

**Research Bank**

Journal article

**Inclusion of features derived from a mixture of time window sizes improved classification accuracy of machine learning algorithms for sheep grazing behaviours**

**Hu, Shuwen, Ingham, Aaron, Schmoelzl, Sabine, McNally, Jody, Little, Bryce, Smith, Daniel, Bishop-Hurley, Greg, Wang, You-Gan and Li, Yutao**

This is the accepted manuscript version. For the publisher's version please see:

Hu, S., Ingham, A., Schmoelzl, S., McNally, J., Little, B., Smith, D., Bishop-Hurley, G., Wang, Y.-G. and Li, Y. (2020). Inclusion of features derived from a mixture of time window sizes improved classification accuracy of machine learning algorithms for sheep grazing behaviours. *Computers and Electronics in Agriculture*, 179, Article 105857. <https://doi.org/10.1016/j.compag.2020.105857>

This work © 2020 is licensed under [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International](https://creativecommons.org/licenses/by-nc-nd/4.0/).

1     **Inclusion of features derived from a mixture of time window sizes improved**  
2     **classification accuracy of machine learning algorithms for sheep grazing**  
3                     **behaviours**

4     Shuwen Hu<sup>1</sup>, Aaron Ingham<sup>2</sup>, Sabine Schmoelzl<sup>3</sup>, Jody McNally<sup>3</sup>, Bryce Little<sup>2</sup>, Daniel  
5             Smith<sup>4</sup>, Greg Bishop-Hurley<sup>2</sup>, You-Gan Wang<sup>1</sup> and Yutao Li<sup>2\*</sup>

6     <sup>1</sup>School of Mathematical Sciences, Queensland University of Technology, Gardens Point,  
7     QLD 4001, Australia

8     <sup>2</sup>CSIRO Agriculture & Food, 306 Carmody Road, St Lucia, QLD 4067, Australia

9     <sup>3</sup>CSIRO Agriculture & Food, New England Highway, Armidale, NSW 2350, Australia

10    <sup>4</sup>CSIRO DATA61, Private Bag 12, Hobart, Tasmania 7001, Australia

11    \*Corresponding author: Yutao Li. CSIRO Agriculture & Food, 306 Carmody Road, St Lucia,  
12    QLD 4067, Australia. E-mail address: [Yutao.Li@csiro.au](mailto:Yutao.Li@csiro.au).

13  
14    **Highlights**

- 15     • Simultaneous inclusion of features derived from mixed time window sizes of sensor  
16         signal data significantly improved the sheep behaviour classification accuracy, in  
17         comparison to those from a single unique time window size.
- 18     • Using features derived from time windows of different lengths provided key  
19         information needed to accurately identify different behaviours that involve multiple  
20         movements of unequal duration.
- 21     • Using Random Forest and a mixed window size approach significantly improved the  
22         ability of identifying the walking behaviour, only accounted for 1% of the ground truth  
23         data.

24 **Abstract:**

25 Inertial motion sensors located on the animal have been used to study the behaviour of  
26 ruminant livestock. The time window size of segmented signal data can significantly affect the  
27 classification accuracy of animal behaviours. To date, there have been no studies evaluating  
28 the impact of a mixture of time window size features on the accuracy of animal behaviour  
29 classification. In this study, data was collected from accelerometers attached to the neck of  
30 17 Merino sheep over a period of two days. We also recorded a ground truth dataset of  
31 behaviour recordings (grazing, ruminating, walking, and standing) over the same time period,  
32 We then investigated the ability of three machine learning approaches, Random Forest (RF),  
33 Support Vector Machine (SVM) and linear discriminant analysis (LDA), to accurately classify  
34 sheep behaviour. Our results clearly show that simultaneous inclusion of features derived from  
35 time windows of mixed sizes, ranging from 2-15 seconds, significantly improved the behaviour  
36 classification accuracy, in comparison to those determined from a single unique time window  
37 size. Of the three ML methods applied here, the Random Forest approach yielded the best  
38 results. Together our results show that including features obtained from mixed window sizes  
39 improved the classification accuracy of sheep behaviours.

40 **Keywords:** Mixture of time window sizes; Features ranking; Classification algorithm; Machine  
41 learning; Sheep behaviour; Accelerometer;

42

43 **1. Introduction**

44

45 Animal behaviour can be used to provide a mechanism for the early detection and  
46 quantitative assessment of animal health status (Martiskainen et al., 2009). Grazing and  
47 ruminating are two important behaviours for ruminants and continuous monitoring of animal  
48 eating behaviour provides vital information about ruminant health, productivity and welfare  
49 (Mansbridge et al., 2018). Traditional methods of animal monitoring are based on direct

50 observation by human operators or assessment of video recording, both are labour-intensive,  
51 time-consuming, and prone to human error (Alvarenga et al., 2016). The rapid development  
52 of sensor technologies provides great opportunities for remotely monitoring animals in a range  
53 of applications (Brown et al., 2013; Schmoelzl et al., 2016). Sensor devices, especially those  
54 using an accelerometer that measures inertial acceleration associated with movement  
55 (usually on three different axes), can give a good insight into individual animal behaviour  
56 patterns (Fogarty et al., 2020).

57 To date, several studies have used sensors to study the behaviour of ruminant animals,  
58 especially in cattle. Greenwood et al. (2017) investigated the possibility of predicting pasture  
59 intake based on behaviour classification. They developed a simple algorithm to predict pasture  
60 intake of individual cattle. Rahman et al. (2018) compared cattle behaviour classification using  
61 the information from sensors located on different parts of the body (collar, halter and ear) and  
62 found that different sensor placement can still achieve good classification accuracy providing  
63 that the feature variation between the training and testing animals is very small. For sheep  
64 behaviour classification, Guo et al. (2018) compared grazing behaviour of sheep on pasture  
65 with different sward surface heights using an inertial measurement unit sensor. They found  
66 that a high accuracy (> 95%) of identifying grazing behaviour from non-grazing behaviour  
67 could be achieved for all epochs (5s, 10s and 15s) with 10s being the best, regardless of  
68 sward surface heights. There are also commercially available monitoring systems that may be  
69 used to capture feeding behaviours for dairy cattle, such as Lely (Bar and Solomon, 2010) and  
70 MooMonitors (Verdon et al., 2018). However, these automatic systems cannot be directly  
71 applied to other species such as sheep as there are likely differences in accelerometer signal  
72 patterns between species. Further, the sensor may be impractical or unsafe for deployment  
73 on sheep due to limitations involving sensor size, shape, weight or method of attachment  
74 (Mansbridge et al., 2018)

75 Recently machine learning algorithms have become very popular and offer great potential  
76 in animal behaviour classification, largely because of their abilities in dealing with high  
77 dimension datasets (such as sensor data) and provide high prediction accuracy for complex

78 phenotypes. For example, Dutta et al. (2015) applied six supervised machine learning  
79 methods, binary tree, LDA, naive Bayes, k nearest neighbour (kNN) and adaptive neuro fuzzy  
80 inference system (ANFIS), to classify five major cattle behaviours (Grazing, Ruminating,  
81 Resting, Walking and other behaviour). They achieved a high accuracy (96%) of classification  
82 by using the bagging ensemble classification with tree learner. For sheep behaviour  
83 classification, Mansbridge et al. (2018) found that RF performed the best when compared with  
84 SVM, kNN and adaptive boosting (AdaBoost) for sheep data collected by  
85 accelerometer/gyroscope sensor attached to the ear and collar. Guo et al. (2018) also reported  
86 that the LDA classifier was the best performer compared to binary tree, naive Bayes, kNN and  
87 ANFIS on classifying grazing activities.

88 Window size for signal segmentation is one of the crucial factors influencing on activity  
89 recognition. Banos et al. (2014) evaluated the impact of different window sizes (0.25 s – 7 s)  
90 on human activity classification with accelerometer data and found that the interval 1-2 s was  
91 the best trade-off between speed and accuracy of recognition. In sheep, Walton et al. (2018)  
92 investigated the effects of sensor position (ear and collar), sampling frequency (8, 16 and 32  
93 Hz) of triaxial accelerometer and gyroscope sensor, and window size (3, 5 and 7 s) on  
94 behaviour classification, and concluded that the combination of 16Hz with 7 s window would  
95 produce the benefits of energy efficient and reasonable classification accuracy (91-93%) in a  
96 real-time sheep monitoring system. Smith et al. (2016) built a separate classifier for each of  
97 five cattle behaviours (grazing, walking, ruminating, resting and “other”) using a “one vs all”  
98 ensemble on 24 Holstein-Friesian dairy cows. Of nine window sizes evaluated (1.5, 2.5, 5, 7.5,  
99 10, 15, 20, 25, 30 s), they found that “the grazing, resting and rumination behaviours produced  
100 their highest mean F-score for the longest window of the study (30 s)”.

101 Since there is no consensus about which time window size for signal segmentation and  
102 machine learning methods should be used, one obvious question is whether the features  
103 derived from a mixture of different time window sizes can be used to improve classification  
104 accuracy of livestock grazing behaviour. In this study, we aimed to: (1) Determine if animal  
105 behaviour classification could be improved by the simultaneous inclusion of features

106 calculated from time windows of different sizes; (2) Compare the performance of three machine  
107 learning methods, RF, SVM and LDA in behaviour classification; (3) Investigate if new features  
108 from cumulative effects can improve the classification performance.

