



UvA-DARE (Digital Academic Repository)

Data-driven methods in application to flood defence systems monitoring and analysis

Pyayt, A.L.

Publication date
2014

[Link to publication](#)

Citation for published version (APA):

Pyayt, A. L. (2014). *Data-driven methods in application to flood defence systems monitoring and analysis*. [Thesis, fully internal, Universiteit van Amsterdam].

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

Appendix C. Description of the Methods

C.1. FFT, STFT and phase shift

One of the possible time-frequency features to be used for levee behaviour monitoring is the phase difference between oscillating signals of different sensors. This feature can be extracted from any monitored system with some periodic behaviour, such as electrical signals, vibrations or more-complex engineering, microeconomic and socio-dynamic systems. A sudden change in frequencies, amplitudes or phase shifts indicates a potential problem in the system.

In case of levee health monitoring, tidal changes in water levels propagate through the soil inside the dike and cause periodic increases in the water pressure at the sensors (see Figure C.1). Because soil that fills the levee is a porous material with a relatively low permeability, water flow experiences resistance and reaches the sensors with some delay. Areas further from the sea exhibit longer time delays (this effect is called "phase shift" in signal analysis).

If the levee is stable (*i.e.*, the structure is not damaged and the soil layers are not eroded), then the "resistivity" of the porous levee remains constant. Consequently, the phase shift between the sensors also stays constant. A change in the phase difference indicates that the levee integrity might be corrupted. Moreover, it also identifies the exact location of a problem, *i.e.*, between the two sensors that revealed an anomalous phase shift.

Similar to the water pressure dynamics, other sensors also respond to the dynamically changing hydraulic forces caused by tides. Therefore, our methodology can also be applied to sensors that measuring inclination, displacement, stress, strain, and other levee health parameters.

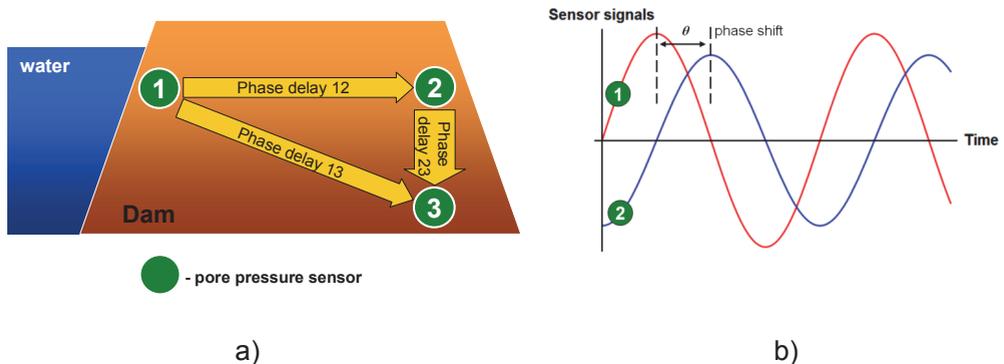


Figure C.1. Levee health monitoring based on pairwise phase shift monitoring. (a) A schematic of the monitored levee (filled with permeable soil) and sensor positions. (b) Illustration of the phase shift concept: a pressure wave caused by sea tides is reaching sensor #2 with a time delay relative to sensor #1, which is located close to the river. This time lag is called a "phase shift" or "phase delay" in signal analysis.

We consider a phase shift between the Fourier transform components calculated for a selected frequency as the time delay metric.

The short-time Fourier transform (STFT) is the correct method to represent the signal in both the time and frequency domains [88]. This property facilitates detection of anomalies by tracking phase changes over time.

Time-frequency representation by STFT is performed using a discrete fast Fourier transform (FFT) algorithm in a sliding window. Each new sliding window overlaps with a previous window to reduce boundary effects. The STFT coefficients have a time delay at each frequency. The STFT for a discrete time series is given by

$$X(n, \omega) = \sum_{m=-\infty}^{\infty} x[m]w[n-m]e^{-j\omega m} \quad (\text{C.1})$$

where $x[m]$ is the analysed signal; w is the window function (e.g., a rectangular window, Hamming, Gaussian), and $x[m]w[n-m]$ is the short-time section of the signal $x[m]$ at time n .

