



Kent Academic Repository

Passfield, Louis and Hopker, James G. (2017) *A mine of information: can sports analytics provide wisdom from your data?* International journal of sports physiology and performance, 12 (7). pp. 851-855. ISSN 1555-0265.

Downloaded from

<https://kar.kent.ac.uk/58658/> The University of Kent's Academic Repository KAR

The version of record is available from

<https://doi.org/10.1123/ijsp.2016-0644>

This document version

Author's Accepted Manuscript

DOI for this version

Licence for this version

UNSPECIFIED

Additional information

Versions of research works

Versions of Record

If this version is the version of record, it is the same as the published version available on the publisher's web site. Cite as the published version.

Author Accepted Manuscripts

If this document is identified as the Author Accepted Manuscript it is the version after peer review but before type setting, copy editing or publisher branding. Cite as Surname, Initial. (Year) 'Title of article'. To be published in *Title of Journal*, Volume and issue numbers [peer-reviewed accepted version]. Available at: DOI or URL (Accessed: date).

Enquiries

If you have questions about this document contact ResearchSupport@kent.ac.uk. Please include the URL of the record in KAR. If you believe that your, or a third party's rights have been compromised through this document please see our [Take Down policy](https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies) (available from <https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies>).



A mine of information: can sports analytics provide wisdom from your data?

Journal:	<i>International Journal of Sports Physiology and Performance</i>
Manuscript ID	IJSPP.2016-0644.R1
Manuscript Type:	Invited Brief Review
Keywords:	training, physical activity, physical performance, exercise training

SCHOLARONE™
Manuscripts

Brief Review

1 Title: A mine of information: can sports analytics provide wisdom from your data?

2

3 Louis Passfield and James G. Hopker

4 Endurance Research Group, School of Sport and Exercise Sciences

5 University of Kent. Chatham Maritime. UK. ME4 4AG

6

7 Corresponding Author

8 Louis Passfield

9 Endurance Research Group, School of Sport and Exercise Sciences

10 University of Kent. Chatham Maritime. UK. ME4 4AG

11

12 Email: l.passfield@kent.ac.uk

13 Twitter: trainalytics

14

15

16 | Word Count: ~~41774299~~

17

18 Abstract: 222

19

20 Figures: 4 Tables: 0

21

22 Abstract

23 This paper explores the notion that the availability and analysis of large datasets has
24 the capacity to improve practice and change the nature of science in the sport and
25 exercise setting. The increasing use of data and information technology in sport is
26 giving rise to this change. Websites hold large data repositories and the development
27 of wearable technology, mobile phone applications and related instruments for
28 monitoring physical activity, training and competition, provide large data sets of
29 extensive and detailed measurements. Innovative approaches conceived to exploit
30 more fully these large datasets could provide a basis for more objective evaluation of
31 coaching strategies and new approaches to how science is conducted. The emergence
32 of a new discipline, sports analytics, could help overcome some of the challenges
33 involved in obtaining knowledge and wisdom from these large datasets. Examples of
34 where large datasets have been analyzed, to evaluate the career development of elite
35 cyclists, and to characterize and optimize the training load of well-trained runners are
36 discussed. Careful verification of large datasets is time consuming and imperative
37 before useful conclusions can be drawn. Consequently, it is recommended that
38 prospective studies are preferred to retrospective analyses of data. It is concluded that
39 rigorous analysis of large datasets could enhance our knowledge in the sport and
40 exercise sciences, inform competitive strategies, and allow innovative new research
41 and findings.

42 In recent years there has been an explosion in the use of information technology
43 within the sport and exercise fields. The data and ~~thus~~ information derived from these
44 advances has long been recognized to have the potential for a profound impact¹.
45 Websites ~~now~~ accumulate large repositories of primary and secondary data that
46 previously would have been impossible for sport and exercise scientists to access and
47 collate by hand. The ~~instrumentation of equipment and~~ invention of wearable
48 technology enables extensive measurements to be gathered during exercise, training
49 and competition. Increasingly, athletes and coaches recognize that such detailed, high
50 quality data can be used to inform objective decision making on aspects of training
51 and performance. In this paper we discuss how rigorous analysis of large datasets may
52 hold the potential to change ~~not only~~ sport, ~~but and~~ the nature of its related sciences
53 too.

54 “Moneyball”², and “Big Data” style stories in high performance sport readily capture
55 the public interest, but ~~there remains a question as to whether it’s not clear that~~
56 scientists are making the most of their available data. There is a risk that the
57 unprecedented ~~capacity for obtaining volume of~~ data ~~is~~ overwhelming and ~~prevents us~~
58 ~~from not~~ ~~used fully ing it to obtain insight and~~ inform practice. Consequently, ~~it seems~~
59 ~~appropriate to ask if we suggest~~ there is scope to advance by following other
60 disciplines (such as business and economics), in developing methods to analyze more
61 rigorously the extensive data sources available to us.³ Rowley³, ~~suggests proposes~~
62 ~~that~~ a wisdom hierarchy of data processing ~~exists~~. This hierarchy ~~sees describes how a~~
63 ~~mass of~~ raw data ~~is~~ converted into information, ~~the~~ information into knowledge, and
64 ~~the~~ knowledge into wisdom. Gaining ~~this~~ knowledge and wisdom from data is
65 challenging, but could spawn a new discipline in the sports sciences, that of sports
66 analytics.