109

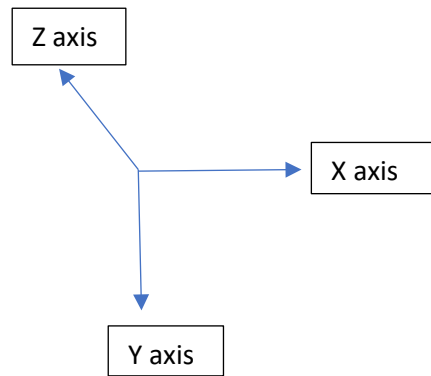
## 110 **2. Materials and methods**

111

### 112 *2.1. Experiment design and data collection*

113

114 Data collection was conducted on animals enrolled in a grazing trial, according to the  
115 Australian Code for the Use and Care of Animals in Research and Teaching, and approved  
116 protocols were approved by the CSIRO Armidale Animal Ethics Committee (Animal Research  
117 Authority 18-13). A total of 20 Merino ewes, habituated to human presence, were kept in a  
118 square mixed sward pasture paddock of 70 m x 70 m. A subgroup of 17 animals was randomly  
119 chosen for device deployment and behavioural annotation. Devices were attached around the  
120 neck of the animals with an elasticated strap (Fig 1) for a period of 48 hours. The sensor  
121 datasets were collected from Actigraph wGT3X-BT fitted with collars around the neck of  
122 Merino ewes, each containing a triaxial microelectromechanical systems (MEMS)  
123 accelerometer. The accelerometer sampled at a frequency of 30 Hz. The X-axis aligned  
124 approximately with the vertical or dorsoventral direction, the Y-axis with the craniocaudal  
125 direction, and the Z axis with the transverse or mediolateral direction (see Fig. 1 for illustration).



126

127

Fig. 1. Location of sensor and its orientation on sheep.

128

129

130

131

132

133

134

135

136

137

Annotation of behaviours was performed by direct observation of animals within the paddock environment. Trained operators were equipped with tablet devices with a custom-designed annotation application (CSIRO AnnoLOG v 1.0.23, (Little, 2018)) installed on a Samsung Galaxy Tab A 7.0" (Samsung, Seoul, Korea). The application allowed users to record time stamped behaviours during the recording period (Fig. 2a), and the output of the behaviour log was presented in tabular form (see Fig. 2b). Four behaviours (Grazing, Ruminating, Walking and Standing) were recorded. Recording time differed between animals and included approximately 30 minutes of annotated behaviour information for each sheep.



b)

Clock date	Clock start time	Animal ID	Behaviour	Clock end time
21/02/2019	8:43:54 AM	1	Graze	8:49:10 AM
21/02/2019	8:49:10 AM	1	Walk	8:49:16 AM
21/02/2019	8:49:16 AM	1	Stand	8:49:47 AM
21/02/2019	8:49:47 AM	1	Walk	8:50:08 AM
21/02/2019	8:50:08 AM	1	Stand	8:51:01 AM
21/02/2019	8:51:01 AM	1	Walk	8:51:05 AM
21/02/2019	8:51:05 AM	1	Graze	8:56:15 AM
21/02/2019	8:56:15 AM	1	Walk	8:56:26 AM
21/02/2019	8:56:26 AM	1	Stand	8:58:16 AM
21/02/2019	8:58:16 AM	1	Walk	8:58:28 AM
21/02/2019	8:58:28 AM	1	Stand	8:59:08 AM
21/02/2019	8:59:08 AM	1	Stand	9:00:26 AM

138

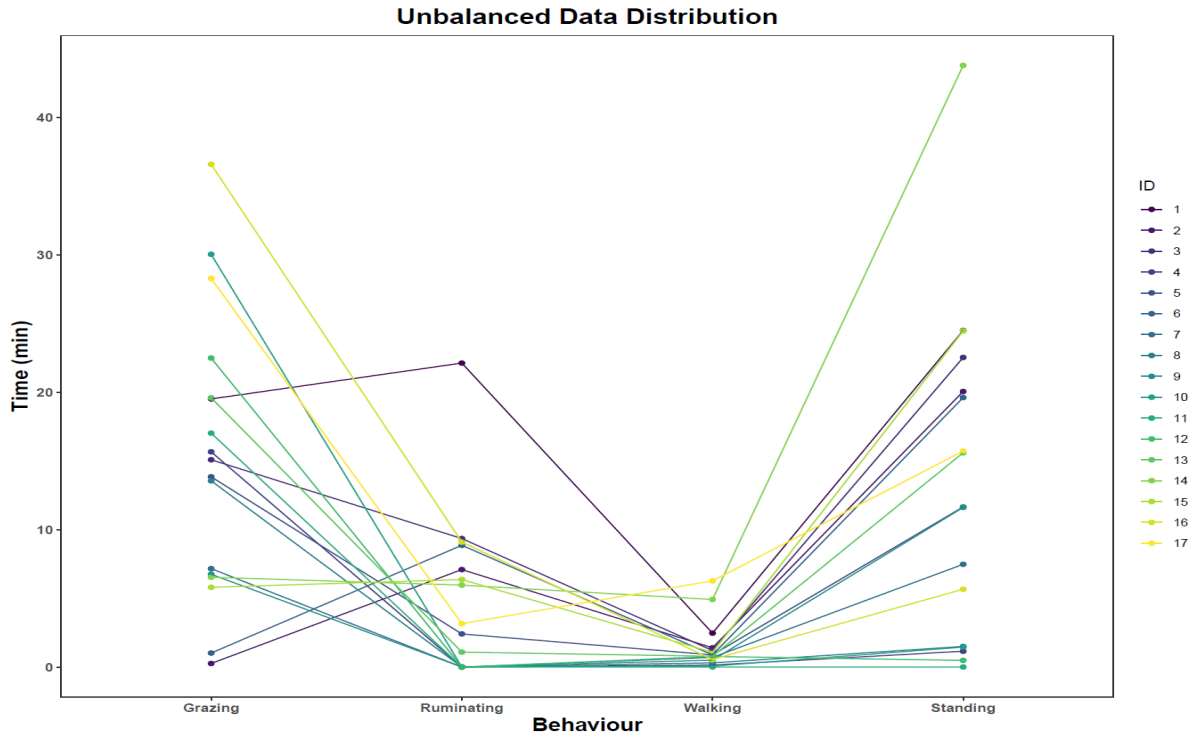
139 Fig. 2. a) Operator interface of the annotation tool CSIRO AnnoLOG v. 1.0.23. b) Tabular  
 140 output of annotated behaviours for animal ID1.

141

142 *2.2 Consolidation of sensor and ground truth datasets*

143 We aligned the raw sensor data and the behaviour observations together via the time  
 144 stamps (i.e. windows, every 1/30 second). A total of 1,052,475 data points was obtained. Fig.  
 145 3 demonstrates the distribution of four different behaviour classes from 17 Merino sheep. Note  
 146 that Walking has a small representation.





147

148

Figure 3. Unbalanced Behaviour Dataset.

149

150 **2.3 Feature extraction from sensor data**

151 In this study, five commonly used time window sizes, 1s, 2s, 5s, 10s and 15s, were applied  
 152 for signal segmentation (Banous et al., 2014; Smith et al., 2016; Walton et al., 2018). For each  
 153 time window, six basic statistical features, minimum, maximum, mean, standard deviation,  
 154 skewness and kurtosis, were computed for each of three axes (X, Y and Z) acceleration data.

155

156 Instead of applying the features from five different time window sizes in isolation, we  
 157 developed a new method that enabled us to conduct the classification analysis with all features  
 158 from different time windows together in the same dataset (see Table 1). In brief, a time window  
 159 of 1s was used as the basis. A total number of 28,425 intervals (bins) were generated with the  
 160 1s window. Using the average of acceleration magnitude values for X-axis ( $a_i$ ) at 1s as an  
 161 example, the corresponding average values for all 1s intervals were denoted as  $a_1, a_2, \dots, a_{28425}$ .  
 162 When using 2s time window, a total number of intervals will be 14,213 with the average values

163 being  $b_1, b_2, \dots, b_{14213}$ . To combine the features of average values from 1s and 2s windows  
 164 together in the same dataset, individual average values of 2s window were used twice to meet  
 165 1s window requirement (see the second column in Table 1). By doing so, the average value  
 166 from a 2s window remained unchanged for the 1<sup>st</sup> and 2<sup>nd</sup> 1s bins. The same concept applies  
 167 to the average values derived for the time window sizes of 5s, 10s and 15s (See Table 1) as  
 168 well as other statistical features (i.e. minimum, maximum, standard deviation, skewness and  
 169 kurtosis values for X-axis acceleration magnitude measurement).

170

171 Table 1 Illustration of deriving average features from mixed time window sizes.

Number of bins	time window				
	1s	2s	5s	10s	15s
1	$a_1$	$b_1$	$c_1$	$d_1$	$e_1$
2	$a_2$	$b_1$	$c_1$	$d_1$	$e_1$
3	$a_3$	$b_2$	$c_1$	$d_1$	$e_1$
4	$a_4$	$b_2$	$c_1$	$d_1$	$e_1$
5	$a_5$	$b_3$	$c_1$	$d_1$	$e_1$
6	$a_6$	$b_3$	$c_2$	$d_1$	$e_1$
...	...	...	...	...	...
...	...	...	...	...	...
28423	$a_{28423}$	$b_{14212}$	$c_{5685}$	$d_{2843}$	$e_{1895}$
28424	$a_{28424}$	$b_{14212}$	$c_{5685}$	$d_{2843}$	$e_{1895}$
28425	$a_{28425}$	$b_{14213}$	$c_{5685}$	$d_{2843}$	$e_{1895}$

172 a, b, c, d and e are the average values of acceleration magnitude values from X axis.

