

Quantile-based short-range QPF evaluation over Switzerland

JOHANNES JENKNER^{1*}, CHRISTOPH FREI² and CORNELIA SCHWIERZ³

¹Institute for Atmospheric and Climate Science, ETH Zurich, Switzerland

²Federal Office of Meteorology and Climatology, MeteoSwiss, Zurich, Switzerland

³Institute for Climate and Atmospheric Science, University of Leeds, UK

(Manuscript received March 10, 2008; in revised form October 6, 2008; accepted October 6, 2008)

Abstract

Quantitative precipitation forecasts (QPFs) are often verified using categorical statistics. The traditionally used 2×2 contingency table is modified here by applying sample quantiles instead of fixed amplitude thresholds. This calibration is based on the underlying precipitation distribution and has beneficial implications for categorical statistics. The quantile difference and the debiased Peirce skill score split the total error into the complementary components of bias and debiased pixel overlap. It is shown that they provide a complete verification set with the ability to assess the full range of rainfall intensities. The technique enables the potential skill in a calibrated forecast to be estimated without spurious influences from the marginal totals and the problem of hedging is therefore avoided. To exemplify the feasibility of quantile-based contingencies, the method is applied to 6.5 years of operational rainfall forecasts from the Swiss Federal Office of Meteorology and Climatology (MeteoSwiss). Daily accumulations of the COSMO model at 7 km grid size are compared to a high-quality gridded observational record of spatially interpolated rain gauge data. The quantile-based scores are applied to single grid points and to predefined regions. A high-resolution error climatology is then built up and reviewed in terms of typical error characteristics in the model. The seasonal QPF performance exhibits the most severe overestimation over the Northern Alps during winter, indicative of the impact of the model ice phase. The QPF performance related to model updates, such as the introduction of the prognostic precipitation scheme, is also evaluated. It is demonstrated that the potential skill continuously increases for subsequent versions of the COSMO model. Over the entire time period, a strong gradient of the debiased Peirce skill score is evident over the Alps, meaning that the potential skill is much higher on the Alpine south side than on the north side.

Zusammenfassung

Zur Verifikation von Quantitativen Niederschlagsvorhersagen (QNV) werden häufig kategorische Fehlermaße verwendet. Die traditionellerweise benutzte 2×2 Kontingenztafel wird hier durch die Anwendung von Quantilen anstelle von festen Amplitudenschwellwerten modifiziert. Diese Kalibrierung orientiert sich an der zugrundeliegenden Niederschlagsverteilung und beeinflusst kategorische Fehlermaße vorteilhaft. Die Quantildifferenz und der angepasste Peirce Skill Score teilen den Gesamtfehler in die sich gegenseitig ergänzenden Komponenten des Bias und der angepassten Gitterpunktsüberlappung auf. Es wird gezeigt, dass sie eine vollständige Verifikationsbasis bilden, die den gesamten Bereich an Niederschlagsintensitäten abdecken kann. Die Methodik erlaubt es, den potenziellen Skill in einer kalibrierten Vorhersage ohne störende Einflüsse der Randverteilungen zu bestimmen und umgeht dadurch die Problematik des "Hedging". Um die Einsetzbarkeit von quantilsbasierten Fehlermaßen zu demonstrieren, wird die Methode auf 6.5 Jahre an Niederschlagsvorhersagen vom Schweizer Bundesamt für Meteorologie und Klimatologie (MeteoSchweiz) angewendet. Tägliche Summen aus dem COSMO-Modell mit einer horizontalen Auflösung von 7 km werden mit einem hochwertigen gegitterten Beobachtungsdatensatz verglichen. Die quantilsbasierten Fehlermaße werden auf einzelne Gitterpunkte und auf festgelegte Gebiete angewendet. Eine hochaufgelöste Fehlerklimatologie wird erstellt und im Hinblick auf typische Fehlercharakteristika im Modell untersucht. Die jahreszeitlich gemittelten QNV-Fehler zeigen die größte Überschätzung über den nördlichen Alpen im Winter, was auf den Einfluss des Eisschemas hinweist. Die QNV-Fehler während verschiedener Phasen der operationellen COSMO-Modellentwicklung, wie z.B. der Einführung des prognostischen Niederschlagsschemas, werden ebenfalls quantifiziert. Dabei wird gezeigt, dass sich der potenzielle Skill kontinuierlich verbessert hat. Über die ganze Periode fällt ein großer Fehlergradient des angepassten Peirce Skill Scores auf, was bedeutet, dass der potenzielle Skill auf der Alpensüdseite viel höher ist als auf der Nordseite.