The phase ϕ at timestamp n of frequency ω can be calculated as an argument of each STFT component [118]:

$$\phi(n, \omega) = \arg(X(n, \omega)) \quad (\text{C.2})$$

The phase delay τ_ϕ is then calculated as:

$$\tau_\phi = -\frac{\phi(n, \omega)}{\omega} \quad (\text{C.3})$$

The minus in the formula (C.3) is required for the positive presentation of the delay in the time domain.

The phase shift $\Delta\phi$ (in radians) and time delay $\Delta\tau_\phi(n, \omega)$ in seconds between two selected sensors is calculated as:

$$\Delta\phi(n, \omega) = \phi_1(n, \omega) - \phi_2(n, \omega) \quad (\text{C.4})$$

$$\Delta\tau_\phi(n, \omega) = -\frac{\Delta\phi(n, \omega)}{\omega} \quad (\text{C.5})$$

where $\phi_1(n, \omega)$ is the phase of sensor #1 and $\phi_2(n, \omega)$ is the phase of sensor #2.

A pairwise phase shift analysis can be extended to the analysis for sensor triplets (Figure C.1(a)). The phase delay within the triangle should not change in time. Application of each individual sensor in a phase-shift analysis of several pairs guarantees redundancy. Thus, detection and localization of the anomaly are possible. A phase delay is calculated from the time-frequency components that are related to two sensors; therefore, we classify this feature as a ‘‘spatial time-frequency feature’’.

C.2. Maximum overlap discrete wavelet transform (MODWT)

The maximum overlap discrete wavelet transform (MODWT) technique is a method for time-frequency representation of time series. Unlike CWT, MODWT does not require high computational costs.

We chose the maximum overlap discrete wavelet transform (MODWT) [93] for feature selection. MODWT is a computationally efficient method for time-frequency representation of time series. The MODWT transform is similar to the discrete wavelet transform (DWT), but it does not produce downsampling of wavelet coefficients [93], which allows it to overcome the lack of translation-invariance present in DWT and does not require the length of the signal to be a power of two. In contrast with CWT, MODWT calculates coefficients at scales 2^j (where j is the level of the transform) without loss of information. This property provides faster computation of MODWT coefficients than CWT computation. The procedure of MODWT coefficient calculation can be described as an application of linear filters (wavelet and scaling filters) via a cascade algorithm (see Figure C.2).

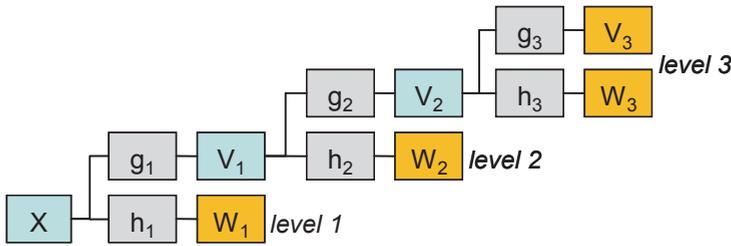


Figure C.2. Maximum overlap discrete wavelet transform (MODWT) cascade algorithm. X is the analysed signal; g_i and h_i are the scaling (low-pass) and wavelet (high-pass) filters, respectively; and v_i and w_i are the approximation and detail MODWT coefficients of the i -th level of decomposition, respectively.

The frequency bands of the MODWT levels are illustrated in Figure C.3.

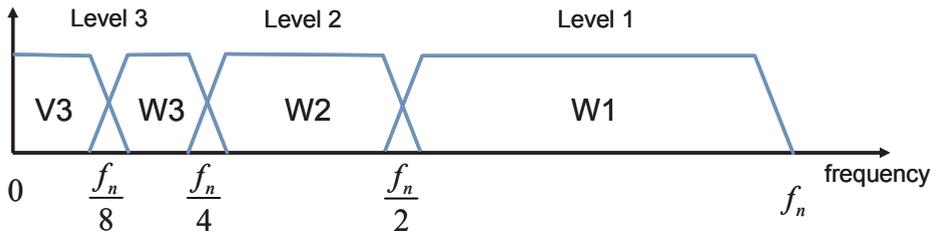


Figure C.3. Frequency-domain representation of the MODWT levels. f_n is the upper limit of the frequency range of the analysed signal.