67 Thornton et al.³⁴ note that ~~the ubiquity of~~ mobile phones and wearable technology
68 present simple methods to assess and promote physical activity but this area is still
69 underdeveloped. Excellence in the nascent field of sports analytics ~~promises will need~~
70 to ~~help~~ sieve ~~the deluge of~~ data from repositories and ~~these~~ devices in order to filter
71 out meaningful information. The benefits of this work could be wide-ranging ~~for the~~
72 ~~coach and scientist~~, such as identifying new talent, optimizing training programs,
73 informing team selection, and deriving and evaluating competition tactics. The

74 success of sports analytics will be governed by whether its findings can be translated
75 clearly and for the benefit of ~~its~~ users, such as exercisers, athletes, and ~~their~~ coaches.
76 A further challenge for sports analytics is that ~~in order to conduct~~ effective data
77 analysis requires; a fusion of diverse expert knowledge ~~has to occur~~; for example, in
78 training theory, sports psychology, data handling and analysis, statistics and
79 mathematical modelling, determinants of performance, and competition strategies. ~~At~~
80 ~~the moment~~ [This presents a genuine interdisciplinary challenge as few, if any,
81 individuals are sufficiently well versed in such disparate areas. Thus for sports
82 analytics to fully mature as a discipline, new opportunities for the development of its
83 practitioners are needed ~~to be conceived~~. This will likely require universities to
84 develop new courses that ~~enable students to combine~~ and acquire a deep
85 understanding of the science of sport, alongside extensive skills for data handling and
86 analysis.

87 ~~In this paper~~ Next we provide two examples of the kind of opportunities that can be
88 found in tackling this challenge, and discuss some consequent issues. We present two
89 preliminary studies from our endurance research group that illustrate different ways of
90 mining and modelling data to look at talent development and optimization of training.
91 Our aim is to promote wider ~~recognition and~~ discussion ~~of the evolving discipline~~ of
92 sports analytics and its potential to influence research and practice in the sport and
93 exercise sciences.

94 Obtaining large datasets for analysis

95 Once a research question has been established, one way of addressing it can be to
96 evaluate existing data ~~that has already been gathered~~. Data mining is a method where
97 raw data is translated into information by analyzing and interpreting its patterns
98 ~~within the data set~~. Data mining may also involve mathematical or statistical
99 modelling, particularly where some kind of predictive capacity is required. The
100 ~~Information information might be~~ obtained ~~from data~~ can be used to help coaches
101 predict changes in sports performance⁵, ~~find events that co-occur or their sequence of~~
102 ~~occurrence, and divide data into similar groups~~⁶. ~~Data mining techniques have been~~
103 ~~used to obtain information by examining~~ examine the relationship between
104 performance and its determinants ~~attributes~~⁶, and to interrogate athletes' existing
105 performance ~~related~~ data to identify new strategies.^{7,8,9} Ofoghi et al.⁷⁻⁸ ~~show how used~~

106 | data mining ~~could be used~~ to inform strategic planning for rider selection and training
107 | prioritization in the multi-discipline ~~events such as the~~ omnium in track cycling, and
108 | Moffatt et al.⁹ for identifying sprint race tactics. There is a cost though, as ~~in many~~
109 | ~~instances~~ the amount or complexity of the data, and preparing it for analysis can
110 | challenge even the most determined, especially where each athlete, team, game or
111 | event, across a season is modelled. It is also very important that the research question
112 | and methods are established before analysis is begun¹⁰. The evaluation of an
113 | hypothesis formed *a priori* helps to reduce the chance of bias and false positives
114 | arising from the analytic process. Otherwise, data fishing or P-hacking in large
115 | datasets ~~is likely to may~~ result in ~~many~~ spurious but statistically significant results.