173

174 Apart from six basic statistics features, we also explored new features of cumulative effects  
 175 of raw data  $X_t, Y_t$  and  $Z_t$ . Table 2 illustrates how the cumulative effects are derived for the  
 176 corresponding features named Xsum, Xvelocity, Xsummean, Xsum2, Xdis and Xsum2mean,

177 using the time series data from X-axis acceleration magnitude measurements. The same  
 178 methods were also applied for the computation of the features for the Y and Z axes  
 179 acceleration magnitude measurements. In addition, the squared acceleration magnitude (acc)  
 180 (Rahman et al. 2018) that considers the joint effects of X, Y and Z measurements were also  
 181 included in this study. As the Y-axis detected the motion from front to back and the Z-axis  
 182 detected the motion from the side to side, the interaction between Y and Z axis measurements  
 183 were further examined using the features named dyz and interyz. For each window size and  
 184 each statistics feature, there were 24 different metrics (including X, Y and Z axis data). A total  
 185 of 720 features were generated.

186

187 Table 2: Illustration of calculation of additional features of cumulative effects of X-axis  
 188 Acceleration magnitude measurement.  $T$ : the total number of intervals for a given time  
 189 window;  $t$ : a particular interval.

Time	X	Xsum	Xvelocity	Xsummean	Xsum2	Xdis
1	$x_1$	$x_1$	$x_1/30$	$x_1/1$	$x_1$	$x_1/30$
2	$x_2$	$x_1 + x_2$	$(x_1 + x_2)/30$	$(x_1 + x_2)/2$	$2x_1 + x_2$	$(2x_1 + x_2)/30$
...	...	...	...	...	...	...
...	...	...	...	...	...	...
T	$x_T$	$(\sum_{t=1}^T x_t)$	$(\sum_{t=1}^T x_t)/30$	$(\sum_{t=1}^T x_t)/T$	$Tx_1 + (T - 1)x_2 + \dots$ $+ x_T$	$(Tx_1 + (T - 1)x_2$ $+ \dots + x_T)/30$

190

191 Table Cont'd

Time	X	Xsum2mean	acc	dyz	interyz
1	$x_1$	$x_1/1$	$\sqrt{x_1^2 + y_1^2 + z_1^2}$	$\sqrt{y_1^2 + z_1^2}$	$y_1z_1$
2	$x_2$	$(2x_1 + x_2)/2$	$\sqrt{x_2^2 + y_2^2 + z_2^2}$	$\sqrt{y_2^2 + z_2^2}$	$y_2z_2$

---

...	...	...	...	...	...
...	...	...	...	...	...
T	$x_T$	$(Tx_1 + (T - 1)x_2 + \dots$ $+ x_T)/T$	$\sqrt{x_T^2 + y_T^2 + z_T^2}$	$\sqrt{y_T^2 + z_T^2}$	$y_T z_T$

---

192

193 **2.4 Machine learning (ML) algorithms for classification**

194

195 **2.4.1 RF**

196

197 RF is a tree-based ensemble method that builds a large collection of decision trees using  
 198 training datasets, and validates predictions using testing datasets (Breiman, 2001). The library  
 199 ranger in R (Wright and Ziegler, 2017) was applied for determining hyperparameters in RF.  
 200 The final parameters were: mtry = 27, Ntree = 12 and default values for all other data.

201

202 **2.4.2 SVM**

203 SVM constructs a linear partition of the high-dimensional space into two subspaces for  
 204 classification or regression (James et al., 2013). Intuitively, as the larger margins tend to  
 205 provide lower classification errors, a good separation is obtained by a hyperplane that has the  
 206 largest distance to the nearest training data. In this study, both linear and radial kernel  
 207 functions were applied. The caret function in R (Kuhn, 2008) was applied. The final parameters  
 208 for the analysis were: the average cost and sigma being 115 and 0.0014 for a radial function,  
 209 and the average cost of 1 for a linear function.

210

211 **2.4.3 LDA**

212

213 LDA is a discriminant analysis that can separate a dataset into two or more classes. LDA  
 214 assumes that the data within each class are drawn from a multivariate Gaussian distribution

215 with a class-specific mean vector and a covariance matrix that is common to all classes  
216 (James et al., 2013). In this study, we used the LDA classifier as a benchmark to compare the  
217 classification performance of RF and SVM.

218

219 For all three methods, a five-fold stratified cross-validation scheme was applied. That is,  
220 the dataset from 17 sheep was randomly partitioned into 5 subsets. Each subset was in turn  
221 used as a test dataset while the other 4 subsets were used as the training dataset. .

222

### 223 *2.5 Performance of the classification*

224

225 Similar to the approach previously used for cattle behaviour classification (Rahman et al.,  
226 2018), we chose four metrics, namely overall accuracy, precision, recall (also called sensitivity)  
227 and F1-score, to assess the classification performance of individual ML algorithms. For each  
228 target behaviour (e.g. Grazing, Ruminating, Walking, or Standing), a binary classifier was  
229 defined as the target behaviour (e.g. Grazing class) against a combined class of all remaining  
230 behaviour classes (Non-Grazing class). The calculation of four metrics can be found in  
231 Rahman et al. (2018).

232

233

## 234 **3. Results**

235

### 236 *3.1 Behaviour classification performance using individual unique time window sizes*

237

238 Table 3 presents the effects of different window sizes on the classification performance of  
239 three ML methods, RF, SVM (with linear kernel), SVM (with radial kernel) and LDA classifiers,  
240 when ignoring cumulative effects, squared acceleration magnitude and interaction effects  
241 between Y and Z axes. When the window size increased from 2s, 5s, 10s to 15s, the

242 classification performance of ML classifiers for grazing behaviour showed a continued  
 243 improvement except in SVM (radial kernel). The 15 second window size gave the best  
 244 classification performance. Among all methods, RF showed the highest F1-score (0.876 –  
 245 0.889). However, when considering ruminating, increasing the window size resulted in a  
 246 reduction in performance in both RF and SVM (radial kernel) with 2s being the best window  
 247 size. SVM (linear kernel) and LDA had no ability or a weak power to identify the ruminating  
 248 behaviour (Table 3). In all cases, none of the ML methods had the ability to recognize walking  
 249 behaviour regardless of the time window applied, except a very weak power identified by RF  
 250 (0.173) at 2s window.

251

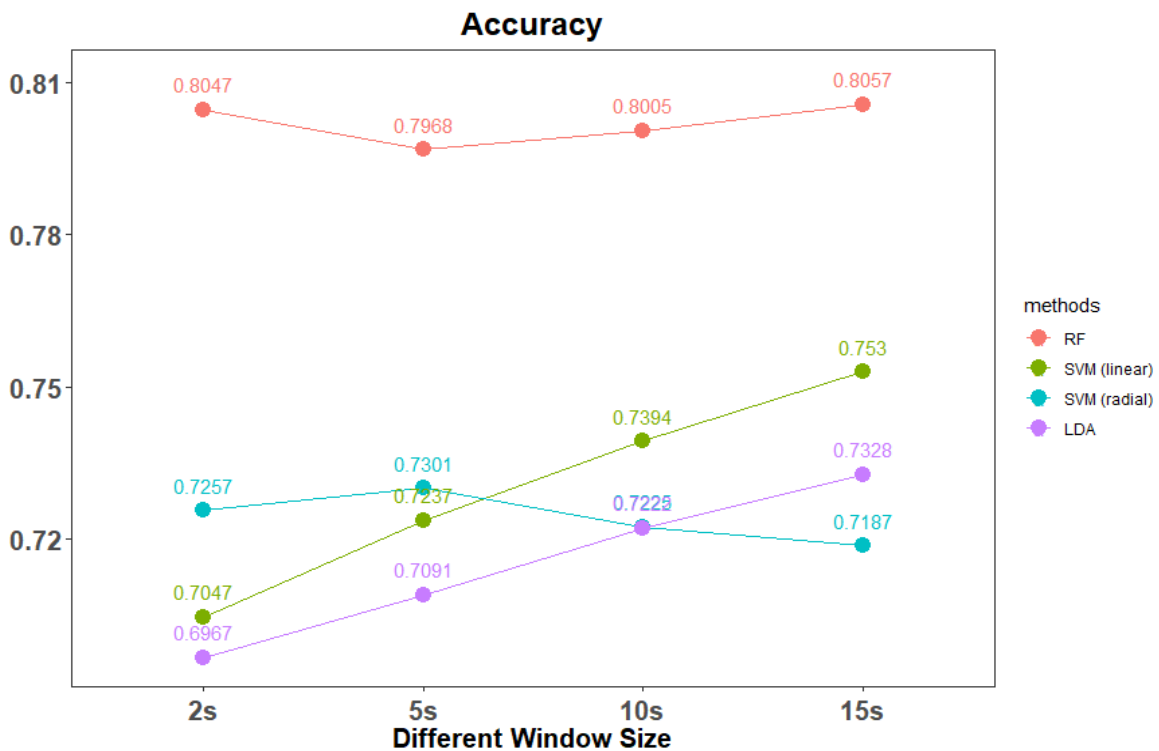
252 Table 3. Effects of individual time window sizes on the behaviour recognition performance  
 253 F1-score, when using RF, SVM and LDA. NA – not available.

		F1 score			
Method	Behaviour	2 sec	5 sec	10 sec	15 sec
RF	Grazing	0.876	0.879	0.886	0.889
	Ruminating	0.655	0.582	0.545	0.550
	Walking	0.173	NA	NA	NA
	Standing	0.794	0.782	0.783	0.781
SVM (linear kernel)	Grazing	0.832	0.850	0.864	0.875
	Ruminating	NA	NA	NA	NA
	Walking	NA	NA	NA	NA
	Standing	0.706	0.721	0.731	0.741
SVM (radial kernel)	Grazing	0.839	0.844	0.842	0.846
	Ruminating	0.286	0.289	0.254	0.224
	Walking	NA	NA	NA	NA

	Standing	0.720	0.719	0.705	0.692
LDA	Grazing	0.823	0.838	0.851	0.863
	Ruminating	0.129	0.057	NA	NA
	Walking	NA	NA	NA	NA
	Standing	0.697	0.709	0.719	0.724

254

255 When comparing the overall accuracies of behaviour classification for each of the three  
 256 ML methods (Fig. 4), the RF performed best of all classifiers regardless of time window  
 257 size. However, the window size did impact on the performance of SVM and LDA classifiers. The  
 258 smaller window size ( $\leq 5s$ ), SVM (radial kernel) performed better than SVM (linear kernel) and  
 259 LDA. For both SVM (linear kernel) and LDA, increasing the window size improved the  
 260 accuracy value, with 15s giving the highest value (Fig. 4). In general, SVM (linear kernel)  
 261 produced higher accuracy than LDA, and it outperformed SVM (radial kernel) when the  
 262 window size was  $\geq 10s$ .