It is easier to work with the equivalent MODWT filters, which are analogous to the wavelet and scaling functions, because the MODWT filters calculate coefficients by direct convolution with a signal [93]:

$$V_{j,t} = \sum_{l=0}^{L_j-1} g_{j,l} \cdot X_{t-l} \quad (C.6)$$

where X is the analysed signal; j is the level of decomposition; $V_{j,t}$ is the vector of approximation coefficients at level j ; g_j is the MODWT equivalent scaling filter; and L_j is defined by Equation (C.8).

$$W_{j,t} = \sum_{l=0}^{L_j-1} h_{j,l} \cdot X_{t-l} \quad (C.7)$$

where X is the analysed signal; j is the level of decomposition; $W_{j,t}$ is the vector of wavelet coefficients at level j ; h_j is the MODWT equivalent wavelet filter; and L_j is defined by Equation (C.8).

The length L of the equivalent MODWT filters for level j can be calculated using the following equation:

$$L_j = (2^j - 1)(L - 1) + 1 \quad (C.8)$$

C.3. Daily-seasonal-annual (DSA) transform

The daily-seasonal-annual (DSA) transform is based on the MODWT transform. This transform assumes a combination of MODWT levels in daily, seasonal (monthly) and annual components (Figure 3.15).

The DSA transform allows one to correlate deviations of dam parameters with natural fluctuations, which have base frequencies at daily, seasonal (monthly) and annual frequencies (e.g., air temperature). The DSA transform represents a MODWT as a set of physically motivated components.

For computation of the DSA transform, time delays in MODWT levels must be compensated. This compensation can be performed by implementing MODWT wavelet and scaling filters for reconstruction of each wavelet level in the following manner [94]:

$$D_j = (W_j * \tilde{h}_j * \tilde{g}_j * \dots * \tilde{g}_1), \quad S_{J_0} = (V_{J_0} * \tilde{h}_j * \tilde{g}_j * \dots * \tilde{g}_1)$$

$$x = \sum_{j=1}^{J_0} D_j + S_{J_0} \quad (C.9)$$

where the variables are as follows:

x – time series;

D_j is the “detail” coefficient of level j ;

S_{J_0} indicates smoothing at the last level J_0 ;

V_{J_0} is the last approximation level of the MODWT;

\tilde{g}_j is the MODWT scaling filter for reconstruction; and

\tilde{h}_j is the MODWT wavelet filter for reconstruction.

Combination of “details”, which correspond to the bands of frequencies, produces a DSA transform [94]. Suppose that j_d is the level for which the frequency band includes the frequency $\frac{1}{24}$ (1/hours). Then, the daily component is

$$Daily = \sum_{j=1}^{j_d} D_j \tag{C.10}$$

If j_s is the level for which the frequency band includes the frequency $\frac{1}{720}$ (1/hours), then the seasonal (monthly) component is

$$Seasonal = \sum_{j=j_d+1}^{j_s} D_j \tag{C.11}$$

Summation of all other detail levels and smoothing at the last level produces an annual component:

$$Annual = \sum_{j=j_m+1}^{J_0} D_j + S_{J_0} \tag{C.12}$$

C.4. Universal threshold

Donoho and Johnstone describe a "universal threshold" λ [70]:

$$\lambda = \hat{\sigma} \sqrt{2 \log n} \tag{C.13}$$

$$\hat{\sigma} = 1.4826 * median \left[|d^{j-1} - median(d^{j-1})| \right]$$

where d^{j-1} is the vector of the finest wavelet coefficients of the wavelet transform and $\hat{\sigma}$ is the MAD estimate.

C.5. Polynomial autoregressive model

The autoregressive (AR) linear model is defined as follows [36]:

$$y(t) + a_1 y(t-1) + \dots + a_{n_a} y(t-n_a) = b u(t-n_k) + \dots + b_{n_b} u(t-n_b-n_k+1) + e(t) \tag{C.14}$$

where $y(t)$ is the output at time t ; a_i and b_i ($i = 1:n$) are the parameter of the model that must be estimated; n_a is the number of poles in the system; n_b is the number of zeroes in the system; n_k is the number of input samples that occur before the inputs affect the output current; $y(t-1)$ and $y(t-n_a)$ are the previous outputs on which the output current depends; $u(t-n_k)$ and $u(t-n_b-n_k+1)$ are the previous inputs on which the current output depends; and $e(t)$ is the white-noise. The n_a and n_b parameters are referred to as model orders.

