

Supplement of Hydrol. Earth Syst. Sci., 21, 2463–2481, 2017
<http://www.hydrol-earth-syst-sci.net/21/2463/2017/>
doi:10.5194/hess-21-2463-2017-supplement
© Author(s) 2017. CC Attribution 3.0 License.



Supplement of

Regionalizing nonparametric models of precipitation amounts on different temporal scales

Tobias Mosthaf and András Bárdossy

Correspondence to: Tobias Mosthaf (tobias.mosthaf@iws.uni-stuttgart.de)

The copyright of individual parts of the supplement might differ from the CC-BY 3.0 licence.

S1 Regionalization example

Two general possibilities to obtain precipitation amount distributions at ungauged locations exist. The first approach is the interpolation of rainfall values for every time step to the target location, followed by an estimation of the distribution function with the interpolated values ($values_{inter}$). The second approach is first fitting a distribution function to all control locations, which is followed by an interpolation of these distributions to the target location (cdf_{inter}). In the following, these possibilities will be compared with each other to motivate the use of the cdf_{inter} approach, which is the method used within our investigations. Although it is commonly accepted to follow the cdf_{inter} approach to obtain precipitation amount distributions at ungauged locations for stochastic rainfall models, we still want to illustrate the deficiencies of the $values_{inter}$ method to motivate the cdf_{inter} approach empirically. Additionally, the resulting estimation errors also appear when rainfall values are interpolated without considering the CDF explicitly. For example the use of interpolated rainfall values for hydrological models may introduce a bias in the discharge estimation caused by poor interpolation results.

In order to ensure equal interpolation weights ϵ_i of the control gauges i for both possibilities, a simple inverse distance weighting (IDW) is used as interpolation technique in this example, which is based on the following Eq. S1:

$$\epsilon_i = \frac{\frac{1}{d_i^2}}{\sum_{i=1}^{30} \frac{1}{d_i^2}} \quad (S1)$$

where d_i is the distance between control gauge i and the respective target gauge. For this interpolation example IDW is preferred over OK for the following reasons: (i) Using OK with daily precipitation values ($values_{inter}$) would lead to the additional challenge of including zero rainfall values within the estimation of the variogram and the kriging itself. The focus of this paper, however, does not lie on interpolating rainfall values, therefore, the simpler IDW method is used for interpolating rainfall values. (ii) IDW leads to the same interpolation weights for both approaches and therefore assures that the better performance of one of the approaches does not originate from the calculation of the weights, but from the chosen interpolation scheme (cdf_{inter} or $values_{inter}$). In the research article, OK is preferred over IDW, because OK is considered as a better interpolation method than the simpler IDW. The nonparametric KDE using SRT for the bandwidth selection is applied for estimating distribution functions at the control gauges.

Another exception within this regionalization example is the inclusion of zero values to show the advantages of interpolating distributions instead of precipitation values regarding P_0 . Zero values can be included within the interpolation of nonparametric distributions by applying the following steps. (i) Fitting a distribution to all precipitation values at each gauge. (ii) Estimate the quantile values for certain quantiles (non-exceedance probabilities) over the whole probability range (0-1) with the inverse of the fitted distributions at each gauge. (iii) Use the interpolation weights from IDW to interpolate the quantile values of different gauges for each chosen quantile. (iv) If the quantile is below P_0 for some (or all)

Table S1: Regionalization example: Basic daily rainfall statistics of observed values at the validation gauge (*data*), interpolated rainfall values (*values_{inter}*), randomly sampled rainfall values of the interpolated nonparametric distribution function (*cdf_{inter}*) and the respective ranges of the calibration gauges. The rainfall statistics are the arithmetic mean (\bar{x}), the standard deviation (s_x) of all rainfall values, the arithmetic mean ($\bar{x}_{>0}$) of non zero values, the probability of zero rainfall P_0 and the maximum value (*max*).

	<i>Data</i>	<i>values_{inter}</i>	<i>cdf_{inter}</i>	<i>Range calibration set</i>
\bar{x}	2.18	2.17	2.27	1.77 - 3.18
s_x	4.56	4.04	4.93	3.88 - 6.47
$\bar{x}_{>0}$	4.39	2.97	4.20	3.73 - 4.47
P_0	0.50	0.27	0.46	0.46 - 0.54
<i>max</i>	56.0	49.12	62.29	42.5 - 102.3

gauges, the quantile value at these gauges will be 0 mm, which are then just included in the interpolation. (v) The highest quantile with 0 mm at the target gauge defines P_0 at the target.

