



Evaluation of the satellite-based Global Flood Detection System for measuring river discharge: influence of local factors

B. Revilla-Romero^{1,2}, J. Thielen¹, P. Salamon¹, T. De Groeve¹, and G. R. Brakenridge³

¹European Commission Joint Research Centre, Ispra, Italy

²Utrecht University, Faculty of Geosciences, Utrecht, the Netherlands

³University of Colorado Boulder, Boulder, USA

Correspondence to: B. Revilla-Romero (beatriz.revilla-romero@jrc.ec.europa.eu)

Received: 4 June 2014 – Published in Hydrol. Earth Syst. Sci. Discuss.: 3 July 2014

Revised: 29 September 2014 – Accepted: 30 September 2014 – Published: 7 November 2014

Abstract. One of the main challenges for global hydrological modelling is the limited availability of observational data for calibration and model verification. This is particularly the case for real-time applications. This problem could potentially be overcome if discharge measurements based on satellite data were sufficiently accurate to substitute for ground-based measurements. The aim of this study is to test the potentials and constraints of the remote sensing signal of the Global Flood Detection System for converting the flood detection signal into river discharge values.

The study uses data for 322 river measurement locations in Africa, Asia, Europe, North America and South America. Satellite discharge measurements were calibrated for these sites and a validation analysis with in situ discharge was performed. The locations with very good performance will be used in a future project where satellite discharge measurements are obtained on a daily basis to fill the gaps where real-time ground observations are not available. These include several international river locations in Africa: the Niger, Volta and Zambezi rivers.

Analysis of the potential factors affecting the satellite signal was based on a classification decision tree (random forest) and showed that mean discharge, climatic region, land cover and upstream catchment area are the dominant variables which determine good or poor performance of the measurement sites. In general terms, higher skill scores were obtained for locations with one or more of the following characteristics: a river width higher than 1 km; a large floodplain area and in flooded forest, a potential flooded area greater than 40%; sparse vegetation, croplands or grasslands and closed to open and open forest; leaf area index > 2; tropi-

cal climatic area; and without hydraulic infrastructures. Also, locations where river ice cover is seasonally present obtained higher skill scores. This work provides guidance on the best locations and limitations for estimating discharge values from these daily satellite signals.

1 Introduction

Flooding is the most prevalent natural hazard at the global scale, often with dire humanitarian and economic effects. According to the International Disaster Database (EM-DAT), an average of 175 flood events per year occurred globally between 2002 and 2011, affecting an average of 116.5 million people, and causing economic losses of USD 25.5 billion. According to MunichRe (2014), the costliest natural catastrophe worldwide in terms of overall economic losses in 2013 was the flooding in southern and eastern Germany and neighbouring states in May and June, with estimated damages of USD 15.2 billion. In June of the same year, flooding in India claimed 5000 lives, with a further 2 million affected (MunichRe, 2014; EM-DAT).

The Global Assessment Report (UNISDR, 2011) states that the proportion of world population living in flood-prone river basins increased by 114% over four decades from 1970 to 2010. Additionally, while economic losses due to river floods have increased over the last 50 years, the number of casualties has decreased. The reduction in loss of life has been associated with the integration of early warning systems with emergency preparedness and planning at local and

national levels (Golnaraghi et al., 2009; Kundzewicz et al., 2012).

Global early warning systems are needed to improve international disaster management. These systems can be used for both early forecasting (for better preparedness) and early detection, as well as for an effective response and crisis management. Their necessity was emphasised in 2005, and since then it has been a key element of international initiatives such as the “Hyogo Framework for Action 2005–2015” and, on a continental level, the European Commission Flood Action Programme. After the 2002 flooding of the Elbe and Danube rivers, the European Commission supported the development of the European Flood Awareness System (EFAS) (Bartholmes et al., 2009; Thielen et al., 2009) by the Joint Research Centre to increase preparedness for riverine floods across Europe. Currently, a number of organisations are involved in rapid mapping activities after major (flood) disasters, such as UNOSAT (2013), GDACS (2013), “Space and Major Disasters” (Disaster Charter, 2014), the Committee on Earth Observation Satellites (CEOS) Flood Pilot and the on-line Dartmouth Flood Observatory (<http://floodobservatory.colorado.edu/>). In Europe, Copernicus is the Earth observation programme which actively supports the use of satellite technology in disaster management and early warning systems for improved emergency management.

Flood warning systems typically rely on forecasts from national meteorological services and in situ observations from hydrological gauging stations. However, this capacity is not equally developed across the globe, and is highly limited in flood-prone, developing countries. Ground-based hydro-meteorological observations are often either scarce or, in cases of transboundary rivers, data sharing among the riparian nations can be limited or absent. Therefore, satellite monitoring systems and global flood forecasting systems are a needed alternative source of information for national flood authorities not in the position to build up an adequate measuring network and early warning system. In recent years, there has been a notable development in the monitoring of floods using satellite remote sensing and meteorological and hydrological modelling (Schumann et al., 2009).

A variety of satellite-based monitoring systems measure characteristics of the Earth’s surface, including terrestrial surface water, over large areas on a regular basis (van Westen, 2013). Such remote sensing is based on surface electromagnetic reflectance or radiance in the optical, infrared and microwave bands. Some key advantages of microwave sensors is that they provide near-daily basis global coverage and, at selected frequencies, relatively little interference from cloud cover. Two presently operating microwave remote sensors with near-global coverage are the Tropical Rainfall Measuring Mission¹ (TRMM), operational from 1998 to present, and the Advanced Microwave Scan-

ning Radiometer for Earth Observation System² (AMSR-E) which was active from June 2002 to October 2011, succeeded by AMSR2, which was launched in May 2012 and is onboard the Japanese satellite GCOM-W1³, and from which brightness temperature data are being distributed from January 2013 onwards. For future work, the European Space Agency (ESA) and NASA have other missions to put similar instruments in orbit, capturing passive microwave energy at 36.5 GHz, such as ESA’s Sentinel-3 satellites (planned launch in 2015 and 2016) and NASA’s Global Precipitation Mission (GPM) (launched in February 2014) to replace TRMM.

Using AMSR-E data initially, De Groeve et al. (2006) implemented a method for detecting major floods on a global scale, based on the surface water extent measured using passive microwave sensing. Also, Brakenridge et al. (2005, 2007) demonstrated that orbital remote sensing can be used to monitor river discharge changes. However, as underlined by Brakenridge et al. (2012, 2013), extracting the microwave signal and converting it into discharge measurements is not straightforward and depends on factors such as sensor calibration characteristics and perturbation of the signal by land surface changes. These changes can be found, for example, in irrigated agricultural zones and in areas where rivers flow along forested floodplains (Brakenridge et al., 2013). As rivers discharge increases, river level (stage), river width, and river flow velocity all increase as well, and the challenge is to measure one or more of these accurately enough to provide a reliable discharge estimator, and compare against a background of other surface changes that may affect what is measured from orbit.

There also remains the need to convert such discharge estimators to actual discharge units. Using ground discharge data or climate-driven runoff models for calibration and validation, methods to convert the remote sensing signal to river discharge have been previously tested at particular stations with output from the Global Flood Detection System (GFDS, <http://www.gdacs.org/flooddetection/>) and by different investigators (Brakenridge et al., 2007, 2012; Khan et al., 2012; Kugler and De Groeve, 2007; Moffitt et al., 2011; Hirpa et al., 2013; Zhang et al., 2013). Yet the results are from different approaches and not easily comparable, making an assessment of the potential performance on a global scale difficult. Furthermore, definite conclusions about the influence of various environmental factors on the signal performance have not been reached. Therefore, in this study, a rigorous broad assessment of the method is undertaken with a systematic evaluation of the relationship between skills obtained between ground- and satellite-based discharges and the local characteristics of the stations. Specifically, this study addresses mean observed discharges, river widths, land cover

¹<http://trmm.gsfc.nasa.gov>.

²http://aqua.nasa.gov/about/instrument_amsr.php.

³http://suzaku.eorc.jaxa.jp/GCOM_W/w_amsr2/whats_amsr2.html.

types, leaf area indices, climatic regions and flood hazard maps, as well as the presence or absence of large floodplains, wetlands, river ice and hydraulic control infrastructure.

Our goal is to assess the potentials and limitations of the satellite-based surface water extent signal data for river discharge measurements with a large number of stations. Moreover, the relationship between ground and satellite sets of discharge measurements and the local surface characteristics is examined in order to provide guidelines for selection of observation sites. For this purpose, river catchments located in a range of different climatic and land cover types were selected in Africa, Asia, Europe, North America and South America. The remainder of the paper is structured as follows: Sect. 2 presents the study regions and data, Sect. 3 describes the analysis methodologies, and the results are discussed in Sect. 4.

2 Study regions and data

2.1 Study regions and in situ discharge data

Figure 1 shows the study basins and in situ discharge locations. The selected stations are all located near major rivers of the world (Global Runoff Data Centre, 2007). The continental distribution and the upstream catchment area of the stations are summarised in Table 1. We selected the locations to be representative of a broad variety of local conditions: they belong to nine different main land cover classes (aggregated from GlobCover, 2009) and five main types of climate (Peel et al., 2007). The characteristics are listed in Table 2.

