

2 **The pragmatic classification of upper extremity motion in**
3 **neurological patients: a primer**

4 **Avinash Parnandi^{1,*}, Jasim Uddin², Dawn M. Nilsen², and Heidi Schambra¹**

5 ¹ NYU School of Medicine; (avinash.parnandi and heidi.schambra)¹@nyulangone.org

6 ² Columbia University Medical Center; (ju2189 and dmn12)²@cumc.columbia.edu

7 * Correspondence: Avinash Parnandi, avinash.parnandi@nyulangone.org

8 Received: date; Accepted: date; Published: date

9 **Keywords:** Machine learning algorithms; wearable sensors; inertial measurement unit;
10 accelerometers, functional primitives; stroke rehabilitation

11 **Abstract**

12 Recent advances in wearable sensor technology and machine learning (ML) have allowed for the
13 seamless and objective study of human motion in clinical applications, including Parkinson's
14 disease and stroke. Using ML to identify salient patterns in sensor data has the potential for
15 widespread application in neurological disorders, so understanding how to develop this approach
16 for one's area of inquiry is vital. We previously proposed an approach that combined wearable
17 inertial measurement units (IMUs) and ML to classify motions made by stroke patients. However,
18 our approach had computational and practical limitations. We address these limitations here in the
19 form of a primer, presenting how to optimize a sensor-ML approach for clinical implementation.
20 First, we demonstrate how to identify the ML algorithm that maximizes classification performance
21 and pragmatic implementation. Second, we demonstrate how to identify the motion capture
22 approach that maximizes classification performance but reduces cost. We used previously
23 collected motion data from chronic stroke patients wearing off-the-shelf IMUs during a
24 rehabilitation-like activity. To identify the optimal ML algorithm, we compared the classification
25 performance, computational complexity, and tuning requirements of four off-the-shelf algorithms.
26 To identify the optimal motion capture approach, we compared the classification performance of
27 various sensor configurations (number and location on the body) and sensor type (IMUs versus
28 accelerometers). Of the algorithms tested, linear discriminant analysis had the highest
29 classification performance, low computational complexity, and modest tuning requirements. Of
30 the sensor configurations tested, seven sensors on the paretic arm and trunk led to the highest
31 classification performance, and IMUs outperformed accelerometers. Overall, we present a refined
32 sensor-ML approach that maximizes both classification performance and pragmatic
33 implementation. In addition, with this primer, we showcase important considerations for
34 appraising off-the-shelf algorithms and sensors for quantitative motion assessment.

35

36 **1 Introduction**

37 Wearable sensors, such as inertial measurement units (IMUs) and accelerometers, provide an
38 opportunity for the objective and seamless capture of human motion. Machine learning (ML)
39 enables computers to learn without being explicitly programmed, and provides an opportunity to
40 rapidly identify patterns in data. Given recent technological and computational advances,

41 combining wearable sensor data with ML algorithms has the potential for rapid, automated, and
42 accurate classification of motion.

43 Researchers have begun using this combined sensor-ML approach in a number of applications.
44 These include human activity recognition [1-3], gesture analysis [4], assessment of bradykinesia
45 in Parkinson’s disease [5], motor function assessment in multiple sclerosis [6], and differentiating
46 between functional and non-functional arm usage in stroke patients [7, 8]. While many of these
47 studies showcase the application of sensors and ML in clinical populations, no previous work has
48 detailed the various hardware and software considerations for using the sensor-ML approach.
49 Furthermore, no guide currently exists to advise investigators in building and troubleshooting this
50 approach, which sits at the intersection of human movement science, data science, and neurology.
51 With the potential for the sensor-ML approach to have widespread applicability to neurological
52 disorders, understanding how to develop this approach for one’s own area of inquiry is paramount.

53 One possible application of the combined sensor-ML approach is the monitoring of rehabilitation
54 dose in stroke patients. Quantifying the dose of rehabilitation entails classifying units of
55 measurement, which are subsequently tallied. In our previous proof-of-principle study, we used
56 IMUs worn by stroke subjects performing a structured tabletop activity to capture motion data.
57 Our units of measurement were functional primitives, elemental motions that cannot be further
58 decomposed by a human observer. We applied an ML algorithm (hidden Markov model with
59 logistic regression) to the IMU motion data to recognize functional primitives embedded in this
60 activity, achieving an overall classification performance of 79% [9]. While promising, this sensor-
61 ML approach had variable classification performance among the primitives (62-87% accuracy). It
62 also did not address research implementation challenges such as the computational complexity and
63 computational costs of the ML approach, or clinical implementation challenges such as the
64 expense [10] and electromagnetic intolerance of the IMUs.

65 In the present study, we address these limitations in the form of a primer, outlining deliberations
66 that researchers developing their own sensor-ML approach would need to consider. We describe
67 our rationale and steps for identifying (1) an algorithm that is highly accurate but computationally
68 tractable, and (2) the type and array of sensors that minimize cost but maximize accuracy. We use
69 functional primitives as the motion type to be classified, and describe our approach for both
70 capturing and identifying these motions. We also use off-the-shelf algorithms and sensors,
71 providing an accessible framework for investigators seeking to address new research and clinical
72 questions with the sensor-ML approach.

73 **2 Methods**

74 To demonstrate the steps in identifying the optimal ML algorithm and sensor array, we use data
75 collected from previous work [9]. Briefly, six mild-to-moderately impaired stroke patients (Table
76 1) moved a toilet paper roll and aluminum can over a horizontal array of targets (Fig. 1).

77 Subjects performed 5 trials moving the object between a center target and eight radially arrayed
78 targets (20 cm away). The task generates the following functional primitives: *reach* (to move into
79 contact with a target object); *transport* (to convey a target object); *reposition* (to move proximate
80 to a target object); and *idle* (to stand at the ready near target object). Functional primitives are
81 discrete, object-oriented motions with a single goal. Functional primitives are non-divisible and
82 are largely invariant across individuals [11], may be represented cortically [12-14], and provide a
83 finer-grained capture of performance in stroke patients who may be unable to accomplish a full
84 activity. Akin to words, functional primitives are combined to make a functional movement [15]

85 (analogous to a sentence), which in turn are combined to make an activity (analogous to a
86 paragraph) [16]. For example, a series of *reach-transport-reposition* primitives could constitute a
87 functional movement for zipping up a jacket, within the activity of dressing.

