence on the number of basis functions in the dictionary, the idea of unrestricted degree (i.e. dictionary of infinite size) is hardly applicable.

The computational problem could be somewhat reduced by choosing a sufficiently small but useful value of p before the actual model building is performed. However, generally the required maximal degree is not known beforehand and needs to be either guessed or found by additional search loop over the whole model building process, since it will differ from one regression task to another, which means either a non-trivial and long trial-and-error process or acceptance of a possibly inadequate model.

4. Adaptive Basis Function Construction

This section introduces Adaptive Basis Function Construction – a possible alternative to the classical subset selection approach. The goal of the ABFC approach is to overcome some of the limitations associated with the subset selection, outlined in the previous section. The ABFC approach is developed for sparse polynomial regression model building without restrictions on model’s degree, enables model building in polynomial time, and does not require repetition of the building process (in contrast to the subset selection approach). The required basis functions are automatically adaptively constructed specifically for data at hand, without using a restricted fixed finite user-defined dictionary. The dictionary in the ABFC is infinite and polynomials of arbitrary complexity can be constructed.

4.1 The models and the basis functions

Generally, a basis function in a polynomial regression model can be defined as a product of original input variables each with an individual exponent:

$$f_i(x) = \prod_{j=1}^d x_j^{r_{ij}}, \quad (7)$$

where \mathbf{r} is a $k \times d$ matrix of nonnegative integer exponents such that r_{ij} is the exponent of the j th variable in the i th basis function. Note that, when for a particular i th basis function $r_{ij} = 0$ for all j , the basis function is the intercept term.

Given a number of input variables d , matrix \mathbf{r} , with a specified number of rows k and with specified values for each of its elements, completely defines the structure of a polynomial model with all its basis functions. The set of basis functions, included in a model, is then

$$f = \left\{ \prod_{j=1}^d x_j^{r_{ij}} \mid i = 1, 2, \dots, k \right\}. \quad (8)$$

For example, if $d = 3$ and $k = 4$, then the matrix

$$\mathbf{r} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 3 \\ 1 & 1 & 1 \end{bmatrix}, \quad (9)$$

corresponds to the set

$$f = \{x_1^0 x_2^0 x_3^0, x_1^1 x_2^0 x_3^0, x_1^0 x_2^1 x_3^3, x_1^1 x_2^1 x_3^1\} = \{1, x_1, x_2 x_3^3, x_1 x_2 x_3\}, \quad (10)$$

which in turn corresponds to the model

$$F(x) = a_1 + a_2 x_1 + a_3 x_2 x_3^3 + a_4 x_1 x_2 x_3. \quad (11)$$

Formally, the problem of finding the best set of basis functions can be defined as finding the best matrix \mathbf{r} with the best combination of nonnegative integer values of its elements:

$$\mathbf{r}^* = \arg \min_{\mathbf{r}} J \left(\left\{ \prod_{j=1}^d x_j^{r_{ij}} \mid i = 1, 2, \dots, k \right\} \right), \quad (12)$$

where $J(\cdot)$ is an evaluation criterion that evaluates the predictive performance of the regression model which corresponds to the set of basis functions.

As neither the upper bounds of \mathbf{r} elements' values nor the upper bound of k are defined, it is possible to generate sparse polynomials of arbitrary complexity, i.e. of arbitrary number of basis functions each with an arbitrary exponent for each input variable. This also means that the searchable state space is infinite.

4.2 The search process

Finding the "best" structure of matrix \mathbf{r} requires search. In this section the five components (outlined in Section 3) of a heuristic search problem are analyzed in the context of the ABFC approach.

Initial state. In ABFC, the state space is infinite therefore a natural initial state of the search is the state that corresponds to the simplest model located in the space. In the current study it is assumed that the simplest model is the one with a single basis function corresponding to

the intercept term. However, also other models could be used as initial states, e.g., an empty model (without any basis functions), a first degree “full” polynomial, or a small randomly generated model. Note that in the current study the basis function corresponding to the intercept term stays in the model at all times and is not allowed to be modified or deleted.

State-transition operators. Using efficient state-transition operators is vital for the search process to be efficient. The employed state-transition operators are the main methodological difference between the subset selection approach and the ABFC approach. Generally, there are two different basic types of modifications to an existing polynomial model: complication and simplification (Jekabsons & Lavendels, 2008a). In the subset selection approach, these are the addition and deletion operators. The addition operator makes the model more complex (by adding a new basis function) but the deletion operator makes it simpler (by deleting an existing basis function).