116 | Analyzing race results

117 | Some websites provide the potential to exploit large datasets by analyzing ~~their~~
118 | ~~information they hold~~ data. With the website's permission it is possible to use web-
119 | spider or web-crawler software to extract data ~~from its databases~~ for subsequent
120 | analysis. We examined the career progression and success of elite cyclists by using
121 | this approach to conduct a retrospective analysis of ~~their~~ race results.¹⁰⁺¹¹ ~~It is Coaches~~
122 | ~~and scientists~~ generally accepted that athletes ~~have to~~ undertake many years of
123 | training to achieve elite status in endurance sports. Yet the development profile of the
124 | most successful senior athletes and ~~the likelihood that whether~~ this involves
125 | performing well in ~~elite~~ junior competitions remains unclear.¹¹⁺² To explore this issue
126 | we extracted race results for major junior and senior elite cycling races from 1980 to
127 | 2014 from one of the freely accessible online databases documenting race results
128 | (www.procyclingstats.com). For the purposes of the study we focused upon 25 major
129 | races and were able to obtain 67,503 results for 5,561 cyclists from 75 countries. This
130 | data included the name, date of birth, nationality, race, finishing position (including
131 | general classification and individual stage results from multi-stage races) of all the
132 | cyclists competing. From this data we were able to establish that the cyclists' average
133 | career length for competing in these most prestigious races was 3 seasons. However,
134 | as the data was heavily skewed by a few highly prolific cyclists, we also used the
135 | semi-interquartile range (SIQR) as an alternative way of depicting cyclists' typical
136 | career length. The SIQR comprises of the 50% of data between 25th and 75th
137 | percentile and it showed that half of all cyclists' careers ranged between 1 and 7

138 years. Notably, a large proportion of cyclists (86%), never achieve a top 10 placing in
139 the major races we studied in their career. Our data mining also revealed ~~findings with~~
140 ~~implications for long-term development of cyclists, and team selection. As shown in~~
141 ~~Figure 1, we identify~~ evidence of a relative age effect¹²⁴³, sometimes referred to as the
142 Matthew effect, within the population of world-class cyclists.

143

144 ~~**** Figure 1 near here ****~~

145

146 There ~~appears to be~~ an over-representation of cyclists at the World Tour level who
147 were born early in the calendar year (January-March). This ~~analysis raises the issue of~~
148 ~~whether observation suggests there is an inappropriate bias in how~~ cyclists are ~~being~~
149 identified and developed ~~by their coaches. To avoid this coaches should encourage a~~
150 ~~later specialization and prematurely, or on an inappropriate basis e.g. more~~ focus
151 upon technical skills, rather than physiological parameters ~~be better for~~ developing
152 young cyclists. ~~Varying~~ the youth cyclists' age group cut-off dates within the
153 competition year (e.g. ~~should~~ 9 or 15 months ~~be used~~ rather than 12 months) could
154 ~~also~~ be considered. ~~Or alternatively, y~~outh teams could ~~also~~ have quotas based upon
155 chronological age within a year. ~~This Only~~ interrogation of a large volume of race
156 data allowed us to describe ~~thi~~ evolution of successful cyclists' and ~~substantiate the~~
157 ~~presence of~~ identify the "Matthew effect" within elite cycling.

158 There are several challenges with establishing the validity and reliability of large
159 datasets, especially where the analysis is retrospective ~~that need to be considered prior~~
160 ~~to conducting a study~~. For this reason, a prospective study design is often preferable
161 ~~in order so~~ that the integrity of the data can be overseen as it is gathered. ~~Trying to~~
162 ~~verify~~ Establishing the veracity of large numbers of observations retrospectively is
163 often impractical. For example, in our study above¹⁰⁴⁴ the collection of retrospective
164 race results from 3rd party websites using web-crawlers assumed these were
165 accurately reported to reflect the "official" finishing positions. Moreover, collecting
166 data in this way brought with it ethical considerations when deciding where, and how
167 fast to crawl. Prior permission was always obtained from the data or website owner.

168 Nonetheless, fast crawlers can have a crippling impact on the performance of a
169 website as the server deals with multiple simultaneous requests. ~~Once the web crawler~~
170 ~~finished gathering data, p~~Pre-processing of the data was imperative to check for errors
171 in its structure, and for subsequent filtering and cleaning. Within the cycling results
172 database ~~for example, some~~ race names ~~had~~ changed over the years, or were listed in
173 both native and English languages across various ~~editions e.g. Tour de Pologne/Ronde~~
174 ~~van Polen/Polen Rundfahrt/Tour of Poland. In some instances, there w~~Where ~~results~~
175 ~~were~~ missing ~~results we that~~ needed ~~to~~ verification ~~of y~~ whether the race took place, ~~or~~
176 ~~if its results were just absent from the database.~~ Similarly, ~~where misspelt~~ cyclists'
177 names were ~~misspelt they needed to be~~ corrected prior to analysis, ~~otherwise their to~~
178 ~~ensure their~~ results ~~would have been inere~~ correctly assigned. In ~~shortsummary, the~~
179 ~~opportunity to analyzinge~~ large data sets ~~can provide as~~ a means of answering ~~to pre-~~
180 specified research questions ~~provides the chance to extractwith~~ novel findings. It does
181 require substantial meticulous and time-consuming work though, and the approach
182 should not be regarded as a surrogate for prospectively conducted studies.
183 Furthermore, conducting prospectively designed studies will help reduce the chance
184 of bias and false positives¹⁰ as mentioned previously.