88 Motion data were recorded with 11 IMUs (XSens Technology) worn on the head, sternum, pelvis,
89 and bilateral hands, forearms, arms, and scapulae. 3D linear accelerations, 3D angular velocities,
90 and quaternions were generated at 240 Hz. To segment and label the motion data as constituent
91 primitives, we synchronously recorded motion (30 Hz) with a single video camera. Trained coders
92 used the video recording to label the beginning and end of each functional primitive, which also
93 labeled the corresponding IMU data. These labels served as the ground truth. This step enabled us
94 to train ML algorithms on motion data and test their classification performance against the ground-
95 truth labels. IMU data were z-score normalized and statistical features were extracted. The
96 statistical features were the following: mean, standard deviation, minimum, maximum, entropy,
97 skewness, energy, and root mean square. These statistical features have been shown to capture
98 human motion efficiently, reducing the computational burden [17-19]. Following prior work, we
99 selected a window size of 0.25 s sliding by 0.1 s [9], from which to derive the statistical features.
100 The statistical feature data were fed to the ML algorithms.

101 The dataset consisted of 2881 primitives, consisting of 810 *reaches*, 708 *transports*, 781
102 *repositions*, and 582 *idles*. It is important to note that this is the sample size of interest (not the
103 number of subjects). Accounting for repeated measures within-subject and at each target, and using
104 this dataset of 2881 primitives with $\alpha = 0.05$, we have 81% power to detect a classification
105 performance of at least 79% (positive predictive value, section 3.2.1 below). We used 79%
106 accuracy as the benchmark for sufficient classification performance as achieved in our previous
107 study [9].

108 **3 Computational details**

109 **3.1 ML methods for classification**

110 In the present study, we sought to identify an ML algorithm that performs well for identifying
111 primitives, i.e. has a high classification performance, but that also is practical, i.e. has low
112 computational overhead and minimal tuning requirements. Supervised ML algorithms work in two
113 phases: training and testing. During training, ML algorithms learn the relationship between a
114 pattern of data characteristics (here, the statistical features) and its class (here, its primitive label).
115 During testing, the trained ML algorithm uses the pattern of data characteristics to identify a new
116 data sample as one of the primitives. This identification is checked against the ground-truth human
117 label, thus reading out classification performance.

118 We considered both generative and discriminative algorithms. Generative algorithms model the
119 underlying distribution of data for each class, seeking to identify data characteristics that enable
120 matching of new data samples to a given class. In contrast, discriminative algorithms model the
121 boundaries between classes and not the data themselves. They seek to identify the plane separating
122 the classes so that, based on location relative to the plane, a new data sample is assigned to the
123 appropriate class.

124 We selected four algorithms that have been found to provide high classification performance in
125 human activity recognition: linear discriminant analysis (LDA) [18], Naïve Bayes classifier (NBC)
126 [17], support vector machine (SVM) [20], and k-nearest neighbors (KNN) [19]. LDA and NBC
127 are generative algorithms, whereas SVM and KNN are discriminative algorithms. We used off-

128 the-shelf versions of these algorithms without any special permutations; in other words, the
129 algorithms are widely available in most machine learning libraries such as scikit-learn [21, 22].

130 **3.2 Algorithm performance metrics**

131 **3.2.1 Classification performance of algorithms**

132 We first evaluated how well the algorithms could classify primitives, measuring classification
133 performance by comparing algorithm-chosen labels against ground-truth human labels. Primitives
134 were classified as true positive (*TP*, labels agreed) and false positive (*FP*, labels disagreed). We
135 used 60% of the data to train the algorithm and 40% to test it, repeating the process 10 times. A
136 validation dataset was not used because we were not optimizing algorithm architectures, and the
137 test dataset provides an unbiased estimate of algorithm performance. Data were randomly selected
138 for each primitive proportional to its prevalence in the complete dataset (i.e., stratified proportional
139 sampling). This ensured that each dataset adequately represented the entire sample population. In
140 addition, to examine the possibility that within-subject dependencies in the training and testing
141 sets leads to an overestimation of classification performance, we also performed a leave-one-
142 subject-out analysis i.e., training the algorithms using data from all but one subject and testing its
143 performance on the data from the remaining subjects. This process was repeated 6 times, once for
144 each subject, and classification performances were averaged.

145 The first metric for classification performance was positive predictive value (PPV; $TP/(TP+FP)$).
146 PPV reflects how often a primitive was actually performed when the algorithm labeled it as such;
147 in other words, PPV is how often a primitive was correctly classified. We generated primitive-
148 level PPVs in a one-versus-all analysis (e.g., *reach* vs. *transport* + *reposition* + *idle* combined).
149 We also generated an overall PPV by combining data for all primitives and tallying all true and
150 false positives. We prefer PPV because it takes into account the prevalence of the primitive in the
151 dataset [23].

152 The second metric for classification performance was the receiver operating characteristic (ROC)
153 curve. ROC curves depict the relative tradeoff between true positive rate (sensitivity; y-axis) and
154 false positive rate (1-specificity; x-axis) and identify the optimal operating point of an algorithm
155 [24]. Perfect classification would lead to a ROC curve that passes through the upper left corner,
156 with an area under the ROC curve (AUC) equal to 1 and an operating point at 100% sensitivity
157 and 100% specificity [24].

158 **3.2.2 Practical performance of algorithms**

159 We next considered the computational complexity of the algorithms in terms of their training and
160 testing times and their tuning requirements. Having a high computational complexity means that
161 specialized computing hardware and advanced expertise would be needed, potentially hindering
162 widespread implementation in research.

163 **3.2.2.1 Training and testing times of the algorithms**

164 The time required to train and test the algorithms was measured for datasets of different sizes. If
165 training time is fast, rapid appraisal and optimization of the algorithm are possible, favoring rapid
166 development and deployment. If the testing time is fast, real-time classification and online
167 feedback are possible, favoring clinical implementation.

168 We first used 20-100% of the dataset (n=2881 primitives) in randomly selected 10% increments.
169 At each increment, we measured (1) the time required to train the algorithm (training time), and

170 (2) the time required for a trained algorithm to classify a primitive (testing time). At each 10%
171 increment, the algorithms were trained *de novo* to avoid overfitting and to provide unbiased
172 estimates.

173 Given the modest size of our dataset, we next used a simulated dataset that could be expected from
174 a typical sample size of 50 subjects performing a variety of activities. The simulated dataset had
175 300,000 primitives with same proportion, mean, and variance as our original dataset. We used 25-
176 100% of the dataset in randomly selected 25% increments. At each increment, we measured the
177 training and testing times, training the algorithm *de novo* as above. Of note, the simulated dataset
178 was used only to generate training and testing times, and was not used for classification
179 performance assessments.