In the ABFC, the two standard operators from subset selection are replaced with other operators that not only add or delete basis functions but also work on the level of individual exponents, modifying the existing basis functions and creating modified copies of them. The basic idea is to use an operator that adds only the simplest (i.e. linear) basis functions which serve as a basic material for further construction of more complex functions using other operators. In this manner there is no need for an operator that explicitly tries to add basis functions of each possible combination of exponent values (as the addition operator in the subset selection). Hence the branching factor of the state space stays not only finite but also relatively small while the state space itself is infinite.

In this study, a set of the following four state-transition operators for the polynomial regression model building are proposed. Operator1: Addition of a new linear basis function with one of its exponents set to one and all the others set to zero. Operator2: Addition of an exact copy of an already existing (in the current model) basis function with one of its exponents increased by 1. Operator3: Decreasing of one of the exponents in one of the existing basis functions by 1. Operator4: Deleting of one of the existing basis functions. Figure 1 gives examples of the operators operating on a simple matrix.

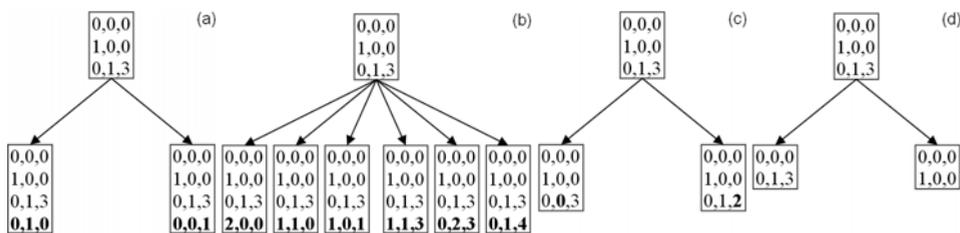


Fig. 1. Example of the four state-transition operators operating on a simple matrix: (a) Operator1; (b) Operator2; (c) Operator3; (d) Operator4

The set of the four state-transition operators is sufficient to generate any polynomial model definable by the matrix r . Their use can also be viewed as a piece of application-domain knowledge. While starting the search from the simplest model, the complication operators (the first two) do the main job - they “grow” the model. The simplification operators (the last two), on the other hand, work as “purifiers” - they decrease the unnecessarily high exponents and delete the unnecessary basis functions. Without the use of simplification operators, a regression model may contain unnecessarily high exponents and include too

many unnecessary basis functions, at the same time preventing truly necessary modifications (this is also known as the nesting effect (Pudil et al., 1994)) and increasing overfitting. Additionally, for all the state-transition operators a special care is taken to prevent basis function duplicates in the resulting model as well as to preserve the intercept term.

The initial state and the state-transition operators together form a state space. Figure 2 shows a small example of a state space in ABFC when the number of input variables is three and all the four state-transition operators are used. Each state represents a set of basis functions included in the regression model. The ordering of the states in the space is such that the simplest models and the simplest basis functions are reached first and, as the search goes on, increasingly complex models and basis functions can be reached.

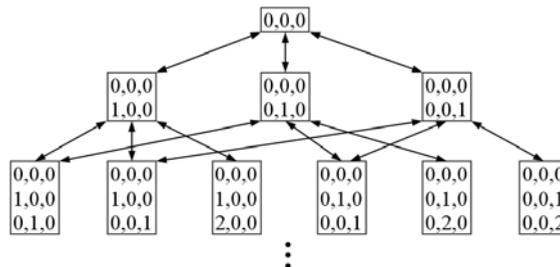


Fig. 2. A small example of the first three layers of a state space in ABFC when $d = 3$ (the space is infinite in the direction of more complex models)

In the Section 3, it is stated that in the subset selection approach the branching factor of a state in the state space increases exponentially with respect to the number of input variables d and pre-specified maximal degree p . In ABFC, the branching factor of the current state in the state space depends on d and on the number of basis functions k , already included in the current model. The upper bound of the number of possible modifications to a model using Operator1 is equal to d ; using Operator2 and Operator3 it is equal to dk ; and using Operator4 it is equal to k . So the upper bound of the branching factor is of order $O(d + 2dk + k) = O(dk)$ that is linear in respect to both d and k .