For Africa, Asia, Europe, North America and South America, daily in situ discharge values were used from the Global Runoff Data Centre (GRDC) database. In addition, for the South African stations, the discharge data were provided by the South African Water Affairs (DWA, <http://www.dwa.gov.za/>). The selected stations for all these continents include daily data between 1998 and 2010; however not all stations have continuous data during this time period. From 1998, the length of the time series was required to be above 6 years. The longest time series available was of 13 years, with a median value of 8.5 years. In situ discharge information may itself be affected by large and variable uncertainty, mostly on the measurement of the cross-sectional area of the channel and mean flow velocity at the gauge or control site (Pelletier, 1988). Although generally unknown, these values are typically between 5 and 20 % at the 95 % confidence level as highlighted in studies such as Hirsch and Costa (2004), Di Baldassarre and Montanari (2009), Le Coz et al. (2014) and Tomkins (2014). However, the uncertainty in river discharge is even higher during floods events when the stage–discharge relationship, the so-called rating curve, is used. As evaluated by Pappenberger et al. (2006), the analysis of rating curve uncertainties leads to an uncertainty of the input of 18–25 % at peak discharge. Di Baldassarre and Mon-

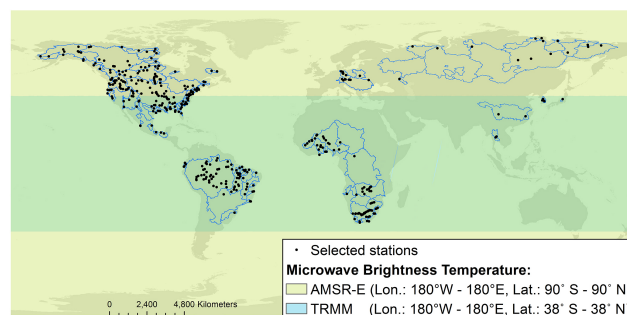


Figure 1. Location of selected stations (398) and corresponding river basins (109). TRMM and AMSR-E brightness temperature product extents are also provided.

tanari (2009) showed that the total rating curve errors increase when the river discharge increases and varies from 1.8 to 38.4 % with a mean value of 21.2 %. For the purposes here, these data are, however, regarded as “ground truth”. We acknowledge the possible errors, however, and note that, for some river reaches, satellite-based methods may actually track discharge changes more accurately than ground-based measurements using stage; however, the extent to which this is true needs to be fully investigated.

2.2 Satellite-derived data

The Global Flood Detection System (GFDS) produces near-real-time maps and alerts for major floods using satellite-based passive microwave observations of surface water extent and floodplains. It was developed and is maintained at the European Commission Joint Research Centre (JRC) in collaboration with the Dartmouth Flood Observatory (DFO). The surface water extent detection methodology using satellite-based microwave data is explained in Brakenridge et al. (2007) and Kugler and De Groeve (2007). Here, only the basic principles are recalled.

At each pixel, the method uses the difference in brightness temperature, at a frequency of 36.5 GHz, between water and land surface to detect the proportion of within-pixel water and land. The retrieved brightness temperature data are first gridded into a product with a pixel size of (near the Equator) $10 \text{ km} \times 10 \text{ km}$ ($0.09^\circ \times 0.09^\circ$), and the system provides a daily output. For our work, the merged TRMM/AMRS-E product was used (<http://www.gdacs.org/flooddetection/download.aspx>); the gridded data are provided in the GCS WGS 1984 projection. For our period of study, 1998–2010, the merged data product was employed for the time period of its availability (June 2002–2010), whereas stand-alone TRMM data was used for the remaining time period (1998 to June 2002) and available latitudes. Note that from 2013 the system has been providing the merged product TRMM/AMSR2; however, this period is out of our scope.

Table 1. Number of catchments by continent and range of upstream areas for the located stations.

Continent	Number of satellite locations forextraction ($n = 398$)	Number of stations for calibration ($n = 322$)	Number of catchments ^a	Upstream catchment areas, (km ^b) approx. range
Africa	75	51	21	46 990–850 500 ^b
Asia	23	3	4	7150–11 000
Europe	13	7	3	9000–132 000
North America	207	183	86	5300–1 850 000
South America	80	78	38	1400–4 680 000

^a Stations used for calibration and validation. ^b South African upstream catchment areas are not available.

Table 2. Climate and land cover type of the 322 sites selected for the calibration and validation, aggregated by continent, climate and land cover.

Climate	Africa	Asia	Europe	North America	South America	Total
Arid	30			25		55
Tropical	10				75	85
Temperate	11		3	51	3	68
Cold		3	4	104		111
Polar ^b				3		3
Total	51	3	7	183	78	322
Land cover	Africa	Asia	Europe	North America	South America	Total
Open forest	4			23		27
Closed to open forest	16	1	1	16	41	75
Closed forest				33		33
Mosaic vegetation predominant ^a	19	2		47	24	92
Mosaic cropland or grassland predominant	5		1	26	9	41
Rainfed crop			4	5	4	13
Sparse vegetation	2			14		16
Sparse vegetation + crops	5			8		13
Urban			1	10		11
Bare areas ^b				1		1
Total	51	3	7	183	78	322

^a Vegetation means a combination of grassland, shrubland and forest. ^b Types of land cover and climate where the number of locations in each type was very low (e.g. 3) were excluded for their respective variables analysis as they will not be representative on a global scale.

In the GFDS system, the microwave signal (s) is defined as the ratio between the measurement over wet pixel (M) and the measurement over a 7 pixel \times 7 pixel array of background calibration (C) pixel, known as the M / C ratio (Brakenridge et al., 2012; De Groeve, 2010). Better discharge signal values may be achieved when the measurement pixel is centred over a river reach and no hydraulic structures are present (Moffitt et al., 2011). However, this is sometimes difficult to achieve due to the desired co-location with gauging stations (Brakenridge et al., 2012) or because the potential measurement pixels within the raster are fixed geographically.

2.3 Other important data sets and maps

The quality of the microwave signal detected by the satellite sensors can be influenced by local ground conditions, including extreme rainfall, snow/ice, land cover/use and topography (Brakenridge et al., 2012). For example, forest is a type of land cover which influences the microwave emission properties due to the biometric features of vegetation such as crown water content and the shape and size of leaves (Chukhlantsev, 2006). In this study, the effects of the local ground conditions on the performance of the satellite signal were analysed as a function of the following factors:

1. *River width*: channel width from Yamazaki et al. (2014), estimation based on SRTM Water Body Database and the HydroSHEDS flow direction map and for which the map was upscaled from 0.025 to 0.1°, taking the mean of the river grid values in the 4 × 4 area.
2. *Mean observed discharge*: for each station, a mean discharge value for the study period was calculated from daily ground data (mainly from the GRDC data set).
3. *Upstream catchment area* (GRDC 2007) data: the GRDC river network was used to visually select those stations located close to the “main rivers” classified by GRDC, and to use the values of the upstream catchment area for each station. Note that upstream catchment area values are missing from all South African stations from DWA data provider.
4. *Presence of floodplains, flooded forest and wetlands*: this was obtained from the Global Lakes and Wetlands Database level 3, a global raster map at 30 s resolution which comprises lakes, reservoirs, rivers and different wetland types (Lehner and Döll, 2004).
5. *Flood extent*: we used the fractional coverage of potential flooding of 25 km by 25 km cells for a 100-year return period from the Global Flood Hazard Map derived using a model grid (HTESSEL + CaMa-Flood) (Papenberger et al., 2012).
6. *Land cover*: we used land cover data from the Global Land Cover 2009 (Bontemps et al., 2010). The 19 labels were aggregated into 8 types of land cover depending on the vegetation type and density to synthesise the outputs (see Appendix Table A1). Further visual category checking was performed using Google Maps display for the sites, and, where necessary, land cover classes changed accordingly. An additional category was added for sparse vegetation areas where crops are grown along or near the river channels.
7. *Leaf area index*: a global reprocessed leaf area index (LAI) from SPOT-VGT is available for a period of 1999–2007 (http://wdc.dlr.de/data_products/SURFACE/LAI/). This LAI product is a global data set of 36 ten-day composites at a spatial resolution of the CYCLOPES products (1 km). For our analysis, a modified version of this product was used, which was up-scaled to a spatial resolution of 10 km.
8. *Climatic areas*: we used the Köppen–Geiger climate map of the world (Peel et al., 2007) to distinguish the main climate areas: tropical, arid, temperate, cold and polar (see Table 2).
9. *Presence of river ice*: through the signal, the presence of river ice cover can also be detected in cold land regions. The Circum-Arctic map of permafrost and ground-ice

conditions (Brown et al., 2002) map was used here. Examples of these rivers are the Yukon and Mackenzie rivers in North America and the Lena River in Russia. As is the case on the ground, discharge under ice cover is left largely unmeasured as both water area and stage no longer are responsive to discharge variation.

10. *Dam location*: hydraulic structures can disrupt the natural flow of water, and therefore may alter the expected performance of the satellite signal on that location. For this analysis the Global Reservoir and Dam (GRanD) (Lehner et al., 2008) data set was used.

3 Methodology

3.1 Satellite signal extraction

In total, 398 locations for satellite-based measurement were selected which overlap spatially and temporally with available in situ stations providing daily measurements. Since satellites never pass directly over the same track at exactly the same time, the operational GFDS applies a 4-day forward-running mean to systematically calculate the signal; this also commonly fills between any missing days (Kugler and De Groeve, 2007). Furthermore, for each observation site, the signal on the GFDS system is calculated as the average signal of all measurement pixels under observation for each location (which can be one or more pixels) (GDACS, 2013). Thus, in some cases, even a 10 km pixel is not large enough as a measurement site, and would entirely saturate with water during flooding. An array of measurement pixels is instead used. In this analysis, we used the signal values from the single pixels which contain the ground station, as well as a multiple pixels selection. This includes, for each location, the pixel itself and also the three nearest neighbours of the 10 km × 10 km grid. In the case of multiple pixels, the signal value was calculated for the spatial median, average and maxima. Similar results were obtained globally when comparing the extracted signals (single or multiple pixels) with the in situ discharge observations. Therefore, we used the temporal and spatial averaging on the multiple pixel array as in the operational GFDS. For each site, a visual check with Google Maps was carried out to assure that the largest river section was included within the finalised measurement sites (see Fig. 2).

3.2 Satellite signal calibration and validation

For those co-located ground stations and satellite measurement sites where both sets of data (signal and in situ discharge) were above 6 years in length, calibration and validation was performed using the ground information as reference. Several stations, mainly in North America, located close to man-made infrastructures such as weirs and generating stations were excluded from this analysis due to

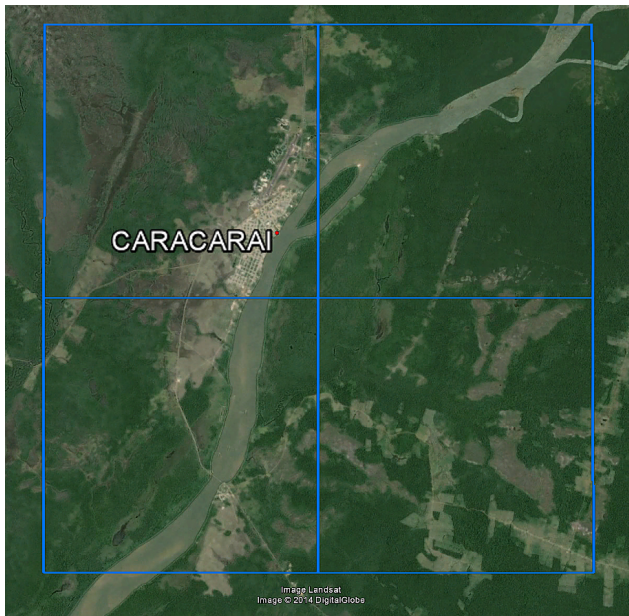


Figure 2. Example of a measurement site: Caracarai station (Rio Branco catchment, Brazil). The blue rectangles outline the measurement pixels and the background image is from 2014 (Google; Landsat, DigitalGlobe).

the rapidly changing behaviour of the in situ-observed discharge. Also, in a satellite-based approach to measure river discharge, the local river characteristics and floodplain channel geometry control the accuracy of rating curves, as is the case for gauging stations on the ground (Brakenridge et al., 2012; Khan et al., 2012; Moffitt et al., 2011). Thus we expect some measurement sites to exhibit a more robust response to discharge changes, and a higher signal-to-noise ratio, than others.