1 Introduction

Precipitation forecasts are of societal, economic, and social interest and decision making often relies on accurate rainfall predictions. Hence, there is a great deal

of research activity to improve Quantitative Precipitation Forecasting (QPF) and weather centers continuously evaluate their operational high-resolution limited-area models (LAM) to trace error sources.

QPF is particularly challenging over complex terrain. The Mesoscale Alpine Programme (MAP, BENOIT et al., 2002) provided many new insights into the dynamics and challenges of predicting orographic precipitation (e.g. RICHARD et al., 2007; ROTUNNO and HOUZE,

*Corresponding author: Johannes Jenkner, Institute for Atmospheric and Climate Science, ETH Zurich, Universitätsstrasse 16, ETH Zentrum, CH-8092 Zurich, Switzerland, e-mail: johannes.jenkner@alumini.ethz.ch

2007, and references therein). A variety of factors determine the formation of heavy precipitation and influence QPF quality in Numerical Weather Prediction (NWP) models. They include: (a) the accuracy of the synoptic-scale upper-level triggers of precipitation and their correct mesoscale interaction with steep orography (e.g. FEHLMANN and QUADRI, 2000; MARTIUS et al., 2006), (b) the representation of the orography (ACCADIA et al., 2005) and model resolution (BUZZI et al., 2004; ZÄNGL, 2007), (c) the low-level moisture field (MARTIUS et al., 2006; MAHONEY and LACKMANN, 2007) and the boundary layer structure (ROTUNNO and HOUZE, 2007), (d) the enhancement of precipitation by microphysical processes (e.g. ZENG et al., 2001; PUJOL et al., 2005) and turbulence (HOUZE and MEDINA, 2005).

Also the formulation and numerics of the NWP system itself can influence QPF. KAUFMANN et al. (2003) evaluated the former hydrostatic LAM of the COSMO consortium¹ and found that the model error is exceedingly sensitive to the activation of convection in parameterized precipitation and the sloping topography in resolved precipitation. KÅLLBERG and MONTANI (2006) compared a non-hydrostatic versus a hydrostatic model which also differed in the data-assimilation schemes and numerics. Above all, they found more intense precipitation extremes in the non-hydrostatic formulation due to dissimilarities of the convection schemes. The prognostic treatment of precipitation markedly improved the wet (dry) bias on the windward (downwind) side of orography which is typically noted for QPF over complex terrain (e.g. ELEMENTI et al., 2005).

From a more generic point of view, HOHENEGGER and SCHÄR (2007) investigated the dynamics of error growth. They used a range of different initial perturbation procedures of a high-resolution ensemble and found a rapid radiation of the initial uncertainties throughout the computational domain and a further amplification over moist convectively unstable regions (compare also HOHENEGGER et al., 2006).

It is clear from the above that the problem of identifying sources of QPF errors is highly complex. An additional drawback arises from the fact that most of our knowledge today is based on case studies or on relatively short-term periods (up to 1-2 years) of investigation. However, QPF quality tends to be more case-dependent than model-dependent (RICHARD et al., 2003). A consistent quantification and rigorous investigation of QPF over a longer-term period is highly desirable to trace QPF errors and identify the main model shortcomings. In view of this, novel verification techniques for precipitation forecasts are currently being developed. Most notably, spatial approaches are able to consider different areal error aspects. An example is the intensity-scale technique (CASATI et al., 2004; MITTERMAIER, 2006) which diagnoses skill as a function of “precipitation rate

Table 1: 2x2 Contingency table with hits H, misses M, false alarms F and correct negatives Z.

	observed yes	observed no
predicted yes	H	F
predicted no	M	Z

intensity” (CASATI et al., 2004) and spatial scale of the error. Another example is the object-based quality measure SAL (WERNLI et al., 2008) which assesses the coherence in structure, amplitude and location of precipitation objects. As a matter of fact, the findings of these and comparable methods are only representative for a predefined verification domain. Hence the attribution of individual verification aspects to specific locations or sites is hardly possible. Pertinent information is therefore lost over mountainous terrain where the spatial variability of atmospheric parameters usually is large (FREI and SCHÄR, 1998; SCHMIDLI et al., 2002).