Search strategy. Most of the heuristic search algorithms of the hill-climbing type can be divided in two categories: those that assume the model state-transition operators to be of either or both the forward and the backward type (e.g., SFS, SBS, and Floating Search algorithms) and those that do not distinguish between the two types (e.g., Steepest Descent Hill-Climbing and Simulated Annealing). The four operators proposed in this study are naturally divided in forward (complication) and backward (simplification) operators; therefore in ABFC both categories of the search algorithms can be applied.

On the other hand, non-hill-climbing search algorithms, e.g., Genetic Algorithms and the like, employ completely different kind of operators (i.e. Crossover and Mutation). While they could be adapted to work with the infinite dictionary of basis functions, their major disadvantage is that, in contrast to the simple hill-climbing algorithms, they are not generally biased towards simpler models. In large state spaces they often spend most of the time exploring too complex models while the "best" ones are in fact mostly the relatively simple ones.

Evaluation measure. The proposed state-transition operators allow using the same methods for model evaluation and comparison as those used in subset selection. However, note that the model complexity penalization criteria, in contrast to the resampling techniques, usually require substantially lower computational resources as well as are less noisy creating less local minima in the state space.

Termination condition. Many different termination conditions can be used to terminate the search process. Some of most widely used ones are the following: a) a user pre-specified number of iterations is reached; b) a user pre-specified size of the model is reached; c) using the available state-transition operators the model could not be improved any further (evaluated by the chosen evaluation measure). The first two termination conditions require the user to set a hyperparameter value. This is a non-trivial task as usually the required information is not available. Adjusting such a hyperparameter may also require too large amounts of computational resources. In this study, the termination condition listed here as the last (c) is employed.

4.3 A concrete practical model building method

This section proposes Floating Adaptive Basis Function Construction (F-ABFC) – a concrete practical polynomial regression model building method, which is a special case of the ABFC approach.

The search procedure of the F-ABFC starts with the simplest model (with only the intercept term included) and uses the Floating Search strategy (hence the name of the method), in particular the Sequential Floating Forward Selection algorithm, SFFS (Pudil et al., 1994), together with the set of the four state-transition operators proposed in the previous section.

In SFFS, the search process consists of two phases – the forward phase and the backward phase. In each iteration of the search, the forward phase is done only once but the number of times the backward phase is performed is determined dynamically. In the forward phase, all the models, which can be generated using the complication operators on the current best model, are evaluated and, if there is improvement over the current best model, the best of the new models is chosen as the new current best model and the search proceeds to the second phase. If there is no improvement, the whole search procedure is stopped. In the backward phase, on the other hand, all the models, which can be generated using the simplification operators on the current best model, are evaluated. In this phase ever simpler models are repeatedly generated and the phase is ended only when, using the available simplification operators, it is impossible to generate a model which is better than the current best one. After the second phase, the search process always proceeds to the next iteration (starting again with the first phase).

According to the studies of many researchers, the Floating Search algorithms, including SFFS, are some of the most efficient heuristic search algorithms for deterministic combinatorial optimization in terms of both required computational resources and quality of the results (Ferri et al., 1994; Jain & Zongker, 1997; Jain et al., 2000; Zongker & Jain, 1996; Pudil et al., 1994; Kudo & Sklansky, 2000; Reunanen, 2006). SFFS also does not have any adjustable hyperparameters, has a tendency to generate simpler models than many other algorithms, and is very simple to implement.

As in (Jekabsons & Lavendels, 2008a; Jekabsons, 2008), to evaluate the predictive performance of a newly generated model, to perform model comparisons, and to steer the

search in direction of the most promising models, in F-ABFC the Corrected Akaike's Information Criterion, AICC (Hurvich & Tsai, 1989) is used. AICC is defined as follows:

$$AICC = n \ln(MSE) + 2k + \frac{2k(k+1)}{n-k-1}, \quad (13)$$

where MSE is the Mean Squared Error of the model of interest in the training data. AICC evaluates model's predictive performance as a trade-off between its accuracy in the training data (the first term of (13)) and its complexity (the last two terms of (13)). Calculation of the AICC for a single model requires a single estimation of model's parameters using OLS and calculation of MSE in training data. The "best" model is that whose AICC value is the lowest.

