Automated Data Quality Monitoring with ADWIN2

Abdullah Farhat January 13, 2021

Automated data-quality monitoring and calibrations

"In most challenging data analysis applications, data evolve over time and must be analyzed in near real time. Patterns and relations in such data often evolve over time, thus, models built for analyzing such data quickly become obsolete over time. In machine learning and data mining this phenomenon is referred to as **concept drift**." [1]

To deal with time-changing data, one needs strategies, at least, for the following:

- 1. detecting when a change occurs
- 2. determining which examples to keep and which to drop
- 3. updating models when significant change is detected

Automated data-quality monitoring and calibrations

OUR APPROACH

1. Identify different data-taking periods Use ADWIN to identify the start of distinct data-taking periods based on changes in the mean of the data stream.

2. Calibrate different data-taking periods to a baseline Use Hoeffding's inequality to estimate the mean of each data-taking period and apply a constant shift to each data taking period by the difference between the means of a baseline period and each subsequent period.

3

ADWIN Algorithm

- ADWIN is an ADaptive WINdowing technique used for detecting distribution changes, concept drift, or anomalies in data streams with established guarantees on the rates of false positives and false negatives [2].
- ADWIN Inputs:
 - confidence value $\delta \in (0,1)$
 - data stream { $x_1, x_2, ..., x_t, ...$ } where each x_t is available at time t drawn from some distribution with expected value μ_t
- ADWIN keeps a sliding window W with the most recently read x_i

MAIN IDEA: whenever two sufficiently large subwindows of W have sufficiently different means, then it is likely the corresponding expected values are different, and the older portion of the window is dropped.

• Moreover, the window size is expected to stay large while μ_t remains constant in W, and becomes small when μ_t changes

ADWIN Algorithm

Partion *W* into subwindows W_0 and W_1 .

Let $|W_0| = n_0$, $|W_1| = n_1$, and |W| = n.

Define:

$$m = \frac{1}{1/n_0 + 1/n_1}$$

$$\delta' = \frac{\delta}{n}$$

$$\epsilon_{cut} = \sqrt{\frac{1}{2m} \ln \frac{4}{\delta'}}$$

The probability for both false positive and false negative is at most δ .

ADWIN: ADAPTIVE WINDOWING ALGORITHM

$$\begin{bmatrix}
 W_0 & W_1 \\
 \overline{x_i, x_{i+1}, \dots, x_{i+n_0}, x_{i+n_0+1}, \dots, x_{i+n}} \\
 W
 \end{bmatrix}$$

To represent the data stream we use a sample of 120,000 Inclusive Deep Inelastic Scattering Monte Carlo events

- generated in the context of the ZEUS experiments
- Includes full detector simulation
- Reconstructed kinematics with all detector effects.

We observe a stream of x and Q^2 , reconstructed by the electron method [3] based on the measurement of the (x, y, z) position and energy E of the outgoing lepton in the calorimeter.

We subdivide the stream into 3 data-taking periods of equal parts and apply a constant or gradual shift of two standard deviations to each (x, y, z) position and energy *E* measurements in the second data taking period.







An example data stream (gradual change)



A higher-dimensional extension of ADWIN improves its ability to find changes in the data distribution.

Two cases:

- 1D: only use information from Q^2
- 2D: use information from (x, Q^2)



 After using ADWIN2 to detect different data-taking periods, each period is calibrated to the baseline period.

• The simple calibration we use is to shift each period by a constant value to force its mean to be equal to the baseline mean

Hoeffding's Inequality:

If $X_1, X_2, ..., X_n$ are independent random variables bounded between [0,1] drawn from the same distribution with expected value μ , and define \overline{X} to be the sample mean, then for any t > 0, $\mathbb{P}(\overline{X} - \mu > t) < e^{-2nt^2}$.

Consequently, to estimate the mean of a distribution with $(1-\alpha)$ %-confidence and a margin of error of t, we need at least n observations, where:

$$n = \frac{\log(2/\alpha)}{2t^2}$$

For a confidence level $\alpha = 0.01$ and a margin of error of t = 0.01:

a minimum sample of 26492 observations is needed to estimate of the mean in each data-taking period.





Corrected data stream

An example data stream (gradual change)





Corrected data stream



ADWIN Algorithm

ADVANTAGES:

- Fast algorithm
- No prior assumption on the underlying distribution of data samples
- No a priori determination of a fixed window size
- Can easily be extended to higher-dimensional anomaly detection
- Does not require any training on simulated data sets

DISADVANTAGES:

- Need to store a large window size when the data stream distribution is stable
- Uses only the mean to characterize changes

QUESTIONS?

References

- Žliobaitė I., Pechenizkiy M., Gama J. (2016) An Overview of Concept Drift Applications. In: Japkowicz N., Stefanowski J. (eds) Big Data Analysis: New Algorithms for a New Society. Studies in Big Data, vol 16. Springer, Cham. <u>https://doi.org/10.1007/978-3-319-26989-4_4</u>
- A. Bifet and R. Gavalda, Learning from time-changing data with adaptive windowing, in Proceedings of the 2007 SIAM international conference on data mining, SIAM, 2007, pp. 443– 448.
- 3. S. Bentvelsen, J. Engelen and P. Kooijman, Reconstruction of (x, Q2) and extraction of structure functions in neutral current scattering at HERA, in Workshop on Physics at HERA Hamburg, Germany, October 29-30, 1991, 1992, pp. 23–42.