In this context, the traditional grid-point verification still provides an expedient alternative. As explained by MURPHY and WINKLER (1987) or MCBRIDE and EBERT (1999), categorical statistics are traditionally used for dichotomous verification purposes such as the validation of precipitation. After the selection of a threshold value it is tested whether model and observations exceed the limit. Most standard methods rest upon a conventional 2x2 contingency table² (Tab. 1) consisting of the number of hits H, misses M, false alarms F and correct negatives Z. From these four entries, several error measures (or scores) such as frequency bias, probability of detection (POD), false alarm ratio, probability of false detection (POFD), threat score, equitable threat score, Peirce skill score (PSS=POD-POFD, PEIRCE, 1884), odds ratio skill score and others can be derived (e.g. MASON, 2003).

A desired property of categorical measures is equitability as defined by GANDIN and MURPHY (1992). The definition of equitable measures implies that random and constant forecasts are treated identically. A necessary prerequisite is that the scoring rule extracts the joint distribution, defined by forecasts and observations together, and disregards the marginal distributions, defined separately by forecasts and observations. In other words, equitable scores display unvarying expected values irrespective of varying marginal totals. In this sense, only equitable scores guarantee a fair comparison of samples with different characteristics. A simple example for a non-equitable score is the POD which exhibits higher values for random forecasts of frequent events than for those of rare events. Non-equitable measures can upgrade forecasts which are not consistent i.e. which do not correspond to the forecaster’s judgment (see MURPHY, 1993, for explanation). Consequently, both model developers and operational forecasters are tempted by hedging (MURPHY and EPSTEIN, 1967) and

¹web site: www.cosmo-model.org

²extensions to multiple categories are straightforward

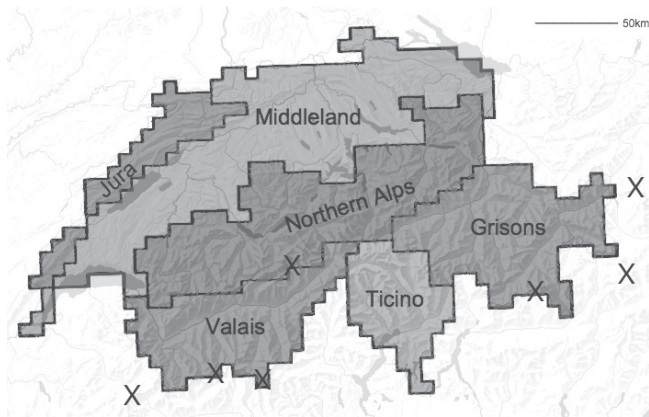


Figure 1: Subdivision of Switzerland into six orographically distinct areas: Jura, Middleland, Northern Alps, Valais, Ticino, Grisons. Square dividing lines indicate the margins of the COSMO grid points. Highest mountain formations are marked with an X.

falsely adjust the number of forecasted events to achieve a better scoring.