The AICC is an improvement over the classical AIC (Akaike, 1974) with the third term in (13) added as a correction term intended for working with small-sized data sets. For problems with relatively small n , AICC is suited better than AIC but converges to AIC as n becomes large (Hurvich & Tsai, 1989). AIC and AICC theoretical justification is based on the relationship between the Kullback-Leibner information and the maximum likelihood principle (Burnham & Anderson, 2002). Note that AIC as well as AICC does not assume that the "true model" (which was presumably used to generate the data) is one of the candidates (Burnham & Anderson, 2002).

In (Jekabsons & Lavendels, 2008b), an issue of the F-ABFC is stated, that, because the branching factor of the ABFC's state space increases very slowly together with d and k , in special cases when the data is of low dimensionality (e.g., $d \leq 4$) and/or the existing structure in the data is very complex (i.e. a very complex model is required) the search algorithm may get stuck in a local minimum too early in the search returning a too simple and underfitted model.

As a remedy for this, here an additional recursion of the state-transition operators is proposed introducing one hyperparameter for the F-ABFC. The idea is to recursively create additional regression models from models already created from the current best model using the same state-transition operators with which they were initially created. This essentially means that if, for example, the recursion depth is set to 2, Operator1 will create not only linear basis functions but also basis functions of the second degree, Operator2 will create not only copies of basis functions with degree increased by 1 but also by 2, and Operator3 will not only try to decrease degrees by 1 but also by 2. However, as still none of the operators add more than one basis function to the model at a time, for the Operator4 the recursion is not used.

The recursion of the operators reduces the number of local minima in the state space which is especially important near the starting-point of the search (the initial model) and enables the search algorithm to find a much better model.

Presence of such a "recursion depth" hyperparameter is a disadvantage as now a user intervention might be required. However, for larger dimensionalities of the input space (when also the increased computational resources are required) it is reasonable to completely disable the recursion (by setting the hyperparameter equal to 1) as with large dimensionalities the branching factor increases sufficiently fast and the problem of too early local minima diminishes.

Figure 3 shows pseudo-code of F-ABFC's search procedure. Note that in practical implementations of F-ABFC maintaining the set of the newly generated models ("MODELS") is not required as a single model can be created, evaluated, and, if it turns out not to be an improvement, immediately discarded.

```

BestModel ← the simplest model
BestModel.PerformOLSandCalculateAICC
loop
    //forward phase
    MODELS ← {all models created from BestModel using Operator1 and Operator2,
with no basis function redundancy}
    if RecursionDepth > 1 then
        for i ← 2 to RecursionDepth do
            MODELS ← MODELS ∪ {all models created from MODELS using the same
operator (with which they were initially created), with no basis function
redundancy}
            foreach Model in MODELS do
                Model.PerformOLSandCalculateAICC
            TestModel ← best of MODELS according to AICC
            if TestModel.AICC < BestModel.AICC then
                BestModel ← TestModel
            else
                break //break the main loop (exit the procedure)
    //backward phase
    loop
        MODELS ← {all models created from BestModel using Operator3 and Operator4,
with no basis function redundancy}
        if RecursionDepth > 1 then
            for i ← 2 to RecursionDepth do
                MODELS ← MODELS ∪ {all models created from MODELS using Operator3
(with which they were initially created), with no basis function redundancy}
                foreach Model in MODELS do
                    Model.PerformOLSandCalculateAICC
                TestModel ← best of MODELS according to AICC
                if TestModel.AICC < BestModel.AICC then
                    BestModel ← TestModel
                else
                    break //break the sub-loop
            end loop
        end loop
    return BestModel

```

Fig. 3. Pseudo-code of F-ABFC's search procedure

In (Jekabsons & Lavendels, 2008a), a version of F-ABFC was developed that slightly differs from the one proposed here in that the method used one additional state-transition operator and the "recursion depth" hyperparameter was not introduced. The paper (Jekabsons & Lavendels, 2008a) empirically demonstrated the computational and predictive performance advantages of F-ABFC comparing to subset selection and a number of other popular regression modelling methods. F-ABFC advantages in real-world practical applications are demonstrated in (Kalnins et al., 2008a; Kalnins et al., 2009b) where it is applied for modelling bending and buckling behaviour of different composite material structures.