185 Analysis of exercise and training data

186 When ~~athletes and coaches~~ monitoring exercise, training and racing, large datasets are
187 now generated routinely. Advances in training technology have resulted in portable
188 devices (such as accelerometers and similar activity monitors, GPS, heart rate
189 monitors, power output meters, and related mobile phone apps), being used habitually
190 to gather data by a wide spectrum of users from recreational exercisers to elite
191 athletes. These devices typically gather data on all the activity of their users with a
192 level of accuracy and detail once unthinkable. Characteristically, this data has been
193 used to describe ~~and recount~~ completed exercise or training bouts and
194 races.^{13+4,14+5,15+6,16+7} However, by exploiting these opportunities more fully, scientists
195 could produce exciting and innovative ~~new~~ findings. With this technology,
196 performance can ~~now~~ be evaluated directly in the field, rather than be inferred from
197 laboratory trials and simulations. Accurate measurements that previously required
198 specialised laboratory equipment ~~are can be now~~ gathered ~~by the coach~~ during normal
199 training and competition (Figure 12). Furthermore, patterns of daily activity and

200 | inactivity can be described to evaluate lifestyle interventions more objectively¹⁷⁺⁸. As
201 | a consequence, more realistic and ecologically valid experimentation can be designed
202 | and questions addressed that were previously beyond the reach of the laboratory-
203 | based scientist. An ~~enticing~~-example of ~~this is these~~ insights ~~that could come from~~
204 | ~~being able to~~ is in accurately quantify prescribing training.

205 | ***** Figure 12 near here *****

206 | To date the process of prescribing training has relied upon the experience and
207 | intuition of those involved (i.e. coaches and athletes), as the necessary research in this
208 | area is lacking¹⁸⁺⁹. Over the past four decades, the scientific basis for prescribing
209 | training programs has advanced little beyond Banister and colleagues' seminal
210 | work^{19+20,21}. This is in marked contrast to the ~~tremendous~~-advances ~~that have been~~
211 | made in our understanding of the adaptations that result from training²². However, this
212 | ~~situation~~ could change with the capability to measure individuals' training and racing
213 | accurately and in detail in the field. The resulting large volumes of field
214 | measurements ~~could present~~ allows the discipline of sports analytics with an early
215 | opportunity to contribute to our understanding of effective training program
216 | prescription²³. Furthermore, ~~this~~-detailed monitoring of training and performance in
217 | the field provides an opportunity to reverse the usual scientific paradigm for research
218 | on this topic. Specifically, instead of conducting experiments to compare the effects
219 | of specific (laboratory-based) training regimens, we can measure study participants'
220 | training, and track their resulting changes in performance. It may then be possible to
221 | determine which aspects of their monitored training is most effective, given sufficient
222 | data. With this scientific paradigm the method of enquiry consists of identifying
223 | which training led to the observed changes in performance, rather than trying to
224 | evaluate how performance changes in response to a ~~carefully~~-restricted laboratory-
225 | based training protocol. Here the bigger the data, the better the insight, as effective
226 | training ~~is likely to~~ may be identified more clearly when the number of participants
227 | involved and the diversity of their training is greater. Exploring a wide range of
228 | training regimes with large numbers of participants is not a viable option for
229 | laboratory-based research, but in a field study it becomes quite plausible. Participants
230 | can be recruited to undertake their usual training program and compete in their

231 preferred competitions, no longer restricted to ~~following~~ scientists' abstract training
232 regimes ~~and or~~ ~~evaluating them~~ with contrived laboratory-based performance trials.

233 ~~Studies involving our endurance research group have demonstrated the potential for~~
234 ~~extracting useful insights from carefully conducted field studies.~~ Galbraith et al.²⁴
235 used GPS devices to record all the training and performances of 14 highly-trained
236 endurance runners for a year-long study. This study resulted in measurements for 2.5
237 million time-points. In ~~our the~~ original analysis we summarized and collapsed this
238 data into 3 training zones, finding total distance, and percent time spent at the highest
239 intensity related to performance. This kind of analysis is difficult to translate into
240 future training prescription for athletes however. Therefore, in order to analyze this
241 dataset more fully Kosmidis and Passfield²⁵ proposed the use of training distribution
242 and training concentration profiles (Figures ~~32 and 3~~ respectively). This training
243 distribution profile is obtained by plotting the amount of time spent above the
244 reference speed during the session. For example, at 0 km·h⁻¹ all the training was
245 completed above this speed and therefore the total number of observations for the
246 session is plotted. In contrast, at 15 km·h⁻¹ only a small fraction of the total
247 observations is seen to occur above this speed. In effect the analysis assumes every
248 possible speed is a training threshold and shows how the pattern of training time
249 changes with speed. The training concentration profile is the derivative of the
250 distribution curve or in statistical terms a concentration curve. It shows the cumulative
251 time spent training at each speed during the session(s) analyzed. By comparing the
252 training distribution profiles with resulting changes in performance, ~~these~~ researchers
253 were able to identify the runners' training speeds that were significantly related to
254 improvement. ~~Not only could they identify these significant speeds for training, but~~
255 ~~†They could also his information was used to~~ model how endurance performance
256 would change in response to training. Notably, ~~the authors observed that~~ the
257 significant training speeds could not be determined from laboratory test data, but only
258 from the analysis of the runners' training and performances. These methods and
259 findings indicate that ~~in the future~~ it may be possible to support the coach by
260 identifying the optimal training sessions for athletes to complete for specific race
261 performances. Perhaps even more importantly, ~~people those promoting exercising~~
262 exercise for health could ~~specify their available training time, and~~ use the same

