



Universiteit  
Leiden  
The Netherlands

## Methods and tools for mining multivariate time series

De Gouveia da Costa Cachucho, R.E.

### Citation

De Gouveia da Costa Cachucho, R. E. (2018, December 10). *Methods and tools for mining multivariate time series*. Retrieved from <https://hdl.handle.net/1887/67130>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/67130>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The following handle holds various files of this Leiden University dissertation:

<http://hdl.handle.net/1887/67130>

**Author:** de Gouveia da Costa Cachucho, R.E.

**Title:** Methods and tools for mining multivariate time series

**Issue Date:** 2018-12-10

## Chapter 6

# Sports Analytics for Professional Speed Skating

Arno Knobbe, Jac Orie, Nico Hofman,  
Benjamin van der Burgh, Ricardo Cachucho

*in journal of the Data Mining and Knowledge Discovery, Volume 31,  
Issue 6, pp 1872–1902, Springer, November 2017*

### Abstract

*In elite sports, training schedules are becoming increasingly complex, and a large number of parameters of such schedules need to be tuned to the specific physique of a given athlete. In this paper, we describe how extensive analysis of historical data can help optimize these parameters, and how possible pitfalls of under- and overtraining in the past can be avoided in future schedules. We treat the series of exercises an athlete undergoes as a discrete sequence of attributed events, that can be aggregated in various ways, to capture the many ways in which an athlete can prepare for an important test event. We report on a cooperation with the elite speed skating team LottoNL-Jumbo, who have recorded detailed training data over the last 15 years. In this project, we analyse this data, and extract actionable and interpretable patterns that can provide input to future improvements in training. We present two alternative techniques to aggregate sequences of exercises into a combined, long-term training effect, one of which based on a sliding window, and one based on a physiological model of how the body responds to exercise. Next, we use both linear modeling and Subgroup Discovery to extract meaningful models.*

## 6.1 Introduction

This paper describes research challenges related to a recent Sports Analytics project between a leading Dutch professional speed skating team and data scientists from Universiteit Leiden and Hogeschool van Amsterdam. During its history, the athletic team, currently called LottoNL-Jumbo, has included numerous Dutch skaters that competed at the European, world, and Olympic level, and presently includes a world record holder and two Olympic medalists. The project involved 15 years of detailed training data kept by the coach of the team (second author of this paper) with the aim of improving the training program and further optimising the performance of current and future skaters of the team. In this paper, we report on the data science techniques required to analyse this non-trivial data, and showcase findings for specific athletes. A number of novel techniques are introduced to deal with the specifics of the recorded data, and to produce interpretable and actionable results that can help fine-tune the training programs.

Speed skating is a winter sport where athletes compete on skates to cover a given distance on an oval (indoor) ice rink. In this paper we focus here on long-track speed skating events, which involves a 400 meter oval track with two lanes. Events over multiple distances exist, ranging from 500 meters to 10 000 meters, with each skater typically specialising in one or two distances, depending on their physiology and training. Although each race in an event involves two skaters, the final standing is determined by the overall ranking of times of all participants. This effectively makes each race a time trial, where the outcome of a given skater is only determined by their own performance. From a data mining point of view, this is attractive since all results can be assumed independent, and one can simply collect all race results of an athlete without having to consider the influence of the ‘opponent’.

The available data, painstakingly collected by the coach, involves primarily descriptions of the daily training activities, partly structured and partly free text. The structured part of the description is very consistent, and captures a classification of the nature of the training (e.g. “*cycling extensive endurance*”), as well as numeric values indicating the duration (in minutes) and intensity (on a subjective scale of 1 to 10) of the session. Training data is specific to individual athletes, and the intensity of the training was obtained from the athlete, post hoc. With six training days per week, and potentially two sessions per day, this amounts to roughly 450 sessions per season, making for a substantial data collection per athlete/season. Next, race results are available that capture test events which will stand as our target. These

were scraped from the internet<sup>1</sup>. As a result, the problem is essentially a *regression task*, since target and predictors are both numeric variables.

While evidently being in the regression domain, it is not immediately clear what the independent variables of our task might be. Clearly, the independent variables should capture aspects of the training program prior to the events in question. However, the available data is a long sequence of events with a small number of characteristics (type, duration, intensity, ...), so some form of transformation is necessary to arrive at an attribute-value representation that is amenable to main-stream regression analysis. In this paper, we take an *aggregation-based feature construction* approach, inspired by earlier work in [14], in order to derive a fairly extensive set of features that capture the preparation (training, but also absence thereof) from various angles, for example focussing on specific periods prior to the test event, or on specific intensity zones. In its basic form, the aggregation will take place over windows of varying lengths (ranging from one day to several weeks) using different aggregate functions and variables, with specifiers such as training type and intensity zones. In a more elaborate approach, developed for this specific purpose, the aggregation will take the form of convolution with a physiology-inspired kernel consisting of several exponential decay functions of varying half times. This kernel is inspired by the so-called *Fitness-Fatigue model* [16], that tries to capture how the human body responds to a specific training impulse over the course of time.

After having obtained a suitable attribute-value representation with potentially predictive features, the next challenge is to produce meaningful models from this dataset. The overall aim of this project is to provide the coaching staff with easy-to-understand, actionable pointers as to how to fine-tune the training routines, and avoid pitfalls of under and over-training. Therefore, we specifically intend to discover interpretable patterns, that are relatively easy to understand by the domain experts, and ideally do not involve a great many variables. We will be working with two types of regression techniques. Assuming mostly linear dependencies between the aggregated features and the target variable, *regularized linear regression* methods such as LASSO [27, 90] are attractive since they select features and produce relatively concise models of the data. However, with the physiological domain at hand, it is likely that non-linear dependencies will also exist, and rather, one expects thresholds to exist on the features, where too large or too low a value (e.g. training load) will produce sub-optimal results. For such phenomena, we expect *Subgroup Discovery* [50, 30, 5, 78] to produce more useful results.

---

<sup>1</sup><http://www.osta.nl>

This Sports Analytics paper has two sides. On the Sports side, we present some interesting findings that are of practical relevance to the team, with the following key contributions:

- The application of the Fitness-Fatigue model and the fitting of this model to individual skaters, where the parameters of the model convey key properties of each skater.
- Various demonstrative, interpretable patterns concerning improved training practices.
- The presentation of results relating to competitions.
- The capability to produce detailed findings for other skaters.

On the Analytics side, we introduce a number of new ways to exploit detailed training data, of relevance not just in the speed skating discipline, and in some cases applicable to other analytics domains also. In this domain the key contributions are:

- Introduction of (conditional) aggregation as a way of aggregating discrete sequences of events, and producing a range of features that capture various aspects of those sequences.
- Aggregation by means of two options: one that is easy to compute and interpret (uniform window), one that is more physiologically plausible, and at the same time harder to compute (the Fitness-Fatigue model).
- Application of linear modeling and Subgroup Discovery in order to select key features and produce interpretable models. 5) Evaluation of models in terms of  $R^2$  and  $p$ -values, that makes linear models and subgroups immediately comparable.