4.4 Computational considerations

Assuming that the “best” model found by the F-ABFC search procedure includes a total of k_* basis functions and in each iteration the number of basis functions in the current model is increased by 1, the total number of evaluated models is of order

$$O\left(\sum_{i=1}^{k_*} di\right) = O\left(dk_* \sum_{i=1}^{k_*} i\right) = O\left(dk_* \frac{k_*(k_*+1)}{2}\right) = O(dk_*^3 + dk_*^2) = O(dk_*^3). \quad (14)$$

Consequently, relatively to the typical subset selection methods, the efficiency of the F-ABFC increases together with the increase in the number of input variables and in the required nonlinearity of the model (the value of p) but decreases together with the increase in the complexity k_* of the “best” found model. Moreover, the relative efficiency of the subset selection additionally substantially decreases in the common case when the required value of p is unknown and needs to be found by trying different values.

Using F-ABFC together with OLS, the associated linear least-squares fitting, required for a single model to be evaluated, demand computations of order $O(nk^2 + k^3)$, where nk^2 operations are required for filling a $k \times k$ matrix and k^3 operations are required for solving a linear equation system (Hastie et al., 2003). However, none of the proposed state-transition operators operate on more than one basis function of a model at a time meaning that, each time the parameters of a newly created model are calculated, only one row and one column of the $k \times k$ matrix will change. Recalculating only the elements of the corresponding row and column, reduces the order of the computations to $O(nk + k^3)$. Moreover, as the Operator4 does not modify any basis function (only deletes one), the order of the computations for this particular operator reduces further to $O(k^3)$.

Yet it must be noted that the F-ABFC can still become computationally rather demanding, especially when the number of input variables and/or the number of samples in the training data gets very large. This is the price to pay for the high flexibility of the method.

4.5 Convergence of the search process

The F-ABFC’s search algorithm is cycle-free because a new model is allocated to “BestModel” (Figure 3) only if it is better than the old one (according to AICC). Moreover, as the AICC criterion tries to estimate model’s true predictive performance, the algorithm will seek for the best trade-off between too simple and too complex models and will stop somewhere in-between them. Additionally there is also a hard bound – the number of basis functions in a model will never exceed the number of samples in the training data as otherwise the OLS cannot estimate model’s parameters.

It should also be noted that, although the state space of F-ABFC is infinite, in practice the models of the best predictive performance are normally located in the part of the space that is relatively near to the initial state where all the models (and their basis functions) are relatively simple and do not yet neither overfit the data nor have basis functions more than samples in the training data. This also means that really only a small finite fraction of the whole infinite state space must be explored.

4.6 Selection bias, selection instability, and model averaging

There are two issues that to some extent plague all the methods of model building (including subset selection and ABFC), especially when working with relatively little data – selection bias and selection instability (also called selection variance). While the issues are attributable to virtually any model building method, they are commonly ignored frequently resulting in models of lower predictive performance.

Selection bias occurs when in the search procedure one uses the same data to compute model's parameters, to perform model building (i.e. evaluation of candidate models, selection of the best one, and steering the search in direction of the most promising models), and to select the final "best" model which will be returned as the result of the model building process (Reunanen, 2003; Reunanen, 2006, Loughrey & Cunningham, 2004; Jekabsons, 2008). The problem is that the more candidates are visited during the search, the greater the likelihood of finding a model that has high accuracy in the training set while having a very low predictive performance (accuracy in the test set) (Reunanen, 2003; Reunanen, 2006; Kohavi & John, 1997; Loughrey & Cunningham, 2004). The random fluctuations in the data will improve the evaluations of some models more than others.

The problem is relevant regardless of the model evaluation measure used – statistical significance tests, complexity penalization criteria, or resampling techniques. In addition, the selection bias occurs even when performing model evaluation using completely independent validation data set (Kohavi & John, 1997; Reunanen, 2006). In any case, the more intensive (relative to the number of samples) is the search process, the larger is the selection bias, and, the larger is the noise in the data, the potentially larger is the harm (in terms of overfitting) done by the selection bias.