Formatted: Superscript

263 method to calculate the most efficient exercise regime that provides ~~the~~ maximum
264 benefits.

265

266 **** Figures 2 and 3 near here ****

267

268 There were some theoretical issues that the training analysis highlighted. Kosmidis
269 and Passfield²⁵ set out with the ambition to retain all of the available data, to minimize
270 the number of assumptions they made, and still utilize a parsimonious model with as
271 few predictor variables as possible. When a data set is summarized, ~~whether such as~~
272 with a mean and standard deviation ~~or something more complex~~, much of the
273 information in the original dataset is compressed ~~in the process~~ too. An advantage of
274 the training distribution and concentration profiles is that they retain all the available
275 data from every session for analysis. Furthermore, ~~relatively assumption-less~~
276 ~~approach to modelling their data meant~~ the authors did not rely on existing models of
277 physiology to make sense of the data. ~~Rather they made the data “talk” and checked~~
278 ~~subsequently to see if their analysis supported traditional physiological models of~~
279 ~~training~~. As mentioned above, their findings did not support existing models used for
280 training, as their traditional laboratory tests results could not be used to identify the
281 training speeds that were related significantly to the changes in performance. If the
282 training data had been described with reference to the laboratory test data (i.e. as
283 percentages of maximum or lactate threshold) ~~at the outset~~, the analysis would not
284 have succeeded. Finally, as with most modeling work, a key challenge is ensuring
285 parsimony to keep the model as simple as is reasonable. The training distribution and
286 concentration curves help this process by reducing the complexity of the underlying
287 dataset whilst still retaining a simpler, yet comprehensive representation of it.

288 ~~There are many challenges to be overcome before it will be possible to introduce a~~
289 ~~rigorous scientific method into the process of prescribing training. Nonetheless some~~
290 ~~important lessons were learned from the studies above~~. Data cleaning and checking
291 was an arduous process, as with the study of cyclists' development profiles discussed
292 earlier¹¹. Every training session was plotted and manually inspected for obvious

293 | errors. This process ~~quickly~~ highlighted ~~that~~ the ~~subsequent analysis would have need~~
294 | to deal with unrealistic “spikes” in the recorded values, and calculations where the
295 | training speed was at, or close to, zero. In addition to ~~clear~~ visibuale data spikes, we
296 | also had to identify unreasonable values e.g. where the apparent speed was ~~clearly~~
297 | above world record pace for the observed distance. These observations were due to
298 | problems with the GPS signal, or runners forgetting to switch off their GPS ~~when~~
299 | ~~cycling or driving home~~ after a training session or race. Most of these issues could be
300 | addressed within the analysis, but a particular ~~challenging~~ issue was how to proceed
301 | in the absence of data. ~~All the runners were asked to submit their training programs,~~
302 | ~~as these were not specified by the research team.~~ By matching the observed training
303 | data to the runners’ training program ~~provided~~, gaps caused by missing training data
304 | were identified. ~~The athletes’ training record could also be used to determine~~
305 | ~~wh~~eOther missing observations in a training session implied a rest period, a gap
306 | between successive sessions, or a runner moving ~~very~~ slowly, ~~or simply missing data.~~
307 | However, as this was a retrospective analysis ~~of the data~~ of data from an earlier
308 | study²⁴, ~~it was not always possible to confirm~~ these assumptions ~~were not always~~
309 | ~~possible to verify.~~ ~~As discussed earlier in this paper v~~erifying the dataset was a
310 | time-consuming but critical part of the analysis. This ~~re-emphasiz~~underlines our
311 | earlier recommendation that scientists prefer conducting prospective studies, ~~as~~
312 | ~~opposed~~ to retrospective analyses of large ~~training~~ data sets when ~~ever~~ possible.