263

264 Fig. 4. The change of overall accuracy for individual ML methods with different time window

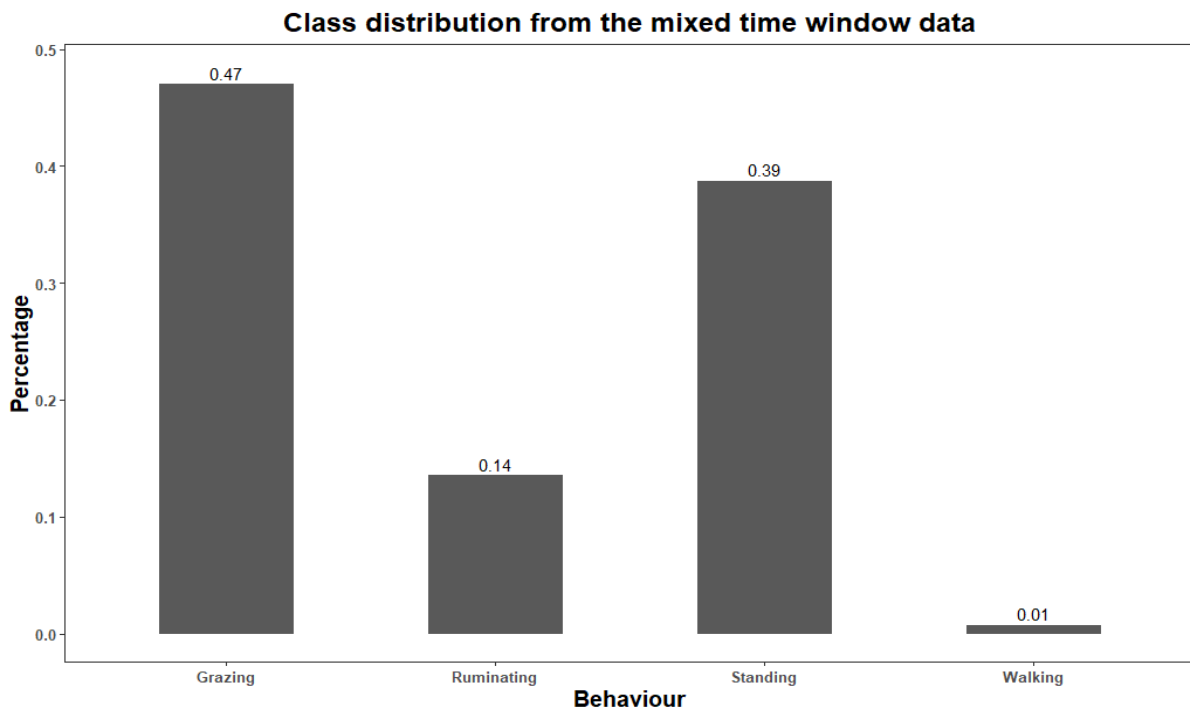
265

sizes.

266

267 **3.2 Classification of behaviours using mixed time window sizes**

268 Fig. 5 illustrates the composition of different behaviours in the newly formed mixed time  
269 window dataset. Grazing behaviour accounted for the largest percentage (47.0%) of the data  
270 and the walking behaviour had the least representation at 1% of data.



271

272 Fig. 5. The percentage distribution of four behaviours in the mixed time window data.

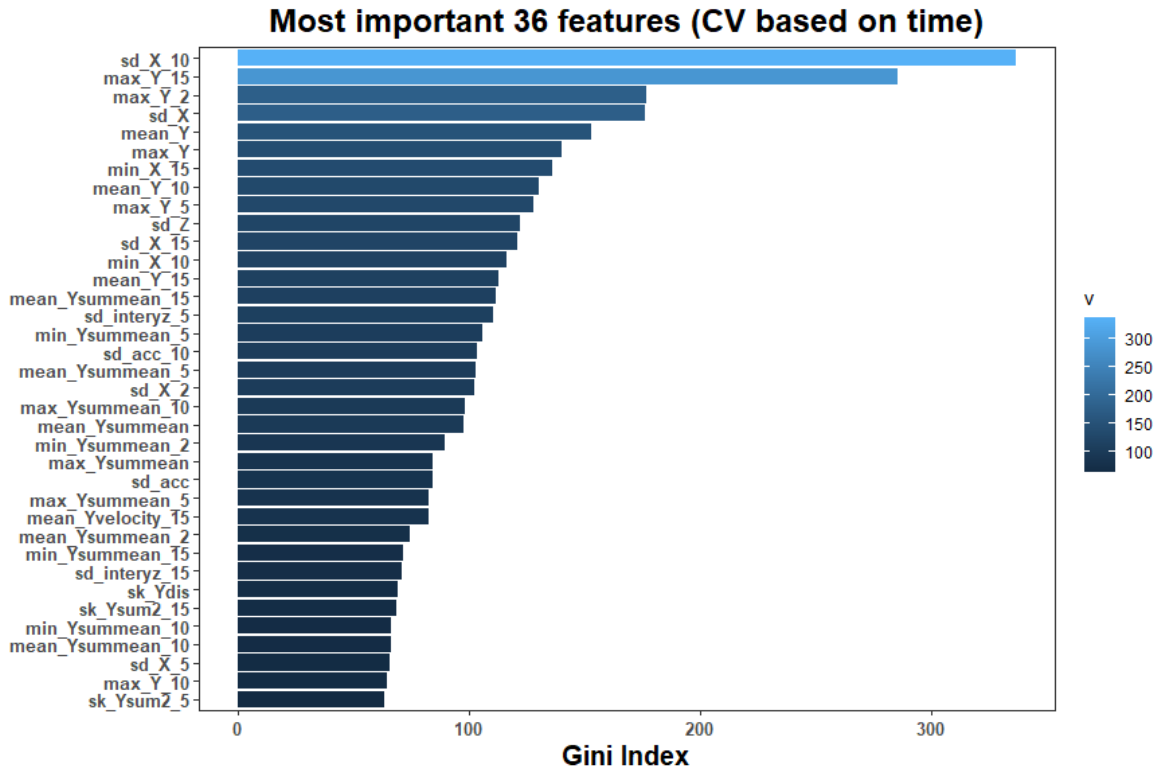
273

274 **3.2.1 RF**

275 Of the 720 features examined with the mixed window size dataset, the top 36 features  
276 (with the highest Gini index values) identified by the RF are shown in Fig. 6. Among the top  
277 36 features, the features derived from different window sizes (1s, 2s, 5s, 10s and 15s) all  
278 contributed to the classification accuracy of different behaviours. One other noteworthy



279 observation was that the features derived from using the cumulative effects (i.e. with “sum” in  
 280 the labels) were also among the top contributing features.  
 281