So far, opinions about equitability substantially diverge in the literature. As proven by GANDIN and MURPHY (1992), the PSS (which is equivalent to the Hanssen-Kuipers discriminant, HANSEN and KUIPERS, 1965), constitutes the unique solution (i.e. the only one possible) which satisfies the strict conditions for equitability, namely an invariable rating of both random and constant forecasts. In fact, the uniqueness of the PSS still applies to conditions with an unequal penalty for misses and false alarms (MANZATO, 2005). However, the uniqueness is lost, if the conditions for equitability are slightly relaxed, i.e. constant forecasts are allowed to score unlike random forecasts (MARZBAN and LAKSHMANAN, 1999). In theory, scores can be equitable or nearly equitable (in the sense that random/constant forecasts do not vary or only slightly vary) and provide a verification free from any computational complications. In practice, the usefulness of measures essentially varies in different constructed and real cases (e.g. MARZBAN, 1998; HAMILL and JURAS, 2006). On the one hand, it is not feasible to compare forecasts for different event frequencies or base rates (e.g. MASON, 2003). The so-called base-rate error is usually neglected, but it lowers the displayed performance of finite-accuracy forecasts for rare events (MATTHEWS, 1996). Most scores are deeply affected by this dilemma, but WOODCOCK (1976) and MASON (1989) find the PSS to be unaffected. Indeed, THORNES and STEPHENSON (2001) only grant the odds ratio skill score and the PSS to cope with small base rates. On the other hand, common scores depend on the bias which HILLIKER (2004) exemplifies by means of the threat score. The PSS approaches the POD in rare-event situations and then is spuriously altered by the bias. Consequently, the PSS exhibits a distinct optimal threshold probability which maximizes its expected value (MASON, 2003) and the PSS is prone to hedging.

MESINGER (2008) points out that the impact of the bias customarily is estimated in a subjective manner, because objective approaches are still missing. He introduces a method to debias the threat score and likewise the equitable threat score. To this end, he converts the standard contingency table to a setting with a unit bias and recommends it to use for an objective verification of the placing of precipitation systems. Although Mesinger’s approach alleviates some of the basic problems of verifying biased distributions, a drawback remains, in that assumptions need to be made which cannot be derived from the underlying distributions alone.

The purpose of our study is twofold. At first, we introduce a refined grid-point verification measure which fulfills the requirements of equitability without making any additional assumptions. Thereby, we make use of the definition of a quantile which is equivalent to the terms of percentile and fractile (WILKS, 2006). Note that the 90 % quantile for example is exceeded every tenth time, meaning that it cuts off the tenth most extreme part from the rest of a dataset. Then, we apply our method to quantify bias and potential skill in a long-term climatology of operational high-resolution precipitation forecasts over complex terrain. The resulting error climatology allows for an extensive model diagnosis to help to identify possible error sources. The dataset is long enough to investigate data subsets, such as seasonal error variations and chronological error evolutions caused by different operational model versions.

The structure of the paper is the following: Some background information about the observational analysis, the verified COSMO model and the geographic setting is provided in Section 2. The refined verification methodology is derived and discussed in Section 3. Then we turn to present the verification results. Section 4 pinpoints seasonal error variations and Section 5 highlights characteristics of different model versions. Finally results are discussed and synthesized in Sections 6 and 7. The mathematical details of the verification methodology and related derivations and discussions are presented in the appendices.

2 Verification data, model and domain

2.1 Observational analysis

The reference dataset used for the model evaluation in this study is a gridded mesoscale analysis for Switzerland, which is derived from rain gauge observations by spatial interpolation. The construction of gridded values is done identically to the daily analysis for Germany used in PAULAT et al. (2008). The underlying observation network encompasses typically 450 stations in Switzerland, corresponding to an average station distance of 10–15 km (KONZELMANN and WEINGARTNER, 2007). The Alpine in-situ observations are among the densest world-wide in high-altitude topography (FREI and SCHÄR, 1998). Despite slight variations

in the reading times across the network, the analyses can be considered as representing 24-hour totals from 06:00 until 06:00 UTC (FREI and SCHÄR, 1998).

The spatial analysis of rain gauge observations is conducted with a modified version of the SYMAP algorithm, a distance and direction weighting scheme by SHEPARD (1984) (see also WILLMOTT et al., 1985). In deviation from the traditional scheme, our procedure encompasses an antecedent climatological scaling of station observations and subsequent re-scaling of the gridded anomalies with a high-resolution climatology (SCHWARB et al., 2001). The procedure is similar to that of WIDMANN and BRETHERTON (2000) and is applied to reduce systematic errors due to biases in the distribution of stations with height. In our procedure, the SYMAP algorithm is applied with a different distance weighting scheme. The purpose is to represent, in the analysis, regional area mean values rather than point values (see FREI and SCHÄR, 1998). The adopted analysis method is similar to that applied in FREI et al. (2006), except that the analysis is undertaken originally on a 2 km grid and is subsequently aggregated to the grid of the NWP model.

