# Dataset Description

January 24, 2019

This section outlines characteristics of nine datasets made use of our paper titled "Self Construction of Multi-layer Perceptron Networks in Lifelong Environments". The nine datasets are detailed as follows:

- *SUSY Problem*: SUSY dataset [Baldi, Sadowski, and WhitesonBaldi et al.2014] is a very popular dataset in the area of big data with five million samples. It is commonly used for classification task with two classes problem. The first eight features represent the kinematic properties of signal process whereas the remaining features are generated from the function of the first eight features. This dataset is not categorized as a nonstationary dataset. However, due to the big size of the dataset, this dataset enables the algorithm to demonstrate the ability to handle the lifelong learning environments.

- *Hepmass Problem*: Like Susy, Hepmass dataset is also a prominent in the big data area. It features a two-classes classification task which aims to separate particle-producing collisions from a background source. The dataset consists of 28 input attributes, where the first 22nd features are low-level features. The next 5th features are high-level features and the last feature is a mass feature. There are 10500000 samples in total contained in the dataset. In our experiment, we utilize only 2 million samples (around 19 percent of the total data) with 2000 time stamps.

- *RLCPS Problem*: the Record Linkage Comparison Patterns (RLCPS) problems consists of five million samples and twelve input features which represent patient particular information (e.g. first name, last name, etc). This dataset is taken from epidemicological cancer registry of the German state at North Rhine-Westphalia. This dataset was collected from the period of 2005 to 2008 and aims to determine whether the underlying record belong to one person [Sariyar, Borg, and PommereningSariyar et al.2011]. In RLCPS, the imbalance data problem becomes the main dataset's characteristics where there are only 20931 out of 5 millions samples are categorized as matches (specific records belong to specific person), whereas the rest (most) of the records are the opposite (unmatches). In this experiement, we conduct the experiment under test-then-train prequential procedure with 5000 time stamps.

- *Indoor RFID Localization Problem*: the RFID localization is the four classes classification problem which aims to identify the object's location in the manufacturing shopfloor utilizing the signal from the RFID tag. From this tag, there are three features extracted from the signals captured by the RFID readers by placing the readers in different four zone manufacturing locations. Each location indicates the class of the sample. The dataset comprises 281.3 K data samples. The experiment is conducted under the test-then-train prequential procedure using 280 time stamps.

- *Rotated MNIST Problem*: The rotated MNIST [Lopez-Paz and RanzatoLopez-Paz and Ranzato2017] is a popular continual learning problem developed from the original MNIST problem [LeCun and CortesLeCun and Cortes2010]. It applies rotation of original MNIST problem with randomly generated angels in between $-\pi$ to $\pi$ inducing abrupt drift. To realize challenging simulation environment, overall data samples are grouped into 65 tasks containing 1000 samples and simulated under the prequential test-then-train protocol examining the generalization performance without feeding a model with relevant samples. Unlike benchmark setting, this procedure allows to obscure the task's boundary and time points of concept changes.

- *Permuted MNIST Problem*: The permutted MNIST [Srivastava, Masci, Kazerounian, Gomez, and SchmidhuberSrivastava et al.2013] is another benchmark problem of CL problem derived from the original MNIST dataset [LeCun and CortesLeCun and Cortes2010]. Random permuttation is implemented in the input pixels resulting uncorrelated input distributions across different data concepts. As with the rotated MNIST problem, it is formed as 70 tasks containing 1000 samples. Three random permutations are realized and lead to three different concepts with unknown task boundaries - each task may be drawn from the same or different concepts. Our simulation is carried out by following the prequential test-then-train simulation protocol.

- *KDDCup Problem*: KDDCup dataset [Stolfo, Fan, Lee, Prodromidis, and ChanStolfo et al.2000] introduces the intrusion detection problem in the form of two classes classification tasks. The task aims to recognize whether the network connection is under the attack condition or not. This problem presents the non-stationary components due to the existence of the various type of intrusions held in military network environment. This dataset was very popular dataset as it was used for many machine learning competitions (e.g. Third International Knowledge Discovery and Data Mining Tools Competition). In this experiment, we use only 10% of the total data (500K) for our numerical study with 100 time stamps under the test-then-train procedure.