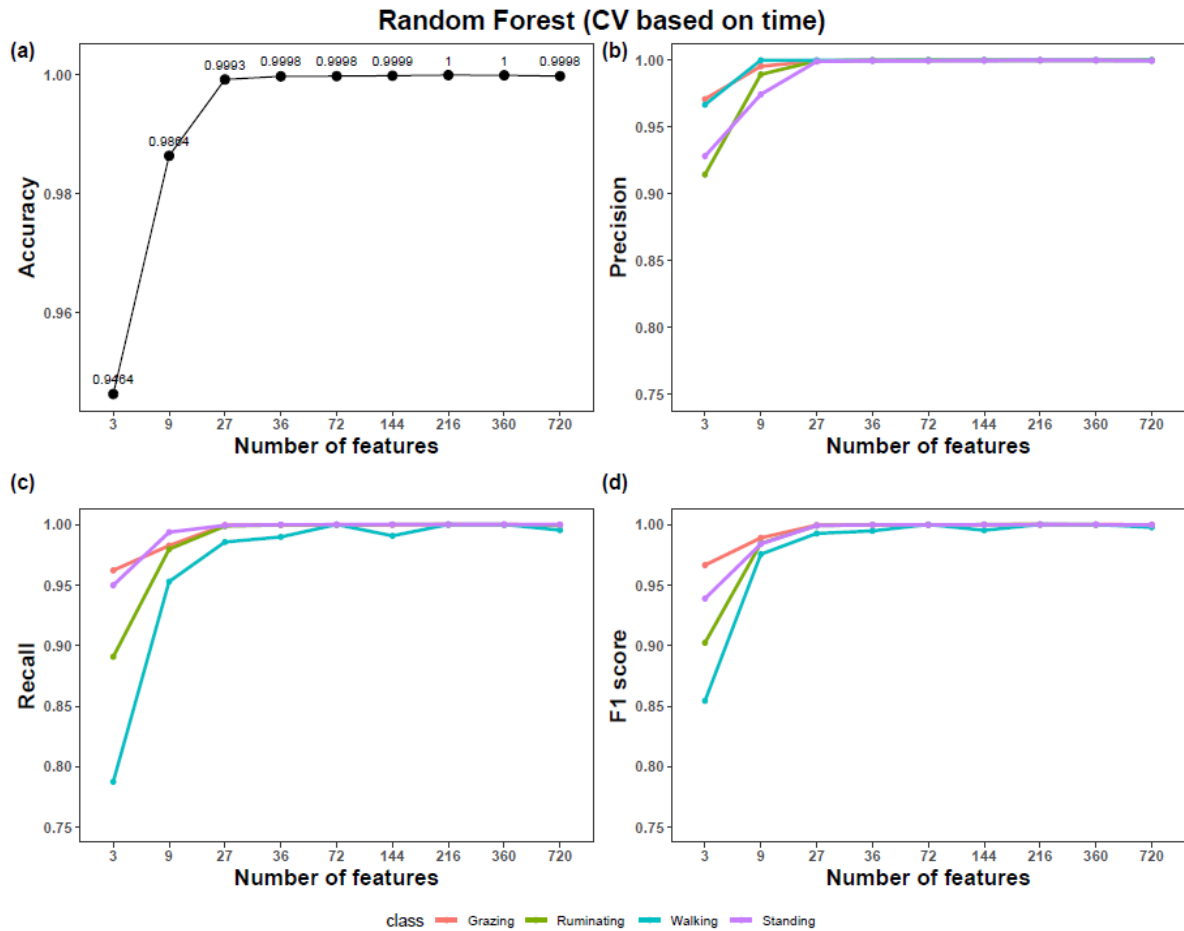


282

283 Fig. 6. The list of top 36 features selected by RF.

284 Next, we compared the behaviour classification performance of different subsets of the  
 285 top features from RF with that of 720 features. The top 3 features produced an overall accuracy  
 286 of 0.946 (Fig. 7(a)). The accuracy value increased to 0.986 with the top 9 features and to  
 287 almost 1 (0.999) with the top 27 features. Similar trends were observed irrespective of the  
 288 performance metric used (See Figs 7b, 7c and 7d). For all four behaviour classes, the RF  
 289 identified the individual behaviour classes with high precision (>0.970), sensitivity (Recall >  
 290 0.950) and F1-score (>0.970) when either the top 9 or 27 features were applied (Fig. 7). Even

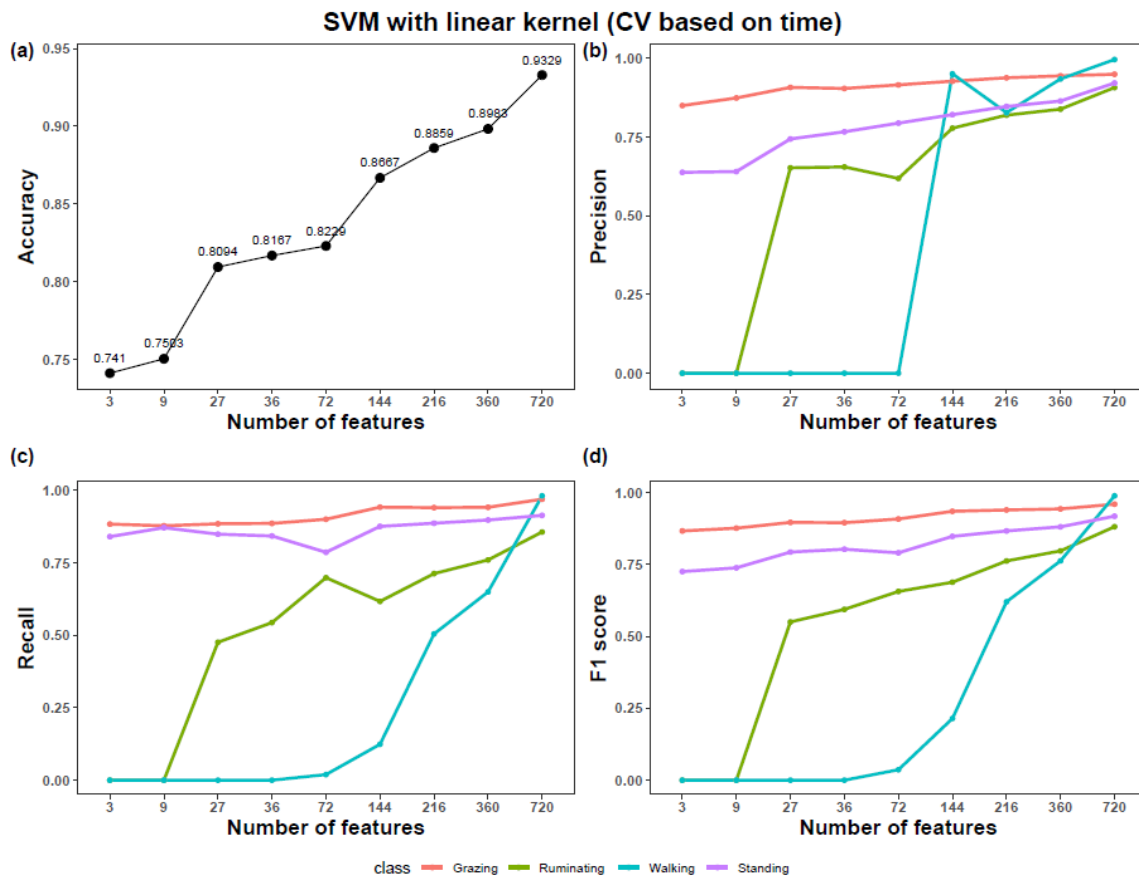
291 with the top three features, the results show that for the most difficult behaviour to classify –  
 292 walking, there was 0.910 (precision), 0.780 (sensitivity) or 0.850 (F1-score) achieved.



293  
 294 Fig. 7. The overall accuracy, precision, recall and F1 score values from using the different  
 295 number of top ranked features chosen from RF for behaviour classification.

296  
 297 **3.2.2 SVM**

298 The SVM classifier was evaluated with both linear and radial kernels. The overall accuracy  
 299 of SVM (linear kernel) was largely dependent on the number of features applied for  
 300 classification (Fig. 8a). While the lowest accuracy was 0.741 with three top features, the  
 301 highest accuracy (0.933) was achieved with all the 720 features.

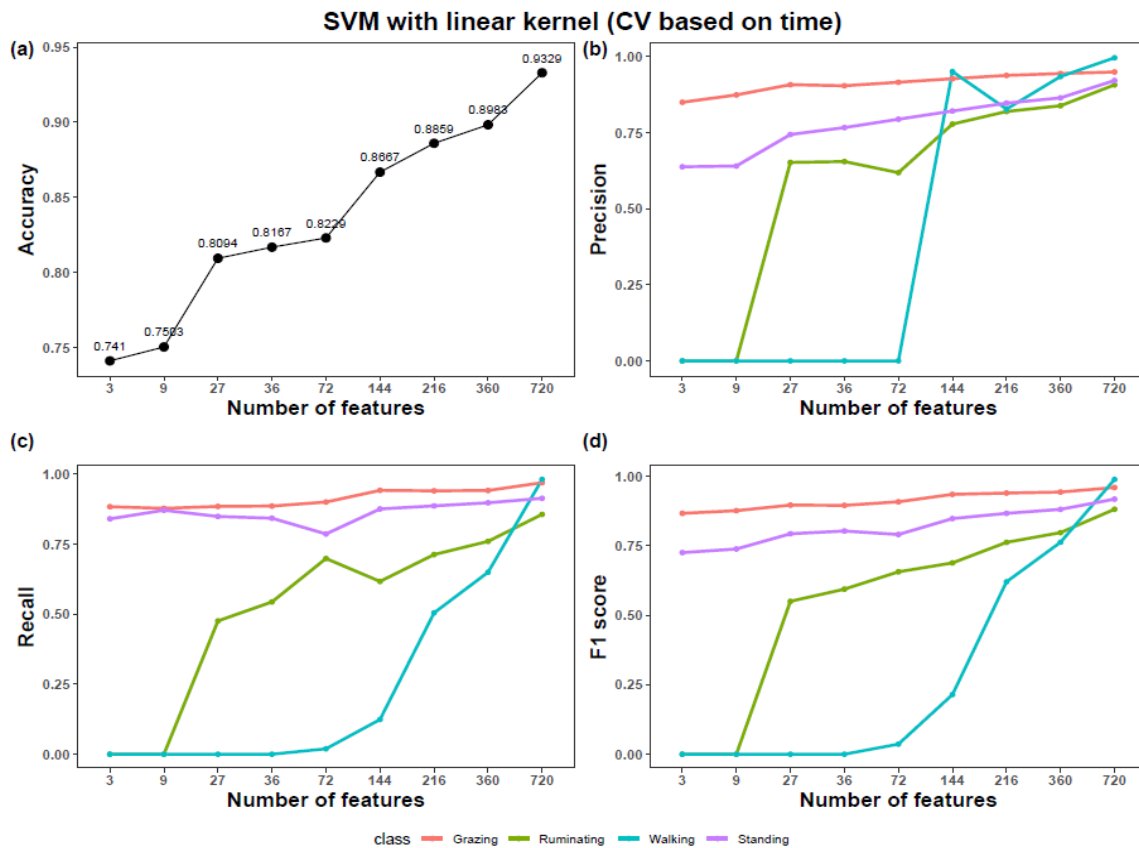


302

303 Fig. 8. The over accuracy, precision, recall and F1 score values from the different number of  
 304 top ranked features chosen from SVM with a linear kernel function and the mixed time  
 305 window sizes.

306 When investigating the classification performance of SVM (linear kernel) for individual  
 307 behaviour classes Fig. 8), increasing the number of features slightly improved the  
 308 classification performance for grazing and standing behaviours(Figs. 8b-8d) . However, for  
 309 walking and ruminating, the change in the number of features had a significant impact on the  
 310 classification performance (Fig. 8d). For example, SVM (linear kernel) had no or little power  
 311 to correctly classify the walking behaviour until the number of features reached more than 144  
 312 (Fig. 8d, F1 score = 0.210).