It should be noted that the verification dataset is affected by systematic biases in the rain gauge measurements. In Switzerland rain-gauge under-catch is expected to range from about 4 % at low altitudes in summer to more than 40 % above 1500 m MSL in winter (SEVRUK, 1985) and has to be considered while interpreting the rainfall bias later on. Moreover, there may be systematic inconsistencies between the model and verification grids as a result of differences in effective resolution. From the observation network we expect an effective resolution of the verification dataset in the order of 15 km in the flatlands and 15–25 km in the mountains. This is close but slightly coarser than the model resolution when taking two nominal grid pixels, i.e. $2 \times 7 \text{ km} = 14 \text{ km}$, as the effective model resolution.

The error climatology derived in this study is based on the non-hydrostatic model of the COSMO consortium in operation at the Swiss Federal Office of Meteorology and Climatology (MeteoSwiss). The model is the Swiss counterpart of the former German Lokal Modell (LM) described briefly in STEPPELER et al. (2003). It has been developed for the purpose of high-resolution weather forecasts with a preferable representation of the meso- γ scale. A detailed description currently can be found on www.cosmo-model.org. The horizontal mesh size of the model is 7 km and the vertical resolution constitutes 45 hybrid height-based levels. The grid structure is based on an Arakawa C setup with Lorenz vertical grid staggering. The basic equations are solved in a fully elastic manner with a dry reference state at rest. Advection is treated by a split explicit scheme based on a filtered leapfrog time integration (ASSELIN, 1972; SKAMAROCK and KLEMP, 1992) with a main time step of 40 s. Moist convection is treated by the mass flux con-

vergence scheme of TIEDTKE (1989). The warm-rain regime refers to a bulk water-continuity model proposed by KESSLER (1969) whereas the ice-cloud regime is based on an extended saturation adjustment technique (LORD et al., 1984).

The preoperational phase of the model at MeteoSwiss started in July 2000 and the model became operational in April 2001. The geographical domain comprises central Europe and receives boundary conditions from the GME model of the German Meteorological Service and later from the Integrated Forecast System (IFS) of the ECMWF. The updating of the boundaries is treated with an adjusted Davies relaxation scheme (DAVIES, 1976). In the present study, daily precipitation accumulations of the period between July 2000 and December 2006 are analyzed over Switzerland. Investigations are confined to operational 00 UTC model runs, from which daily sums are constructed using lead times between 6 and 30 hours.

Since both preoperational and operational forecasts are evaluated, there are various model updates within the considered period. Altogether, there are roughly 60 changes of the model setup including bug corrections. Most of them are not expected to have a significant impact on daily rainfall fields. However, few upgrades are considered to influence QPF quality decisively. We refer to the most important changes in Section 5.

2.2 Geographic aspects

Altogether 859 grid points of the COSMO model are evaluated within the borders of Switzerland ($\sim 41000 \text{ km}^2$). To highlight regional disparities with respect to QPF quality, the whole area is subdivided into six orographically distinct parts (Fig.1) using contour lines of the model topography as well as other landmarks. The area of the Jura comprises the Swiss part of the Jura mountains as well as some adjacent pixels in the canton Jura. The low elevations in the Middleland reach from Lake Geneva in the southwest to Lake Constance in the northeast. The hilly relief only varies slightly in height here. The northern Alpine crest as well as the approximate canton borders between the Valais, the Ticino and the Grisons split up the Alps into four smaller mountainous domains. The Valais and the Grisons can be regarded as central Alpine domains, whereas the Northern Alps and the Ticino lie on opposing sides of the Alpine crest. Overall highest elevations (around 3000 m MSL in the model topography) are found in the middle of the northern Alpine crest, along the southern border of the Valais and in the south of the Grisons. They are marked with an X in Fig.1. Individual regions contain the following numbers of COSMO grid points: Jura 55, Middleland 271, Northern Alps 208, Valais 114, Ticino 71, Grisons 140.