In our example the distribution of daily rainfall values (1D) for the gauge Esslingen / Neckar is estimated from rainfall values of 30 neighboring gauges (see Fig. S1 (a)). In Fig. S1 (b) and (c), parts of the distribution functions resulting from both methods and the original EDF are shown. Clear disadvantages of the *values_{inter}* method are the overestimation of days with rainfall and thus an underestimation of the probability of no rainfall (Fig. S1 (b)) and a clear underestimation of the CDF for higher quantiles (Fig. S1 (c)).

As the *cdf_{inter}* method does not provide rainfall values automatically, which are needed to calculate basic statistical measures, random rainfall values are generated with the inverse of the interpolated CDF. The number of these random values is equivalent to the number of observed daily rainfall values of the validation gauge. In Table S1 basic statistics of precipitation amounts are listed for both methods and observations. Looking at the mean values of all rainfall (\bar{x}) values, the *values_{inter}* method seems to reproduce this statistic very well. Considering the other statistics in Table S1 and Fig. S1 this is most probably caused by two disadvantages of this method: an overestimation of days with small rainfall amounts (see P_0) and a simultaneous underestimation of higher rainfall intensities (see $\bar{x}_{>0}$ and *max*). This argument is reaffirmed by the smaller standard deviation of *values_{inter}* and the illustrations of the precipitation amount distributions in Fig. S1. The *cdf_{inter}* method mainly provides better results summarizing the listed statistics. Only a tendency of overestimating high rainfall intensities can be observed.

As the *values_{inter}* method has great problems in reproducing probabilities of zero rainfall and the shape of the distribution function, this method is not recommended to be used with rainfall over a great range of aggregations. For higher aggregations these disadvantages may have no noticeable effect, but for smaller aggregations with a greater skewness the problems might even increase. This would lead to a more pronounced underestimation of high quantile values, which are mostly the decisive ones for subsequent applications. As

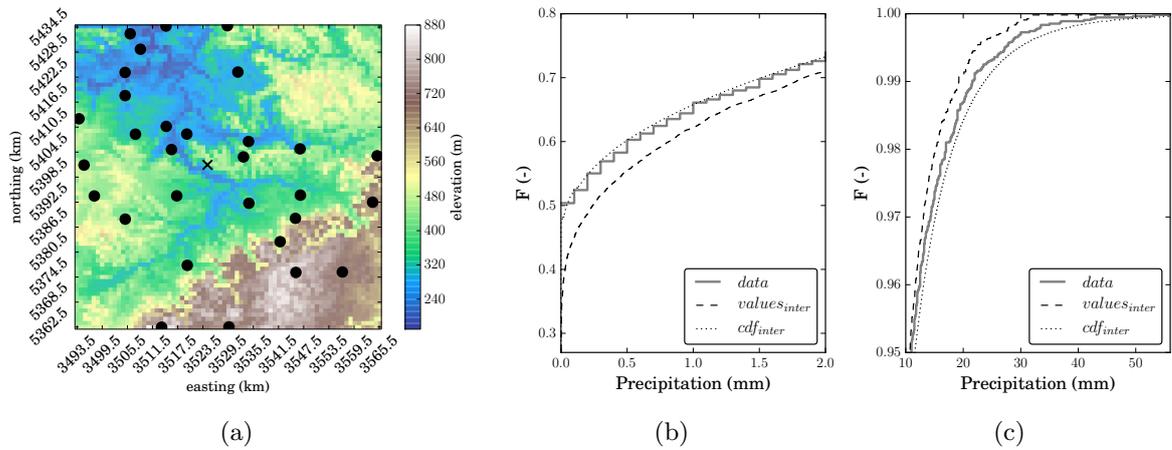


Figure S1: Regionalization example: (a) shows the daily target gauge (black cross) and the 30 neighboring daily gauges (black dots) of the regionalization example. In (b) and (c) parts of the EDF of the target gauge (*data*), the EDF of the interpolated rainfall values (*values_{inter}*) and the interpolated nonparametric estimate (*cdf_{inter}*) of the target CDF are depicted.

the *cdf_{inter}* method exhibits better results concerning the basic rainfall volume statistics, it seems to be the better choice for the purpose of interpolating precipitation amount models.