## 6.2 Speed Skating and Sports Analytics

The (long-track) speed skating takes place on an oval track 400 meters in length. Races are typically held with two participants at each time (skating in separate lanes), but each participant is ranked on their individual time. Both men and women compete, in separate competitions. Races come in various distances, but the most common distances at major events are 500 m, 1000 m, 1500 m, 3000 m, 5000 m and 10000 m, of which the last distance only applies to men, who in turn do not skate the 3000 m. These disciplines are usually divided into sprint, medium, and long distance, and skaters typically

specialize in one of these, and compete in only a few of these disciplines, although participation is facultative. Each category requires a specific type of physiology, which explains the specialization of athletes. Furthermore, each distance requires a specific type of training, and exercises for one distance may actually harm the performance on other distances [64, 45]. This implies that our analysis will often be specific to a small number of similar distances, or even be specific to individual athletes. Since the training programs are well-developed, and the senior athletes will have several years of experience working with the coach, the produced findings may be subtle, which will often call for an athlete-specific approach.

Even though races for a specific event can be considered time trials, there will be a level of variance in the race results that cannot be explained by differences in race preparation and training. It is a well-know fact within speed skating that times are determined to a reasonable degree by the ice rink. First of all, the ice properties will differ from one venue to the next, and top venues will have better ice maintenance techniques and experience. Another factor that influences the times, besides ice quality, is the altitude of the venue, with higher ice rinks tending to be faster due to the lower air resistance. In order to compensate for the difference in rink speeds, we have opted to work with *relative times* rather than recorded times.

**Definition 9** (Relative time). *For a race of distance  $d$ , by a skater of gender  $g$ , at ice rink  $r$ , finishing in time  $t$  (in seconds), the relative time  $t_{rel}$  is defined as*

$$t_{rel} = t/t_{rec}(d, g, r)$$

where  $t_{rec}(d, g, r)$  is the record for a specific discipline and ice rink.

Relative time will produce race results slightly above 1.0 or exactly 1.0 if a race either produced or replicated a rink record. The use of relative time not only allows comparisons between results at different venues, but theoretically also comparisons between results in different disciplines or even between men and women, although one might be comparing apples and oranges here on other aspects of the data. The rink records were scraped from the Dutch Wikipedia page<sup>2</sup> that collects local records. Note that the use of records to estimate the speed of a rink is not flawless. First of all, records are continuously subject to improvement, such that the definition of relative time is sometimes problematic. Next, some venues rarely host international events, such that their records do not fairly represent the theoretical speed of the rink, and actually produce under-estimates of the relative times athletes set (meaning they appear faster than they are). In order to avoid constantly

---

<sup>2</sup>[https://nl.wikipedia.org/wiki/Lijst\\_van\\_snelste\\_ijsbanen\\_ter\\_wereld](https://nl.wikipedia.org/wiki/Lijst_van_snelste_ijsbanen_ter_wereld)

having to update the list of records and subsequent analyses, which is rather time-consuming, we choose to fix the records at a particular point in time. A minor side effect of this is that some newer results may actually have a relative time below 1.0. Our collected list of international rink records will be made available<sup>3</sup> as part of this publication.

### 6.3 Feature Construction by Aggregation

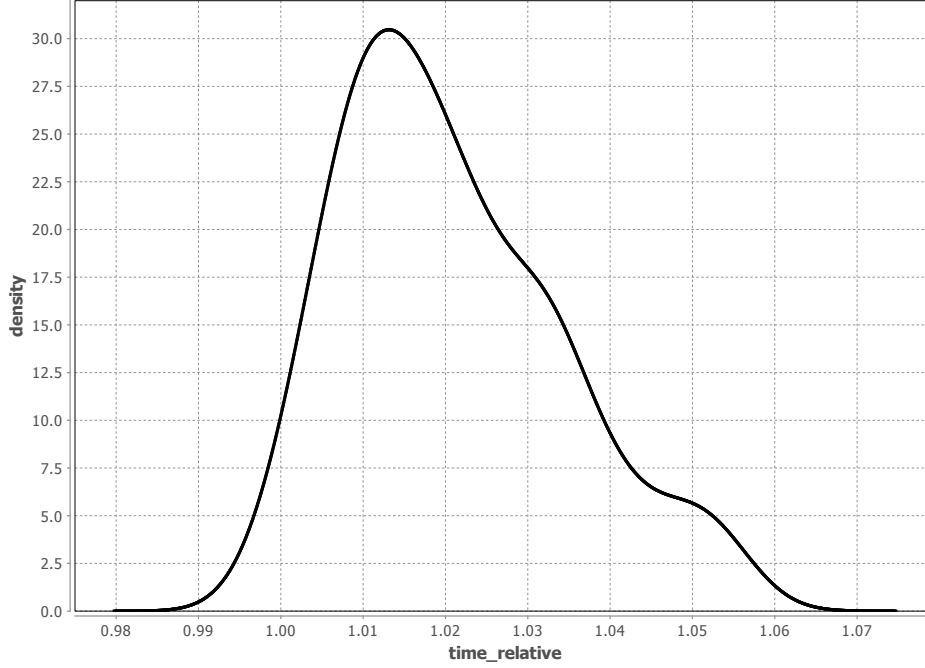
As explained in the introduction, the training data takes the form of a sequence of annotated events, corresponding to the individual exercises an athlete performs. While being valuable information, this sequential representation will require certain transformations in order to elicit general characteristics of an athlete's preparation for a race. Individual exercises will generally not play a deciding role in the outcome (unless of very extreme nature), and it is the combined nature of exercises that determines the effect of the training program. Therefore, some form of aggregation is required to draw out the various aspects of training that potentially play a role. Although there is a large body of knowledge about the effect of certain types of exercise in the sports physiology literature, it is still uncertain what aspects of training and preparation determine the variance in race results that still remains, as for example exemplified in Figure 6.1. For this reason, our feature construction approach will include a rather large collection of features, with the aim of including many angles and leaving room for discovery in later stages of the analysis. Furthermore, it is not quite clear how long the effect of specific exercises lasts for individual athletes, and thus what period prior to the test event should actually be included in the aggregation. Our set of constructed features will therefore involve time windows of various lengths, ranging from one day to several weeks.

In this paper, we will consider two general aggregation approaches, the first of which involves uniform aggregation over the various windows. The second based on convolution using a kernel that is based on the exponential decay function.

---

<sup>3</sup><http://datamining.liacs.nl/rink-records.txt>





**Figure 6.1:** Distribution of relative times over 1000 m as estimated by kernel density estimation (based on 75 results, with a Gaussian kernel of  $\sigma = 0.0051$ ). The probability density estimate continuing below 1.0 is an artefact of the KDE. The best time in the list is actually a rink record in the Hague (the Netherlands) at 1.0.

### 6.3.1 Uniform Aggregation

Before defining the notion of a (time) window, we first formalize the events to be aggregated, as they appear in the data of our application.

**Definition 10** (Exercise). *An exercise is defined as a tuple  $e = (t, ampm, dur, int, load)$ , where  $t$  is the date of the exercise,  $ampm$  is a binary variable indicating the morning or afternoon session. Numeric values  $dur$ ,  $int$ , and  $load$ , indicate the duration (in minutes), the subjective intensity (on a scale of 1 to 10), and the load (in intensity-minutes) of the exercise, respectively.*

The three crucial numeric attributes of an exercise specify the following:

- The *duration* simply specifies the length of the exercise. Durations tend to be rounded to quarters of hours (especially for longer exercises), but this is deliberate, and athletes generally adhere to the required duration.
- The *intensity* indicates how heavy the exercise was, as perceived by the

athlete, with 10 being the intensity of a race itself. During training, values of 9 or 10 are rare. Although intensity is a subjective measure, the athletes are very used to it, and will rate specific trainings fairly consistently.