It has been acknowledged that, for large rivers, using the daily GFDS signal as a floodplain flow surface area indicator of discharge might result in a few days of lag when comparing with ground-based discharge (Brakenridge et al., 2013). Thus, stage may immediately rise at a gauging station as a flood wave approaches, but flow expansion out into the floodplain requires some increment of time. This time lag may introduce error into the scatterplots used to calculate the rating equations and therefore lower skill scores obtained when analysing both data sets. In addition, in previous studies (Khan et al., 2012; Zhang et al., 2013), it was observed that, in some cases, an overestimation of satellite-measured discharge existed during low-flow periods when using a single rating equation for the full period to calibrate signal into discharge units. For this reason, we decided to use a rating equation for each month individually. In this case the time series data for a fixed month can be treated as stationary and the derived daily discharge values also adjusted better during low-flow periods.

To calibrate satellite signal into discharge measurements, the first 5 years of data were used for both satellite signal and ground discharge for each location. Regression equations were obtained using monthly means from daily values and GFDS-measured discharge was derived from this.

$$Q_{\text{GFDSmeasured of X month}} = a_{\text{month}} + b_{\text{month}} \cdot \text{signal} \quad (1)$$

For the sake of simplicity, for this paper, the equations were restricted to linear equations. However, as the relation is purely empirical, we leave further research into a flexible way to fit these relations as follow-on work. Note that fitting straight lines to curves will reduce goodness of fit and predictive accuracy. Power law fitting was also tested to calibrate the signal into discharge units, yielding similar results (see open discussion author's response no. 2).

The validation of the satellite-derived daily discharge data was carried out with daily in situ data on a 2-year period, and skills scores were calculated to quantify the agreement between both satellite- and ground-measured discharge. We are aware of the limited number of years (data) with available time series for both variables, which might influence the robustness of the calibration. In some cases there were longer time series available, but, in order to standardise the analysis for all the stations, we used 5 years (1998–2002 or 2003–2008 for northern stations with AMSR-E signal) and the following 2 years for validation purposes (2003–2004 and 2009–2010, respectively). Note that, for 36 out of the 322 stations available, data length was between 6 years and 3 months to almost 7 years. Validation was still carried out for the same period, but the data used for calibration were slightly reduced. As an example, Fig. 3a presents the scatterplot for the month of March for the Senanga station (long 23.25 degree, lat –16.116 degree) in the Zambezi River (Africa) with mean values derived from the period 1998 to 2002. For the same location, Fig. 3b shows the in situ-observed and the GFDS-measured discharge derived from the GFDS signal for the period 2003–2004.

3.3 Skill scores

The initial analysis of the correlation of the remote sensing signal to in situ discharge was assessed for each station and site pair through the Pearson correlation coefficient (R). For the validation, the performance of the satellite-measured discharge was also assessed using the Nash–Sutcliffe efficiency (NSE) statistic in addition to the R skill score. Spearman's rank correlation coefficient (ρ) was also calculated to assess the validation performance.

One of the advantages of the R coefficient is that it is independent of the units of measurement, which permits the comparison of dimensionless GFDS signal data. A small value indicates a weak or non-linear relationship between the satellite signal and discharge. For this study, we grouped the computed R values into three ranges as follows: < 0.3 , $[0.3-0.7]$, and > 0.7 . While Pearson benchmarks linear relationship,

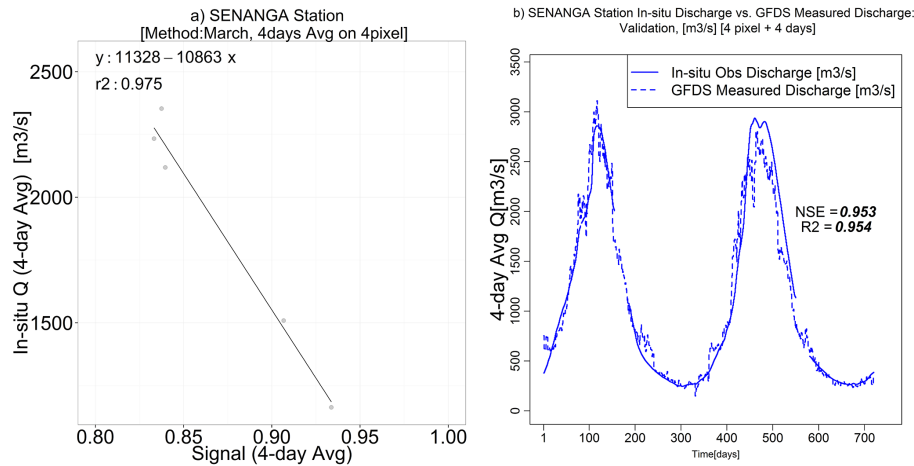


Figure 3. (a) Scatterplot for the Senanga station (long 23.25 degree, lat -16.116 degree) in the Zambezi River (Africa). Monthly mean for March from 1998 up to 2002. (b) Validation hydrograph for 2003–2004 and skill scores for Senanga. The (monthly) rating equations were used to calibrate the signal into discharge units. Different rating equations were used for different months.

Spearman benchmarks monotonic relationship. Spearman's validation scores just obtained a mean value 6 % higher than Pearson mean score (see open discussion author's response no. 2). In this manuscript, results are analysed based on the scores obtained using Pearson correlation coefficient.

Nash–Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970) is typically used to assess the predictive power of hydrological models and was calculated here to describe the accuracy of satellite-derived discharge in comparison to gauge-observed discharge values. Higher values of the Nash–Sutcliffe statistic should indicate more correlated results, without other factors taken into account, such as autocorrelation (Brakenridge et al., 2012). However, the degree of correlation of these variables does not verify the discharge magnitudes (Brakenridge et al., 2013). An NSE value of 1 corresponds to a perfect match of modelled to observed data, whereas $NSE=0$ indicates that the model predictions are as accurate as the mean of the observed data. The resulting scores will be classified as in Zaraj et al. (2013): < 0 , $[0.2-0.5]$, $[0.5-0.75]$, and > 0.75 .

3.4 Factors affecting the satellite signal

Understanding the influence of local factors on the accuracy of the satellite flood detection is critical for practical use of the remotely sensed signal. We analysed the accuracy effects of river width, mean daily discharge, upstream catchment area, presence of large floodplain, flooded forest and wetlands, potential flood extent, land cover type, LAI, climatic areas, presence of river ice and hydraulic structures. To assess their influence, the fractional coverage over the measurement site was retrieved for variables with spatial coverage.

First, we use the skill scores (R and NSE) obtained from a simple analysis for each individual factor or variable. Second, we seek to understand which of the surface variables

have the greatest importance in determining sites with a good or poor performance. For this purpose, we use a decision tree technique called random forest (RF). Among other features, this allows for ranking of the relative importance of each variable. The technique is described by Breiman (2001) and implemented in *R* by Liaw and Wiener (2002), where the reader is referred for a more detailed explanation. As a summary of the RF algorithm, *n* bootstrap samples are randomly selected from the data set; a different subset is used for each bootstrap; and for each sample a tree is grown, obtaining *n* trees. RF is called an ensemble method because it applies the method for a number of decision trees, in this case 500, in order to improve the classification rate. Some stations are left out of the sample (out of bag – oob) and used to gain an internal unbiased estimate of the generalisation error (oob errors) and to obtain estimates of the importance of the variables (Breiman, 2001). These values are averaged over the *n* trees. For the variables classification, the node impurity is measured by the Gini index. Gini's mean difference was first introduced by Corrado Gini in 1912 as an alternative measure of variability. One of the parameters derived from it, the Gini index, is also referred to as the concentration ratio (Yitzhaki and Schechtman, 2013). The Gini index is mostly popular in economics; however it is also used in other areas, such as building decision trees in statistics to measure the purity of possible child nodes, and it has been compared with other equality measures (Gonzalez et al., 2010). The variables with larger decreases in Gini values (lower Gini) are those with higher importance in the classification analysis.

Although the information is hidden inside the model structure for “black-box models” such as RF, the prediction power is high (Palczewska et al., 2013). This method is relatively robust given outliers and noise because it uses randomly chosen subsets of variables at each split of each tree (Breiman, 2001;

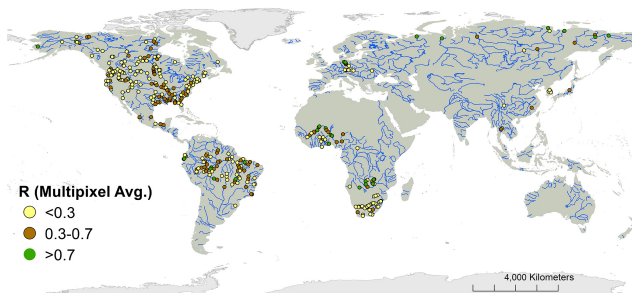


Figure 4. Location of stations and R skill score between in situ-observed discharge and satellite signal (4 days and 4 pixels average). Globally, 169 sites have $R > 0.3$, of which 42 have $R > 0.5$.