S2 Probability distributions of precipitation amounts in a spatial context (corresponds to section 5 in the research article)

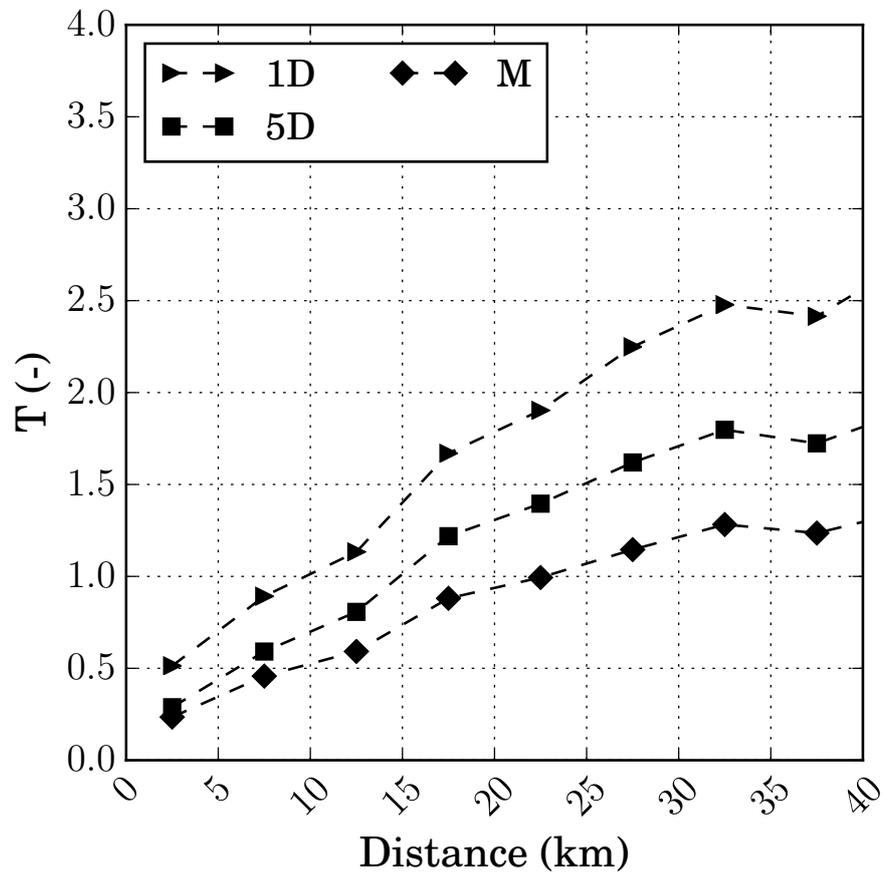


Figure S2: T statistic over distance: The graph shows the mean values of the T statistic for temporal resolutions from 1D to M.

S3 Usage of daily values for sub-daily values - Empirical cross validation (corresponds to section 9 in the research article)

To estimate the usage of daily observations for sub-daily distribution functions with the rescaling procedure described in section 9 of the research article, a cross validation is applied based on the high resolution gauges only, which are used as daily gauges one after another. The resulting sub-daily statistics of scaled values for these pseudo daily gauges are compared to their original sub-daily values by calculating the mean squared errors over all gauges. The scaled nearest neighbor values are compared to nearest neighbor values and to interpolated rainfall values. The interpolation is done by OK with ten neighbors using a single variogram model. During the cross validation a nearest neighbor gauge is defined as the gauge with the closest distance and at least 50 % of data overlapping. For the interpolation of the rainfall values with OK then again only this data overlapping period is chosen.

In Fig. S3 the results are shown for quantile values, but the standard deviation, the mean values and $Q_{V_{th}}$ were also investigated. The cross validation of the different statistical variables are very similar. For all of them the scaled nearest neighbor values (NNS) lead to the best results in summer and winter. Therefore daily gauges seem to be useful for the interpolation of sub-daily nonparametric and parametric models.