- The *load* is defined as the product of duration and intensity, with the intention of capturing the total energy expenditure during the exercise. Although load is actually a derived attribute (it does not appear in our normalized database, for that reason), we include it in the definition of an exercise because it plays a crucial role both in the modeling as in the reasoning behind the training program<sup>4</sup>.

Note that the races themselves also appear as exercises in the database, since it is crucial to include the training load produced by such intense events, when considering the preparation for other races. In speed skating, several races often take place in a single weekend, such that later races are influenced by earlier ones.

**Definition 11** (Time Window). *A (time) window  $w_{t,m}$  is a set of exercises  $e_1, \dots, e_n$ , such that all dates  $e_i.t$  are before  $t$ , and not more than  $2m - 1$  days before  $t$ :  $t - 2m + 1 \leq e_i.t < t$ .*

Note that day  $t$  itself is not part of the window. For reasons that will become clear in later sections, we have opted to define the length of a window in terms of its middle  $m$ , essentially indicating the ‘centre of mass’ of the window. A window  $w_{t,1}$  will thus include the one day prior to  $t$ ,  $w_{t,2}$  indicates the three days prior to  $t$ , and so on.

For a window  $w$ , the following primitive aggregates will be considered:

**Count** Simply the number of exercises in  $w$ :  $|w|$ .

**Sum(duration)** The sum of durations of the exercises in  $w$ :  $\sum_i e_i.dur$ .

**Sum(intensity)** The sum of intensities of the exercises in  $w$ :  $\sum_i e_i.int$ .

**Sum(load)** The sum of loads of the exercises in  $w$ :  $\sum_i e_i.load$ .

**Avg(duration)** The average duration of the exercises in  $w$ :  $\sum_i e_i.dur / |w|$ .

**Avg(intensity)** The average intensity of the exercises in  $w$ :  $\sum_i e_i.int / |w|$ .

---

<sup>4</sup>Note that the definition in terms of a product of duration and intensity might be too simplistic, since neither duration nor intensity may be a linear scale. Doubling the length of an exercise may have an exaggerated effect if the intensity is (too) high, and doubling the intensity makes the exercise entirely different in nature, addressing different metabolic systems.

**Avg(load)** The average load of the exercises in  $w$ :  $\sum_i e_i.load/|w|$ .

**Max(duration)** The maximum duration of the exercises in  $w$ :  $\max_i e_i.dur$ .

**Max(intensity)** The maximum intensity of the exercises in  $w$ :  $\max_i e_i.int$ .

**Max(load)** The maximum load of the exercises in  $w$ :  $\max_i e_i.load$ .

Aggregation using the minimum was deemed senseless, since a very light training has little effect, and one could interpret daily rest periods as very light exercises anyway.

### Specifiers

Each primitive aggregate listed above can be applied to all the exercises in a given window, or just to subsets of exercises from specific categories. These subsets are specified by what we will refer to as specifiers. We apply the following specifiers:

**Morning/afternoon sessions** By specifying *am*, *pm* or no specifier at all, the aggregate can include only the morning sessions, only the afternoon sessions, or all sessions, respectively. Note that during the winter, the coach will plan exercises on the ice in the morning, and alternative training in the afternoon, so distinguishing between those may be fruitful.

**Intensity intervals** Exercises at different intensities will trigger different parts of the system, and hence will produce a different training stimulus. As specifier, we optionally select only exercises within specific intensity intervals  $[l, u]$ , where  $l \in [1, 10]$  and  $u \in [l, 10]$ .

Note that each type of specifier will introduce multiple variants of the primitive aggregates. For *ampm*, adding specifiers will raise the number of aggregates (per window size) from 10 to 30. For the intensity-intervals, there will be  $\frac{1}{2} \cdot 10 \cdot (10 + 1) = 55$  variants of each primitive. However they are only applied to the 4 primitive aggregates that do not involve intensity and load, producing  $55 \cdot 4 = 220$  aggregates. In order to avoid combinatorial explosion of the aggregate collection, we do not include combinations of specifiers (such as low intensities in morning sessions). In total, there will be 250 aggregates per window.

### Aggregation and Convolution

Observe that such uniform aggregation over a window can be seen as a form of convolution with a rectangular kernel [85]. The convolution of a time series  $x(t)$  (in this case any of the training parameters that are aggregated) applies a kernel to the series to obtain a new series  $y(t)$  as follows:

$$y(t) = h * x(t) = \sum_{i=-\infty}^{\infty} h(i)x(t-i)$$

Here,  $h$  refers to the kernel, which is required to sum to 1 over its domain. In the case of a uniform window, the kernel is essentially a rectangular function (remember that  $2m - 1$  is the length of the window):

$$h(t) = \begin{cases} 1/(2m-1) & \text{if } 0 \leq t \leq 2m-1 \\ 0 & \text{otherwise} \end{cases}$$

Since the rectangular kernel is zero over a large part of its domain, the convoluted function  $y(t)$  can be computed much faster. In the next section, we will introduce a kernel that is both more natural and more expensive to compute.

#### 6.3.2 The Fitness-Fatigue Model

Although uniform aggregation is intuitive and straightforward to implement (and as we will see, provides fairly good models), it is somewhat unnatural. First of all, it is unlikely that all exercises over a period of, say, four weeks will have the same influence on the level of fitness at a race. Rather, one would expect that exercises several weeks ago have a much smaller influence than more recent ones. Second, the hard distinction between an exercise 28 days ago, and one 29 days ago seems unnatural, and may introduce minor artefacts in the constructed features. Finally, there is a general pattern where the initial effect of an exercise is negative, while after a short period of rest and recuperation, the effect is positive. Ideally, the aggregated features should exhibit such behaviour.

In this section, we introduce a type of aggregation based on convolution with a more natural gradually progressing kernel. We will use a multi-component kernel that is taken from the physiology literature [16] and aims to model the complex way in which a human body responds to exercise by initial fatigue,

followed by a slight improvement in performance, the effects of which die out gradually over the course of several week, returning to a state of fitness comparable to that prior to the exercise.