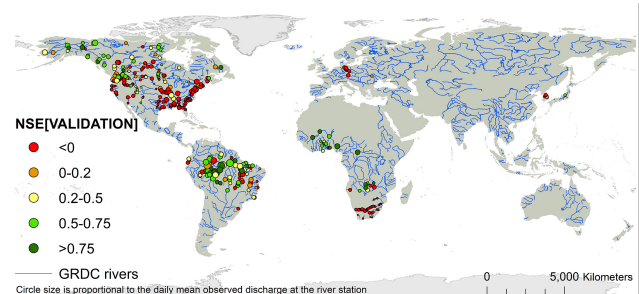


Figure 5. Nash–Sutcliffe efficiency of the validation ($n = 332$ stations). Globally, 154 stations have $NSE > 0$, of which 80 stations have $NSE > 0.50$.

Chan and Paelinckx, 2008). To further increase robustness, Strobl et al. (2009) state that results from the RF and conditional variable importance should always be tested by doing multiple RF runs using different seeds and sufficiently large *n*tree values to obtain robust and stable results.

The quality index chosen to rank variable importance and classify good or poor locations, in the RF analysis, was the NSE score. A threshold of $NSE = 0$ splits the data into two groups, obtaining about 50 % of the data above (true or good predictive) and below (false or poor predictive) that value of NSE. The results presented here are the average of 200 runs. Furthermore, four different training sets were used by a random 70 %/75 %/80 %/90 % of the stations and were validated with the remaining 30 %/25 %/20 %/10 % of stations, respectively.

4 Results and discussion

As a first step we analysed the relationship between the satellite signal and the in situ-observed discharge to have an initial understanding of the performance between the two data sets (Sect. 4.1). Then we calibrated the satellite signal with in situ discharge data. With the regression equations obtained, we calculated satellite discharge measurements. A 2-year validation period was carried out for each station using the skill scores as described in Sect. 3.3 (Sect. 4.2). This was followed by an assessment of how different variables contribute in a positive or negative way to the overall skill (Sect. 4.3). Variables included in the analysis are daily mean river discharge, river width, upstream catchment area, potential flood hazard area, land cover, LAI, climatic zones, presence of large floodplains, flooded forest and wetlands, river ice and hydrologic structure. Finally, the relative importance of all variables in comparison to each other has been assessed (Sect. 4.4).

Before analysing the validation results, it is important to highlight two possible different sources of error which might influence the outputs. Firstly, the signal-to-noise ratio might be low for a site or have intermittent instrument noise occasionally producing positive spikes in discharge. Secondly, the

rating curve may be offset, which will result in a consistent bias on the discharge values for that location even though the time series are strongly correlated.

4.1 Correlation of raw satellite data vs. gauge observations

The first step was to look at the “raw” correlation between daily ground-station-measured water discharge and the satellite signal and to calculate the empirical linear relation between these two variables for each site. The full time series, including low flows, were used for the calculation and executed for 398 stations. Figure 4 shows the R skills obtained. Of a total of 398 sites, 169 have an $R > 0.3$ and 42 of them have $R > 0.5$. Correlations might have perhaps been higher if regression had not been restricted to linear equations (Brakenridge et al., 2007, 2012).

4.2 Satellite signal calibration, validation and evaluation through skill scores

For the stations with over 6 years of contemporary data for both in situ discharge and satellite signal, we obtained regression equations for each month of the year and station using the first 5 years of data. Next, using these equations, we carry out a calibration of the daily signal into discharge units. Afterwards, the validation of the GFDS-measured discharge was implemented for the following 2 years. In some regions, such as northern Asia, the lack of available recent long time series (after 2002) meant that the number of stations available for calibrating the satellite into discharge measurements was reduced. Stations where the number of years matching observed discharge and satellite signal was shorter than 6 years were excluded from the validation exercise despite performing well. Finally, out of 398, a total of 332 stations remained for calibration and validation.

For NSE score, Fig. 5 shows that 154 out of 332 stations are larger than 0: 13 located in Africa, 77 in North America, 62 in South America, 1 in Asia and 1 in Europe. Nevertheless, it needs to be noted that, in arid regions, results calculated with the skill scores such as NSE are penalised by

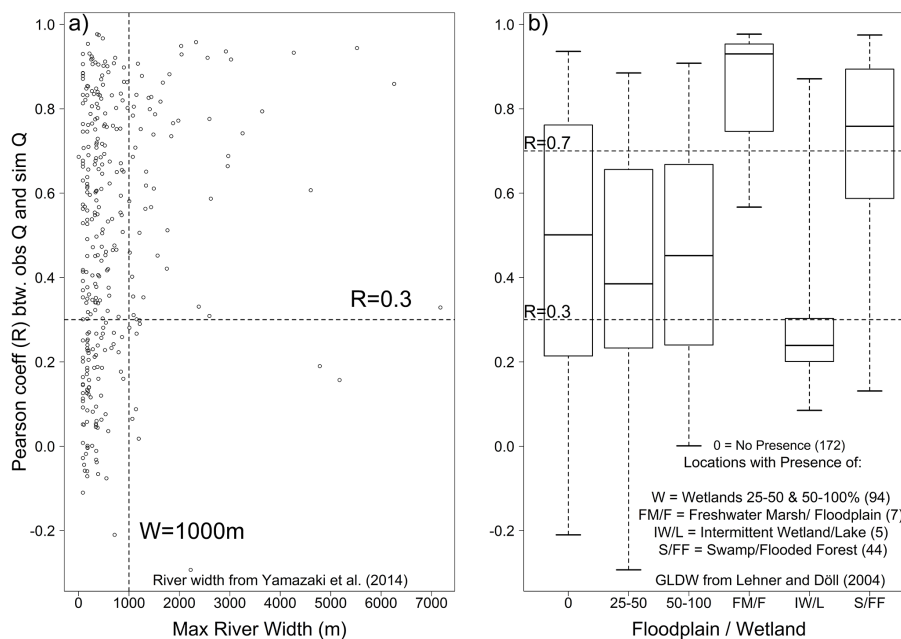


Figure 6. (a) Relationship between R obtained from the validation of satellite-measured discharge and the maximum river width for each location. (b) Relationship between the same R score and the presence of significant floodplains, flooded forest and wetlands. The horizontal dotted line shows the $R = 0.3$ and $R = 0.7$ threshold, and the vertical line is the river width equal to 1 km.

low average discharge compared to high-flow conditions. If, instead of using all the available time series, a “dry stream” threshold had been applied, the scores obtained for these sites could have been higher when analysing the remaining data set period where flow is present.

4.3 Analysis of the factors affecting the satellite signal

4.3.1 River width and presence of floodplain and wetlands

As a first step to analyse the potential relationship between the individual local characteristics and the performance of the locations in global terms, we study the R score of the validation for the 322 stations in relation with the maximum river width value at each location (Fig. 6a). Results indicate that locations with a river width higher than 1 km are more likely to score an R larger than 0.3. In fact, the mean R score is 0.60 and 26 out of 64 ($\sim 41\%$) have $R > 0.75$. However, there are a number of stations with lower river width that also obtained high scores. As the retrieval of the satellite signal also depends on the floodplain geometry. As soon as the river floods and water goes over-bank, the proportion of water in the wet pixel greatly increases. Thus the score should also be high for small rivers with a proportionally large floodplain. Figure 6b shows the R scores by location, where the majority of the area belongs to floodplain, flooded forest and wetlands category, or their absence. In our study, higher median scores were obtained for those located in large freshwater marsh and floodplains, followed by those on swamps and flooded forest.

These results give a first indication on the characteristics of the locations with better performance.

4.3.2 River discharge and potential flooding

Flooding is determined by the discharge as well as the potential flood hazard. Figure 7a shows that 84 out of 95 stations with $R < 0.3$ also have mean discharge values lower than $500 \text{ m}^3 \text{ s}^{-1}$ ($\log_{10}(500) \approx 2.7$), of which 55 stations had a mean discharge lower than $200 \text{ m}^3 \text{ s}^{-1}$. These stations are mainly located in South Africa, and in some areas of North America. Therefore, it can be concluded that the mean discharge can be considered a key variable that determines the appropriateness of locations for which satellite discharges can be derived: as 77% of the stations with $Q < 500 \text{ m}^3 \text{ s}^{-1}$ have $R < 0.3$, while 91.5% of the stations with $Q > 500 \text{ m}^3 \text{ s}^{-1}$ have $R > 0.3$, locations with discharge of less than $500 \text{ m}^3 \text{ s}^{-1}$ might not provide reliable results for a global satellite-based monitoring system. Alternatively, non-permanent rivers and streams exhibiting only seasonal or ephemeral flow (typical for dry regions) may require a different monitoring approach, wherein a “dry” threshold is established for the signal data.

After excluding the global stations with low skill score due to low flows and studying the remaining stations, we can better understand the performance of the system in relation to other local characteristics. Figure 7b shows for each location the relationship between the validation R and the percentage of area in each pixel covered by potential flooding during a 100-year return period flood event, obtained with the

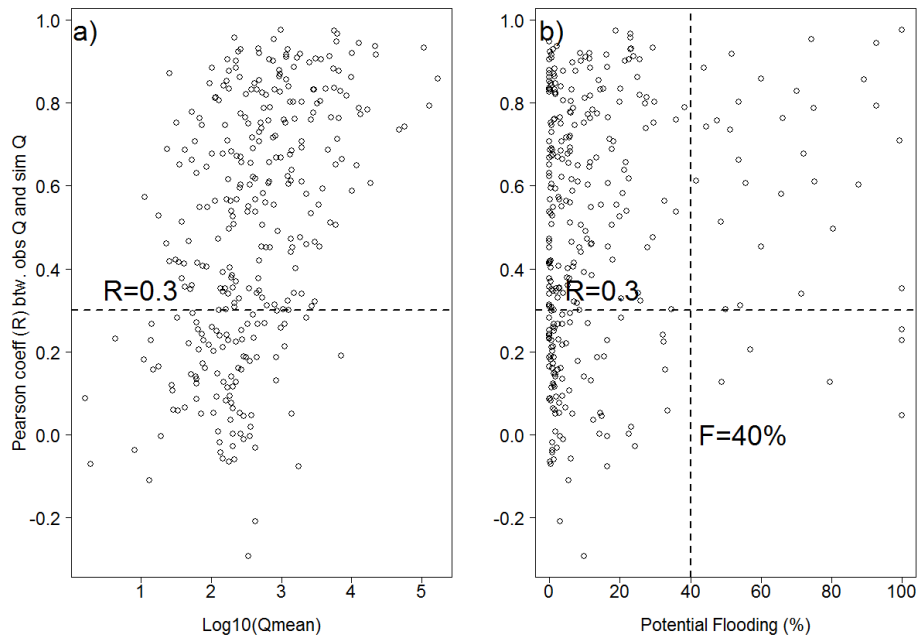


Figure 7. (a) Relationship between R obtained from the validation of satellite-measured discharge and the mean in situ-observed discharge (\log_{10} displayed) for each station. (b) Relationship between the same R score and the potential percentage of flooded area per pixel for a 100-year return period flood event (Pappenberger et al., 2012). The horizontal dotted line shows the $R = 0.3$ threshold, and the vertical line is the 40% potential flooding threshold.

model grid (HTESSEL + CaMa-Flood) (downscaled from a $25 \text{ km} \times 25 \text{ km}$ pixel; Pappenberger et al., 2012). A value of 100 means completely flooded across its area, 50 means 50% of the area within the cells is flooded, and 0 means that the area is not flooded. Although there is no clear trend for all the points, results indicate that locations with a percentage of potential flooding larger than 40% are expected to score an R larger than 0.3.