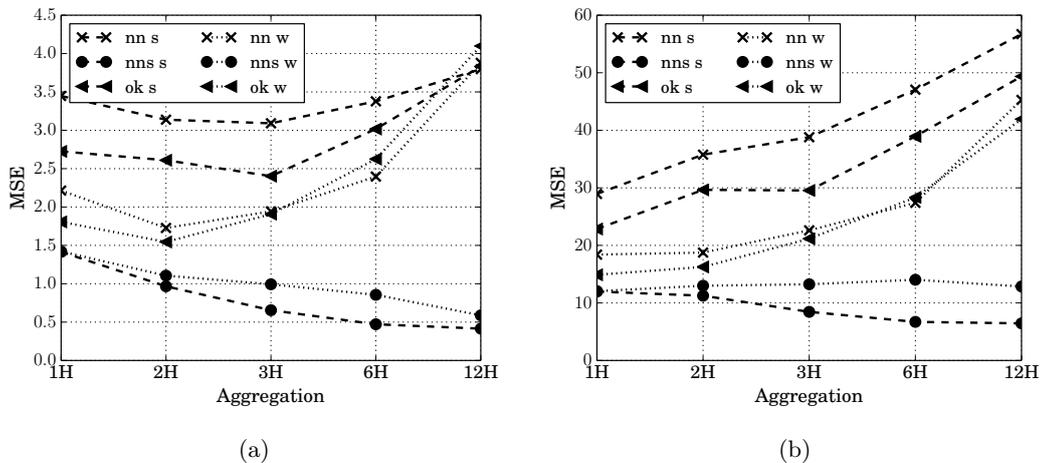


Figure S3: The mean squared errors (mse) for quantile values of discrete quantiles (in 0.001 steps) greater than Q_{th} (see Table 1 in the research article) (a) and greater than 0.995 (b) in winter (dotted) and summer (dashed) for nearest neighbor (NN), nearest neighbor scaled (NNS) and OK of rainfall values over different aggregations. At first the mean squared error over discrete quantiles is calculated for each gauge which is followed by calculating the mean of these over the whole study region.

S4 Performance of point models (corresponds to section 10.2 in the research article)

Table S2: Ranking numbers including only median and mean of the W^2 criterion for point wise estimation. The bold numbers indicate the overall best model for each temporal resolution.

Winter season								
	1H	2H	3H	6H	12H	1D	5D	M
P-Exp-MLM	69.36	55.29	41.37	23.09	15.01	15.55	6.10	79.80
P-Gamma-MLM	18.53	12.17	8.30	4.96	3.94	5.02	4.15	4.24
P-Mixed-Exp-MLM	4.50	3.61	3.38	3.09	2.76	3.15	6.31	79.80
P-Pareto-MLM	3.60	4.04	4.15	3.95	3.85	3.77	3.70	724.63
P-Weibull-MLM	10.67	6.56	4.51	3.22	2.91	3.60	3.44	12.92
P-Gamma-MOM	185.72	89.81	52.09	20.97	11.07	11.54	4.39	2.82
P-Pareto-MOM	5.98	6.12	6.17	5.24	4.85	4.71	3.51	21.92
P-Weibull-MOM	83.01	38.01	22.44	9.72	6.48	6.72	3.82	3.11
NP-SRT	2.10	2.03	2.01	2.00	2.00	2.11	2.09	2.00
NP-SJ	2.00	2.03	2.04	2.08	2.09	2.00	2.00	2.34
Summer season								
	1H	2H	3H	6H	12H	1D	5D	M
P-Exp-MLM	167.98	98.67	57.06	20.66	9.58	10.70	7.57	66.59
P-Gamma-MLM	34.10	17.41	10.24	5.48	3.67	3.98	4.17	3.39
P-Mixed-Exp-MLM	7.25	3.93	3.31	2.71	2.75	2.79	7.78	66.59
P-Pareto-MLM	3.41	4.30	3.88	3.12	2.96	3.28	3.76	484.25
P-Weibull-MLM	14.01	6.91	4.82	3.51	3.03	3.20	3.32	14.47
P-Gamma-MOM	262.32	103.55	48.55	15.54	7.41	7.20	4.12	3.43
P-Pareto-MOM	14.48	11.02	7.39	4.26	3.63	3.85	3.44	22.37
P-Weibull-MOM	87.51	33.70	17.78	7.34	4.77	4.56	3.46	3.05
NP-SRT	2.00	2.00	2.00	2.00	2.00	2.00	2.07	2.00
NP-SJ	2.15	2.19	2.14	2.20	2.14	2.05	2.00	2.21

Table S3: Ranking numbers including only median and mean of the L_d criterion for point wise estimation. The bold numbers indicate the overall best model for each temporal resolution.