180 **3.2.2.2 Tuning requirements of the algorithms**

181 We also assessed the algorithm’s need for tuning, the adjustment of algorithm parameters to
182 maximize classification performance. A high tuning requirement requires the extensive analysis
183 of the algorithm to identify its optimal parameters, potentially limiting implementation in settings
184 that lack domain expertise. Of note, tuning requirements were only used to index complexity, but
185 we did not tune the algorithms themselves in the assessment of classification performance.

186 We operationalized the algorithms’ tuning requirements as the number of parameters that can be
187 adjusted. We also qualitatively classified the level of domain knowledge required to implement
188 and tune the algorithms. Based on typical US educational programs, “low” domain expertise
189 indicates a basic knowledge of statistics, “medium” indicates undergraduate-level knowledge of
190 machine learning, and “high” indicates graduate-level knowledge of machine learning.

191 **3.2.3 Optimal sensor characteristics**

192 We then focused on the hardware side, seeking the best balance between ease of motion capture
193 and high classification performance. We first considered the use of IMUs compared to
194 accelerometers alone. IMUs are a combination of sensors, including accelerometers, gyroscopes,
195 and magnetometers. Many IMU hardware-software systems generate 3D linear accelerations, 3D
196 angular velocities, 3D magnetic heading, and 4D quaternions, resulting in 10 data dimensions per
197 sensor. We used accelerations, angular velocities, and quaternions for derivation of statistical
198 features (section 2), as these data types have been used previously for human activity recognition
199 [19, 25, 26]. In contrast, 3D accelerometers generate only 3D linear accelerations, resulting in 3
200 data dimensions per sensor.

201 While IMUs are data-rich, they are challenged by electromagnetic drift. Magnetic environments
202 lead to potentially inaccurate gyroscopic measurements and therefore necessitate frequent
203 recalibration. While accelerometers are data-sparse, they are largely unaffected by a magnetic
204 environment.

205 Another practical consideration for sensor choice is system expense. IMU systems can cost on the
206 order of thousands of dollars [10] whereas accelerometry systems cost in the hundreds [27]. It is
207 possible that cost and set-up time could be optimized by reducing the number of sensors or by
208 using accelerometers alone. Although simplified and less expensive motion capture would favor
209 clinical implementation, it may come at the cost of reduced classification performance.

210 In this analysis, we subsampled data from the IMUs to extract accelerometry data, ensuring that
211 comparisons were based on identical sensor locations and primitive motions. LDA was trained and
212 tested on the separate datasets to read out effects on classification performance.

213 3.2.3.1 Optimal sensor number and configuration for classification

214 We first evaluated how the number of sensors and their location on the body affects classification
215 performance. We used exhaustive search to systematically test all possible sensor configurations
216 [28]. This approach provides an unbiased appraisal of all sensor combinations for each incremental
217 reduction in sensor number.

218 3.2.3.2 Optimal sensor type for classification

219 We also evaluated how sensor type affected classification performance. We compared
220 classification accuracies using IMU data versus accelerometry-only data. This allowed us to
221 determine whether accelerometers, with their reduced dimensionality, could enable sufficient
222 accuracy to warrant their use in lieu of IMUs.

223 4 Results

224 4.1 Classification performance of algorithms

225 We first determined the classification performance of multiple ML algorithms using PPVs (Table
226 2). LDA and SVM had high classification performance for all primitives (overall PPV 92.5% and
227 92%, respectively). KNN had intermediate performance (PPV 87.5%) and NBC had the lowest
228 performance (PPV 80.2%), particularly for *reaches* (PPV 77%) and *transports* (PPV 71%). In the
229 leave-one-subject-out analysis, which addressed the possibility of within-subject dependencies,
230 similar overall classification performances were identified (PPVs of 89% for LDA, 90% for SVM,
231 83% for KNN, and 75% for NBC).

232 To further characterize classification performance, we generated ROC curves for each primitive
233 (Fig. 2). All algorithms detected *idle* with high accuracy (AUC > 0.87). For the other primitives,
234 LDA and SVM had AUCs 0.95-0.99, indicating very high classification performance. KNN also
235 had high classification performance for *reach* (AUC 0.94) and *transport* (AUC 0.90) and
236 intermediate classification performance for *reposition* (AUC 0.87). In contrast, NBC had the
237 lowest classification performance on the remaining primitives (AUC 0.80-0.85). We also
238 identified the optimal operating point, indicating the best tradeoff between sensitivity and
239 specificity, for each algorithm (Fig. 2). At their respective optimal operating points, LDA and
240 SVM achieved high sensitivities (0.83-0.95) and specificities (0.83-0.95) for all primitives. KNN
241 achieved a high sensitivity (0.91) and specificity (0.86) for *transport*, but had moderate
242 sensitivities (0.80-0.88) and specificities (0.79-0.86) for other primitives. NBC had the lowest
243 sensitivities (0.74-0.81) and specificities (0.74-0.79) for all primitives. In sum, these findings
244 indicate that LDA and SVM have the highest classification performance of the algorithms tested.

245 4.2 Training and testing times of the algorithms

246 We next evaluated the pragmatic aspects of implementing the algorithm to gauge real-world
247 applicability. We first calculated the time required to train and test the algorithm on increasing
248 quantities of data (Fig. 3) from our dataset of 2880 primitives. In terms of training times, NBC and
249 LDA were on the order of seconds (12 s and 26 s, respectively), with training times growing
250 linearly with increasing data quantity. SVM was on the order of minutes (5.6 min), with training
251 times growing quadratically with increasing data quantity. KNN required no time to train as an
252 inherent property of the model. In terms of testing, LDA, NBC, and SVM required sub-millisecond

253 times (approximately 0.03 ms), whereas KNN required the longest time (1.5 ms) with testing times
254 growing linearly with increasing dataset size.

255 To investigate the real-world ramifications of training and testing requirements, we generated a
256 dataset with 300,000 primitives (Fig. 4). Training times became prohibitively long for SVM (up
257 to 23 h) but were manageable for the other algorithms (up to 13 min). Testing time was relatively
258 high for KNN (up to 2.3 min), whereas LDA, NBC, and SVM required nominal testing times
259 (<0.03 ms). Given their consistently low training and testing times, LDA and NBC have the best
260 practical performance of the algorithms tested.

261 **4.3 Tuning requirements of the algorithms**

262 To gauge the difficulty of algorithm implementation, we characterized their tuning requirements
263 (Table 3). NBC has the lowest number of parameters (1) and requires a low amount of domain
264 knowledge in machine learning to optimize it. KNN has a moderate number of parameters (5), but
265 their optimization is reasonably intuitive and requires a low level of domain knowledge. LDA has
266 fewer parameters (3), but they require a medium level of domain knowledge. SVM has many
267 parameters (9) and requires a high level of domain knowledge to build an accurate and efficient
268 model. In sum, these findings indicate that NBC and KNN are the easiest to implement, and LDA
269 implementation requires a modestly higher skillset.