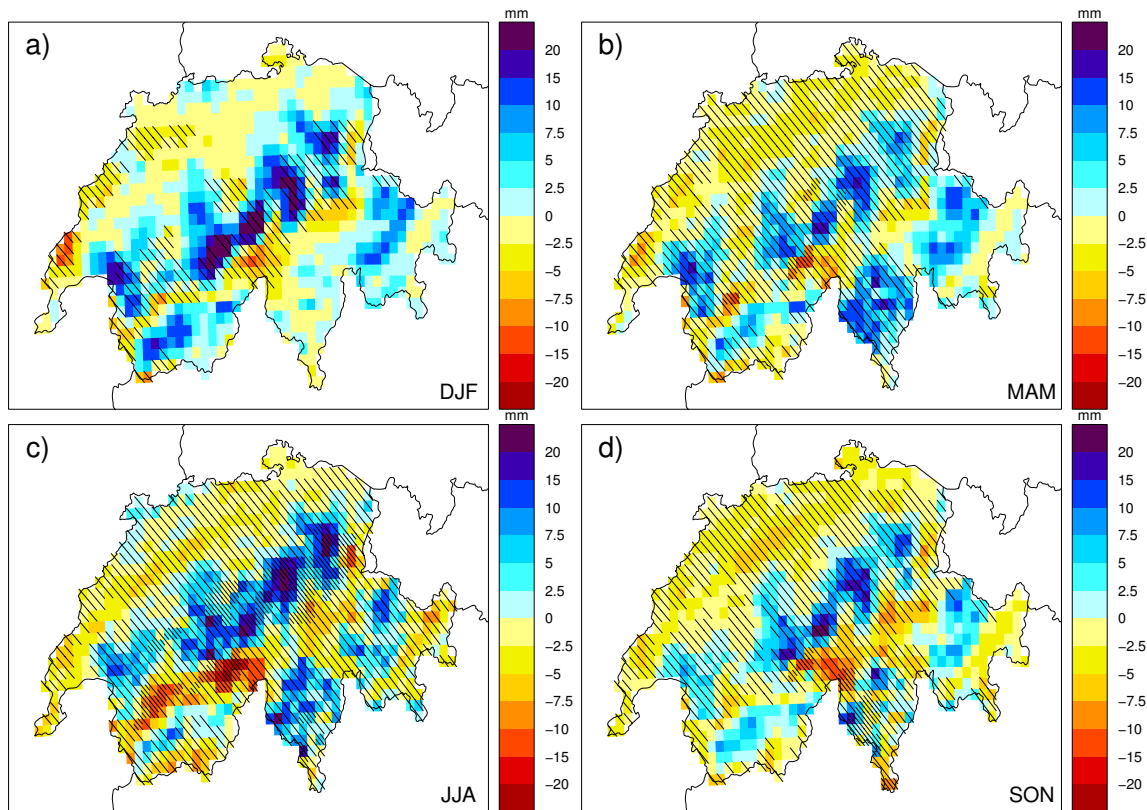


Figure 2: Difference [mm/day] of 90 % quantile values (COSMO–observations) for different seasons: a) winter, b) spring, c) summer and d) fall. Lightly hatched grid points indicate observed 90 % quantiles over 10 mm/day and densely hatched grid points indicate observed 90 % quantiles over 20 mm/day.

3 Methodology of verification

The standard 2x2 contingency table (Tab. 1) for dichotomous events (relating to a predefined precipitation threshold) contains the following entries: hits H, misses M, false alarms F and correct negatives Z. The four numbers are used to compute several different categorical scores (see Section 1 for examples). Varying the threshold from low (representing weak precipitation rates) to high values (representing strong precipitation rates) renders possible a comprehensive investigation of different intensities.

However, there are some shortcomings, if standard categorical statistics are applied directly. Firstly, a predetermined amplitude threshold splits the precipitation distributions at an unknown percentile, i.e. it is not obvious a priori whether the threshold represents common or rare events within the considered sample. In an extreme case, single cells of the contingency table can become zero. Then some scores cannot be computed (due to a division by zero) and statements about model behavior are hard to make. Secondly, the distributions under comparison usually differ considerably with respect to their range of values and exhibit a distinct offset/bias. Customary scores do not fulfill the requirements of equitability (GANDIN and MURPHY, 1992) or fail to be firm with respect to hedging (STEPHENSON, 2000). Thirdly,

the joint distribution comprises three degrees of freedom, if the four entries are only linked to the sample size. Three scores are required to display all verification characteristics. STEPHENSON (2000) proposes the triplet of odds ratio skill score, PSS and frequency bias. According to his comments, it is possible to draw complementary information out of the considered datasets, if concurrent scores are applied simultaneously. But it remains ambiguous, how to attribute individual verification aspects to these measures which project onto each other. Fourthly, we argue that it is not possible to integrate amplitude-based scores over a range of intensities and condense forecast performance in such a way. In principle, it is not meaningful to average scores for different thresholds, because it is not obvious how many data points fall within a certain range of thresholds.