The core of this kernel is based on the *exponential decay* function, as follows:

$$h_e(m) = e^{-\lambda m}, \quad m \geq 0$$

The parameter  $\lambda$  here determines the speed with which this kernel decays towards zero, in other words, the speed with which the effects of exercise diminish over time. Although the exponential decay is defined in terms of  $\lambda$  (with unit  $s^{-1}$ ), we will primarily define a specific kernel in terms of the parameter  $\tau$  (in units  $s$ , or more conveniently, in days), which corresponds to the ‘mean lifetime’ of the kernel, and as such can be interpreted as the centre of mass of the kernel. This makes this parameter immediately comparable to parameter  $m$  of a uniform window, which is also the centre of mass of the kernel. The simple relationship between  $\tau$  and  $\lambda$  is as follows:

$$\tau = 1/\lambda$$

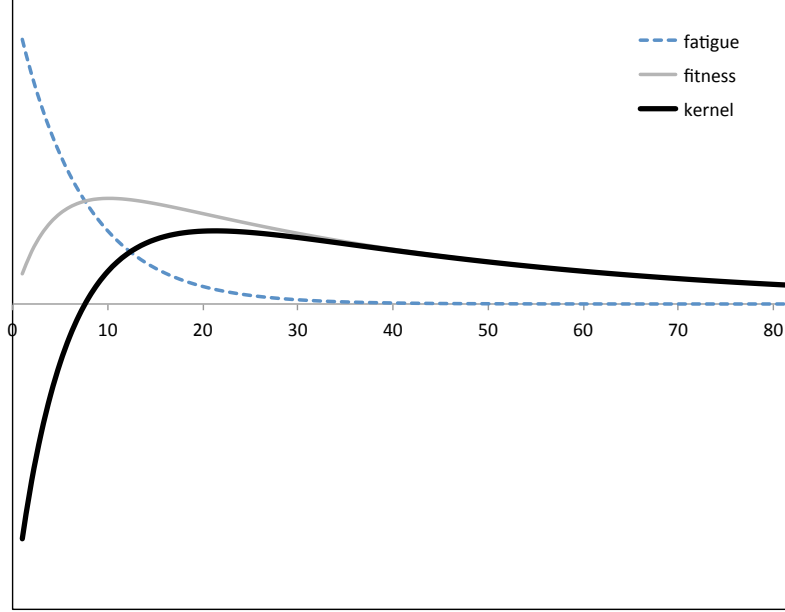
The exponential decay function effectively models the diminishing positive effect of an exercise as time passes. However, it does not include the tiring effects of exercise in the few days after training, which may outweigh the positive influence of training. For this reason, [16] introduced the so-called *Fitness-Fatigue* model, which models these two effects as two components of a kernel with different weights and different decay factors, as follows:

$$h_2(m) = e^{-\lambda_{fit}m} - Ke^{-\lambda_{fat}m}, \quad m \geq 0$$

$\lambda_{fit}$  determines the speed with which the positive effects of training (the *fitness*) diminish, and typically corresponds to a  $\tau_{fit}$  in the order of months, while  $\lambda_{fat}$  determines the shape of the fatigue curve immediately after the exercise. The associated  $\tau_{fat}$  is typically in the range of two weeks. Initially, the influence of fatigue is about twice as big as that of the improved fitness (determined by the value of  $K$ ).

The fitness in the above two-component model is assumed to be immediately improved after the exercise, which in practice is not the case. The desired adaptation in various metabolic systems will not take effect until several days after the exercise, such that the fitness will need to be modeled with an additional component [16], producing the following three-component kernel:

$$h_{ff}(m) = (e^{-\lambda_{fit}m} - e^{-\lambda_{del}m}) - Ke^{-\lambda_{fat}m}, \quad m \geq 0$$



**Figure 6.2:** The three-component Fitness-Fatigue kernel (in solid black) as a function of time after the exercise (in days). The *fitness* and *fatigue* parts are also shown, in solid grey and dashed, respectively.

where  $\lambda_{del}$  affects the exponential function that reduces the initial fitness. In Figure 6.2, the combination of fitness and fatigue into this kernel is demonstrated. In [16], values are given for the associated parameters, obtained by fitting the convolved function to athletic data, producing the values below. Although these values seem reasonable, they will be athlete- and specialism-specific, such that we will fit these values to specific datasets collected, in the experimental section. The published values for the parameters are as follows:

$$\tau_{fit} = 50 \text{ days}, \quad \tau_{del} = 5 \text{ days}, \quad \tau_{fat} = 15 \text{ days}, \quad K = 2.0$$

## 6.4 Modeling Approaches

### 6.4.1 Regularized Linear Regression

In the previous section, we explained the procedures to build large sets of interpretable features about the training, that might be able to explain the

target variables of performance. These target variables might be a linear function of a subset of the aggregate features, but we do not know which ones beforehand. In order to find a good subset of aggregate features for each target variable, we use LASSO [90], a method for estimating generalized linear models using convex penalties ( $l_1$ ) to achieve sparse solutions [27].

Assume  $\bar{t}$  is the mean of the target variable:

$$\bar{t} = \frac{1}{n} \sum_{i=1}^n t_i$$

The coefficient of determination  $R^2$  is now defined as:

$$R^2(t, f) = 1 - \frac{\sum_i (t_i - f_i)^2}{\sum_i (t_i - \bar{t})^2} \quad (6.1)$$

where  $\sum_i (t_i - f_i)^2$  is the sum of squared differences between the actual and predicted target value, and  $\sum_i (t_i - \bar{t})^2$  is the sum of squared differences between the target value and the constant function  $t = \bar{t}$ . Note that  $R^2$  is between 0 and 1 whenever the model  $f$  is produced using the Ordinary Least Squares method, but may be lower than 0 for functions derived differently.  $R^2$  is often interpreted as the *explained variance*, where a value of 0 means that no variance in the dependent variable can be explained by variance in the independent variable(s), and a value of 1 means that all variance can thus be explained (a perfect fit of the data).

### 6.4.2 Subgroup Discovery

The previous section focussed on linear models, assuming that the dependencies we hope to discover are indeed linear in nature. Unfortunately, in the domain we are focussing on, it is quite likely that the relationship between (extent of) training and performance is non-linear. When doing a certain training routine, it can be expected that the relationship is in fact curved, with peak performance being achieved at a certain optimal load on the human body. Doing too little will not achieve the right effect, but doing too much of the training also produces sub-optimal performance. Specifically, one can expect thresholds in the training load above (or below) which performance will rapidly diminish. Therefore, we will also experiment with modeling techniques that are more local in nature, and find subsets of the data where performance was surprisingly low, as well as finding variables and thresholds that will identify such sub-optimal subsets.

Our paradigm of choice for such (potentially) non-linear data is that of *Subgroup Discovery* [50, 54]. It is a data mining framework that aims to find interesting *subgroups* satisfying certain user-specified constraints. In this process, we explore a large search space to find subsets of the data that have a relatively high value for a user-defined quality measure. We consider constraints on attributes, and determine which records satisfy these constraints. These records then form a subgroup. The constraints on the attributes (the *description*) form an intensional specification of a part of our dataset, and the subgroup forms its extension (that is, an exhaustive enumeration of the members of the subgroup).

A number of papers in the literature discuss SD variants for regression tasks, which to some extent are applicable to our case. One group of techniques focusses on finding subgroups where the target shows a surprisingly high (or conversely, low) average value [30, 5, 78, 59]. Typical proposed quality measures use statistical tests to capture the level of deviation within the subgroup, often weighted by the size of the subgroup, for example the *mean test* or *z-score* [78],

$$\varphi_z(S) = \sqrt{|S|} \frac{\mu_S - \mu_0}{\sigma_S}$$

where  $\mu_S$  and  $\mu_0$  stand for the subgroup and database means of the target, respectively, and  $\sigma_S$  denotes the standard deviation within the subgroup  $S$ . Other works consider the distribution of target values within the subgroup [41], and use statistical measures for assessing distributional differences.