The AR model can be written in a compact way using the following notation:

$$A(q)y(t) = B(q)u(t) + e(t) \tag{C.15}$$

where:

$$A(q) = 1 + a_1 q^{-1} + \dots + a_{n_a} q^{-n_a} \quad (C.16)$$

$$B(q) = b_1 q^{-n_k} + \dots + b_{n_b} q^{-n_b - n_k + 1}$$

q^{-1} is the backward shift operator, which is defined as

$$q^{-1}u(t) = u(t-1) \quad (C.17)$$

The following block diagram shows the AR model structure (Figure C.4).

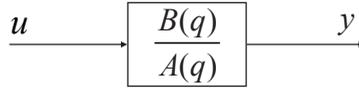


Figure C.4. Autoregressive block diagram, where u is the model input, y is the model output, and $B(q)/A(q)$ is a transfer function.

The *training* procedure of the AR model is based on finding the best order (n_a, n_b). The stopping criterion corresponds to the root mean square error (RMSE) in Equation (C.18), where RMSE is defined in Equation (C.19).

$$\sigma_e \leq \text{threshold} \quad (C.18)$$

$$\sigma_e = \left(\frac{1}{T} \sum_{t=1}^T (y(t) - y'(t))^2 \right)^{1/2} \quad (C.19)$$

where y is the model output and y' is the expected output. The RMSE is 0 when $y(t) = y'(t)$. In addition, the RMSE must be less than the threshold that is selected by the user.

Another approach for the best model selection is to use the Akaike Information Criterion (AIC) [23] as follows:

$$AIC = \ln \sigma^2 + \frac{2d}{N} \quad (C.20)$$

where σ^2 is the variance of the model error (mean squared error), d is the model length, and N is the length of the training set. The minimal value of AIC corresponds to the best model.

C.6. Artificial neural networks

An artificial neural network is a non-linear approximation of an input signal to output time series. In the current study, we use the feed-forward neural network (FFNN). The structure of FFNN is presented in Figure C.5.

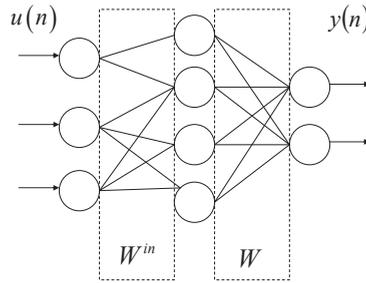


Figure C.5. The structure of the three-layered (input, hidden and output layers) feed forward neural network that is used for the static nonlinear transformation of the input $u(n)$ values into the output values $y(n)$. W^{in} and W represent the matrices of the connections (weights) between the layers.

The NN training is conducted by optimizing the weights of the network so that the input signal for the NN produces the right output signal. The inputs of the NN propagate through the input, hidden and output layers. The following transfer function is used: $f(x)=\tanh(u)$.

The gradient descent is used as the optimization algorithm. The learning error function is defined as follows:

$$E = \frac{1}{2} \sum_{t=1}^T (y_t - d_t)^2 \quad (\text{C.21})$$

where y_t is the model output, d_t is the target output, and T is the number of model and target instances.

For more information regarding FFNN, please refer to [55].

The following types of models are considered and approximated using FFNN:

$$y_t = F(u_{t-k}, \dots, u_t, y_{t-m}, \dots, y_{t-1}), \quad (\text{C.22})$$

where F is a nonlinear function, u is the input signal, y is the model output signal, k is the time-delay for the input signal, and m is time delay for the output signal that was used as an input value in the model.

C.7. Model quality assessment

The quality of modelling (for the models described in Sections C.5 and C.6) can be characterized using the following quality measures described below.

1) The coefficient of determination R^2 is calculated as follows:

$$R^2 = 1 - \frac{\sum_{t=1}^T (y_t - d_t)^2}{\sum_{t=1}^T (y_t - \frac{1}{T} \sum_{t=1}^T y_t)^2} \quad (\text{C.23})$$

where y_t is the model output, d_t is the target output, and T is the number of the model and target instances. This metric lies in the range of $[1, -\infty)$, where 1 corresponds to the best model fit.

2) Root-mean-square error, see Equation (C.19).