313

314 | Summary

315 | Technological advances in recent years have enabled large datasets to be gathered in
316 | sport and exercise settings. Examples of these large datasets are information held by
317 | websites, and data generated by people monitoring their regular exercise, training or
318 | competitions. Careful analysis of these large datasets can enhance our knowledge in
319 | the sport and exercise sciences, support the coach by informing competitive strategies,
320 | and allow innovative new research and findings. The interest in making more from the
321 | data in sport and exercise sciences appears to be spawning a new discipline of sports
322 | analytics. This discipline necessitates the fusion of a diverse range of knowledge in
323 | computing, mathematics, statistics and sports sciences, that may require new
324 | development opportunities before the discipline can develop fully. Examples of

325 preliminary work exploring large datasets from websites and GPS devices have been
326 discussed along with some of the issues that this work presents. A common theme for
327 this kind of work is that careful quality checking of the large dataset is imperative and
328 time-consuming. Identification of missing data and strategies for dealing with it is
329 also critical. Accordingly, it is recommended that prospective studies are preferred to
330 retrospective analyses of data.

331

332 References

- 333 1. Liebermann DG, Katz L, Hughes MD, Bartlett RM, McClements J, and
334 Franks IM. Advances in the application of information technology to sport
335 performance. *J Sports Sci.* 2002;20(10):755-769. PubMed doi:
336 10.1080/026404102320675611
- 337 2. Lewis M, *Moneyball: The Art of Winning an Unfair Game.* New York. Norton
338 & Company; 2004.
- 339 3. Rowley J, The wisdom hierarchy: representations of the DIKW hierarchy.
340 *Journal of Information Science.* 2007; 33:163–180. doi:
341 10.1177/0165551506070706
- 342 4. Thornton JS, Frémont P, Khan K, Poirier P, Fowles J, Wells GD, Frankovich
343 RJ. Physical activity prescription: a critical opportunity to address a
344 modifiable risk factor for the prevention and management of chronic disease: a
345 position statement by the Canadian Academy of Sport and Exercise Medicine.
346 *Br J Sports Med.* 2016;50(18):1109-14. PubMed doi: 10.1136/bjsports-2016-
347 096291.
- 348 5. Cangley P, Passfield L, Carter H, and Bailey M. A model for performance
349 enhancement in competitive cycling. *Movement & Sport Sciences.* 2012;1:59-
350 71.

- 351 6. Ofoghi B, Zeleznikow J, MacMahon C, Raab M. Data mining in elite sports:
352 A review and a framework. *Meas Phys Educ Exerc Sci*. 2013;17:171-186. doi:
353 10.1080/1091367X.2013.805137
- 354 7. Chen I, Homma H, Jin C, Yan HH. Identification of elite swimmers' race
355 patterns using cluster analysis. *Int J Sports Sci Coach*. 2007;2, 293–303. doi:
356 10.1260/174795407782233083
- 357 8. Ofoghi B, Zeleznikow J, MacMahon C, Dwyer D. A machine learning
358 approach to predicting winning patterns in track cycling omnium. In M.
359 Bramer (Ed.) *Proceedings of the International Federation for Information*
360 *Processing (IFIP)*. Conference on Advances in Information and
361 Communication Technology 2010; pp. 67–76. Brisbane, Australia: Springer
362 Berlin Heidelberg.
- 363 9. Moffatt J, Scarf P, Passfield L, McHale IG, Zhang K. To lead or not to lead:
364 analysis of the sprint in track cycling. *J Quant Anal Sports*. 2014;10(2): 161-
365 172. doi: 10.1515/jqas-2013-0112
- 366 10. Head ML, Holman L, Lanfear R, Kahn AT, Jennions MD. The extent and
367 consequences of P-hacking in science, *PLoS Biol*. 2015; 13(3):1-15. PubMed
368 doi:10.1371/journal.pbio.1002106
- 369 11. Hopker J, Dietz KC, Schumacher YO, Passfield L, Using retrospective
370 analysis of race results to determine success in elite cycling. *Journal of*
371 *Science and Cycling*. 2015;4:29.
- 372 12. Schumacher YO, Mroz R, Mueller P, Schmid A, Ruecker G. Success in elite
373 cycling: A prospective and retrospective analysis of race results. *J Sports Sci*.
374 2006;24(11):1149-56. PubMed doi. 10.1080/02640410500457299
- 375 13. Armstrong N, Welsman J. Physiology of the child athlete, *Lancet*,
376 2005;366:s44-45.
- 377 14. Esteve-Lanao JO, San Juan AF, Earnest CP, Foster CA, Lucia AL. How do
378 endurance runners actually train? Relationship with competition performance.