270 **4.4 Optimal sensor characteristics**

271 **4.4.1 Optimal sensor number and configuration**

272 To evaluate the effect of the sensor number and configuration on classification performance, we
273 used an exhaustive search process, which evaluated all combinations of sensor number and
274 location. We note that exhaustive search arrived at the same optimal configurations for IMUs as
275 for accelerometers. Seven sensors on the head, sternum, pelvis, and UE of the active side resulted
276 in the highest classification performance (IMU PPV 92.5%; accelerometer PPV 84%). In
277 comparison, when exhaustive search progressively added sensors to the non-active forearm, then
278 hand, then upper arm, then scapula, classification performance worsened (IMU PPV 88%;
279 accelerometer PPV 80%) (Fig. 5). When exhaustive search progressively removed sensors on the
280 trunk and then head, performance also worsened. Subsequent removal of sensors from the scapula,
281 then arm, and then hand further worsened performance, arriving at PPVs of 71% and 62% for
282 IMUs and accelerometers, respectively, for the remaining forearm sensor.

283 **4.4.2 Optimal sensor type for classification**

284 To finish, we evaluated classification performance using IMU versus accelerometry data only.
285 Classification performance using accelerometry data was consistently lower than for IMU data for
286 all sensor configurations (Fig. 5; Table 4). Classification performance with accelerometers was
287 lower especially for *reaches* (PPV 77% vs. 93%; Table 4), which include different arm
288 configurations to grasp the objects (e.g. supinating to side-grasp the aluminum can versus
289 pronating to overhand grasp the toilet paper roll). These findings indicate that IMU data enable a
290 superior level of classification, particularly with more variable motions involving forearm
291 rotations.

292

293

295 The combination of wearable sensors and machine learning offers exciting opportunities in
296 numerous applications, including human activity recognition [1-3] and assessment of impaired
297 motion [5, 7, 8]. We recently proposed an approach that uses wearable sensors and ML algorithms
298 to classify functional primitives, which could be summed to quantify rehabilitation dose. In this
299 study, we aimed to address limitations in this previous work, including a modest computational
300 performance, high computational complexity, and hardware drawbacks. We present our analyses
301 as a primer for considering software and hardware variables in the capture and classification of
302 motion data. We sought to identify—from both performance and practical standpoints—the best
303 machine learning algorithm, sensor configuration, and sensor type to classify functional primitives
304 in stroke patients.

305 Among the ML algorithms, LDA represented the best balance of classification performance and
306 pragmatic implementation. Among sensor configurations, seven sensors on the paretic arm and
307 trunk enabled better classification performance than more or fewer sensors on the body. Among
308 sensor types, IMU data enabled better classification performance than accelerometers. To our
309 knowledge, this is the first study to systematically outline the steps of identifying optimal ML
310 algorithms, sensor configurations, and sensor types to automatically classify motion patterns of
311 neurological patients.

312 **Optimal performer in classification.** Evaluating the ability of the ML algorithms to classifying
313 functional primitives, we found that LDA and SVM had the highest classification performance.
314 LDA performs well because it aims to reduce dimensionality while preserving as much
315 discriminatory information as possible. This approach leads to tight clusters and high separation
316 between the classes [29]. SVM performs well because it projects training data to a high-
317 dimensional space. This approach leads to maximal separation between classes that may not be
318 possible in the original feature space [30]. Overall, LDA aims to find commonalities within and
319 differences between data classes, whereas SVM aims to find a classification boundary that is
320 furthest from the data classes. Importantly, these algorithms maximize rigor in the training phase
321 by being less susceptible to noisy or outlier data [31, 32]. LDA accomplishes this by using the
322 clusters' centers and ignoring outlier samples to classify [31], while SVM uses the most closely
323 spaced data (i.e., the most difficult to discriminate) to define class boundaries [32]. It is worth
324 noting that LDA assumes that the underlying classes are normally distributed (unimodal
325 Gaussians) with the same covariance matrix [29]. If real-world motion data are significantly non-
326 Gaussian, LDA may not capture the complex data structures required for accurate classification.
327 In this case, classification performance can be tuned by allowing the covariance matrices among
328 classes to vary, resulting in a regularized discriminant analysis [33].

329 By comparison, KNN showed a marginally lower classification performance, likely due to its
330 susceptibility to noise [34]. KNN relies on the assumption that samples from the same class exist
331 in close proximity. Given a new sample, KNN assigns it to the class with the majority of closest
332 neighbors [35]. In our current setup of KNN, all nearest sample points are given the same
333 weighting. Therefore when assigning a class label, a noisy sample will be weighted the same as
334 other statistically important samples. KNN classification performance can be tuned by choosing
335 an appropriate weighting metric (e.g., inverse squared weighing) [36], which ensures that samples
336 closer to the test sample contribute more to classifying it. Performance may also be tuned by using
337 mutual nearest neighbors, where noisy samples are detected using pseudo-neighbors (neighbors of
338 neighbors) and are assigned lower weights [37].

339 Finally, NBC had the lowest performance compared to other algorithms. NBC uses Bayes' rule
340 and prior information to classify a new sample, using the posterior probability of it belonging to a
341 class [38]. Its lower performance may be attributed to its underlying assumption of conditional
342 independence between data features [39]. This assumption is violated for data streams that are
343 correlated, such as data from adjacent sensors on the body, like the hand and wrist. The
344 performance of NBC could be improved by applying principal components analysis to the dataset
345 as a pre-processing step, and then training the NBC [40].

346 Comparing these results with our prior work [9], we found that the four algorithms outperformed
347 the hidden Markov model-logistic regression (HMM-LR) classifier for identifying the functional
348 primitives in stroke patients. The improved performance may be due in part to differences in the
349 training datasets. Our previous study trained the algorithm on healthy control data and tested on
350 stroke patient data to examine the generalizability of the model. It is conceivable that if the HMM-
351 LR classifier been trained and tested in stroke patients only, its performance would have been
352 higher.