4.3.3 Land cover types and climatic areas

Figure 8 presents a global evaluation of the R score obtained during the validation and its classification by the land cover type of the stations. The bare land cover category was excluded from this study as only one of the selected locations belongs to that class. Looking at the median of the box plot (see Fig. 8), we found that some of the locations with higher density of vegetation such as those located on “closed forest” and “mosaic with predominant vegetation” (including forest, scrublands and grasslands) obtained lower median scores values. In contrast, the locations with lower vegetation density such as “sparse vegetation”, “mosaics with predominant cropland/grasslands”, “open forest” and “closed to open forest” land cover types obtained larger median R scores, around 0.6–0.8. Similar results can be observed when looking at the interquartile range or spread of the box plots: “closed to open forest” and “mosaics with predominant cropland/grasslands” obtained better results. At the same time, “closed forest” and “mosaic with predominant vegetation”

had lower scores. In addition, those sites with a combination of sparse vegetation and crops growing near the river channel had a lower median value when comparing with those on sparse or mosaic crop land cover. Note that the sites denoted “sparse with crops” are located in arid climatic areas, whereas most of the “sparse” sites are in cold or polar regions and are therefore run by different processes. In addition, sites with a majority of artificial/urban land cover (not shown) obtained a low median value of 0.267.

The relationship between locations by main Köppen–Geiger climatic areas (Peel et al., 2007) and R score obtained is shown in Fig. 9. Globally the tropical regions (Africa and South America) obtained the highest median scores ($R \approx 0.8$), followed by cold regions ($R \approx 0.6$). Lower median score values ($R \approx 0.3$) were obtained for arid and temperate regions. It is important to clarify that these results are not only due to direct climate characteristics but also, for example, due to the characteristics of the rivers in those areas. In the case of the arid regions, it is mainly related to reduced daily average discharges, a characteristic of many of these stations. Note that polar climate was excluded from this evaluation as only three locations belong to that class.

4.3.4 Leaf area index (LAI)

LAI values typically range from 0 for bare ground to 6 or above for a dense forest; however CYCLOPES underestimates over dense vegetation (forest) (Zhu et al., 2013). Therefore, for this product LAI range is limited to [0–4], as

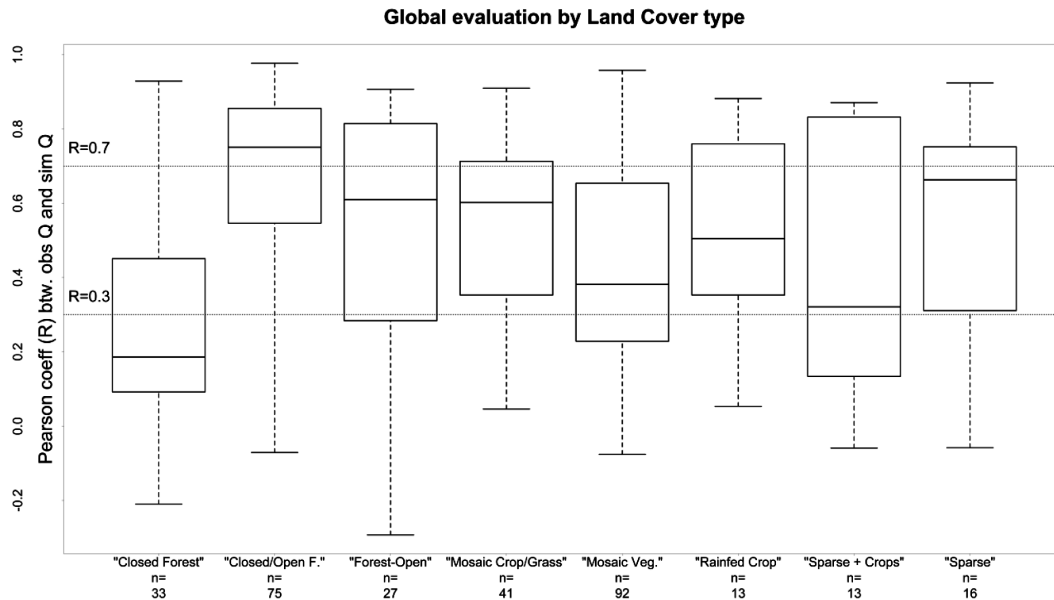


Figure 8. Global evaluation of the R score obtained during the validation and its classification by the land cover type of the stations. Land cover types were aggregated from the GlobCover (2009) and modified by means of a visual check with Google Maps. Note that artificial and bare land cover were excluded in this figure.

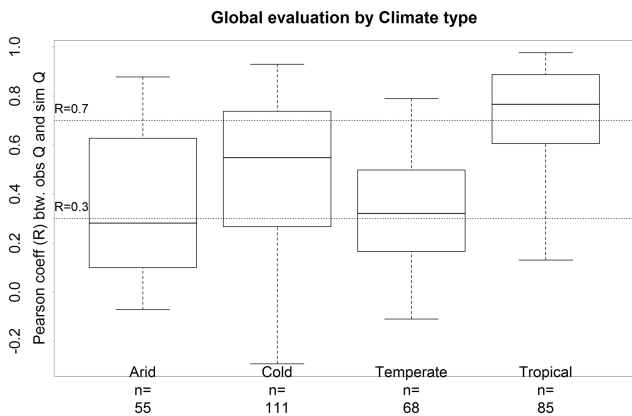


Figure 9. Global evaluation of the R score obtained during the validation and its classification – 2 only main types by the Köppen–Geiger climate area (Peel et al., 2007). Note that polar climate was excluded from this analysis as only three stations fell into this category.

seen in our analysis. Despite this, CYCLOPES is the most similar product to LAI references map (ibid.). According to the study carry out by Zhu et al. (2013) monthly CYCLOPES LAI values for the period 1999 to 2007 by four main groups of vegetation are predominantly as follows: bare ground [0], forest [0–3.5], other woody vegetation [0–1.5], herbaceous vegetation [0–2], and cropland/natural vegetation mosaics [0–3]. The highest annual mean LAI values are obtained by evergreen broadleaf forest (3.16), included in our “closed to open forest” class.

We decided to study the relationship between the mean LAI and the skill obtained in the validation for each location, also looking at complementary variables such as the land cover and the geographical region which the stations belong to. Figure 10 shows that locations with a mean [LAI > 2] predominantly have a “closed to open forest” type in South America (31 stations), of which 29 have an R score higher than 0.6. For [LAI > 2], there are also 12 North American locations with “closed forest” land cover, but in general scores are poorer for those locations. Additionally, 18 stations with mosaic vegetation from North and South America obtained [LAI > 2], and 16 of those obtained [$R > 0.6$]. For [LAI < 2], both the land cover and geographical locations are distributed along the scatterplots, from poor to high correlations.

4.3.5 River ice

Figure 11a shows the scores obtained for the locations with presence or absence of river ice, including a range from continuous to sporadic (Brown et al., 2002). It can be seen that stations located in areas with river ice tend to have a good correlation between in situ- and satellite-measured discharge (based on 33 stations), as the system tends to capture the annual spring ice break-up and freezing well, as indicated in the studies by Brakenridge et al. (2007) and Kugler (2012). At these locations, once ice-covered, the system has no sensing capability and the retrieved signal may seem analogous to low-flow conditions. However, there is an important difference when analysing time series of signal between ice-covered high-latitude river and all-year-round low-flow rivers. When an ice-melting process takes place, an

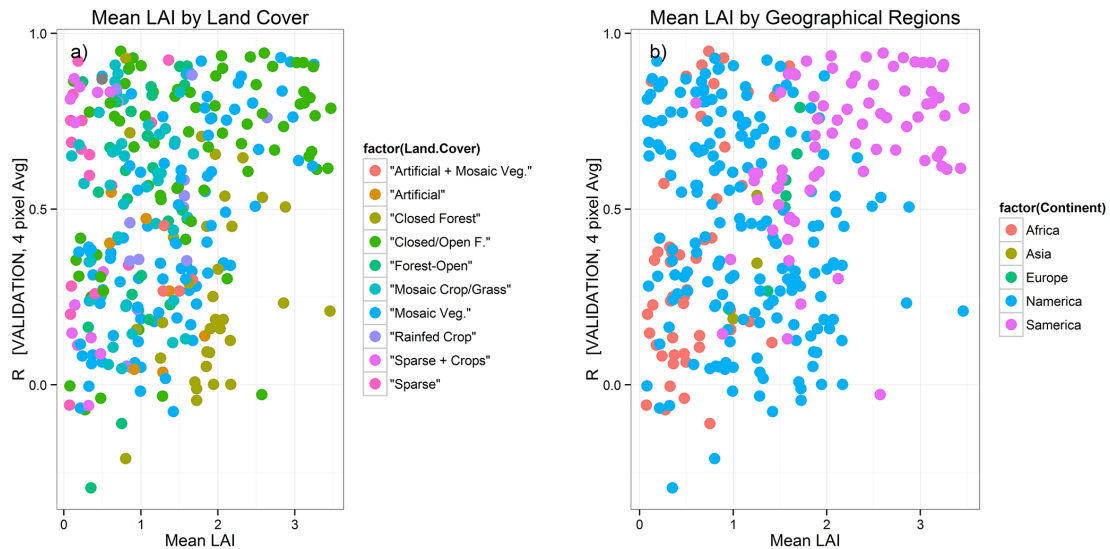


Figure 10. Evaluation of the R score obtained during the validation and its classification by LAI according to factors of (a) land cover and (b) geographical regions (continent).