- 379 *Med Sci Sports Exerc.* 2005; 37(3):496-504.
- 380 15. Esteve-Lanao J, Foster C, Seiler S, Lucia A. Impact of training intensity
381 distribution on performance in endurance athletes. *J Strength Cond Res.* 2007;
382 21(3):943-949. PubMed doi: 10.1519/R-19725.1
- 383 16. Muñoz I, Seiler S, Bautista J, España J, Larumbe E, Esteve-Lanao J. Does
384 polarized training improve performance in recreational runners. *Int. J. Sports*
385 *Physiol. Perform.* 2014;9(2):265-272. PubMed doi: 10.1123/ijsp.2012-0350
- 386 17. Neal CM, Hunter AM, Brennan L, O'Sullivan A, Hamilton DL, DeVito G,
387 Galloway SD. Six weeks of a polarized training-intensity distribution leads to
388 greater physiological and performance adaptations than a threshold model in
389 trained cyclists. *J Appl Physiol.* 2013;114(4):461-471. PubMed doi:
390 10.1152/jappphysiol.00652.2012
- 391 18. Strath SJ, Kaminsky LA, Ainsworth BE, Ekelund U, Freedson PS, Gary RA,
392 Richardson CR, Smith DT, Swartz AM. *Circulation.* 2013;128(20):2259-79.
393 PubMed doi:10.1161/01.cir.0000435708.67487
- 394 19. Borresen J, Lambert MI. The quantification of training load, the training
395 response and the effect on performance, *Sports Med.* 39(9):779-795. PubMed
396 doi: 10.2165/11317780-000000000-00000
- 397 20. Banister EW, Calvert TW, Savage MV, Bach TM. A system model of training
398 for athletic performance. *Australian Journal of Sports Medicine*, 1975;7: 57-
399 61.
- 400 21. Calvert TW, Banister EW, Savage MV, Bach T. A systems model of the
401 effects of training on physical performance. *IEEE Trans. Syst. Man Cybern.*
402 1976;6:94-102.

- 403 22. Mann T, Lamberts RP, Lambert MI. Methods of prescribing relative exercise
404 intensity: Physiological and practical considerations. *Sports Med.* 2013;43(7):
405 613-625. PubMed doi: 10.1007/s40279-013-0045-x
- 406 23. Jobson SA, Passfield L, Atkinson G, Barton G, Scarf P. The analysis and
407 utilization of cycling training data. *Sports Med.* 2009;39(10):833-844.
408 PubMed doi: 10.2165/11317840-000000000-00000
- 409 24. Galbraith A, Hopker JG, Cardinale M, Cunniffe B, Passfield L. A 1-year study
410 of endurance runners: Training, laboratory tests, and field tests. *Int J Sport*
411 *Phys Perf.* 2014; 9:1019-25. PubMed doi: 10.1123/ijsp.2013-0508
- 412 25. Kosmidis I, Passfield L. Linking the performance of endurance runners to
413 training and physiological effects via multi-resolution elastic net. *arXiv*
414 *preprint 2015;arXiv:1506.01388.*

415

416 Figure Legends

417

418 Figure 1: ~~The percentage of riders placing in the top 10 of World Tour cycling races~~
419 ~~by birth month. Data is percentage normalized for month length. The horizontal line~~
420 ~~at 8.33 represents the uniform distribution over the 12-month period.~~¹¹

421

422 ~~Figure 2:~~ A training session for an endurance runner, showing running speed over
423 time. Data were gathered by wrist-worn GPS recording every second for each variable
424 measured.

425

426 Figure ~~23~~: A training distribution profile for the training session shown in Figure ~~12~~
427 as proposed by Kosmidis and Passfield²⁴²⁵²⁵ for analyzing large training datasets. The
428 distribution profile shows the total session time spent training above the
429 corresponding speed.

430

431 Figure ~~34~~: A training concentration profile for the training session shown in Figure ~~12~~
432 as proposed by Kosmidis and Passfield²⁴²⁵²⁵. The concentration profile shows the
433 session time spent training at the corresponding speed.

