Developing a Typology for Publicly Available Activity Sensing Data

Fiona Draxler, Yannik Peters, Vanessa Lux







MASS Workshop 2025 5th of June, 2025

Working with Sensor Data...





- What sensor datasets are out there already?
 - For secondary analyses
 - For testing algorithms
 - ...
- How "good" are they?
 - Extrinsic quality, e.g., variable documentation and clarity about licenses
 - Intrinsic quality, e.g., signal-to-noise ratio, sampling rates

\rightarrow We collected and analyzed publicly available datasets on activity.

Example Dataset: WISDM (#21)





- 51 participants performing 18 different activities
- Watch + smartphone (different models)
- Accelerometer + gyroscope
- Developed for a classification scenario

Donated on 10/5/2019			CITE
Contains accelerometer and gyros subjects perform 18 activities for 3	st 99 0 citations 9 11808 views		
Dataset Characteristics	Subject Area	Associated Tasks	Creators
Multivariate, Time-Series	Computer Science	Classification	💄 Gary Weiss
Feature Type	# Instances	# Features	
Real	15630426	6	DOI
			10.24432/C5HK59
Dataset Information			^ Liconso
Additional Information			License
For a detailed description of the da	ataset, please see the following pdf file	that is stored with the data: WISDM-dataset	Commons Attribution 4.0 International (CC
description.pdf. The raw accelero rate of 20Hz. It is collected from 5	meter and gyroscope sensor data is co test subjects as they perform 18 active	ollected from the smartphone and smartwatch wities for 3 minutes apiece. The sensor data for	at a BY 4.0) license.
each device (phone, watch) and ty	pe of sensor (accelerometer, gyroscor	be) is stored in a different directory (so there a	This allows for the sharing and adaptation of
data directories). In each directory	/ there are 51 files corresponding to the	e 51 test subjects. The format of every entry is	the datasets for any purpose, provided that the appropriate credit is given.
same: <subject-id, activity-code,="" t<="" th=""><th>time stamp, x, y, z>. The descriptions of</th><th>f these attributes are provided with the attribu</th><th>ite</th></subject-id,>	time stamp, x, y, z>. The descriptions of	f these attributes are provided with the attribu	ite
10-second window. See the datase	et description document for details. Alt	though this data can most naturally be used for	ing a ir
activity recognition, it can also be	used to build behavioral biometric mod	dels since each sensor reading is associated w	/ith a
specific subject.			
SHOW LESS A			
Has Missing Values?			
No			
Variable Information			^
subject-id: value from 1600- 1650	that identifies one of the 51 test subje	cts	
activity-code: character between	'A' and 'S' (no 'N') that identifies the ac	tivity. The mapping from code to activity is	
provided in the activity_key.txt file	and in our dataset description docume	ent	
SHOW MORE V			
Dataset Files			~
Eile		Siza	
rite		017.6	

Extrinsic Data Quality







Core questions:

- What are the data?
- What can/may I do with the data?
- First insights into usability

MASS Workshop 2025 5th of June, 2025

Intrinsic Data Quality







Core aspects:

- Completeness
- Correctness
- Consistency

Across sensors, participants, (datasets ← later)

Outlook: Across-Dataset Classification



GESIS Leibniz Institute for the Social Sciences

- Features derived from sensor axes in datasets 5 and 12

 → Activity prediction (driving, running, sitting, sleeping, standing, walking)
- Random forest model

true										
		driving	running	sitting	sleeping	standing	walking			
predicted	driving	8470	0	144	31	0	22			
	running	0	1338	0	0	0	6			
	sitting	0	1	5354	3	10	44			
	sleeping	0	0	8	5147	0	0			
	standing	0	0	17	0	3964	0			
	walking	6	56	61	0	70	9128			

As researchers interested in sensors...

- How can we make the dataset collection (even) more useful to you?
- What data quality aspects are particularly important?
- What are your favorite datasets that we should include?





Contact us!

Fiona.Draxler@uni-mannheim.de

Vanessa.Lux@gesis.org