Winter season								
	1H	2H	3H	6H	12H	1D	5D	M
P-Exp-MLM	155.82	90.56	88.07	46.07	22.66	26.09	8.50	662.51
P-Gamma-MLM	41.46	18.49	14.49	5.47	3.05	6.49	4.88	9.00
P-Mixed-Exp-MLM	8.42	2.11	2.00	2.00	2.03	2.00	9.34	662.86
P-Pareto-MLM	2.79	2.08	3.14	3.14	2.53	2.63	2.49	254.29
P-Weibull-MLM	20.16	7.45	5.23	2.52	2.02	3.91	3.01	11.41
P-Gamma-MOM	80.17	48.33	42.68	20.49	9.17	11.93	2.25	2.68
P-Pareto-MOM	2.13	2.69	4.06	4.01	3.12	3.37	2.20	29.24
P-Weibull-MOM	40.71	23.72	22.21	11.56	5.89	6.99	2.00	2.00
NP-SRT	11.44	20.03	33.43	40.16	34.22	29.10	15.24	10.40
NP-SJ	9.23	20.09	31.78	41.13	36.07	27.15	14.77	15.01
Summer season								
	1H	2H	3H	6H	12H	1D	5D	M
P-Exp-MLM	77.26	134.82	131.76	49.97	16.56	15.13	14.10	783.93
P-Gamma-MLM	17.41	23.12	19.89	7.34	3.75	4.00	5.83	3.83
P-Mixed-Exp-MLM	3.11	2.00	2.00	2.00	2.02	2.00	15.32	783.93
P-Pareto-MLM	4.17	9.14	8.91	4.01	2.59	2.30	2.90	225.44
P-Weibull-MLM	7.07	7.01	5.60	3.13	2.45	2.82	3.28	11.15
P-Gamma-MOM	26.83	41.72	39.18	17.55	8.10	6.24	2.33	3.83
P-Pareto-MOM	2.00	3.44	4.16	3.02	2.36	2.23	2.38	25.44
P-Weibull-MOM	10.89	17.32	17.46	9.03	5.14	3.92	2.00	2.00
NP-SRT	4.15	17.05	29.54	35.41	34.15	28.76	16.93	7.27
NP-SJ	4.76	19.90	33.97	44.22	39.03	31.09	15.64	10.11

S5 Variogram estimation (corresponds to section 10.3.1 in the research article)