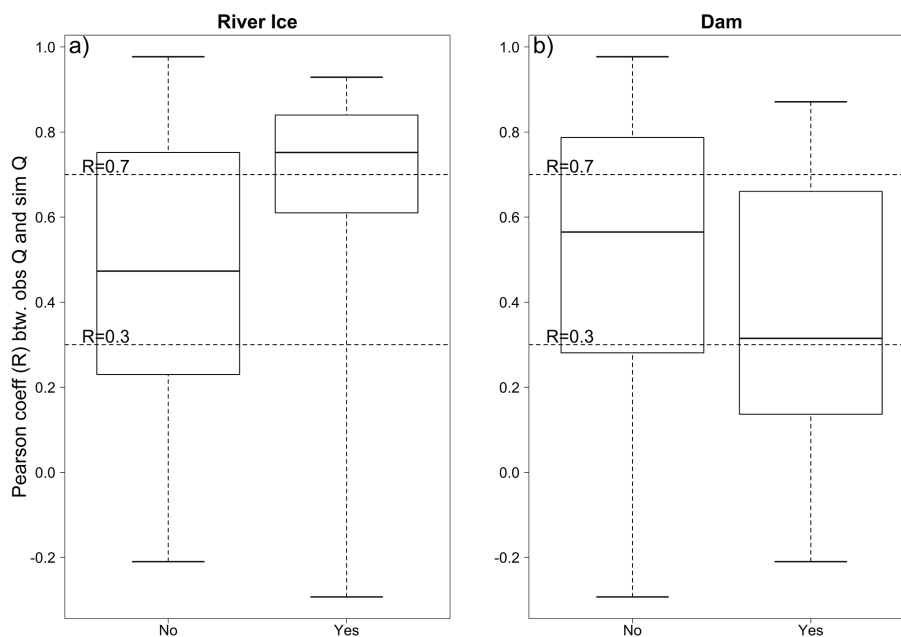


Figure 11. Evaluation of the R score obtained during the validation and its classification by (a) presence or absence of a river ice (Brown et al., 2002), and (b) presence or absence of a nearby dam or hydraulic control infrastructure using the Global Reservoir and Dam (GRanD) database (Lehner et al., 2008) and a visual check with Google Maps. For the validated locations, it is worth nothing that all stations with river ice (33) and most of them with dams (34 out of 48) are located in North America.

increase in river runoff occurs, and for many places this is translated into a strong change on the signal values. For the other types of rivers, low flows are generally a characteristic for most of the year, and if the signal-to-noise ratio is low, the signal retrieved is very noisy, which is one motivation for setting a “dry” threshold for such sites.

4.3.6 Hydraulic structures

The correlation between satellite and discharge data depends on both variables. Typically it is assumed that observed discharges are “ground truth”; however, when influenced by structures and dams, the ground discharge may not be well monitored with regard to flow area/flow width variation. For

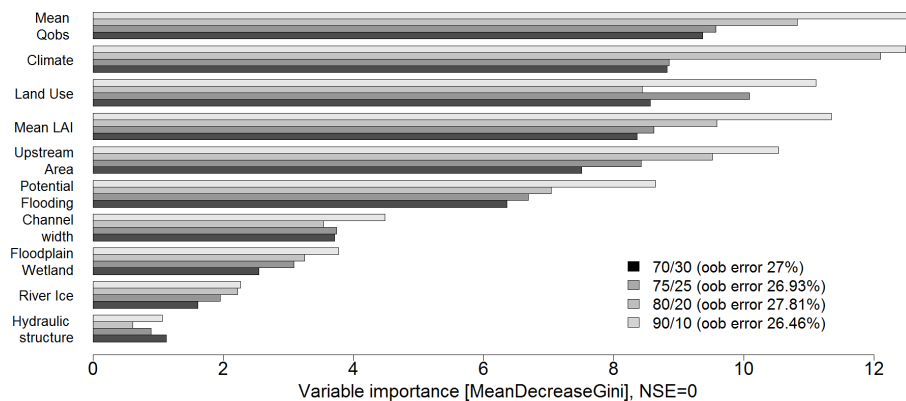


Figure 12. Average variable importance of 200 runs using the RF methodology. The Nash–Sutcliffe score was chosen as a quality index for categorising the stations as true (good predictive) or false (poor predictive). With a threshold of $NSE = 0$, we have about 50 % of the stations above and 50 % below that value. Results are shown for the different training and test groups. For all the test groups and runs, the average highest variable importance was obtained for mean observed discharge, climatic region, land cover/mean LAI and upstream catchment area, and the lowest for dam/hydraulic structure presence and river ice.

example, when there is a major increase in river discharge but a flood is avoided by artificial levees, we cannot expect the satellite signal to accurately capture the flood hydrograph; moreover, downstream flooding may be attenuated by an upstream flood control dam and reservoir, and thus the gauge location is critical. Figure 11b shows the influence of the presence or absence of a nearby dam using the Global Reservoir and Dam (GRanD) database (Lehner et al., 2008) or visually identified hydraulic control infrastructure. Locations where the dam or other element was present (48 stations) obtained lower median R scores. Therefore, ideally, observation sites should be located in areas without hydraulic control infrastructures.

4.4 Variable importance

Based on the individual analysis of the signal potential influence factors, we found that, in order to understand site performances, on some occasions multiple variables need to be analysed in a simultaneous way. For example, the generally low scores obtained at the eastern USA stations might be due to a number of factors: $\sim 64\%$ of these stations have a mean discharge value lower than $500 \text{ m}^3 \text{ s}^{-1}$ and $\sim 88\%$ of the stations are located at river width lower than 1 km. In addition, $\sim 59\%$ of the stations are located in wetland areas. Another example of the importance of analysing several factors can be seen with the locations (11 stations) which obtained low R but their mean observed discharge is higher than $500 \text{ m}^3 \text{ s}^{-1}$. All of them have a potential probability of flooding lower than 21 %, the land cover for 10 out of 11 is forest, 5 of them are located in wetlands, and 2 of them have a nearby hydraulic structure. Despite exhibiting a mean discharge greater than $500 \text{ m}^3 \text{ s}^{-1}$, these other local characteristics may be the cause of the poor performance. Therefore, we decided to use a classification decision tree technique (RF), which split the

data set at each node according to the value of one variable at a time (the best split) from a selected set of variables so as to understand the importance of each variable. RF is called an ensemble method because it is performed for a number of decision trees, in this case 500 trees, in order to improve the classification rate.

The result presented here is the rank of the importance of variables to classify a location with a good or poor performance. These values are obtained as an output of the RF analysis and are, in addition, the average of 200 independent runs. As explained in Sect. 3.4 the variable importance based on the mean decrease in Gini index was calculated for the $NSE = 0$ to distinguish between the sites with a good (above 0) from poor performance (below 0), and we also tested it with a threshold NSE of 0.50.

Figure 12 presents the variable importance for the four test groups. Features which produced large values of the “mean decrease in Gini” are ranked as more important than features which produced small values. For our locations and data available, the mean daily observed discharge has the highest importance, followed by the climatic region, land cover/mean LAI and upstream catchment area. At the same time, the presence of hydraulic structures (mainly dams) and of river ice has the lowest importance to classify a location as good or poor performance. However, this does not mean that it has no influence. Although discharge is correlated with upstream catchment area and to some degree also LAI with land cover type, both were included in this case to understand which variable might help us most to classify the sites.

Although the effect of the correlations on these measures has been studied recently (see Archer and Kimes, 2008; Strobl et al., 2009; Nicodemus and Malley, 2009; Nicodemus et al., 2010; Nicodemus, 2011; Auret and Aldrich, 2011; Tolosi and Lengauer, 2011; Grömping, 2009; Gregorutti et

al., 2013), there is not yet a consensus on the interpretation of the importance measures when the predictors are correlated and on what the effect of this correlation is on the importance measure.

In order to test the effect on the results when correlated variables were included in our analysis, an independent RF analysis was carried out (not shown in this paper) for the same variables but excluding the river width and the presence of floodplains and wetlands variables. Results also showed that the mean daily observed discharge had the highest importance and that the presence of hydraulic structures (mainly dams) and river ice had the lowest importance for classifying a location as good or poor performance.

5 Conclusions and future research

In this article we presented an evaluation of the skill of the Global Flood Detection System to measure river discharge from remote sensing signal. From the 322 stations validated, the average continental R skills are as follows: Africa 0.382, Asia 0.358, Europe 0.508, North America 0.451 and South America 0.694. Approximately 48 % of these stations have an NSE score higher than zero: 13 located in Africa, 77 in North America, 62 in South America, 1 in Asia and 1 in Europe. Results showed that the low skills scores received by stations were, for the majority of cases, due to low-flow conditions. For example, 84 out of 95 stations with $R < 0.3$ have mean discharge values lower than $500 \text{ m}^3 \text{ s}^{-1}$. These are located mainly in South Africa (25 cases) and North America (53 cases), which penalised their average continental skills. Note that our focus was on factors affecting the method globally, and that these skill values do not directly indicate measurement accuracy at a site (which could be improved, for example, by use of non-linear rating equations and/or accommodation of any phase shift or timing differences in flow-area- versus state-based discharge monitoring).

In order to better understand the impact of the local conditions on the performance of the sites, we first looked at specific factors individually. In general terms, higher skill scores were obtained for locations with one, or more than one, of the following characteristics: a river width higher than 1 km; a large floodplain area; in flooded forest; a potential flooded area per pixel greater than 40 % during a 100-year return period flood event; a land cover type of sparse vegetation, croplands or grasslands and closed to open and open forest; LAI above 2; location in a tropical climatic area; and a location where no dams or hydraulic infrastructures are present. Also, out of our locations, high-latitude rivers with seasonal ice cover tend to exhibit good performance.

Secondly, we performed a classification decision tree analysis, based on RF, to obtain the variable importance when classifying a site as good or poor. The output of this analysis showed that mean observed discharge, climatic region, land cover and mean LAI and upstream catchment area and were

the variables with higher importance, whereas river ice and dam obtained the lowest importance. Both the individual and the combined classification analysis of these local characteristics give us critical evidence of the relationship between the ground and satellite discharge measurements and when it is expected to perform well. Furthermore, it provides a guideline for future selection of measuring sites.