353 **Optimal performer in practicality.** We also determined the most pragmatic algorithms with
354 respect to their training and testing times and their tuning requirements. In terms of training times,
355 KNN did not have any computational overhead. This is expected, since KNN requires no training
356 and shifts its computations to the testing phase. Training times for LDA and NBC grew gradually
357 with dataset size, but took at most minutes with a real world-sized dataset. LDA had lower training
358 times than NBC on a smaller dataset, but required more training time as the dataset increased. This
359 is explained by the scatter matrix computations and optimization of LDA, which become
360 computationally expensive as the dataset size increases [18]. By contrast, SVM training time
361 increased quadratically with dataset size, because finding an optimal hyperplane between classes
362 entails solving a quadratic programming problem [20]. Complex algorithms such as SVM thus
363 require more processing time for large datasets, which limits real-world application. For example,
364 for a modestly sized study, training times for SVM may be on the order of days. This lag would
365 be prohibitive for rapid tuning, significantly delaying algorithm optimizations. Conversely,
366 performance of LDA and NBC could be rapidly appraised after training, alerting an investigator
367 to further tune the algorithm or to move on from it.

368 In terms of testing times, SVM, LDA, and NBC required sub-milliseconds to classify primitives,
369 whereas KNN took seconds-minutes and testing times grew linearly with dataset size. This can be
370 explained by the exhaustive and computationally expensive search performed by KNN [41].
371 During testing, the KNN algorithm searches for the k nearest neighbors that have similar data
372 characteristics as the test sample. With increasing samples and dimensionality of the data, the
373 search broadens and takes more time. If an investigator wishes to classify primitives offline, KNN
374 testing times may be acceptable. For applications requiring near- or real-time classification (e.g.
375 for online feedback), the other algorithms should be considered instead. Alternatively, the
376 computational complexity of KNN can be reduced by selecting an efficient search algorithm (e.g.,
377 KD tree) [42], which limits the search space during testing.

378 In terms of ease of tuning to increase classification performance and reduce training/testing time,
379 we determined that NBC had the lowest parameter complexity and requirement for domain
380 knowledge in machine learning. KNN has a moderate number of tuning parameters, but they are
381 relatively straightforward to understand and address. LDA has fewer tuning parameters than KNN,
382 but moderate domain knowledge is required to select the amount of regularization allowing the
383 covariance among classes to vary [33]. SVM requires the highest amount of parameter tuning, and

384 necessitates a deep understanding of statistics, optimization, probability theory, and machine
385 learning [43]. This level of domain knowledge is prohibitive for SVM use in an unsupported
386 research setting.

387 Weighing classification performance and pragmatic implementation, we judged LDA to be the
388 best choice for our application. Investigators will similarly need to weigh their performance goals,
389 time resources, and available level of expertise for ML implementation in their own motion
390 classification questions.

391 **Optimal IMU configuration.** On the hardware side, we determined the optimal sensor location
392 and configuration to facilitate data capture while maintaining high classification performance.
393 Seven sensors (not more or fewer) enabled optimal classification performance, and the best sensor
394 configuration was placed on the active limb and trunk. This result is expected, given that the
395 participants performed a unimanual task. Interestingly, accuracy worsened with more sensors,
396 likely because of the increased dimensionality of the dataset. This may cause the ML algorithm to
397 overfit the training data, resulting in lower classification performance during the testing phase [44].
398 Finally, we found that if only one sensor was available, the forearm location was the most
399 informative, although classification performance was modest. This location is nonetheless
400 appealing, given recent advances in smartwatches capable of capturing motion.

401 **Optimal data characteristics.** Finally, we determined the sensor type that led to the highest
402 classification performance. Accelerometry data consistently generated lower accuracies than IMU
403 data, likely due to its fewer dimensions. Although IMUs enable higher classification performance
404 than accelerometers, they also have some drawbacks: a higher risk of electromagnetic drift leading
405 to inaccurate data estimates and the need for more frequent recalibrations, a higher consumption
406 of energy [45], and a higher cost [10]. Thus there is a tradeoff between robust motion capture and
407 practical motion capture. We believe that the benefits of richer data and better classification of
408 IMUs outweigh their practical limitations. However, there currently exist no benchmarks for the
409 level of classification accuracy needed to justify clinical implementation. If these accuracy
410 benchmarks are lower than those achieved by IMUs, and if investigators are constrained by
411 financial resources or the magnetic noisiness of an environment, accelerometers could be
412 appropriate.

413 **5.1 Limitations and future work**

414 Our study has some limitations to be considered. We importantly do not suggest that we have
415 found the definitive approach for classifying primitives in rehabilitating stroke patients, for two
416 reasons. First, our analysis was performed on a dataset of mild-to-moderately impaired stroke
417 patients, limiting generalization to stroke patients with severe impairment. To achieve high
418 classification performance across the range of stroke impairment, separately trained ML models
419 may be needed for different impairment levels. Second, the activity used in this study was highly
420 structured. The motion characteristics of the resulting primitives were thus more consistent and
421 limited than what would be found in a real-world rehabilitation setting. The training and testing of
422 algorithms on functional primitives with an array of kinematic characteristics is still required, and
423 is ongoing in our laboratory. We use this circumscribed dataset here to showcase the practical
424 deliberations required in the development of a sensor-ML approach for motion classification.

425

426

427 6 Conclusion

428 In summary, we present a primer that details how one can optimize both the software and hardware
429 facets of motion capture. This work outlines computational and practical considerations for
430 implementing a sensor-ML approach in quantitative research. Specific to our application, we
431 demonstrate how to refine a strategy that builds towards the precise and pragmatic classification
432 of functional primitives in stroke patients. We found that LDA had the best combination of
433 classification performance and pragmatic performance. We also found that seven sensors on the
434 paretic UE and trunk maximized classification performance, and that IMUs enabled superior
435 classification compared to accelerometers.

436 7 Declarations

437 **Ethics approval and consent to participate:** Institutional Review Board-approved testing
438 occurred at Columbia University. Subjects gave written informed consent to participate in this
439 study, in accordance with the Declaration of Helsinki.

440 Consent for publication: Not applicable

441 **Availability of data and material:** The dataset analyzed for the current study are available from
442 the corresponding author on reasonable request.

443 **Pre-print:** This manuscript has been released as an online preprint on arXiv
444 (<https://arxiv.org/abs/1902.08697>)

445 **Competing interests:** All authors report no financial or non-financial competing interests.

446 **Funding:** The work was supported by K02NS104207 (HS) and K23NS078052 (HS).

447 **Authors' contributions:** AP analyzed and interpreted the data and wrote the paper. JU collected
448 and coded the data. DN created the activities battery and interpreted the data. HS collected, labeled,
449 and interpreted the data and wrote the paper.

450 **Acknowledgements:** We would like to acknowledge Dr. Sunil Agarwal for provision of the sensor
451 system and Jorge Guerra for his early machine learning analysis. We thank Shatif Hughes and Dr.
452 Hyoungseop Kim for their assistance with data labeling.

453

454 8 Figure legends

455 **Figure 1. Tabletop activity set-up.** Healthy individual wearing the sensors and transporting the
456 object from center to a target.