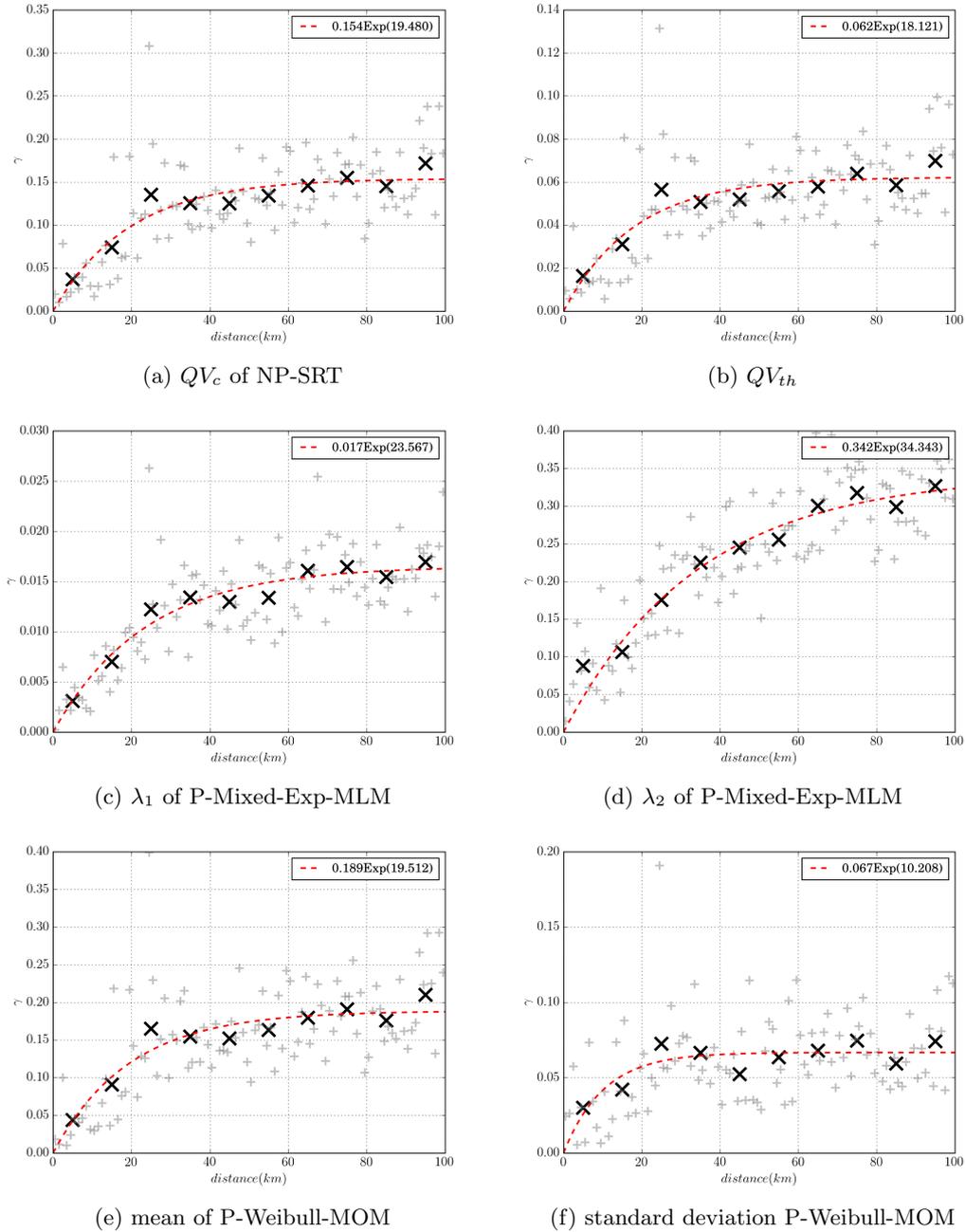


Figure S4: Winter calibration sample 2 for 1H. The black crosses represent the empirical variogram values of the 10 km distance classes, which are used for the least squares fit. The grey crosses represent the empirical variogram values of 1 km distance classes.