The locations with a very good performance will be selected for a potential future project where satellite-measured discharge could be calculated for longer periods and on a daily basis from the remote sensing signal, analogous to the Dartmouth Flood Observatory method. This will represent a major step forward in developing continental and global hydrological monitoring systems as these data can fill the gaps where real-time ground discharge measurements are not available (the case at many locations globally). We found that some of the sites with good performance are located within international river basins such as the Niger, Volta and Zambezi in Africa. In addition, for the studied locations with good signal performance but rather short contemporary time series with in situ-observed discharge (such as the Siberian stations), the calibration of the signal to obtain discharge measurements could be executed at any point when additional ground data are available. This will also be beneficial for all stations, including those with time series longer than 7 years.

Zhang et al. (2013) recently demonstrated the potential of integrating satellite signal provided by the Global Flood Detection System in improving flood forecasting. This first attempt at data assimilation was carried out for a single station (Rundu, northern Namibia – included in this study) with the conceptually simple Hydrological MODel (HyMOD). Hence, a prospective study with the inclusion of all these stations for post-processing through data assimilation and error correction of the stream-flow forecast in hydrological models could be done. For instance, for the pre-operational Global Flood Awareness System (GloFAS) (Alfieri et al., 2013) and the African Flood Forecasting System (AFFS) (Thiemig et al., 2014) in an analogous way as it is already being done with ground-gauge-observed streamflow on the European Flood Awareness System (Bartholmes et al., 2009; Thielen et al., 2009). Hence, work towards the integration of global flood detection and forecasting systems such as GFDS and GloFAS, respectively, can provide more comprehensive information for decision makers.

Appendix A

Table A1. Studied land cover types from GlobCover (2009) aggregated into broader categorical classes by type and vegetation density.

Label	Aggregated classes
Rainfed croplands	Rainfed croplands
Sparse (< 15 %) vegetation	Sparse vegetation
Closed to open (> 15 %) broadleaved evergreen or semi-deciduous forest (> 5 m)	Closed to open forest
Closed to open (> 15 %) mixed broadleaved and needle-leaved forest (> 5 m)	Closed to open forest
Closed to open (> 15 %) (broadleaved or needle-leaved, evergreen or deciduous) shrubland (< 5 m)	Closed to open forest
Closed to open (> 15 %) herbaceous vegetation (grassland, savannahs or lichens/mosses)	Closed to open forest
Closed to open (> 15 %) broadleaved forest regularly flooded (semi-permanently or temporarily) – fresh or brackish water	Closed to open forest
Closed to open (> 15 %) grassland or woody vegetation on regularly flooded or waterlogged soil – fresh, brackish or saline water	Closed to open forest
Open (15–40 %) broadleaved deciduous forest/woodland (> 5m)	Open forest
Open (15–40 %) needle-leaved deciduous or evergreen forest (> 5m)	Open forest
Mosaic cropland (50–70 %)/vegetation (grassland/shrubland/forest) (20–50 %)	Mosaic cropland or grassland
Mosaic grassland (50–70 %)/forest or shrubland (20–50 %)	Mosaic cropland or grassland
Mosaic vegetation (grassland/shrubland/forest) (50–70 %)/cropland (20–50 %)	Mosaic vegetation predominant
Mosaic forest or shrubland (50–70 %)/grassland (20–50 %)	Mosaic vegetation predominant
Closed (> 40 %) broadleaved deciduous forest (> 5m)	Closed forest
Closed (> 40 %) needle-leaved evergreen forest (> 5 m)	Closed forest
Closed (> 40%) broadleaved forest or shrubland permanently flooded – saline or brackish water	Closed forest
Artificial surfaces and associated areas (urban areas > 50 %)	Urban

Acknowledgements. We acknowledge the Global Runoff Data Centre and South African Water Affairs for providing historic discharge measurements. Furthermore, we would like to acknowledge the team from the Joint Research Centre Crisis Monitoring and Response Technologies (CRITECH) for support and access to the Global Flood Detection System signal historical data. Also, Philippe Roudier, Simone Russo, Angel Udias and Feyera Hirpa are thanked for their valuable input and methodology advice; Ad de Roo for PhD supervision; and the editor and the two reviewers are gratefully acknowledged for their valuable feedback. G. R. Brakenridge acknowledges funding support from the NASA Hydrology Program.

Edited by: H. Cloke

References

- Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J., and Pappenberger, F.: GloFAS –global ensemble streamflow forecasting and flood early warning, *Hydrol. Earth Syst. Sci.*, 17, 1161–1175, doi:10.5194/hess-17-1161-2013, 2013.
- Archer, K. J. and Kimes, R. V.: Empirical characterization of random forest variable importance measures, *Comput. Stat. Data Anal.*, 52, 2249–2260, doi:10.1016/j.csda.2007.08.015, 2008.
- Auret, L. and Aldrich, C.: Empirical comparison of tree ensemble variable importance measures, *Chemomet. Intelligent Labor. Syst.*, 105, 157–170, doi:10.1016/j.chemolab.2010.12.004, 2011.
- Bartholmes, J. C., Thielen, J., Ramos, M. H., and Gentilini, S.: The European flood alert system EFAS –Part 2: Statistical skill assessment of probabilistic and deterministic operational forecasts, *Hydrol. Earth Syst. Sci.*, 13, 141–153, doi:10.5194/hess-13-141-2009, 2009.
- Bontemps, S., Defourny, P., Bogaert, E. V., Arino, O., Kalogirou, V., and Perez, J.R.: GLOBCOVER 2009 – Products Description and Validation Report, available at: <http://due.esrin.esa.int/globcover/> (last access: 15 February 2014), 2010.
- Brakenridge, G. R., Nghiem, S. V., Anderson, E., and Chien, S.: Space-based measurement of river runoff, *Eos Trans. AGU*, 86, 185–188, doi:10.1029/2005EO190001, 2005.
- Brakenridge, G. R., Nghiem, S. V., Anderson, E., and Mic, R.: Orbital microwave measurement of river discharge and ice status, *Water Resour. Res.*, 43, W04405, doi:10.1029/2006WR005238, 2007.
- Brakenridge, G. R., Cohen, S., Kettner, A. J., De Groeve, T., Nghiem, S. V., Syvitski, J. P. M., and Fekete, B. M.: Calibration of satellite measurements of river discharge using a global hydrology model, *J. Hydrol.*, 475, 123–136, 2012.
- Brakenridge, G. R., De Groeve, T., Cohen, S., and Nghiem, S. V.: River Watch, Version 2: Satellite River Discharge and Runoff Measurements: Technical Summary, University of Colorado, Boulder, CO, USA, available at: <http://floodobservatory.colorado.edu/SatelliteGaugingSites/technical.html>, last access: 1 December 2013.
- Breiman, L.: Random Forests, *Machine Learning*, 45, 5–32, 2001.
- Brown, J., Ferrians Jr., O. J., Heginbottom, J. A., and Melnikov, E. S.: Circum-Arctic Map of Permafrost and Ground-Ice Conditions. Version 2. [Permafrost], Boulder, Colorado USA: National Snow and Ice Data Center, 2002.
- Committee on Earth Observation Satellites (CEOS) Flood Pilot, available at: <http://www.ceos.org/>, last access: 1 September 2014.
- Chan, J. C.-W. and Paelinckx, D.: Evaluation of Random Forest and Adaboost tree-based ensemble classification and spectral band selection for ecotope mapping using airborne hyperspectral imagery, *Remote Sens. Environ.*, 112, 2999–3011, 2008.
- Chukhlantsev, A. A.: Modeling of microwave emission from vegetation canopies, *Microwave Radiometry of Vegetation Canopies*, Springer Netherlands, Chap. 6, 147–175, 2006.
- De Groeve, T., Brakenridge, G. R., and Kugler, Z.: Near Real Time Flood Alerting for the Global Disaster Alert and Coordination System, edited by: Van de Walle, B., Burghardt, P., and Nieuwenhuis, C., *Proceedings of the 4th International ISCRAM Conference*, 33–40, 2006.
- De Groeve, T. and Riva, P.: Global Real-time Detection of Major Floods Using Passive Microwave Remote Sensing, *Proceedings of the 33rd International Symposium on Remote Sensing of Environment Stresa, Italy, May 2009*.
- De Groeve, T.: Flood monitoring and mapping using passive microwave remote sensing in Namibia, *Geomatics, Nat. Hazards Risk*, 1:1, 19–35, 2010.
- Di Baldassarre, G. and Montanari, A.: Uncertainty in river discharge observations: a quantitative analysis, *Hydrol. Earth Syst. Sci.*, 13, 913–921, doi:10.5194/hess-13-913-2009, 2009.
- Disaster Charter: Space and Major Disasters, available at: <http://www.disasterscharter.org/> (last access: 1 September 2014), 2013.
- EM-DAT: The OFDA/CRED International Disaster Database, Université Catholique de Louvain, Brussels, Belgium, available at: <http://www.emdat.be>, last access: 1 December 2013.
- Fekete, B. M., Vorosmarty, C. J., and Grabs, W.: Global, composite runoff fields based on observed river discharge and simulated water balances, *GRDC Report 22*, Global Runoff Data Center, Koblenz, Germany, 1999.
- GDACS: Global Disaster Alert and Coordination System, Global Floods Detection System available at: <http://www.gdacs.org/>, last access: 1 December 2013.
- Global Runoff Data Centre: Major River Basins of the World, 56068 Koblenz, Germany: Federal Institute of Hydrology (BfG), available at: <http://grdc.bafg.de/> (last access: 20 January, 2013), 2007.
- Global Runoff Data Centre: The River Discharge Time Series, 56068 Koblenz, Germany: Federal Institute of Hydrology (BfG), available at: <http://grdc.bafg.de/>, last access: 20 January 2013.
- Golnaraghi, M., Douris, J., and Migraine, J.-B.: Saving Lives Through Early Warning Systems and Emergency Preparedness, *Risk Wise*, Tudor Rose, 137–141, 2009.
- Gonzalez, L., Velasco Morente, F., Gavilan Ruiz, J. M., Sanchez-Reyes and Fernandez, J. M.: The Similarity between the Square of the Coefficient of Variation and the Gini Index of a General Random Variable, *J. Quant. Methods Econom. Business Admin.*, 10, 5–18, ISSN 1886-516X, 2010.
- Gregorutti, B., Michel, B., and Saint-Pierre, P.: Correlation and variable importance in random forests., *Cornell University Library*, arXiv: 1310.5726 [stat], 2013.
- Grömping, U.: Variable Importance Assessment in Regression: Linear Regression versus Random Forest, *The American Statistician*, 11/2009; 63, 308–319, doi:10.1198/tast.2009.08199, 2009.
- Hirpa, F. A., Hopson, T. M., De Groeve, T., Brakenridge, G. R., Gebremichael, M., and Restrepo, P. J.: Upstream satellite remote sensing for river discharge forecasting: Application to ma-