In the majority of these quality measures, the interestingness is computed from the distribution of the subgroup alone, or when compared to that of the entire dataset. Here, we take a slightly different approach, and consider the subgroup description a dichotomy of the data, where both the distribution of the subgroup as well as of the *complement* play a role in determining the quality of the dichotomy. Therefore, we introduce a new quality measure for numeric targets in SD. This quality measure uses the well-known notion of  $R^2$  to capture how well a subgroup and its corresponding complement describe the data, in comparison to the distribution of the entire dataset, so ignoring the dichotomy. Hence, we treat the subgroup as a model of the data, to be more specific a step function of two parts. The following two averages over the target  $t$  provide the constant prediction for, respectively, the subgroup and its complement:

$$\bar{t}_{subg} = \frac{1}{|S|} \sum_{i \in S} t_i$$



$$\bar{t}_{comp} = \frac{1}{n - |S|} \sum_{i \notin S} t_i$$

These two average values now lead to the following step function:

$$f_S(i) = \begin{cases} \bar{t}_{subg} & \text{if } i \in S \\ \bar{t}_{comp} & \text{if } i \notin S \end{cases}$$

The quality measure *Explained Variance* is now simply defined as follows:

$$\varphi_{EV}(S) = R^2(t, f_S)$$

This quality measure uses the definition of  $R^2$  given in the previous section, which was independent of the nature of the model  $f$ . Note that by using this quality measure, we have a way of directly comparing the discovered subgroups (with corresponding step functions) to the linear models, which is a clear benefit over the traditional quality measures. We furthermore observe that the step functions, despite representing dichotomies, can be based on subgroups of multiple conditions. Therefore, the resulting step functions will be multi-dimensional (involving potentially multiple features). The quality measure introduced here was implemented in the Cortana Subgroup Discovery tool [66].

## 6.5 Experimental Results

In order to demonstrate the kinds of analyses and results of the proposed methods on actual data, and to test the benefits of individual techniques proposed above, we experiment with data from four athletes of the LottoNL-Jumbo team, two male and two female. All experiments were performed using three software components:

1. A relational database that organizes all the different datasets and meta-data concerning exercise and competition: the *Performance Sports Repository* (PSR).
2. A dashboard accessible over the Internet, that provides various views and visualizations of the data, and allows online aggregation and linear modeling of the data.
3. The generic Subgroup Discovery tool Cortana<sup>5</sup>, which was extended

---

<sup>5</sup>Sources in Java and an executable of this tool can be downloaded from <http://datamining.liacs.nl/cortana.html>.

**Table 6.1:** Dataset details for competition results of four skaters.

	Gender	Distance	Competitions	Sessions
M1	Male	1000 m	75	2 758
M2	Male	500 m	142	2 930
F1	Female	1500 m	60	2 230
F2	Female	500 m	22	1 095

for this purpose with a direct database connection to the PSR, and the Explained Variance quality measure [66].

We will demonstrate the results on the datasets listed in Table 6.1, collected from four athletes. Next to competition data, we also have physical test data, for which we provide results for one of the athletes (M1 in the table below), for which we have 146 records.

We will generally describe three types of modeling of the data: 1) univariate models, either using a linear or a step function, where we rank all features by  $R^2$ , 2) multivariate models using LASSO, and 3) Subgroup Discovery using Cortana. Note that univariate step models can also be interpreted as subgroups with a single condition, such that results between settings 1 and 3 overlap to some extent. When mining for subgroups, we use beam search to a fairly shallow depth, typically to a maximum depth of three or less, depending on the experiment. When not further specified, the subgroups (or their step functions) presented involve a single feature ( $d = 1$ ). The width of the beam is set to a default of  $w = 100$  (candidates that proceed to the next level). For the numeric attributes, the Cortana setting ‘best’ is applied, which means that for each attribute, all numeric threshold are considered and the optimal split point is selected. The resulting locally optimal subgroup is added to the result list if of sufficiently quality, and conditionally added to the beam for further refinement.

Before considering more systematic experiments, for example, testing the merits of the Fitness-Fatigue model, we first present results for a single skater, and demonstrate the kinds of input given to the coach concerning possible changes to the training schedule.

### 6.5.1 Demonstration of results

This section discusses results for female skater F1 who specializes in the medium to long distances. We first focus on 1500 m races, for which we have 60 examples over a period of five years. The average relative time for this skater was 1.0391, so 3.91% above the track record. We have training details for the entire five years, such that we can aggregate the preparation for each of these races easily.

**Uniform features** We start with univariate results in the simplest setting: uniform aggregates without specifiers and simple linear models. The best-fitting aggregate that was found was `max_load_1` with the following model:

$$y = -0.000014x + 1.042$$

The explained variance of this model is a mere  $R^2 = 0.0563$ . The model starts with a reasonable intercept, and encourages a high load (the product of duration and intensity) on the day prior to the race. Although the effect is minor, this suggests that a peak right before a race (possibly due to another race in the same weekend) is beneficial. The step function associated with this aggregate function, with a threshold around 360, has a more pronounced  $R^2 = 0.1233$ . The two levels are  $t_{rel} = 1.043$  for low maximum loads vs. 1.031.

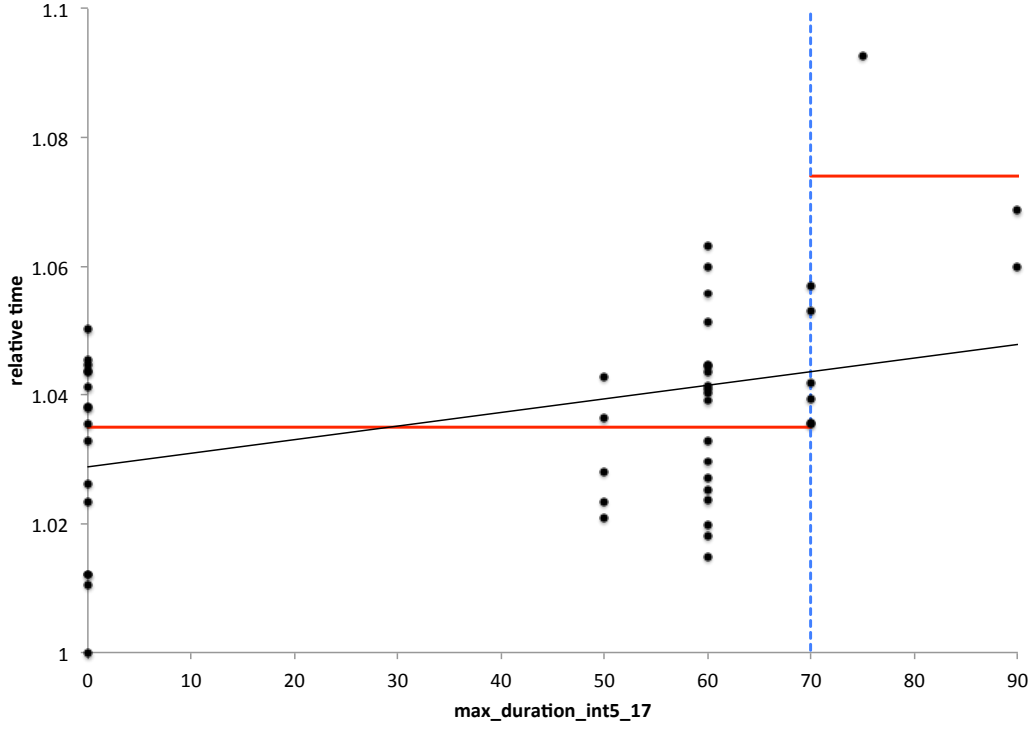
The second-best aggregate is `avg_intensity_19`, with model

$$y = -0.003952x + 1.052$$

which suggests that the intensity should be kept high over a period of almost three weeks to improve race times. This aggregate actually scores highest on explained variance of the step function, with a minimum average intensity of 3.93 (low to moderate intensity). Lower intensities suggest an expected relative time of  $t_{rel} = 1.044$ , whereas higher values on average lead to relative times of  $t_{rel} = 1.028$  ( $R^2 = 0.2231$ ). For clarity, details of these three features are listed in Table 6.2.