313



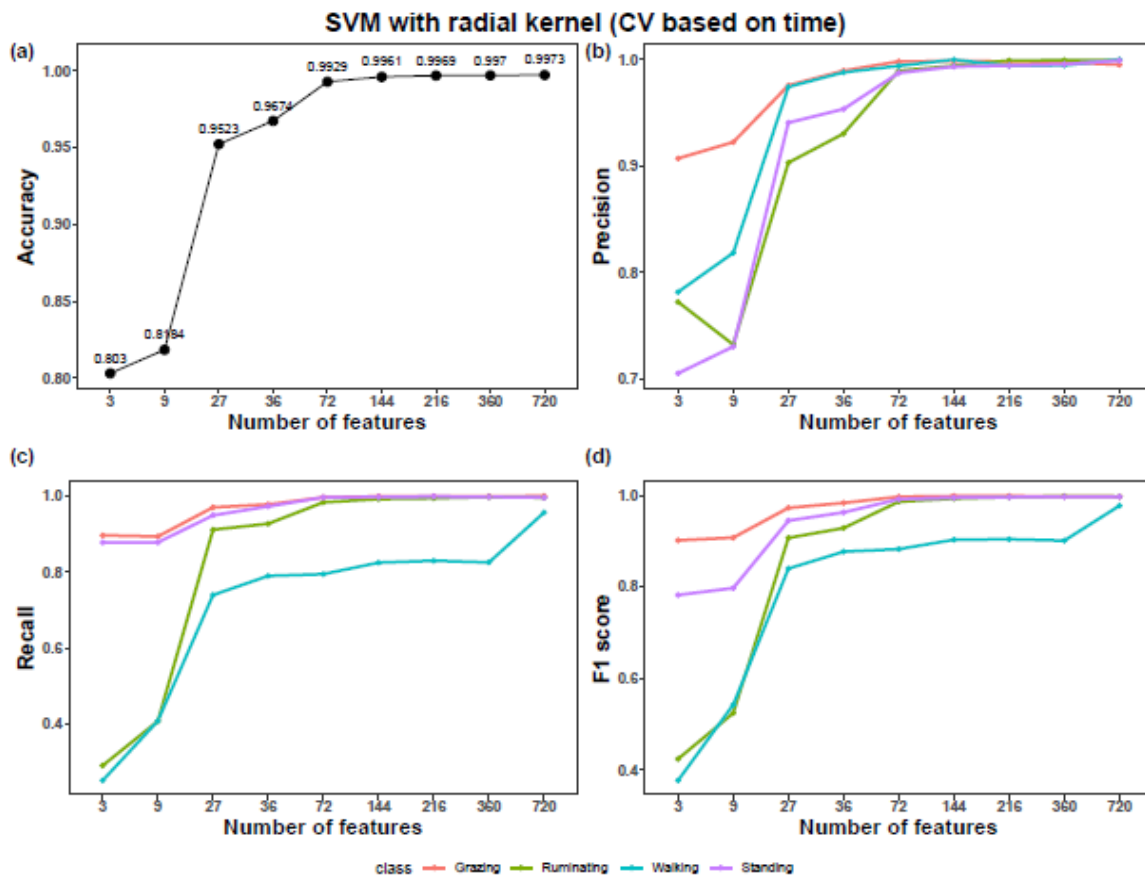
314

315 Fig. 8. The overall accuracy, precision, recall and F1 score values from the different number  
 316 of top ranked features chosen from SVM with a linear kernel function and the mixed time  
 317 window sizes.

318

319 When comparing SVM (radial kernel) with SVM (linear kernel), the overall accuracy was  
 320 significantly improved by 8.40% with top three features, and by 17.70% with 27 features (Fig.  
 321 9a vs Fig. 8a). SVM (radial kernel) performed well in classification of all individual behaviours  
 322 when the number of features was more than 27. This can be demonstrated by high precision  
 323 (> 0.900), reasonable sensitivity (> 0.750) and medium to high F1 score (> 0.850) in Fig 9.  
 324 The most noticeable results are for the walking behaviour.

325



327

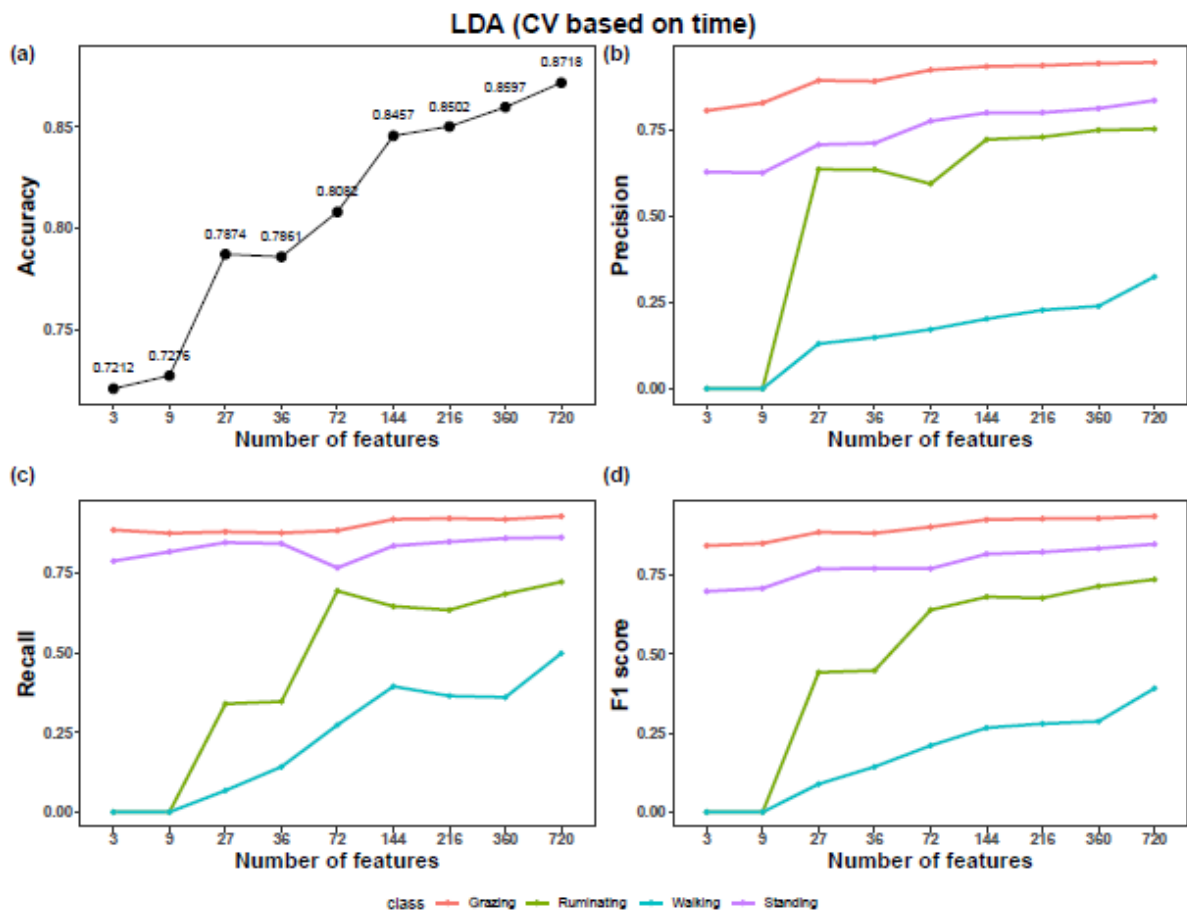
328 Fig. 9. The overall accuracy, precision, recall and F1 score values from the different number  
 329 of top ranked features chosen from SVM with a radial kernel function and the mixed time  
 330 window sizes

331

332 **3.2.3 LDA**

333 Fig. 10 presents the classification results from the LDA classifier. The overall accuracy  
 334 from the LDA classifier followed a similar trend as the SVM classifier with linear kernel (Figure  
 335 10(a) vs Figure 8(a)). However, the SVM (linear kernel) still gave 2.750% (with 3 features) -  
 336 7.010% (with 720 features) better accuracy values than the LDA classifier. When comparing  
 337 the classification performance for individual behaviour, again, LDA had a very similar  
 338 performance to the SVM with linear kernel for the grazing, ruminating, and standing

339 behaviours. For the walking behaviour, although the LDA classifier had an overall low  
 340 precision (<0.300), low recall (<0.500) and low F1 score (<0.375), surprisingly, it did show  
 341 some ability to recognise the walking behaviour when the number of features were less than  
 342 72. This was a stark contrast to the SVM classifier with linear kernel for the walking behaviour.



343  
 344 Fig. 10. The overall accuracy, precision, recall and F1 score values from the different  
 345 number of top ranked features chosen from LDA with the mixed time window sizes.  
 346

### 347 4. Discussion

348  
 349 Accurate classification of animal behaviour from sensor derived data is influenced by a  
 350 number of factors, including experimental design, sensor placement position, data sample rate,  
 351 signal segmentation window size, feature selection and different machine learning methods  
 352 applied (Banos et al., 2014; Rahman et al., 2018; Walton et al., 2018; Mansbridge et al., 2018;

353 Fogarty et al., 2020). Among them, window size, feature selection and analytical methods play  
354 critical roles.

355

#### 356 *4.1. Window size*

357

358 To our knowledge, this is the first study that evaluated the classification of animal  
359 behaviour using features derived from a mixed window size approach. It is a common practice  
360 to spend extensive amount of time to evaluate a range of time window sizes in order to identify  
361 an optimal size that produces reasonable classification accuracy from sensor datasets. In this  
362 study, developing a new combined window size approach enabled us not only to minimise  
363 subjective selection of window sizes, but also to systematically and simultaneously consider  
364 the features from multiple window sizes together to capture the irregular duration of animal  
365 behaviours, that occur both within and across behaviours.