434

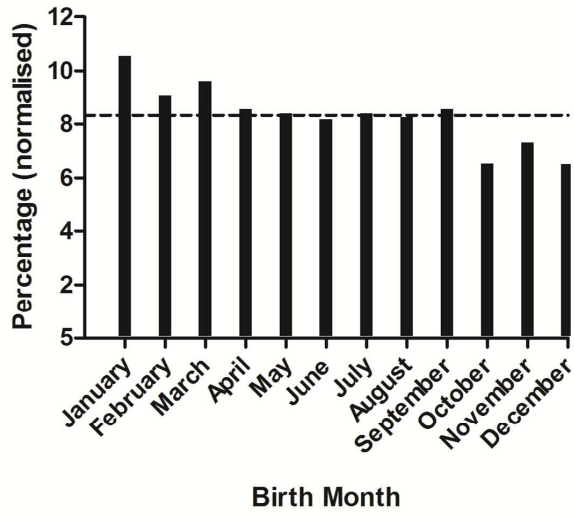
435

Formatted: List Paragraph

436

Figures

437 Figure 1

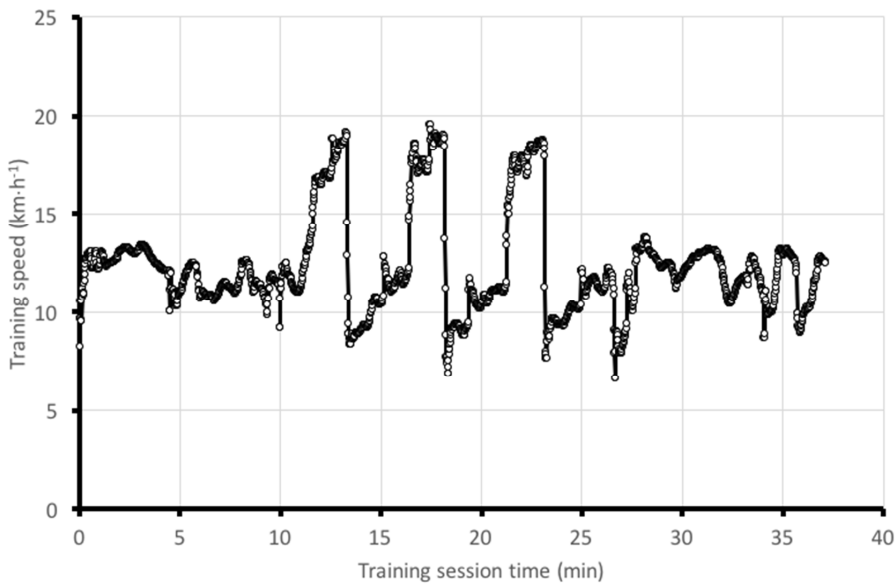


Formatted: Font: (Default) Times New Roman
Formatted: Normal
Formatted: Level 1

438

439 Figure 12

440

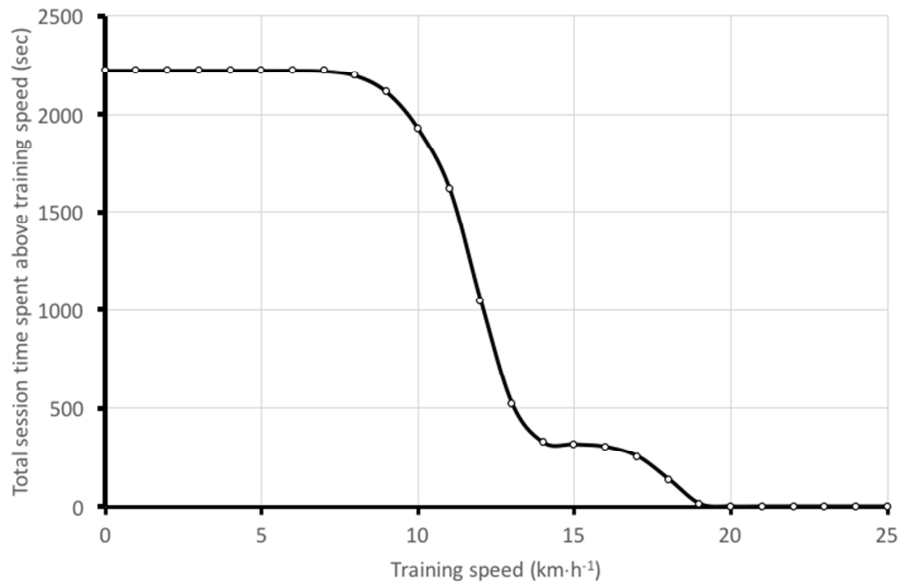


441

442

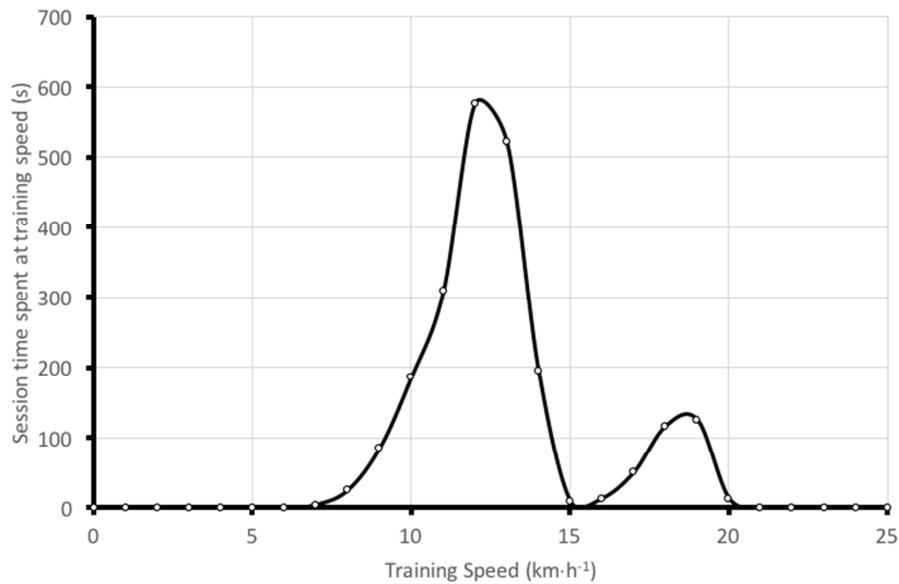
443
444
445
446

Figure 23



447
448

Figure 34



449
450