**Specifiers** The addition of specifiers does a great deal to the quality of the univariate models. The following aggregate scores the best on  $R^2$ :

- `max_duration_int5_17` (the largest period spent at intensity 5 for the last 17 days prior to the race)

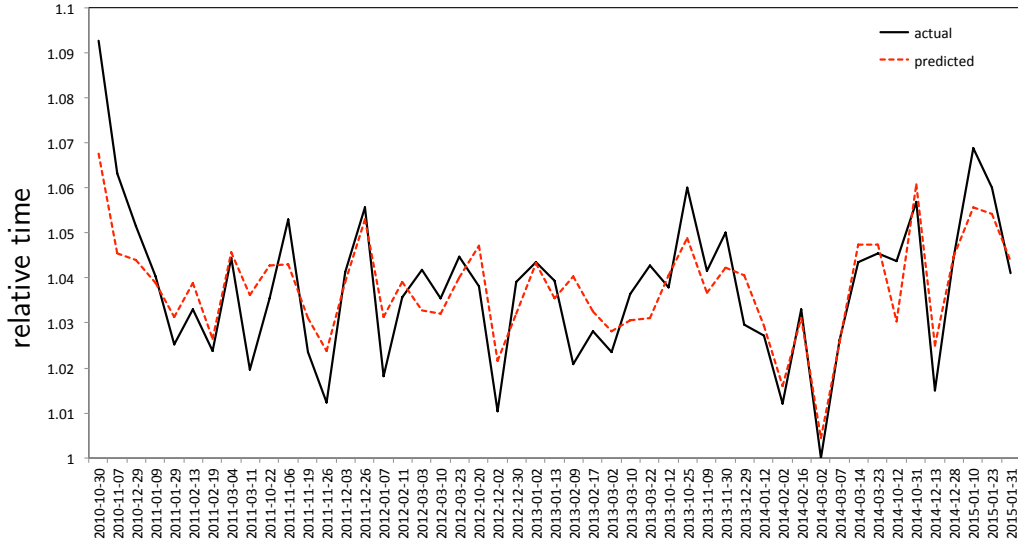


**Figure 6.3:** Graph showing the relation between `max_duration_int5_17` and the relative time of the subsequent race. The vertical dashed line indicates the threshold (inclusive to the left) and the solid horizontal lines the two average times for subgroup and complement. The black line with a slope indicates the best-fitting simple linear equation.

Exaggerating this type of training has a detrimental effect on the race outcome: longer durations of this specific training type lead to higher times. The step function ( $R^2 = 0.3080$ ) caps this value at 70 minutes. With longer durations of intensity 5, relative times of  $t_{rel} = 1.074$  are expected, compared to  $t_{rel} = 1.035$  below this threshold (Fig. 6.3).

**Fitness-Fatigue model** Switching from uniform to FF features, we note that the following parameters provide the best (linear) aggregate, based on the sum of duration:  $\tau_{fit} = 39, \tau_{del} = 4.0, \tau_{fat} = 7.0, K = 2.0$ . The corresponding kernel is the one featured in Fig. 6.2 earlier in this paper. The associated explained variance is  $R^2 = 0.1002$ .

**Multi-variate model** The individual features presented so far do not lead to very well-fitting models, despite their role in informing the coach of ways to



**Figure 6.4:** Obtained results (actual) over 60 historical races, compared to the predicted results of a relatively simple LASSO model based on the training sessions prior to each race.

**Table 6.2:** Overview of selected uniform features.

feature	linear model	$R^2$ linear	$R^2$ step	low	high
max_load_1	$-0.000014x + 1.042$	0.0563	0.1233	1.043	1.031
avg_intensity_19	$-0.003952x + 1.052$	0.0557	0.2231	1.044	1.028
max_duration_int5_17	$0.000211x + 1.029$	0.1491	0.3080	1.035	1.074

optimize the training and avoiding some pitfalls of under- and overtraining. More precise models can of course be obtained by involving multiple features. The graph in Fig. 6.4 presents the 60 results achieved by the skater, as well as the predicted times, by a multi-variate linear model. The model, induced by the LASSO procedure, involves 18 features, selected from the larger pool of uniform aggregates (ignoring the specifiers). The quality of the model is  $R^2 = 0.721$ , obtained on the training set<sup>6</sup>. Although such models (and more accurate models involving more complex aggregates) can be used to predict the outcome of an upcoming race, this particular prediction is of limited value. Rather, the model is more valuable from a knowledge discovery point of view, pointing to the features that matter most. In this case, the 18 features mostly include duration over short windows (one to five days), and intensities over longer windows (approx. two weeks).

<sup>6</sup>No distinction between training and test set was made for these experiments.

**Subgroups** In the above text, we have already presented the following three subgroups<sup>7</sup>:

- $\text{max\_load\_1} \leq 360$  ( $R^2 = 0.1233$ )
- $\text{avg\_intensity\_19} > 3.93$  ( $R^2 = 0.2231$ )
- $\text{max\_duration\_int5\_17} \leq 70$  ( $R^2 = 0.3080$ )

The first two subgroups relate to the experiment with just uniform aggregates, where the second subgroup is the optimal step function found at depth 1. The third subgroup relates to the experiment involving specifiers. While subgroups at depth 1 are interesting since they point to individual predictive features, they capture only shallow effects. We now present subgroups at greater depth, that indeed describe more complex concepts. The best subgroup found on the uniform windows (without specifiers) by Cortana at depth  $d \leq 2$  is

$$\text{avg\_intensity\_20} > 3.94 \vee \text{sum\_duration\_2} > 170 \quad (R^2 = 0.4232)$$

Although Cortana produces subgroups as conjunctions of conditions, for reasons of presentation this was logically inverted<sup>8</sup> in the above subgroup. The subgroup, covering 17 cases, describes races with an average of 1.0299, compared to 1.0526 for the remainder. It specifies that whenever the average intensity over the last 20 days is too low, and the total duration of exercises over the last 2 days is also low, this has a negative effect on the race result. Note how the explained variance has almost doubled at  $d \leq 2$ . Adding a third feature to the subgroups only produced a marginal improvement, which is not uncommon in SD.

The addition of specifiers in combination with deeper subgroups produced slightly better results, with the top subgroup being as follows:

$$\text{avg\_duration\_int5\_17} < 60 \vee \text{sum\_duration\_int6789\_10} > 115 \quad (R^2 = 0.446)$$

Note that compared to the  $d = 1$  result of  $R^2 = 0.3080$ , this is a reasonable improvement. The subgroup specifies that a lower duration of intensity 5 exercises over 17 day, or a higher duration of high-intensity exercises over 10 days, is good.

---

<sup>7</sup>Although subgroups are interpreted here as dichotomies, we present them as either lower or upper thresholds, such that cases meeting the condition(s) specified correspond to the faster races.