457 **Figure 2. Performance characteristics of machine learning algorithms for (A) Reach, (B)**
458 **Transport, (C) Reposition, and (D) Idle.** Receiver operating characteristic (ROC) curves show
459 the trade-off between true positive rate (or sensitivity) and false positive rate (1-specificity).
460 Curves closer to the top-left corner indicate a better classification performance. The optimal
461 operating point for each algorithm (solid circles), reflect the best tradeoff between sensitivity and
462 specificity for an algorithm. The area under the curve (AUC), a measure of classification
463 performance, is shown in parenthesis for each algorithm. AUC=1 represents perfect classification.
464 LDA had the highest AUCs followed closely by SVM, indicating high classification performances.
465 NBC had consistently the lowest AUCs, indicating the weakest classification performance.

467 **Figure 3. Algorithm (A) training times and (B) testing times on sample dataset. The dataset**
 468 **is comprised of 2880 primitives.** We computed times to train and test each algorithm on 20-100%
 469 of the dataset in increments of 10%. To avoid overfitting and compute an unbiased estimate of
 470 training and testing times, ML algorithms were trained and tested de novo with each incremental
 471 increase. For training with the complete sample dataset, SVM required the most time (336 s) while
 472 the other algorithms finished training rapidly (<30 s). For testing, KNN required the most time
 473 (1.5 ms), while the other algorithms finished testing rapidly (<0.03 ms). Please note break in the
 474 y-axis to highlight the difference in the algorithm testing times.

475 **Figure 4. Algorithm (A) training times and (B) testing times on real world-sized dataset. The**
 476 **dataset is comprised of 300,000 simulated primitives.** We evaluated training and testing times
 477 required by each algorithm for quartile increases in dataset size. Please note the break in the y-
 478 axes to highlight differences in training and testing times. To avoid overfitting and compute
 479 unbiased estimates, the algorithms were trained and tested de novo at each quartile. For training
 480 with the entire dataset, SVM required the most time (1380 min) while the other algorithms required
 481 less time (LDA: 13 min; NBC: 2.5 min; KNN: 0 min, as per model property). For testing, KNN
 482 required the most time (2.3 min). The remainder of algorithms (LDA, NBC, and SVM) needed a
 483 testing time of <0.09 ms, which grew marginally with increasing sample sizes.

484 **Figure 5. Classification performance for full and reduced sensor counts.** Performance was
 485 computed using LDA and data from with progressively reduced sensor counts. Seven sensors
 486 (pelvis, sternum, head, and the active shoulder, upper arm, forearm, and hand) gave the best
 487 classification performance, with a performance drop-off at more or fewer sensors. IMU data
 488 consistently supported higher classification than accelerometry data, achieving PPV 92.5% vs.
 489 82% at the seven sensors.

490 9 References

- 491 [1] R. J. Lemmens, Y. J. Janssen-Potten, A. A. Timmermans, R. J. Smeets, and H. A. Seelen, "Recognizing
 492 complex upper extremity activities using body worn sensors," *PLoS One*, vol. 10, no. 3, p. e0118642, 2015.
- 493 [2] D. Biswas *et al.*, "Recognition of elementary arm movements using orientation of a tri-axial accelerometer
 494 located near the wrist," *Physiol Meas*, vol. 35, no. 9, pp. 1751-68, Sep 2014.
- 495 [3] N. A. Capela, E. D. Lemaire, and N. Baddour, "Feature selection for wearable smartphone-based human
 496 activity recognition with able bodied, elderly, and stroke patients," *PLoS One*, vol. 10, no. 4, p. e0124414,
 497 2015.
- 498 [4] H. Junker, O. Amft, P. Lukowicz, and G. J. P. R. Tröster, "Gesture spotting with body-worn inertial sensors
 499 to detect user activities," vol. 41, no. 6, pp. 2010-2024, 2008.
- 500 [5] B. M. Eskofier *et al.*, "Recent machine learning advancements in sensor-based mobility analysis: Deep
 501 learning for Parkinson's disease assessment," in *2016 38th Annual International Conference of the IEEE*
 502 *Engineering in Medicine and Biology Society (EMBC)*, 2016, pp. 655-658: IEEE.
- 503 [6] I. Carpinella, D. Cattaneo, and M. Ferrarin, "Quantitative assessment of upper limb motor function in
 504 Multiple Sclerosis using an instrumented Action Research Arm Test," *Journal of neuroengineering*
 505 *rehabilitation*, vol. 11, no. 1, p. 67, 2014.
- 506 [7] E. M. Bochniewicz, G. Emmer, A. McLeod, J. Barth, A. W. Dromerick, and P. Lum, "Measuring
 507 Functional Arm Movement after Stroke Using a Single Wrist-Worn Sensor and Machine Learning,"
 508 *Journal of Stroke and Cerebrovascular Diseases*, vol. 26, no. 12, pp. 2880-2887, 2017.
- 509 [8] K. Leuenberger *et al.*, "A method to qualitatively assess arm use in stroke survivors in the home
 510 environment," *Medical & Biological Engineering & Computing*, vol. 55, no. 1, pp. 141-150, 2017.
- 511 [9] J. Guerra *et al.*, "Capture, learning, and classification of upper extremity movement primitives in healthy
 512 controls and stroke patients," (in eng), *IEEE Int Conf Rehabil Robot*, vol. 2017, pp. 547-554, Jul 2017.