366 Most importantly, we have demonstrated that using mixed window sizes in combination  
367 with the ML method RF, significantly improved the classification accuracies for all behaviour  
368 classes, especially for walking and ruminating. The overall accuracy, precision, recall and F1  
369 score for individual behaviour, when mixed window sizes were applied, were in strong contrast  
370 to those when individual time window size was applied. This is largely due to the features  
371 derived from different window sizes being inter-related. The correlations between these  
372 features of different window sizes provide additional information for ML methods to correctly  
373 identify individual animal behaviours. The biological basis for this can be explained by the  
374 need to classify specific short duration movements as a component of a longer movement.  
375 For instance, grazing might involve the lowering of the head and several biting events. Only  
376 classifying the lowering of the head would likely lead to inaccurate classification as would the  
377 classification of biting-like behaviours alone. The combination of the two is however likely to  
378 be much more informative. In contrast, when analysing the datasets with the features from  
379 one window size only, many feature correlations were not accounted for, especially in the

380 cases where mixed behaviours occurred in a given time window (i.e. unequal length of animal  
381 behaviour), errors were expected to arise.

#### 382 *4.2. Feature selection*

383 To date an optimum number of features that can be used for classifying animal behaviours  
384 has varied greatly between studies, depending on the nature of sensor and behaviour data  
385 and analytical methods applied. For example, by analyzing the data collected by collar  
386 mounted motion sensors, Guo et al. (2018) found that the top 5 features, mean of  
387 accelerometer Z-axis, entropy of accelerometer Y-axis, entropy of accelerometer Z-axis, mean  
388 of gyroscope X-axis and mean of gyroscope Y-axis, can be used in classifying the grazing  
389 versus non-grazing activities in sheep. Mansbridge et al. (2018) identified 39 being the  
390 optimum number of features that can be used successfully in the classification of eating  
391 behaviors in sheep with a high accuracy (91% for ear and 92% for collar data). These features  
392 ranged from dominant frequency, zero crossings, signal area, spectral entropy, to basic  
393 statistics such as mean, min, max, standard deviation, and kurtosis.

394 In this study for each axis acceleration magnitude measurements, apart from the common  
395 features derived from 6 basic statistics and the squared acceleration magnitude (acc), we also  
396 evaluated the effects of new features of cumulative effects of measurements for a given time  
397 window size on classification performance. The reasons for using these features include: 1)  
398 to properly evaluate the efficiency of the new approach - a mixed time window sizes, it is  
399 crucial to compare the new method with conventional methods using commonly used features.  
400 2) using new features from the accumulative effects was to examine if they could better  
401 capture actual change of motion movements for mixed behaviours. When applying the RF  
402 with the mixed window sizes, of 720 features, we found 9 top ranking features that contributed  
403 the most in the classification accuracy were all related to basic statistics (e.g. SD\_X\_10,  
404 max\_Y\_15 Max\_Y\_2, sd\_X, mean\_Y, max\_Y, min\_X\_15, mean\_Y\_10 and max\_Y\_5, see  
405 Figure 6) of X and Y-axis measurements. The X and Y-axis in this study aligned with upward  
406 and downward, and front to back movements of the neck, respectively hence kinetically related  
407 to grazing and ruminating behaviors. Among the 27 top ranking features, there were also 9



408 features derived from the average of the accumulative effects of Y axis (Ysummean, Fig. 6)  
409 for different window sizes. This indicates that the average cumulative values of Y axis from  
410 mixed window sizes may have better-reflected changes of acceleration, therefore contribute  
411 to additional improvement in the accuracy of behaviour classification.

412 Feature selection can also be impacted by where signal information sensor placement  
413 position **is one of the key factors that impact classification accuracy of sheep behaviour**  
414 (Barwick et al., 2018).

415

416

### 417 *4.3. Machine learning methods*

418

419 Different machine learning methods could yield different outcomes when dealing with  
420 unbalance behaviour classes. Among three machine learning algorithms (RF, SVM and LDA)  
421 evaluated, the RF classifier performed the best with over 99% accuracy, followed by the SVM  
422 (radial kernel), the SVM (linear kernel) and LDA classifiers. RF has been known to produce  
423 good classification accuracies in sheep behaviour (Alvarenga et al., 2016), especially  
424 classification of grazing and rumination behaviour (Walton et al., 2018; Mansbridge et al.,  
425 2018). This is mainly due to its great ability in handling non-linearly correlated data and  
426 robustness to noise (Mansbridge et al., 2018). SVM (radial kernel) performed better than the  
427 SVM (linear kernel) and LDA classifiers, also because its radial kernel function can non-  
428 linearly separate the sensor signals associated with irregular length of individual behaviours.  
429 There are other machine learning methods that can also be applied to provide good  
430 classification accuracy of sheep behaviour, depending on the sources of signal information of  
431 sensor placement (e.g. ear, collar, or leg, Barwick et al., 2018), and trade-off between energy  
432 consumption and classification accuracy (Le Roux et al, 2018). Future work needs to be  
433 carried out with ensemble classifiers in which several different classifiers are trained

434 simultaneously and their classification decisions can be combined at the end (D'Este et al.,  
435 2014; Dutta et al., 2015).

436 All the results presented in this study were obtained using a five-fold stratified cross-  
437 validation scheme based on time, rather than the cross-validation based on individual sheep.  
438 The initial analysis using a five-fold cross-validation approach based on subsets of sheep  
439 produced much worse results than that of a stratified cross-validation approach (see  
440 supplementary results I). The primary reason was due to the small number of animals used in  
441 this study and the big variation between individual sheep behaviours in feature space.  
442 Therefore, it is difficult to obtain the consistent results in test datasets when different subsets  
443 of sheep were used as training datasets. [To minimize the impact of individuality on classification](#)  
444 [accuracy of animal behaviours in future sensor application, it will be crucial to: 1\) obtain results from](#)  
445 [larger numbers of animals; and 2\) explore and validate results using a number of repeated k-fold CV](#)  
446 [to improve the prediction on both population and individual results](#)

447 Limitations of the study that may have influenced the results, include the small number of  
448 animals and limited paddock space (70 m x 70 m) used during the experiment. However, the  
449 study aimed to serve as a proof of concept that incorporation of features calculated across  
450 time windows of different lengths has the potential to improve classification accuracy. We  
451 believe this principle has been demonstrated and the broader applicability of the approach  
452 can be tested in future trials involving larger numbers of animals.

453

## 454 **4. Conclusions**

455

456 This study demonstrated that the sheep behaviours of grazing, ruminating, walking, and  
457 standing can be differentiated with a high accuracy using ML algorithm RF and a mixed  
458 window size approach. One clear benefit of applying the RF, mixed window approach was the  
459 ability to accurately classify walking behaviour, that only accounted for 1% of the ground truth  
460 data, when conventional approaches failed. One possible explanation for this outcome is that

461 behaviour classification requires the information contained in features derived from time  
462 windows of different length to provide the context needed for accurate identification.

463

#### 464 **Acknowledgements**

465

466 We would like to acknowledge the financial support SH received from the Queensland  
467 University of Technology. We thank Dr Peter Hunt for allowing us to deploy devices on  
468 animals enrolled in his experimental work.

469

#### 470 **References**

471

472 Alvarenga, F.A.P., Borges, I., Palkovic, L., Rodina, J., Oddy, V.H., & Dobos, R.C. (2016).

473 Using a three-axis accelerometer to identify and classify sheep behaviour at pasture.

474 *Applied Animal Behaviour Science*, 181, 91–99.

475 Banos, O., Galvez, J.-M., Damas, M., Pomares, H., & Rojas, I. (2014). Window size impact in

476 human activity recognition. *Sensors*, 14, 6474-6499.

477 Bar, D., & Solomon, R. (2010, March 2-5). Rumination collars: What can they tell us. [Paper

478 presentation]. The First North American Conference on Precision Dairy Management,

479 Toronto, ON, Canada.

480 Barwick, J., Lamb, D., Dobos, R., Schneider, D., Welch, M., & Trotter, M. (2018). Predicting

481 Lameness in Sheep Activity Using Tri-Axial Acceleration Signals. *Animals (Basel)*, 8(1),12.

482 Brown, D.D., Kays, R., Wikelski, M., Wilson, R., & Klimley, A.P. (2013). Observing the

483 unwatchable through acceleration logging of animal behavior. *Animal Biotelemetry*, 1, 20.

484 D'Este, C., Timms, G., Turnbull, A., & Rahman, A. (2014). Ensemble aggregation methods for

485 relocating models of rare events. *Engineering Applications of Artificial Intelligence*, 34, 58-

486 65.

487 Dutta, R., Smith, D., Rawnsley, R., Hurley, G.B., Hills, J., Timms, G., & Henry, D. (2015).  
488 Dynamic cattle behavioural classification using supervised ensemble classifiers.  
489 *Computers and Electronics in Agriculture*, 111, 18-28.

490 Fogarty E.S., Swain D.L., Cronin G.M., Moraes L.E., & Trotter M. (2020). Behaviour  
491 classification of extensively grazed sheep using machine learning. *Computers and*  
492 *Electronics in Agriculture*, 169: 105-175.

493 Gonzalez, L.A., Bishop-Hurley, G.J., Handcock, R.N., & Crossman, C. (2015). Behavioural  
494 classification of data from collars containing motion sensors in grazing cattle. *Computers*  
495 *and Electronics in Agriculture*, 110, 91-102.

496 Guo, L., Welch, M., Dobos, R., Kwan, P., & Wang, W. (2018). Comparison of grazing  
497 behaviour of sheep on pasture with different sward surface heights using an inertial  
498 measurement unit sensor. *Computers and Electronics in Agriculture*, 150, 394-401.

499 Greenwood, P.L., Paull, D.R., McNally, J., Kalinowski, T., Ebert, D., Little, B., Smith, D.V.,  
500 Rahman, A., Valencia, P., Ingham, A.B., & Bishop-Hurley, G.J. (2017). Use of sensor-  
501 determined behaviours to develop algorithms for pasture intake by individual grazing cattle.  
502 *Crop and Pasture Science*, 68, 1091-1099.