While the deterministic search algorithms of the hill-climbing type (including the SFFS algorithm of the F-ABFC) are usually less intensive and consequently more robust against overfitting than, for example, Simulated Annealing or Genetic Algorithms (Loughrey & Cunningham, 2004; Guyon & Elisseeff, 2003), the problem of selection bias remains relevant.

The second issue, selection instability, is related to the fact that small perturbations of the data (deleting or adding samples, adding noise, rescaling the values) can lead the model building process to vastly different models. This is because the large variability of estimates of the evaluation methods can lead to different local minima (Breiman, 1996; Kotsiantis & Pintelas, 2004; Guyon & Elisseeff, 2003; Cherkassky & Mulier, 2007). This variance is undesirable because variance is often the symptom of a "bad" model that does not generalize well and because the model may be failing to capture the "whole picture" (Guyon & Elisseeff, 2003).

One of the ways to reduce both the selection bias and the selection instability, is to employ model combining (also called model ensembling or averaging) techniques (Breiman, 1996; Opitz & Maclin, 1999; Cherkassky & Mulier, 2007; Jekabsons, 2008). While a typical model building process usually consists in choosing only one best description for the data discarding the remainder, combining a number of models in some reasonable manner appears more reliably accurate as this can have the effect of smoothing out erratic models that overfit the data and gain more stability in the modelling process.

A typical model combination procedure consists of a two-stage process (Cherkassky & Mulier, 2007). In the first stage, a number of different models are constructed. The parameters of these models are then held fixed. In the second stage, these individual models are linearly combined to produce the final model.

Both stages can be done in different ways. In this study, to increase the predictive performance of models built by the F-ABFC, a CV-type resampling of the training data together with unweighted model averaging (Opitz & Maclin, 1999; Duin, 2002) is employed. As this resampling and model averaging works on top of the F-ABFC, the method is called Ensemble of Floating Adaptive Basis Function Construction (EF-ABFC). During resampling, the whole training data is randomly divided into ν disjoint subsets (ν typically being equal to 10). Then ν overlapping training data sets are constructed by dropping out a different one of these ν subsets. Such procedure is also employed to construct training sets for ν -fold CV, so model ensembles constructed in this way are also called cross-validated committees (Parmanto et al., 1996).

Combining models via simple unweighted averaging requires them to be not too underfitted as well as not too overfitted (Duin, 2002). To lower the overfitting, in each CV iteration the unused 10th data subset is used as a validation data set for “re-evaluation” (using MSE) of the best models of each F-ABFC iteration and for selection of the one “final best” model from any iteration. Note that this validation set is never used for model evaluation during the search. Instead it is used strictly only for the “re-evaluation” and “re-selection” after the F-ABFC search process has already ended. Also note that as an evaluation measure in the search algorithm still the AICC is applied. This “re-evaluation” using the validation data set can detect whether the search process at some iteration may have started to generate overfitted models and select a model of some earlier iteration that is (hopefully) not (or at least less) overfitted (see Figure 4).

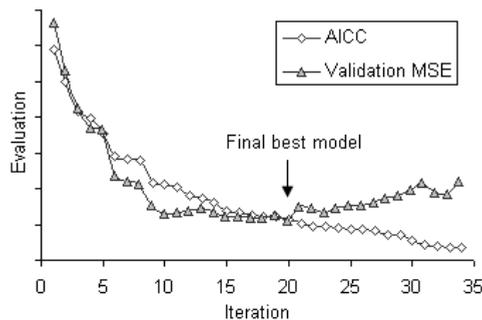


Fig. 4. An example of how a less overfitted model is selected using “re-evaluation” in validation set. Note that here starting from the 35th iteration the AICC values also start to increase (in contrast to the training error which always decreases) however this might be too late due to selection bias

The so far described process produces ν models built by ν independent F-ABFC runs each using a different combination of CV-partitioned data subsets. Next, the ν models from the ν CV iterations are combined using the unweighted model averaging. Note that prior to combining, all the models are re-fitted to the whole training data set (without the CV partitioning). This is done to compensate for the smaller training sets used during the individual model building.

Model combining by unweighted model averaging consists in taking an unweighted average of predictions of all the models:

$$F_{comb} = \frac{1}{v} \sum_{i=1}^v F_i, \tag{15}$$

where F_i is i th individual model from the i th CV iteration and F_{comb} is the combined model. For polynomial regression this simply means summation of all the polynomials and then a division of all the parameters of F_{comb} (that is also a polynomial) by v . Note that the parameter values of F_{comb} will not necessarily be optimal in the sense of the least-squares loss (in fact they will be optimal only in special cases, e.g., when all F_i 's are identical).

