Rapid automatic clean-up toolkit for large corrupted tidal datasets

Vamsi Krishna Sridharan¹

¹University of California Santa Cruz

November 26, 2022

Abstract

Tides are critical to coastal and oceanic processes. While tidal data are available readily, they are often corrupted by various sources of error. An automated, fast MATLAB toolbox is developed to clean-up tidal timeseries data from estuarine and oceanic locations corrupted by errors. This toolbox will immensely speed up delivery of quality-controlled tidal data. It will also reduce errors in quality control, which typically involves several manual tasks. The toolbox corrects poorly interpolated and noisy data, erroneous outliers, and instrumentation bias such as spurious jumps, drifts, spikes, and modulations in the true signal. Signal clean-up involves multiple stages. First, thresholds are imposed on higher order temporal derivatives of the signal to remove gross interpolations and noise saturated signal chunks, followed by a moving median threshold to remove outliers. Then the surviving signal is filtered into tidal, subtidal and long-period components, and the long-period component is subject to a maximal overlap discrete wavelet transformation, in which the transform coefficients corresponding to multi-scale edge features are removed. Subsequently, local information in the subtidal and tidal components is compared relative to the whole signal to correct spurious amplitude modulations and sudden biases. Consequently, these components are added to recover the uncorrupted signal, and large data gaps are filled with short term harmonic reconstruction. For estuarine locations, the correlation in the spectrogram between two nearby stations is initially used to quantify and remove river influence in the signal. Applications to datasets at multiple global locations demonstrate the value of the toolbox.

Rapid automatic clean-up toolkit for large corrupted tidal datasets

Vamsi Krishna Sridharan, University of California, Santa Cruz, Affiliate: Southwest Fisheries Science Center 110 McAllister Way, Santa Cruz, CA 95062 | Email: vamsikrishna.sridharan@ucsc.edu

1. Introduction

- Tides are ubiquitous in estuaries, coasts and oceans around the world
- Tidal harmonic signals are observed in all hydrodynamic and water quality data sets
- More that 30,000 data sets from multiple sources are prone to error and need to be cleaned up to produce useful data
- Although QA/QC protocols exist, they are ad-hoc, quantityspecific, and not easily automatable
- We present a robust, quantity-independent workflow that produces QA/QCed data rapidly for multi-spectral signals

Figure 1.

Typical tide gauge. Image source: Google Earth Pro. Inset image source: https://tidesandcurrents.noaa.gov/ stationhome.html?id=9414290

San Francisco

station 9414290

3. Use of Machine Learning for optimal parameter selection

• For long signals, optimal breakpoints are identified by optimizing ratio of computational cost to short term harmonic analysis efficacy using gradient descent learning



- Thresholding and outlier detection parameters are optimized using a cost function and gradient descent learning
- Other operations could potentially be optimized as using machine learning as well

Microwave tide gauge

2A. Salient features

- Developed in MATLAB
- The toolbox corrects
- Blocky interpolation and noisy data
- o erroneous outliers
- instrumentation bias such as spurious jumps and drifts
- narrow- and broadspectrum spikes
- \circ modulations in the true signal

2B. Methodology

Data is censored to the instrument's operational range ➤Thresholding on signal's temporal derivatives to remove gross errors, followed by moving median threshold to remove outliers. Surviving signal is

- decomposed
- Long-period component is subject to maximal



Method is robust, can deal with different quantities, multispectral signals, and can be improved using machine learning.



Figure 2. Workflow for signal clean-up from all locations

overlap discrete wavelet transformation to remove multi-scale edges > Local information in the subtidal and tidal components is compared relative to the whole signal to correct spurious amplitude modulations and sudden biases > Signal components are added to recover the uncorrupted signal, and large data gaps are

- filled with short term harmonic reconstruction
- > For estuarine locations, the correlation in the spectrogram between two nearby stations is used first to quantify and remove river influence in the signal.

component of the target signal, f.

Figure 4. Results: (A) Estuarine locations, (B) Correction of signal drift at Bournemouth, England, (C) Correction of jump and outliers in salinity, (D) Robustness to good data

5. References

Cho, H.Y., Oh, J.H., Kim, K.O., Shim, J.S., 2013. Cornish, C.R., Percival, D.B., Bretherton, C.S., 2003. Daubechies, I. 1992. Donoho, D.L., Johnstone, J.M. 1994. EuroGOOS DATA-MEQ working group. 2010. Flinchem, E.P., Jay, D.A. 2000. Hastie, T., Tibshirani, R., Friedman, J. 2009. Hoitink, A.J.F., Jay, D.A. 2016. Integrated Ocean Observing System. 2016. Kukulka, T., Jay, D.A., 2003. The Mathworks Inc. 2016. Pawlowicz, R., Beardsley, B. Lentz, S. 2002. Woodworth, P.L., Player, R. 2002.