- *SEA Problem*: The SEA is an artificial dataset [Ditzler and PolikarDitzler and Polikar2013] which features two classes classification problem. It

2

contains three input features, where the first two features are relevant features and the third feature acts as a noise. A sample is classified as a class 1 if the following condition $f_1 + f_2 < \theta$ is satisfied. Otherwise, the sample is classified as a class 2. The changing of the class threshold three times $\theta = 4 \longrightarrow 7 \longrightarrow 4 \longrightarrow 7$ triggers the two types of drift : abrupt and recurring. The original SEA dataset features have values between zero to ten. However, in this case, we utilize the modified version of SEA problem used by [Elwell and PolikarElwell and Polikar2011] which has a class imbalanced problem properties with 5% to 25% class proportion. The total sample of this dataset is 200K. The SEA dataset is very crucial in developing a controlled simulation environment where the type of drift and the time instant when the concept drift occurs are known. These two classes output is determined based on $d$-dimensional random hyperplane $\sum_{j=1}^{d} w_j x_j > w_o$. For this dataset, we conduct the experiment under prequential test-then-train process with 100 time stamps.

- *Hyperplane Problem*: Hyperplane dataset is an artificial dataset generated by the massive online analysis (MOA) - a popular framework in the data stream field [Bifet and GavaldàBifet and Gavaldà2007]. It is categorized as a binary classification problem which aims to separate data points into two classes based on to the position of $d$- dimensional random hyperplane $\sum_{j=1}^{d} w_j x_j > w_o$. This dataset has a gradual drift characteristic where the data samples are generated from one distribution with a probability of one. Then, this probability is reduced gradually until the second distribution completely replaces the first one. The hyperplane dataset consists of 120K sample and we conduct 120 time stamps for the experiment under prequential test-then-train procedure.

# References

[Baldi, Sadowski, and WhitesonBaldi et al.2014] P. Baldi, Paul D. Sadowski, and Daniel Whiteson. 2014. Searching for exotic particles in high-energy physics with deep learning. *Nature communications* 5 (2014), 4308.

[Bifet and GavaldàBifet and Gavaldà2007] Albert Bifet and Ricard Gavaldà. 2007. Learning from time-changing data with adaptive windowing. In *In SIAM International Conference on Data Mining*.

[Ditzler and PolikarDitzler and Polikar2013] G. Ditzler and R. Polikar. 2013. Incremental Learning of Concept Drift from Streaming Imbalanced Data. *IEEE Trans. on Knowl. and Data Eng.* 25, 10 (Oct. 2013), 2283–2301.

[Elwell and PolikarElwell and Polikar2011] R. Elwell and R. Polikar. 2011. Incremental Learning of Concept Drift in Nonstationary Environments. *Trans. Neur. Netw.* 22, 10 (Oct. 2011), 1517–1531.

[LeCun and CortesLeCun and Cortes2010] Yann LeCun and Corinna Cortes. 2010. MNIST handwritten digit database. http://yann.lecun.com/exdb/mnist/. (2010). `http://yann.lecun.com/exdb/mnist/`

[Lopez-Paz and RanzatoLopez-Paz and Ranzato2017] David Lopez-Paz and Marc' Aurelio Ranzato. 2017. Gradient Episodic Memory for Continual Learning. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 6467–6476. `http://papers.nips.cc/paper/7225-gradient-episodic-memory-for-continual-learning.pdf`

[Sariyar, Borg, and PommereningSariyar et al.2011] Murat Sariyar, Andreas Borg, and Klaus Pommerening. 2011. Controlling False Match Rates in Record Linkage Using Extreme Value Theory. *J. of Biomedical Informatics* 44, 4 (Aug. 2011), 648–654. `https://doi.org/10.1016/j.jbi.2011.02.008`

[Srivastava, Masci, Kazerounian, Gomez, and SchmidhuberSrivastava et al.2013] Rupesh K Srivastava, Jonathan Masci, Sohrob Kazerounian, Faustino Gomez, and Jürgen Schmidhuber. 2013. Compete to Compute. In *Advances in Neural Information Processing Systems 26*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Eds.). Curran Associates, Inc., 2310–2318. `http://papers.nips.cc/paper/5059-compete-to-compute.pdf`

[Stolfo, Fan, Lee, Prodromidis, and ChanStolfo et al.2000] Salvatore J. Stolfo, Wei Fan, Wenke Lee, Andreas Prodromidis, and Philip K. Chan. 2000. Cost-based Modeling for Fraud and Intrusion Detection: Results from the JAM Project. In *In Proceedings of the 2000 DARPA Information Survivability Conference and Exposition*. IEEE Computer Press, 130–144.