<sup>8</sup>The original subgroup discovered,  $\text{avg\_intensity\_20} \leq 3.94 \wedge \text{sum\_duration\_2} \leq 170$ , covers the complement.

**Validation of results** For the results above, one could wonder to what extent each result is statistically significant. For individual models, be it linear or step function, it is possible to compute a  $p$ -value that indicates to what extent the model might be due to chance. Such  $p$ -values will be reported in the detailed experiments in the next section. However, we should note that we are generating a substantial number of features, such that we are in fact testing multiple hypotheses. The best ranked result may thus appear to be significant, even though this is just a consequence of the many models considered. [23] presents a method for validating the results of an SD algorithm, by means of a *distribution of false discoveries*. This distribution is obtained by running the algorithm repeatedly on the data after swap-randomising the target attribute, thus capturing what maximum qualities can be obtained from random data (that resembles the original data). Using the distribution, it is possible to set a lower bound on the quality (in this case explained variance) as a function of the desired significance level  $\alpha$ . Assuming a significance level  $\alpha = 0.05$ , this validation method produces a lower bound of  $R_{min}^2 = 0.2907$  for the uniform data without specifiers, searching for subgroups at depth  $d = 1$ . This means that our optimal subgroup

$$\text{avg\_intensity\_19} > 3.93 \quad (R^2 = 0.2231)$$

is not actually significant at  $\alpha = 0.05$ . It is good to note that the lower bound produced by the swap randomization depends on the specific settings of the SD run. Specifically, if the extent of the search is bigger, more hypotheses will be tested, such that the lower bound will increase in order to account for the higher probability of finding a seemingly interesting subgroup by chance. When increasing the search depth to  $d \leq 2$ , the procedure produces a lower bound of  $R_{min}^2 = 0.4212$ . This makes the earlier depth 2 result (without specifiers)

$$\text{avg\_intensity\_20} > 3.94 \vee \text{sum\_duration\_2} > 170 \quad (R^2 = 0.4232)$$

just barely significant.

### 6.5.2 Fitness-Fatigue model

In this section, we analyse the specific merits of the Fitness-Fatigue model introduced in Sec. 6.3.2. We start by considering the four parameters the model involves, in order to fit the kernel to the specific physiological properties of the individual athlete. Rough values for the optimal setting were determined by informal experimentation, after which an extensive grid search was used

**Table 6.3:** Optimal parameters for the Fitness-Fatigue model for four speed skaters.

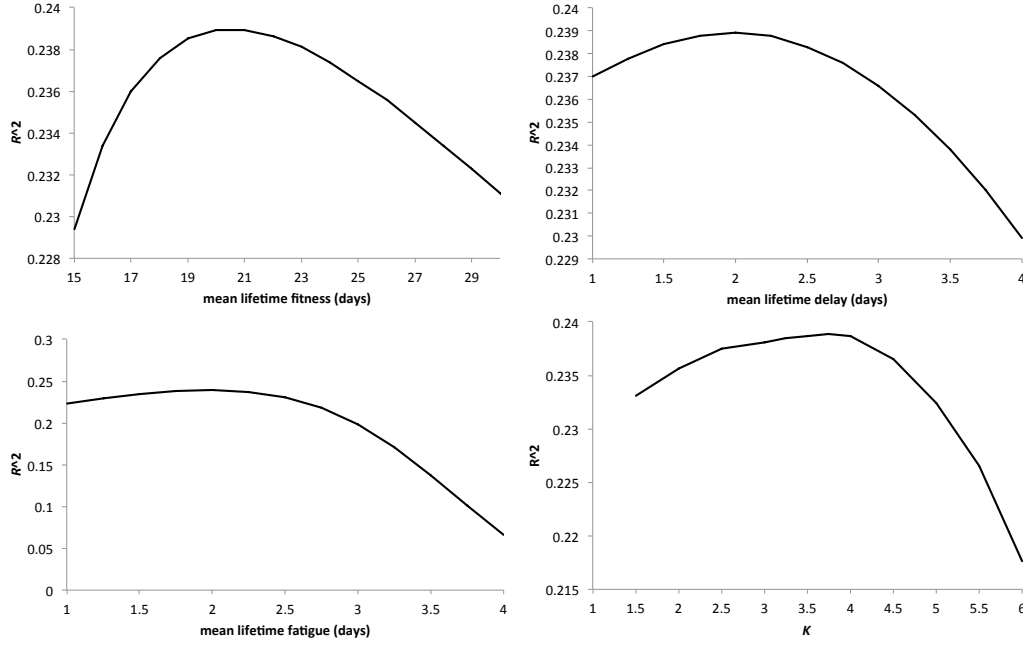
	Input context				Optimal parameters			
	Gender	Distance	Time	$n$	$\tau_{fit}$	$\tau_{del}$	$\tau_{fat}$	$K$
M1	Male	1000 m	1.0207	75	13	2.0	2.0	2.00
M2	Male	500 m	1.0212	142	21	2.0	2.0	3.75
F1	Female	1500 m	1.0391	60	39	4.0	7.0	2.00
F2	Female	500 m	1.0691	22	29	4.0	4.0	2.00

to determine the optimal values for each athlete involved. These results are demonstrated in Table 6.3, for four athletes and their respective distance of speciality. The left columns indicate the gender of the athlete, the preferred distance, the average relative time, and the number of races  $n$  available. The remaining columns indicate the optimal values for  $\tau_{fit}$ ,  $\tau_{del}$ ,  $\tau_{fat}$  and  $K$ . The three decay parameters are in days. As an optimization criterion, we select the  $R^2$  of the (univariate) linear model of the best feature found.

In order to analyse the stability of these parameters, we selected the third athlete (the one with the most available data), and varied each parameter individually, fixing the remaining the parameters to the optimum found earlier. The  $R^2$  of the best feature was recorded for each setting of the parameters. Fig. 6.5 demonstrates for each parameter how sensitive it is to change, in terms of quality of fit of the FF model. We note that all functions are very well-behaved and smooth over the domains considered, with the selected optimum clearly being undisputed. Furthermore, observe that the functions appear to be convex, making them fairly straightforward to optimize. Hence, the relatively simple grid search used in the pragmatic setting can be easily replaced by a more efficient hill-climber.

Let's consider the table of FF parameters in more detail. First of all, the rough numbers are very plausible from a physiological point of view. Clearly, the fatigue and (delayed) gain in fitness should be in the order of a few days, while the prolonged benefit of the exercise remains for a longer period in the order of several weeks. Also, the optimal values clearly differ per athlete, as a function of the different physiology and type of training the athlete is subjected to generally. Table 6.3 also suggests a difference between men and women, with men having a shorter time scale than women, both for the recuperation and how long the benefit lasts, although such conclusions are hard to draw from only four cases. Note also that the values reported





**Figure 6.5:** Analysis of sensitivity of the features to varying the four parameters  $\tau_{fit}$ ,  $\tau_{del}$ ,  $\tau_{fat}$  and  $K$  (from top left to bottom right). Clearly, these functions are smooth and well-behaved.

here are somewhat different from the ones reported in [16], which are:  $\tau_{fit} = 50$ ,  $\tau_{del} = 5$ ,  $\tau_{fat} = 15$ ,  $K = 2.0$ . As a last observation, we note that  $\tau_{del}$  and  $\tau_{fat}$  tend to assume very similar values. What impact this has from a physiological perspective (the fatigue and beneficial adaptation of the body go hand in hand?) is hard to say, but at least from a modeling perspective, it is a good opportunity to dispense with one parameter, and make the fitting of models more efficient.