503 Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data*  
504 *Mining, Inference, and Prediction*. 2nd ed. Springer: Berlin, Germany.

505 Kuhn, M. (2008). Building Predictive Models in R Using the caret Package. *Journal for*  
506 *Statistical Software*, 28(5): 26.

507 Le Roux

508 Little, B. (2018). CSIRO AnnoLOG (1.0.23) [Mobile application software]. Retrieved from  
509 <https://play.google.com/>

510 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical*  
511 *Learning*. Springer: Heidelberg, Germany.

512 Mansbridge, N., Mitsch, J., Bollard, N., Ellis, K., Miguel-Pacheco, G. G., Dottorini, T., & Kaler,  
513 J. (2018). Feature Selection and Comparison of Machine Learning Algorithms in

514 Classification of Grazing and Rumination Behaviour in Sheep. *Sensors*, 18(10), 3532.  
515 doi:10.3390/s18103532

516 Martiskainen, P., Jarvinen, M., Skon, J.-P., Tiirikainen, J., Kolehmainen, & M., Mononen, J.  
517 (2009). Cow behavior pattern recognition using a three-dimensional accelerometer and  
518 support vector machines. *Applied Animal Behaviour Science*, 119, 32–38.

519 Menze, B.H., Kelm, B.M., Masuch, R., Himmelreich, U., Bachert, P., Petrich, W., & Hamprecht,  
520 F.A. (2009). A comparison of random forest and its Gini importance with standard  
521 chemometric methods for the feature selection and classification of spectral data. *BMC*  
522 *Bioinformatics*, 10, 213.

523 Rahman, A., Smith, D.V., Little, B., Ingham, A.B., Greenwood, P.L., & Bishop-Hurley, G.J.  
524 (2018). Cattle behaviour classification from collar, halter, and ear tag sensors. *Information*  
525 *Processing in Agriculture*, 5, 124-133.

526 Smith, D., Rahman, A., Bishop-Hurley, G.J., Hills, J., Shahriar, S., Henry, D., Rawnsley, R.  
527 (2016). Behavior classification of cows fitted with motion collars: decomposing the multi-  
528 class classification problem into a set of binary problems. *Computers and Electronics in*  
529 *Agriculture*, 131, 40-50.

530 Verdon, M., Rawnsley, R., Raedts, P., & Freeman, M. (2018). The Behaviour and Productivity  
531 of Mid-Lactation Dairy Cows Provided Daily Pasture Allowance over 2 or 7 Intensively  
532 Grazed Strips. *Animals*, 8, 115.

533 Walton, E., Casey, C., Mitsch, J., Vazquez-Diosdado, J.A., Yan, J., Dottorini, T., Ellis, K.A.,  
534 Winterlich, A., & Kaler, J. (2018). Evaluation of sampling frequency, window size and  
535 sensor position for classification of sheep behaviour. *Royal Society Open Science*, 5,  
536 171442. <http://dx.doi.org/10.1098/rsos.171442>.

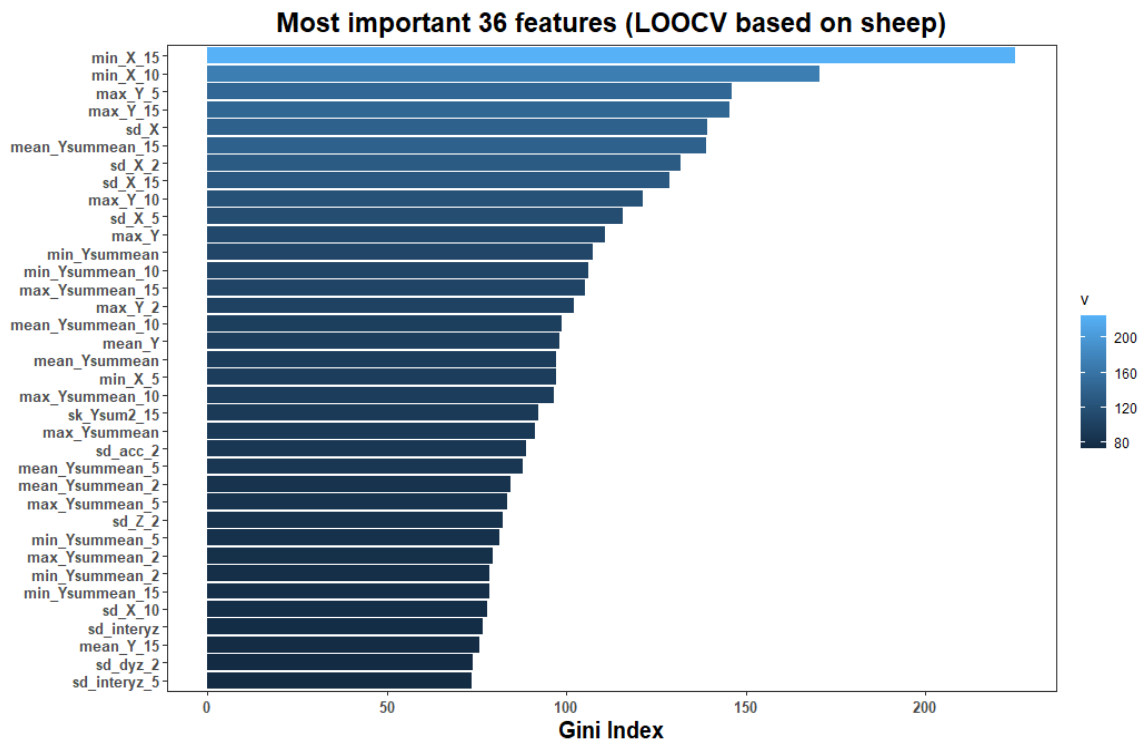
537 Watanabe, N., Sakanoue, S., Kawamura, K., & Kozakai, T. (2008). Development of  
538 an automatic classification system for eating, ruminating and resting behavior of cattle  
539 using an accelerometer. *Japanese Society of Grassland Science*, 54, 231–237.

540 Wright, M. N. & Ziegler, A. (2017). ranger: A Fast Implementation of Random Forests for High  
541 Dimensional Data in C++ and R. *Journal for Statistical Software*, 77(1): 17.

542 Yoshitoshi, R., Watanabe, N., Kawamura, K., Sakanoue, S., Mizoguchi, R., Lee, H.-J., &  
543 Kurokawa, Y. (2013). Distinguishing cattle foraging activities using an accelerometry-  
544 based activity monitor. *Rangeland Ecology & Management*, 66 (3), 382–386.

545

546

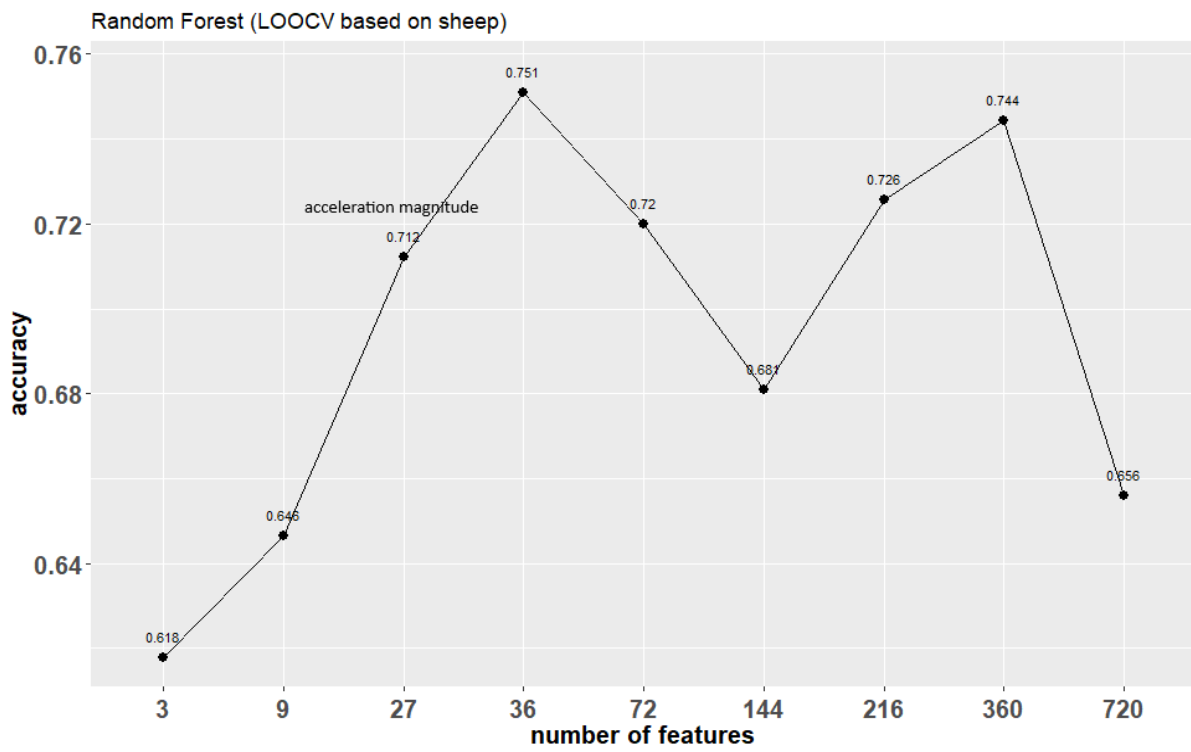


548

549 Fig. 1s. The list of 36 of the most important features selected by RF based on the leaving  
 550 one animal out cross-validation scheme.

551

552



553

554 Fig. 2s. The over accuracy values from the different number of top ranked features chosen  
555 from Random Forest (RF) when using the mixed time window approach and leaving one  
556 animal out cross-validation scheme.