The employed model combining method is similar to Bagging (bootstrap aggregating (Breiman, 1996)) where the training set is bootstrapped (usually to build varied decision trees), and the unweighted average of the resulting models is taken.

Figure 5 gives an outline of the EF-ABFC model building process when the number of CV folds v is three. Note however that for all the practical applications of this study $v = 10$ is used. This is because too small number of models in ensemble will yield too little diversity hindering the models to correct each others errors, but, on the other hand, using too many models will yield no further improvement (Breiman, 1996; Opitz & Maclin, 1999; Kotsiantis & Pintelas, 2004; Parmanto et al., 1996). Moreover, too large number of CV folds can yield unreliable validation MSE estimates for the selection of the individual final best models, as then the individual validation sets may be too small.

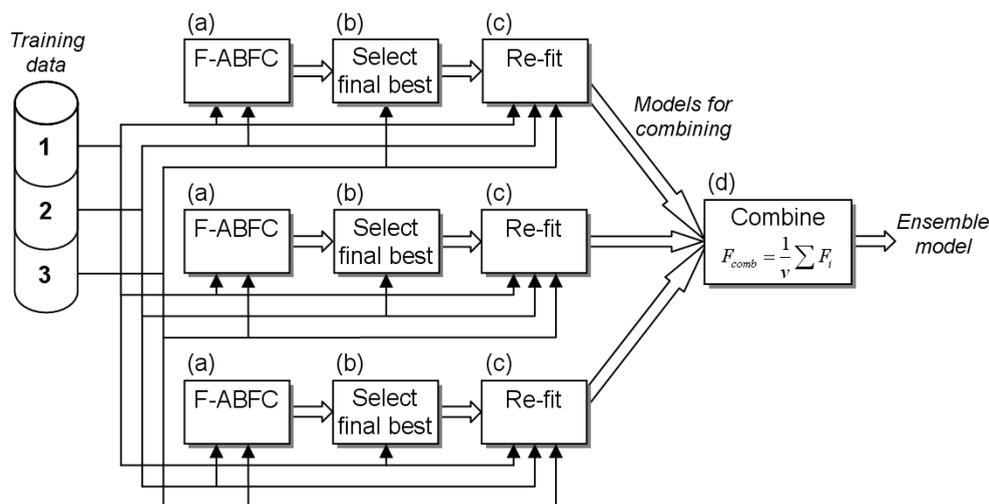


Fig. 5. An outline of the EF-ABFC modelling process when $v = 3$: (a) search for the best model according to AICC using F-ABFC; (b) select the one final best model according to MSE in validation data set; (c) re-fit the model (recalculate its parameters) using the whole training data; (d) combine the models

In recent literature, there is ever growing confidence that model ensembles often perform better than individual models and consistently reduce prediction error (Breiman, 1996; Opitz & Maclin, 1999; Kotsiantis & Pintelas, 2004; Jekabsons, 2008). However, model ensembles are not always the best solutions (Kotsiantis & Pintelas, 2004): if there is too little

data, the gains achieved via an ensemble may not compensate for the decrease in accuracy of individual models, each of which now sees an even smaller training set. On the other end, if the data set is sufficiently large, even a single flexible model can be quite adequate. Using large data sets also substantially decreases potential selection bias, so superiority of EF-ABFC over F-ABFC in such situations is expected to diminish.

The most significant disadvantage of the EF-ABFC compared to F-ABFC is that it requires larger computational resources. However, the fact, that before the model combining the v models are built completely separately, allows for an easy parallelization of the process dividing the execution time by v . In this study however the parallelization is not done.

The paper (Jekabsons, 2008) empirically demonstrated the computational and predictive performance advantages of EF-ABFC comparing to subset selection and a number of other popular regression modelling methods. EF-ABFC advantages in real-world practical applications are demonstrated in (Kalnins et al., 2008b; Kalnins et al., 2009a) where it is applied for modelling bending and buckling behaviour of different composite material structures.

4.7 Remarks

This section covers various aspects (extensions, limitations, etc.) of the ABFC not discussed in the previous sections.