- for rivers in South Asia, *Remote Sens. Environ.*, 131, 140–151, 2013.
- Hirsch R. M. and Costa J. E.: U.S. Stream Flow Measurement and Data Dissemination Improve EOS, *Transactions, American Geophysical Union*, 18 May 2004, 85, 197–203, 2004.
- Khan, S. I., Hong, Y., Vergara, H. J., Gourley, J. J., Robert Brakenridge, G., De Groeve, T., Flamig, Z. L., Policelli, F., and Yong, B.: Microwave satellite data for hydrologic modeling in ungauged basins, *IEEE Geosci. Remote Sens. Lett.*, 9, 663–667, 2012.
- Kugler, Z. and De Groeve, T.: *The Global Flood Detection System*, Office for Official Publications of the European Communities, Luxembourg, 2007.
- Kugler, Z.: Remote sensing for natural hazard mitigation and climate change impact assessment, *Quarterly J. Hungarian Meteorol. Serv.*, January–March 116, 21–38, 2012.
- Kundzewicz, Z. W.: *Changes in Flood Risk in Europe*, Wallingford: IAHS Press. p. 516, IAHS special publication; 10, United Nations: Report of the United Nations Conference on Sustainable Development Rio de Janeiro, Brazil, 20–22 June 2012, A/CONF.216/16, 2012.
- Le Coz, J., Renard, B., Bonnifant, L., Branger, F., and Le Boursicaud, R.: Combining hydraulic knowledge and uncertain gaugings in the estimation of hydrometric rating curves: A Bayesian approach, *J. Hydrol.*, 509, 573–587, 2014.
- Lehner, B. and Döll, P.: Development and validation of a global database of lakes, reservoirs and wetlands, *J. Hydrol.*, 296, 1–22, doi:10.1016/j.jhydrol.2004.03.028, 2004.
- Lehner, B., Reidy Liermann, C., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P., Döll, P., Endejan, M., Frenken, K., Magome, J., Nilsson, C., Robertson, J., Rödel, R., Sindorf, N., and Wisser, D.: High resolution mapping of the world's reservoirs and dams for sustainable river flow management, *Frontiers in Ecology and the Environment*. Source: GWSP Digital Water Atlas. Map 81: GRanD Database (V1.0), available at: http://atlas.gwsp.org/index.php?option=com_content&task=view&id=209&Itemid=1 (last access: 11 March 2014), 2008.
- Liaw, A. and Wiener, M.: Classification and Regression by random Forest, *R News*, 2, 18–22, 2002.
- Moffitt, C. B., Hossain, F., Adler, R. F., Yilmaz, K. K., and Pierce, H. F.: Validation of a TRMM-Based Global Flood Detection System in Bangladesh, *Int. J. Appl. Earth Observ. Geoinform.*, 13, 165–177, doi:10.1016/j.jag.2010.11.003, 2011.
- MunichRe: Munich Reinsurance: January 2014 press release, Münchener Rückversicherungs-Gesellschaft, Geo Risks Research, NatCatSERVICE, available at: http://www.munichre.com/en/media_relations/press_releases/2014/2014_01_07_press_release.aspx, last access: 20 January 2014.
- Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I – A discussion of principles, *J. Hydrol.*, 10, 282–290, 1970.
- Nicodemus, K. K.: Letter to the editor: On the stability and ranking of predictors from random forest variable importance measures, *Briefings in Bioinform.*, 12, 369–373, doi:10.1093/bib/bbr016, 2011.
- Nicodemus, K. K. and Malley, J. D.: Predictor correlation impacts machine learning algorithms: implications for genomic studies, *BCM Bioinform.*, 25, 1884–1890, doi:10.1093/bioinformatics/btp331, 2009.
- Nicodemus, K. K., Malley, J. D., Strobl, C., and Ziegler, A.: The behavior of random forest permutation-based variable importance measures under predictor correlation, *BMC Bioinform.*, 11, 13 pp., doi:10.1186/1471-2105-11-110, 2010.
- Palczewska, A., Palczewski, J., Robinson, R. M., and Neagu, D.: Interpreting random forest models using a feature contribution method, *Proceedings of the 2013 IEEE 14th International Conference on Information Reuse and Integration, IEEE IRI 2013*, 112 pp., 2013.
- Pappenberger, F., Matgen, P., Beven, K. J., Henry, J. B., Pfister, L., and de Fraipont, P.: Influence of uncertain boundary conditions and model structure on flood inundation predictions, *Adv. Water Resour.*, 29, 1430–1449, doi:10.1016/j.advwatres.2005.11.012, 2006.
- Pappenberger, F., Dutra, E., Wetterhall, F., and Cloke, H. L.: Deriving global flood hazard maps of fluvial floods through a physical model cascade, *Hydrol. Earth Syst. Sci.*, 16, 4143–4156, doi:10.5194/hess-16-4143-2012, 2012.
- Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-Geiger climate classification, *Hydrol. Earth Syst. Sci.*, 11, 1633–1644, doi:10.5194/hess-11-1633-2007, 2007.
- Pelletier, P. M.: Uncertainties in the single determination of river discharge: a literature review, *Can. J. Civil Eng.*, 15, 834–850, 1988.
- Rosso, R. A.: linear approach to the influence of discharge measurement error on flood estimates, *Hydrol. Sci. J.*, 30, 137–149, doi:10.1080/02626668509490975, 1998.
- Sandri, M. and Zuccolotto, P.: A bias correlation algorithm for the Gini variable importance measure in classification trees, *J. Comput. Graphical Stat.*, 17, 611–628, doi:10.1198/106186008X344522, 2008.
- Schumann, G., Bates, P. D., Horritt, M. S., Matgen, P., and Pappenberger, F.: Progress in Integration of Remote Sensing-derived Flood Extent and Stage Data and Hydraulic Models, *Rev. Geophys.*, 47, RG4001, doi:10.1029/2008RG000274, 2009.
- South African Water Affairs (DWA) database, available at: <http://www.dwa.gov.za/Hydrology/>, last access: 10 July 2013.
- Strobl, C., Malley, J., and Tutz, G.: An Introduction to Recursive Partitioning: Rationale, Application, and Characteristics of Classification and Regression Trees, Bagging, and Random Forests, *Psychol. Methods*, 14, 323–348, doi:10.1186/1471-2105-9-307, 2009.
- Thielen, J., Bartholmes, J., Ramos, M.-H., and de Roo, A.: The European Flood Alert System –Part 1: Concept and development, *Hydrol. Earth Syst. Sci.*, 13, 125–140, doi:10.5194/hess-13-125-2009, 2009.
- Thiemig, V., Bisselink, B., Pappenberger, F., and Thielen, J.: A pan-African Flood Forecasting System, *Hydrol. Earth Syst. Sci. Discuss.*, 11, 5559–5597, doi:10.5194/hessd-11-5559-2014, 2014.
- Tolosi, L. and Lengauer, T.: Classification with correlated features: unreliability of feature ranking and solutions, *Bioinform.*, 27, 1986–1994, doi:10.1093/bioinformatics/btr300, 2011.
- Tomkins, K. M.: Uncertainty in streamflow rating curves: Methods, controls and consequences, *Hydrol. Process.*, 28, 464–481, 2014.
- UNISDR: *Global Assessment Report: Revealing Risk, Redefining Development*, Chap. 2.2, Global disaster risk trends, United Nations, printed in the UK, ISBN 978-92-1-132030-5, 22–27, 2011.

- UNOSAT, UNITAR Operational Satellite Applications Programme, available at: <http://www.unitar.org/unosat/maps>, last access: 1 December 2013.
- Van Westen, C. J.: Remote sensing and GIS for natural hazards assessment and disaster risk management, in: *Treatise on Geomorphology*, edited by: Shroder, J. and Bishop, M. P., Academic Press, San Diego, CA, Vol. 3, Remote Sensing and GIScience in Geomorphology, 259–298, 2013.
- Yamazaki, D., O'Loughlin, F., Trigg, M. A., Miller, Z. F., Pavel-sky, T. M., and Bates, P. D.: Development of the global width database for large river, *Water Resour. Res.*, 50, 3467–3480, doi:10.1002/2013WR014664, 2014.
- Yitzhaki, S. and Schechtman, E.: *The Gini Methodology. A Primer on a Statistical Methodology*, Springer Series in Statistics, Vol. 272, ISBN: 978-1-4614-4720-7, 2013.
- Zaraj, Z., Zambrano-Bigiarini, M., Salamon, P., Burek, P., Gentile, A., and Bianchi, A.: Calibration of the LISFLOOD hydrological model for Europe. Calibration Round 2013JRC Technical Report, European Commission, Joint Research Centre, Ispra, Italy, 2013.
- Zhang, Y., Hong, Y., Wang, X., Gourley, J. J., Gao, J., Vergara, H. J., and Yong, B.: Assimilation of passive microwave streamflow signals for improving flood forecasting: A first study in Cubango River Basin, Africa, *IEEE J. Selected Topics. Appl. Earth Observ. Remote Sens.*, 6, 2375–2390, 2013.
- Zhu, Z., Bi, J., Pan, Y., Ganguly, S., Anav, A., Xu, L., Samanta, A., Piao, S., Nemani, R. R., and Myneni, R. B.: Global data sets of vegetation leaf area index (LAI)3g and fraction of photosynthetically active radiation (FPAR)3g derived from global inventory modeling and mapping studies (GIMMS) normalized difference vegetation index (NDVI3G) for the period 1981 to 2011, *Remote Sens.*, 5, 927–948, 2013.