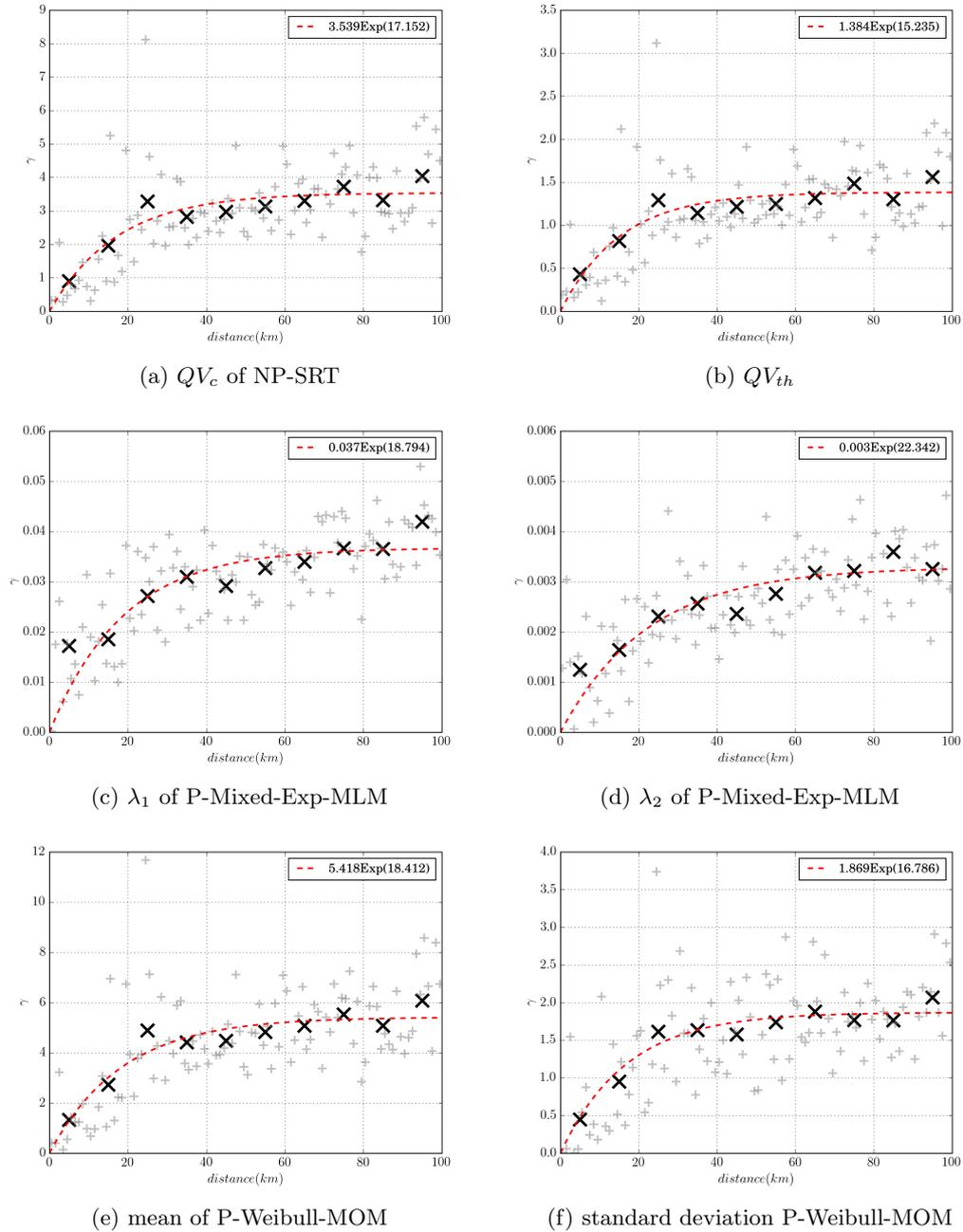


Figure S5: Winter calibration sample 2 for 12H. The black crosses represent the empirical variogram values of the 10 km distance classes, which are used for the least squares fit. The grey crosses represent the empirical variogram values of 1 km distance classes.

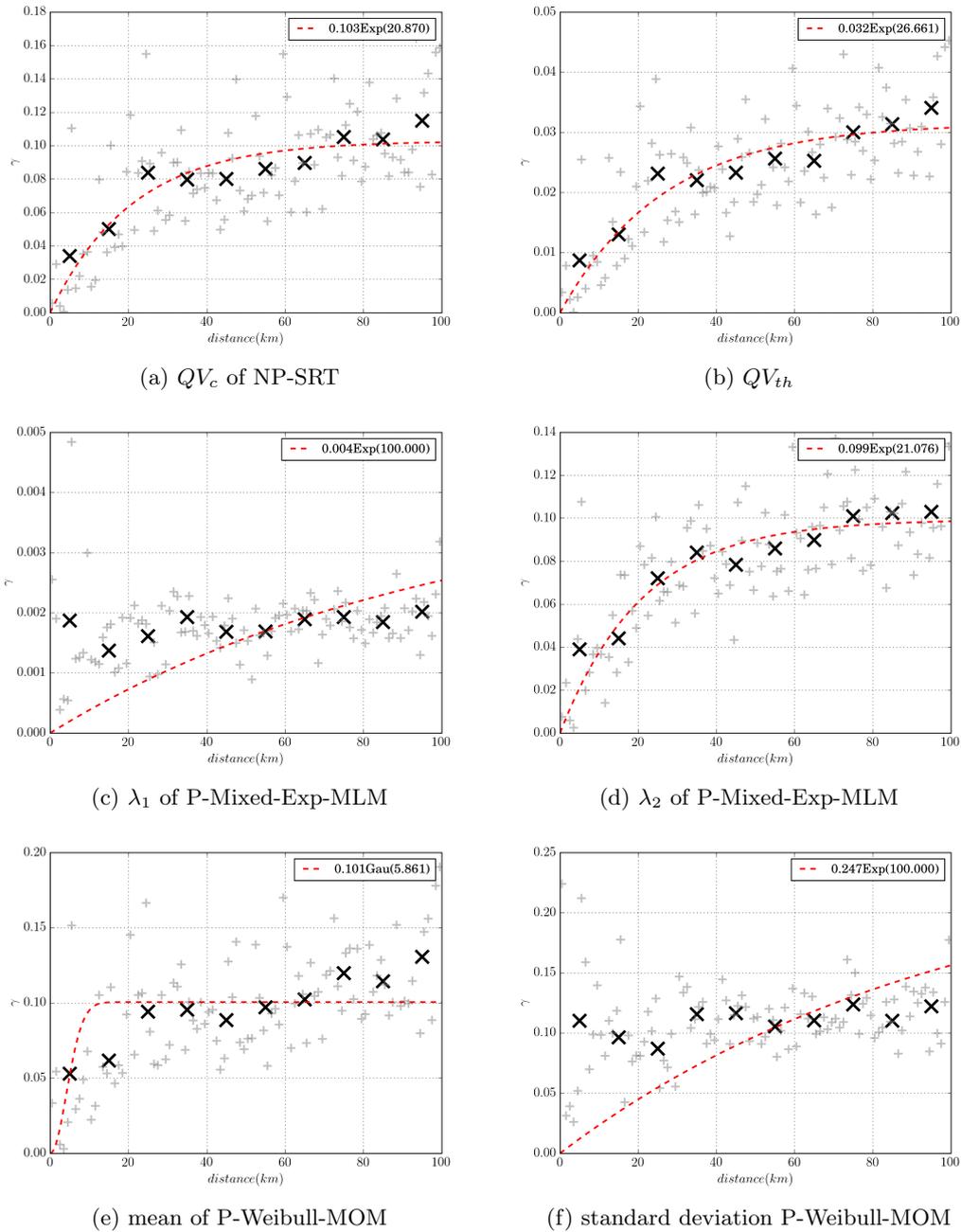


Figure S6: Summer calibration sample 2 for 1H. The black crosses represent the empirical variogram values of the 10 km distance classes, which are used for the least squares fit. The grey crosses represent the empirical variogram values of 1 km distance classes.