4.7.1 Incorporating domain knowledge

The ABFC methods attempt to model arbitrary dependencies in data with little or no knowledge of the system under study. In problems of moderate and large dimensionality the user usually is not required to tune any hyperparameters. However, if there is sufficient additional domain knowledge outside the specific data at hand, it may be appropriate to place some constraints on the final model. If the knowledge is fairly accurate, such constraints can improve the accuracy while saving computational resources.

For example, the constraints might be one or more of the following: 1) limiting the maximal degree of all the basis functions (similarly as in the subset selection), i.e. $0 \leq \sum_{j=1}^d r_{ij} \leq p$ for all i ; 2) limiting the maximal value of exponents for each particular input variable in all the basis functions, i.e. $0 \leq r_{ij} \leq p_j$ for all i , where p_j is maximal exponent of the j th variable; 3) restricting contributions of specific input variables that are not likely to interact with others so that those variables can enter the model in basis functions only solely - with exponents of all other variables fixed to zero. These constraints, as well as far more sophisticated ones, can be easily incorporated in the ABFC. However, note that in all the experiments described in this chapter no constraints are used.

4.7.2 Robustness

The ABFC methods described in this study estimate model parameters via minimization of the squared-error loss, i.e. using OLS. However, while the squared-error loss is the most commonly used, it is known that it loses its robustness against grossly outlying samples as well as in very sparse high-dimensional data sets (Cherkassky & Ma, 2002).

One solution of this problem is to use a more robust loss function. The squared-error loss in ABFC is not fundamental. Any other loss function can be used to estimate the parameters

and to evaluate the models by simply replacing the routine “PerformOLSandCalculateAICC” of the search procedure (Figure 3) with a more robust one. Note that while this would make the methods more robust, the computational advantage of OLS would be lost. In any case, gross outliers (in output variable as well as input variables) that can be detected through a preliminary data analysis should be considered for removal before applying ABFC.

4.7.3 Other types of basis functions

The ABFC methods described in this study can generate regression models with basis functions of only nonnegative integer exponents. However, in principle the exponents can also be allowed to take negative or even fractional values. Appropriate adaptation of the state-transition operators can enable generating such models. Keeping the same initial model as before, the search now could go in direction of both positive and negative exponents.

4.7.4 Integrating ABFC into other modelling methods

The result of running an ABFC procedure is a simple polynomial regression model. Such models are also utilized as “sub-models” in a number of other regression modelling methods. For example, the ABFC methods can be used in Polynomial Neural Networks (usually induced by Group Method of Data Handling) (Nikolaev & Iba, 2006) for adaptation of each individual neuron’s functional form and degree. The methods also can serve for generation of local regression models in Locally-Weighted Regression (also called Moving Least Squares) (Cleveland & Devlin, 1988; Kalnins et al., 2008b; Kalnins et al., 2005) adaptively generating a model each time a query is received. ABFC can also induce piecewise polynomial models for appropriately partitioned data sets.

The polynomial basis functions can also be viewed as nonlinear transformations (or features) of the original input variables. In this manner the ABFC methods can also be viewed as methods for automatic adaptive feature construction. For example, the constructed features can further serve as inputs for Support Vector Machines (Vapnik, 1995; Smola & Scholkopf, 2004) similarly to the features constructed using genetic algorithm in (Ritthoff et al., 2002).

All these applications of ABFC can make the original methods more flexible and therefore, if treated appropriately, produce models of higher predictive performance.

4.7.5 Using ABFC for solving classification problems

The ABFC methods can also be used for solving binary classification problems where the output variable y can take value of only either 0 or 1. This can be done, for example, by constructing basis functions for logistic regression (also called maximum entropy classifier) models. Logistic regression (Hastie et al., 2003; Witten & Frank, 2005) represents log odds of y being equal to 1 as a linear model:

$$\ln(P/(1-P)) = F(x) = \sum_{i=1}^k a_i f_i(x), \quad (16)$$

Thank You for previewing this eBook

You can read the full version of this eBook in different formats:

- HTML (Free /Available to everyone)
- PDF / TXT (Available to V.I.P. members. Free Standard members can access up to 5 PDF/TXT eBooks per month each month)
- Epub & Mobipocket (Exclusive to V.I.P. members)

To download this full book, simply select the format you desire below

