

A key step in data-driven environmental modelling is **input variable selection** to ensure that the least number of variables with minimum relationship between inputs and outputs.

Hydrological predictions in ungauged catchments is one such area where the information on influential predictors of runoff signatures guides in understanding **dominant controls of** analysis of hydrological similarity among ungauged catchments, e.g. to identify reference (donor) catchment(s).

Large-sample hydrology can help to gain useful insights on how these significant predictors change over different



Input variable selection for hydrological predictions in ungauged catchments: with or without clustering?

1 Middle East Technical University, Department of Civil Engineering, Water Resources Laboratory, Ankara, Turkey 2 Middle East Technical University, Department of Statistics, Ankara, Turkey

Nilay Dogulu¹, Inci Batmaz², and Elcin Kentel¹

- predictor of the system.



Input Variable Selection (IVS)

Information (PMI)		Partial Correlation (PCIS)			Iterative Input Selection (IIS)		
)w flow % quantile)	Medium flow (Q _{mean})	High flow (5% quantile)	Low flow (95% quantile)	Medium flow (Q _{mean})	High flow (5% quantile)	Low flow (95% quantile)	Medium flow (Q _{mean})
q mean ry daily vs PET slope	mean daily precip days	mean daily precip long	mean daily precip lat	mean daily precip lat long	aridity water % seaso nality and timing of precip	water % aridity con duct snow %	aridity water %
lat max water conte ev nt	mean daily precip slope aridity	mean daily precip slope mean daily PET	mean daily precip long	slope GVF max	mean daily precip	max water content lat	mean daily precip
permea bility elev	mean daily precip mean daily PET	mean daily precip % forest	ave. dura of dry perio ds elev ave. dura highprecip events	aridity mean daily precip max water content	water % mean daily precip	slope LAI diff	water % mean daily precip
dura recip nts freq	slope sand % mean daily precip	mean daily precip slope LAI max	slope freq of dry days permea bility	mean daily precip slope LAI max	mean % daily precip	slope snow %	ave. dura highprecip events water %
of ays sity sity seaso nality and timing of precip	mean daily PETLAILAIdiffmaxporosity	mean daily precip elev	slope soil depth	elev LAI diff	mean daily precip % forest	slope % forest	mean daily precip % forest
elev mean daily precip	aridity slope mean daily precip	aridity snow % mean daily PET	elev con duct permea bility	mean daily precip slope	aridity soil depth	mean daily precip soil depth	slope elev
porosity q lry /s	long mean freq daily of dry precip days	mean daily precip lat soil depth	con duct of dry periods lat	mean daily precip long lat	mean daily precip soil depth	con duct high precip freq	aridity mean daily PET
sity of dry days % forest	seaso mean nality daily and PET timing mean of daily precip precip	mean daily precip freq of dry days elev	mean high daily precip PET freq % forest	arid ity seasonality and timing of precip	mean daily precip freq of dry days	freq of dry days forest	seaso nality and timing of precip
nduct lat n daily precip	mean daily precip aridity	mean daily precip conduct	Conduct GVF max	mean mean daily daily precip PET slope	aridity GVF diff	ave. duration of dry periods sand %	aridity
nduct silt % elev	aridity sand %	aridity long snow %	permeability long porosity	aridity mean daily precip %	% carbonate rocks	max water content silt %	aridity
ty permea bility y	mean daily precip aridity	mean daily precip snow %	% carbonate rocks area mean daily precip	mean daily precip aridity	% carbonate rocks conduct	slope % carbonate rocks	% carbonate rocks conduct

Summary

• IVS is important for hydrological predictions in ungauged catchments! And clustering of input data

• Climate and topography related variables are more frequently selected compared to soil, geology and vegetation attributes. However, there are cases where geology becomes more important as a

• We found that selected variables for predicting hydrological attributes are not the same for different

• Moreover, selected variables for a given cluster vary depending on the hydrological attribute of interest. • More in-depth analysis of IVS results is in progress • The effect of clustering method choice could be

• Next step: training data-driven models for each cluster to make hydrological predictions on validation

Contact







Input variable selection for hydrological predictions in ungauged catchments: with or without clustering?

Nilay Dogulu (1), Inci Batmaz (2), and Elcin Kentel (3)

(1) Middle East Technical University (METU), Civil Engineering Dept., Water Resources Laboratory, Ankara, Turkey (ndogulu@metu.edu.tr), (2) Middle East Technical University (METU), Statistics Dept., Ankara, Turkey (ibatmaz@metu.edu.tr), (3) Middle East Technical University (METU), Civil Engineering Dept., Water Resources Laboratory, Ankara, Turkey (ekentel@metu.edu.tr)

A key step in data-driven environmental modelling, including for hydrological purposes, is input variable selection (IVS) to ensure that the least number of variables with minimum redundancy are used to characterize the inherent relationship between inputs and outputs. Hydrological predictions in ungauged catchments is one such area where the information on influential predictors of runoff signatures guides in understanding dominant controls of meaningful information transfer from gauged to ungauged locations (i.e. regionalization). This understanding is valuable especially for the analysis of hydrological similarity among ungauged catchments, e.g. to identify reference (donor) catchment(s). Large-sample hydrology can help to gain useful insights on how these significant predictors change over different runoff signatures representing particular hydrological conditions as well as within different groups of similar catchments. This study explores the added value of clustering for input variable selection in the case of catchments across continental USA using the CAMELS dataset (Addor et al., 2017). We employ the method of k-means clustering on the input space consisting of topography, soil, geology, vegetation and climate attributes as a way to deal with heterogeneity and complexity in hydrological processes, and thus, aiming to identify similar groups of catchments. The input variables (for predicting selected hydrological attributes) are determined by three IVS filter algorithms (partial mutual information, partial correlation input selection, and iterative input selection), then evaluated and compared for among different clusters and when no clustering method is applied. We present and discuss the results for three hydrological attributes – 95% flow percentile (low flows), mean daily discharge (medium flows), and 5% flow percentile (high flows) - in order to account for variability in hydrological conditions, and to investigate the dominant catchment and/or climate characteristics controlling low/medium/high flow predictability at ungauged locations. The findings of our study have implications for reference (donor) catchment selection and hydrological model independent regionalization (aka hydrostatistical or data-driven methods), and are particularly relevant to understanding hydrological similarity and catchment classification in the absence of local runoff observations.

Addor, N., Newman, A. J., Mizukami, N., & Clark, M. P. (2017) The CAMELS data set: catchment attributes and meteorology for large-sample studies. Hydrology and Earth System Sciences, 21(10), 5293-5313.