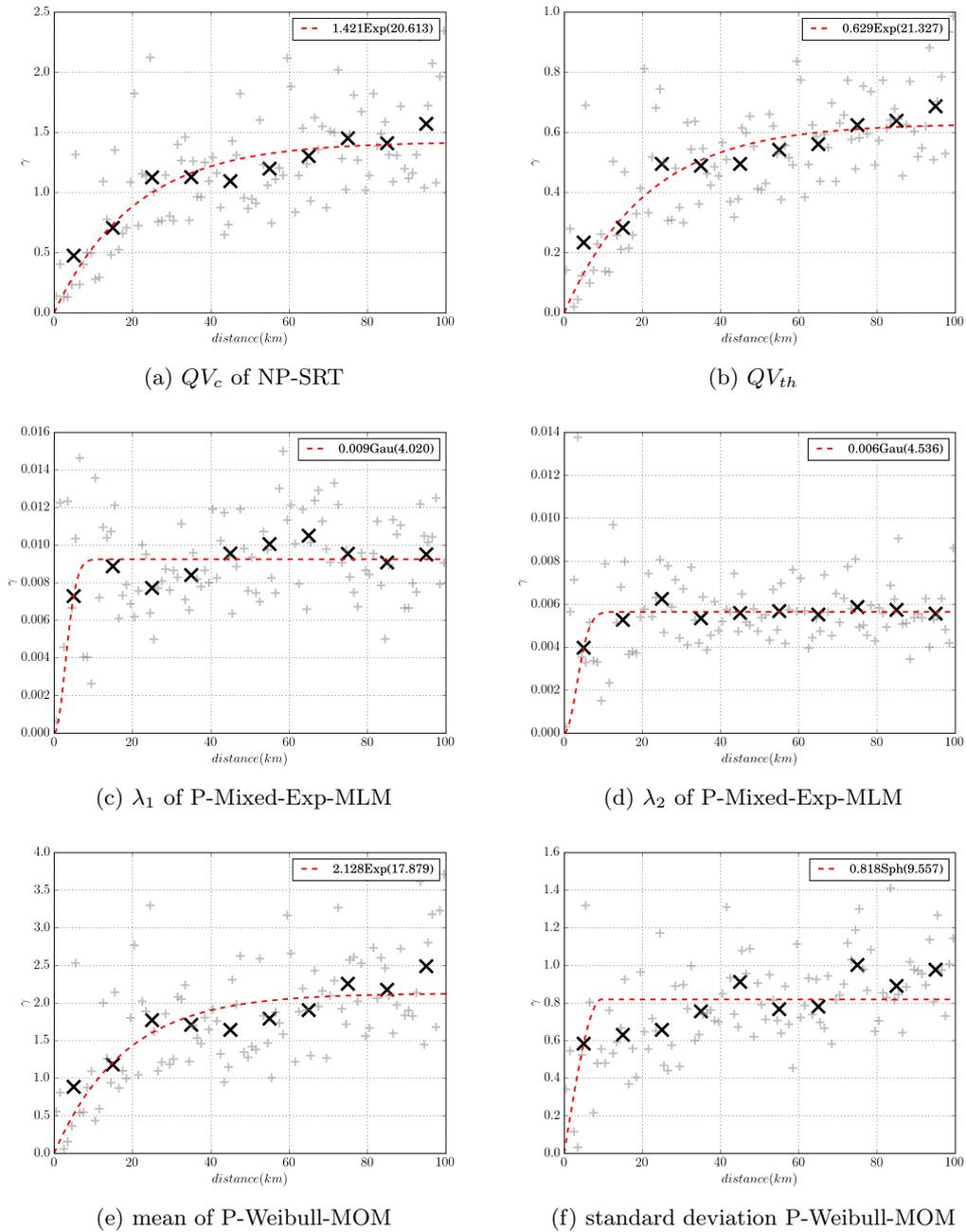


Figure S7: Summer calibration sample 2 for 12H. The black crosses represent the empirical variogram values of the 10 km distance classes, which are used for the least squares fit. The grey crosses represent the empirical variogram values of 1 km distance classes.