- 513 [10] *Analog Devices Inertial Measurement Units*. Available: <https://www.mouser.com/ProductDetail/Analog-Devices/ADIS16497-3BMLZ?qs=sGAEpiMZZMt8xBNXJng%252BZAmSdp7Aftz5BZDziPI1t2Q%2FH2F5IOOUxA%3D%3D>
- 514
- 515
- 516 [11] G. Sumbre, G. Fiorito, T. Flash, and B. Hochner, "Neurobiology: motor control of flexible octopus arms," *Nature*, vol. 433, no. 7026, pp. 595-6, Feb 10 2005.
- 517
- 518 [12] M. Desmurget *et al.*, "Neural representations of ethologically relevant hand/mouth synergies in the human precentral gyrus," *Proc Natl Acad Sci U S A*, vol. 111, no. 15, pp. 5718-22, Apr 15 2014.
- 519
- 520 [13] M. S. A. Graziano, "Ethological Action Maps: A Paradigm Shift for the Motor Cortex," *Trends Cogn Sci*, vol. 20, no. 2, pp. 121-132, Feb 2016.
- 521
- 522 [14] F. A. Mussa-Ivaldi and S. A. Solla, "Neural primitives for motion control," *IEEE Journal of Oceanic Engineering*, vol. 29, no. 3, pp. 640-650, 2004.
- 523
- 524 [15] C. E. Lang, D. F. Edwards, R. L. Birkenmeier, and A. W. Dromerick, "Estimating minimal clinically important differences of upper-extremity measures early after stroke," *Arch Phys Med Rehabil*, vol. 89, no. 9, pp. 1693-700, Sep 2008.
- 525
- 526
- 527 [16] C. Bregler, "Learning and recognizing human dynamics in video sequences," in *IEEE Conference on Computer Vision and Pattern Recognition*, 1997, pp. 568-574.
- 528
- 529 [17] L. Bao and S. Intille, "Activity recognition from user-annotated acceleration data," *Pervasive computing*, pp. 1-17, 2004.
- 530
- 531 [18] A. M. Khan, Y.-K. Lee, S. Y. Lee, and T.-S. Kim, "A triaxial accelerometer-based physical-activity recognition via augmented-signal features and a hierarchical recognizer," *IEEE transactions on information technology in biomedicine*, vol. 14, no. 5, pp. 1166-1172, 2010.
- 532
- 533
- 534 [19] N. Ravi, N. Dandekar, P. Mysore, and M. L. Littman, "Activity recognition from accelerometer data," in *Aaai*, 2005, vol. 5, no. 2005, pp. 1541-1546.
- 535
- 536 [20] C. Schuldt, I. Laptev, and B. Caputo, "Recognizing human actions: a local SVM approach," in *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, 2004, vol. 3, pp. 32-36: IEEE.
- 537
- 538
- 539 [21] Mathworks. (2018). *Analyze and model data using statistics and machine learning*. Available: <https://www.mathworks.com/products/statistics.html>
- 540
- 541 [22] (2018). *scikit-learn Machine Learning in Python*. Available: <http://scikit-learn.org/stable/>
- 542 [23] R. Parikh, A. Mathai, S. Parikh, G. C. Sekhar, and R. J. I. j. o. o. Thomas, "Understanding and using sensitivity, specificity and predictive values," vol. 56, no. 1, p. 45, 2008.
- 543
- 544 [24] B. JC, "Statistical Analysis and Presentation of Data," in *Evidence-Based Laboratory Medicine; Principles, Practice and Outcomes*, C. Price and R. Christenson, Eds. 2 ed. Washington DC, USA: AACC Press, 2007, pp. 113-40.
- 545
- 546
- 547 [25] S. Ha and S. Choi, "Convolutional neural networks for human activity recognition using multiple accelerometer and gyroscope sensors," in *2016 International Joint Conference on Neural Networks (IJCNN)*, 2016, pp. 381-388: IEEE.
- 548
- 549
- 550 [26] Y. Shin, W. Choi, and T. Shin, "Physical activity recognition based on rotated acceleration data using quaternion in sedentary behavior: A preliminary study," in *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, 2014, pp. 4976-4978: IEEE.
- 551
- 552
- 553 [27] *Analog Devices Accelerometers*. Available: <https://www.mouser.com/ProductDetail/Analog-Devices/ADXL357BEZ?qs=sGAEpiMZZMvwE4h8i4g3cg1FWtTpVweBI4PgoJiZDJ%2FRVLPerCg45w%3D%3D>
- 554
- 555
- 556 [28] E. W. Weisstein. *Exhaustive Search*. From *MathWorld--A Wolfram Web Resource*. Available: <http://mathworld.wolfram.com/ExhaustiveSearch.html>
- 557
- 558 [29] A. J. Izenman, "Linear discriminant analysis," in *Modern multivariate statistical techniques*: Springer, 2013, pp. 237-280.
- 559
- 560 [30] C. Burges, "A tutorial on support vector machines for pattern recognition," *Journal Data mining knowledge discovery*, vol. 2, no. 2, pp. 121-167, 1998.
- 561
- 562 [31] W. Zhao, R. Chellappa, and N. Nandhakumar, "Empirical performance analysis of linear discriminant classifiers," in *Proceedings. 1998 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No. 98CB36231)*, 1998, pp. 164-169: IEEE.
- 563
- 564
- 565 [32] T. B. Trafalis and S. A. Alwazzi, "Support vector machine classification with noisy data: a second order cone programming approach," *J International Journal of General Systems*, vol. 39, no. 7, pp. 757-781, 2010.
- 566
- 567
- 568 [33] J. H. Friedman, "Regularized discriminant analysis," *Journal of the American statistical association*, vol. 84, no. 405, pp. 165-175, 1989.
- 569

- 570 [34] C. M. Van der Walt and E. Barnard, "Data characteristics that determine classifier performance," 2006.
- 571 [35] L. E. Peterson, "K-nearest neighbor," *Scholarpedia*, vol. 4, no. 2, p. 1883, 2009.
- 572 [36] D. Bridge, "Classification: k Nearest Neighbours," ed, 2007.
- 573 [37] H. Liu and S. Zhang, "Noisy data elimination using mutual k-nearest neighbor for classification mining,"
- 574 *Journal of Systems Software*, vol. 85, no. 5, pp. 1067-1074, 2012.
- 575 [38] K. P. Murphy, "Naive bayes classifiers," *University of British Columbia*, vol. 18, 2006.
- 576 [39] M. Tom, "Generative and discriminative classifiers: Naive bayes and logistic regression," ed, 2005.
- 577 [40] L. Fan and K. L. Poh, "A comparative study of PCA, ICA and class-conditional ICA for naïve bayes
- 578 classifier," in *International Work-Conference on Artificial Neural Networks*, 2007, pp. 16-22: Springer.
- 579 [41] P. Cunningham and S. J. J. M. C. S. Delany, "k-Nearest neighbour classifiers," vol. 34, no. 8, pp. 1-17,
- 580 2007.
- 581 [42] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Communications of*
- 582 *the ACM*, vol. 18, no. 9, pp. 509-517, 1975.
- 583 [43] B. Scholkopf and A. J. Smola, *Learning with kernels: support vector machines, regularization,*
- 584 *optimization, and beyond*. MIT press, 2001.
- 585 [44] O. M. Giggins, K. T. Sweeney, and B. Caulfield, "Rehabilitation exercise assessment using inertial sensors:
- 586 a cross-sectional analytical study," *J Neuroeng Rehabil*, vol. 11, p. 158, Nov 27 2014.
- 587 [45] Q. Liu *et al.*, "Gazelle: Energy-efficient wearable analysis for running," no. 9, pp. 2531-2544, 2017.