Having a stable and physiologically plausible FF model of the training response, it is now time to turn to the question whether the model indeed produces a better fit, compared to our baseline of uniform aggregates. To this end, we again consider the explained variance  $R^2$  of the linear model on the best feature found, first for uniform features, then for the exponentially decaying features. Since above, the FF model was optimized without specifiers (intensity zones and morning/afternoon distinction), we compare the results with a similar setting for the uniform features. Table 6.4 presents these results. The columns marked “ $R^2$  linear” indicate the explained variance of the simple linear model. The indicated  $p$ -values for each result refer to the statistical significance of a linear regression  $t$ -test: the significance

**Table 6.4:** Comparison of goodness of fit of best univariate linear model for (middle) the uniform aggregates and (right) the Fitness-Fatigue model.

	Gender	Distance	Uniform		Fitness-Fatigue	
			$R^2$ linear	$p$ -value	$R^2$ linear	$p$ -value
M1	Male	1000 m	<b>0.41</b>	$4.16 \times 10^{-10}$	0.30	$3.85 \times 10^{-7}$
M2	Male	500 m	0.17	$2.46 \times 10^{-7}$	<b>0.23</b>	$7.02 \times 10^{-10}$
F1	Female	1500 m	0.06	$9.95 \times 10^{-3}$	<b>0.11</b>	$1.69 \times 10^{-3}$
F2	Female	500 m	<b>0.54</b>	$9.53 \times 10^{-5}$	0.50	$2.58 \times 10^{-5}$

of the best model, testing the hypothesis that the coefficient of the model is not 0 (in other words, testing whether the dependent variable is indeed influenced by the independent variable).

Based on these numbers, there is clearly not a consistent benefit of the FF model, over the less natural uniform features. Especially in the case of the first athlete, the uniform features are in fact more accurate. The feature in question (although there are multiple variants of similar score) concerns the sum of the duration over 9 days, which is an indication of too intense training and hence too high levels of fatigue as a result. Also for the fourth athlete, the uniform features come out on top. Still, for individual athletes the FF model may show a considerable improvement in fit over the unnatural uniform features.

**Alternative aggregate functions** The presentation of the FF model in terms of convolution translates into SQL as the **SUM** aggregate function. If features based on **SUM** have a potential benefit, so might the alternative functions **AVG** and **MAX**. We present an additional experiment in Table 6.5 that investigates the added value of these two aggregates to the plain implementation of the FF model used so far. The last two columns of the table show that in two cases, **AVG** or **MAX** do outperform the standard convolution. This is interesting, since there is clearly the potential to improve the models in this way. However, the features are less intuitive to understand since they are not based on the standard definition of convolution in terms of summation.

**Table 6.5:** Analysis of the benefit of adding **AVG** and **MAX** as aggregates to the Fitness-Fatigue model.

	Gender	Distance	SUM		SUM, AVG, MAX	
			$R^2$ linear	$p$ -value	$R^2$ linear	$p$ -value
M1	Male	1000 m	0.30	$3.85 \times 10^{-7}$	0.30	$3.85 \times 10^{-7}$
M2	Male	500 m	0.24	$7.02 \times 10^{-10}$	<b>0.28</b>	$1.60 \times 10^{-11}$
F1	Female	1500 m	0.11	$1.69 \times 10^{-3}$	<b>0.15</b>	$5.22 \times 10^{-4}$
F2	Female	500 m	0.50	$2.58 \times 10^{-5}$	0.50	$2.58 \times 10^{-5}$

## 6.6 Conclusion

We have presented a general approach to the modeling of training data in elite sports, with a specific application to speed skating in the LottoNL-Jumbo team. Our approach computes the combined effect of a training schedule by aggregating details of the individual training sessions, and thus capturing a considerable number of aspects of training and how one prepares for important test moments, such as physical tests and races. Since it is not entirely clear from the literature what aspects of training contribute the most, and what parameters individual athletes need to tweak in order to optimize the training to their specific physique and specialization, we produce a reasonably large collection of promising features. The most relevant features are then selected by a number of techniques, specifically univariate linear regression, the LASSO regression process and Subgroup Discovery. The linear modeling methods assume that the dependencies of interest are indeed monotonic and linear, that is, adding more load to the exercise will increase the (long-term) fitness of the athlete’s body. Clearly, this is not generally the case, and one would expect there to be certain thresholds, above which training is no longer beneficial. For each aspect of training, there is an optimal volume, above or below which training is ineffective. This suggests that non-linear models, or models that are able to represent thresholds (such as subgroups) will outperform linear models.

As mentioned in our introduction, we aim to discover interpretable and actionable patterns in the data, such that the coach can immediately incorporate the most significant findings in the preparation for upcoming events, as well as in future training schedules. We believe that our presented approach, that deliberately presents simple results, and gives clear guidelines and boundaries on training load, makes this possible. In fact, individual

findings on the athletes of the team have led to (subtle) modifications in training regimens, most notably where sprinters were sometimes subjected to too much aerobic exercise. It is good to stress again the athlete-specific nature of our analyses. Luckily, for a reasonable number of skaters, we have a long enough history to have a substantial database of training-response examples, where natural variation in preparation has produced a productive dataset. Athlete-specific data leads to athlete-specific findings, and one should therefore not interpret any discovered pattern as a general rule of exercise physiology, but rather as an opportunity to optimize training for that athlete.

We have presented a number of anecdotal results for a specific skater, demonstrating that interpretable and actionable results can be found. The best-fitting subgroup suggests that for a good result, this skater should avoid longer exercises at intensity 5 (over a longer window), as well as (slightly) increase the exposure to intensities 6 to 9. Although separate results appear very significant, a more thorough analysis using swap-randomization is necessary to account for the many features and models being considered. For this specific skater, statistically significant results could be obtained, despite there only being 60 available races. With up to 142 race results, other skaters will allow for much more significant findings.

The Fitness-Fatigue model, introduced as a more natural way to aggregate training impulses over time, produced reasonable results. After experimenting with four different skaters, two male and two female, very consistent and realistic values were found for the four parameters of this model. Although slight variations did occur, most notably between the male and female skaters, the general picture did match that of the coach. Knowing these (athlete-specific) parameters in detail allows the coach to mix exercise and recuperation in a more precise manner. From a modeling point of view, we also demonstrated that the optimal values of these parameters can be found efficiently, due to their well-behaved nature.

In a number of detailed experiments, we compared different choices in our modeling approach. A first experiment compared the uniform window to the Fitness-Fatigue kernel. Given the unnatural nature of a rectangular kernel, one would expect the FF model to be superior. Somewhat surprisingly, our data did not support this hypothesis. The FF model did indeed produce superior models for two athletes, but the uniform window performed equally superior on the remaining two, leaving this comparison unresolved.

Finally, we considered an experiment whether using the aggregate functions MAX and AVG, alongside the more obvious SUM, would be beneficial. This was

indeed sometimes the case, although not by a large margin. Whether such slightly more accurate models are in fact attractive is questionable, since the combination with the non-trivial FF kernel does not lead to very interpretable patterns.