588 10 Table legends

589 **Table 1. Demographic and clinical characteristics of patients.** Shown are number of

590 participants, mean age (range), gender, race, hand dominance, paretic side, mean Fugl-Meyer

591 assessment score at first assessment (range; maximum 66), and time since stroke (range). Inclusion

592 criteria were age ≥ 18 years; premorbid right-hand dominance; unilateral motor stroke;

593 contralateral arm weakness with Medical Research Council score $< 5/5$ in a major muscle group.

594 Exclusion criteria were traumatic brain injury; musculoskeletal, medical, or non-stroke

595 neurological condition interfering with assessment of motor function; contracture at shoulder,

596 elbow, or wrist; moderate dysmetria or truncal ataxia; visuospatial neglect; apraxia; global

597 inattention; blindness.

598 **Table 2. Classification performance of machine learning algorithms for functional**

599 **primitives.** Positive predictive values (PPV) with associated 95% confidence intervals are shown.

600 PPV reflects how often a primitive was actually made when the algorithm identified it as such,

601 was calculated for the primitives of reach, transport, reposition, and idle. Primitive-level PPVs

602 were computed in one-versus-all analysis (e.g., reach vs. transport + reposition + idle combined).

603 The overall PPV was assessed by combining data for all primitives and tallying all true and false

604 positives. Overall classification performance was highest for linear discriminant analysis (LDA)

605 and support vector machine (SVM), moderately high for k-nearest neighbors (KNN), and lowest

606 for Naïve Bayes classifier (NBC).

607 **Table 3. Complexity of algorithm implementation.** Algorithm parameter tuning is necessary to

608 achieve optimal classification performance. Shown are algorithm tuning characteristics, as

609 indicated by number and specifics of the tuning parameters. Also shown is a graded estimate of

610 the level of domain knowledge required to tune these parameters. NBC is considered the simplest

611 to tune while SVM is the most difficult. LDA has a handful of parameters that require medium

612 domain knowledge to negotiate. KNN has a moderate number of parameters that are intuitive to

613 tune and require little domain knowledge. Level of domain knowledge: low, basic knowledge of

614 statistics; medium, undergraduate-level knowledge of ML; high, graduate-level knowledge of ML.

615 **Table 4. Primitive-level classification using IMU or accelerometer data.** Classification

616 performance is shown using the 7-sensor configuration (pelvis, sternum, head, and the active

617 shoulder, upper arm, forearm, and hand). Accelerometers provided systematically poorer
 618 classification performance compared to IMUs across all primitives. Classification performance
 619 using accelerometry data was particularly low for reach (PPV 77%) and relatively higher for idle
 620 (PPV 88%).

621 **11 Tables**

622

N	6
Age (years)	61.7 (46.5 -71.0)
Gender (Female/Male)	2F/4M
Dominant arm (Right/Left)	5R/1L
Paretic side (Right/Left)	6R
Impairment (Fugl-Meyer score)	52.8 (45 - 62)
Time since stroke (years)	12.0 (2.0 - 31.1)

Table 1. Demographic and clinical characteristics of patients. Shown are number of participants, mean age (range), gender, race, hand dominance, paretic side, mean Fugl-Meyer assessment score at first assessment (range; maximum 66), and time since stroke (range). Inclusion criteria were age ≥ 18 years; premorbid right-hand dominance; unilateral motor stroke; contralateral arm weakness with Medical Research Council score $< 5/5$ in a major muscle group. Exclusion criteria were traumatic brain injury; musculoskeletal, medical, or non-stroke neurological condition interfering with assessment of motor function; contracture at shoulder, elbow, or wrist; moderate dysmetria or truncal ataxia; visuospatial neglect; apraxia; global inattention; blindness.

623
 624
 625

Algorithm	PPVs for functional primitives				Overall PPV
	Reach	Transport	Reposition	Idle	
LDA	93 \pm 1.47%	91 \pm 1.65%	93 \pm 1.47%	92 \pm 1.56%	92.5 \pm 1.52%
NBC	77 \pm 2.42%	71 \pm 2.61%	83 \pm 2.16%	85 \pm 2.06%	80.2 \pm 2.30%
SVM	92 \pm 1.56%	90 \pm 1.73%	92 \pm 1.56%	93 \pm 1.47%	92 \pm 1.56%
KNN	86 \pm 2.00%	87 \pm 1.94%	85 \pm 2.06%	89 \pm 1.80%	87.5 \pm 1.90%

Table 2. Classification performance of machine learning algorithms for functional primitives. Positive predictive values (PPV) with associated 95% confidence intervals are shown. PPV reflects how often a primitive was actually made when the algorithm identified it as such, was calculated for the primitives of *reach*, *transport*, *reposition*, and *idle*. Primitive-level PPVs were computed in one-versus-all analysis (e.g., *reach* vs. *transport* + *reposition* + *idle* combined). The overall PPV was assessed by combining data for all primitives and tallying all true and false positives. Overall classification performance was highest for linear discriminant analysis (LDA) and support vector machine (SVM), moderately high for k-nearest neighbors (KNN), and lowest for Naïve Bayes classifier (NBC).

626
 627
 628
 629
 630

Algorithm	# tuning parameters	Tuning parameters	Level of domain knowledge
LDA	3	Prior probability, regularization term, optimizer	Medium
NBC	1	selection of prior distribution	Low
SVM	9	Kernel function, kernel parameters (scale, offset), regularization term, # of iterations, Nu, prior probability, convergence parameter, optimizer	High
KNN	5	# of neighbors (K), distance metric, search algorithm, tie breaker, weighing criterion	Low

Table 3. Complexity of algorithm implementation. Algorithm parameter tuning is necessary to achieve optimal classification performance. Shown are algorithm tuning characteristics, as indicated by number and specifics of the tuning parameters. Also shown is a graded estimate of the level of domain knowledge required to tune these parameters. NBC is considered the simplest to tune while SVM is the most difficult. LDA has a handful of parameters that require medium domain knowledge to negotiate. KNN has a moderate number of parameters that are intuitive to tune and require little domain knowledge. Level of domain knowledge: low, basic knowledge of statistics; medium, undergraduate-level knowledge of ML; high, graduate-level knowledge of ML.

631
632

Primitives	Classification performance (PPV)	
	IMU	Accelerometer
Reach	93%	77%
Transport	91%	80%
Reposition	93%	82%
Idle	92%	88%
Average	92.5%	82%

Table 4. Primitive-level classification using IMU or accelerometer data. Classification performance is shown using the 7-sensor configuration (pelvis, sternum, head, and the active shoulder, upper arm, forearm, and hand). Accelerometers provided systematically poorer classification performance compared to IMUs across all primitives. Classification performance using accelerometry data was particularly low for *reach* (PPV 77%) and relatively higher for *idle* (PPV 88%).

633
634
635