In this study, a refined version of categorical statistics is proposed to address the above problems and avoid confusing verification results. The contingency table is defined by means of frequency thresholds instead of amplitude thresholds. To this end, sample quantiles are computed for both data records independently. Thus, the two datasets are compared using the same relative cut-off (according to the definition of a quantile) within each distribution. In other words, a nonlinear calibration is performed similar to CASATI et al. (2004) and the bias is omitted automatically. The full derivation is outlined

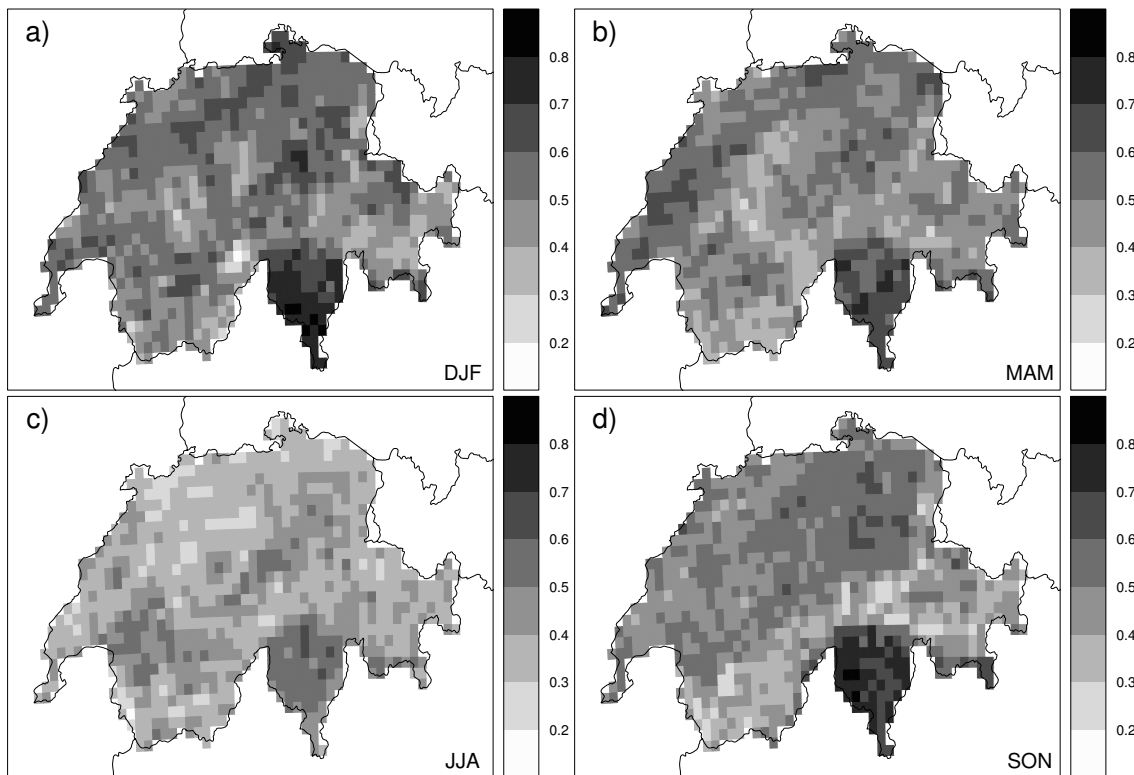


Figure 3: Peirce skill score for 90 % quantiles of daily precipitation sums [1: perfect, 0: random forecast]: a) winter, b) spring, c) summer and d) fall.

in Appendix A. By using quantiles, only one degree of freedom is left within the contingency table (Appendix A). It is hence sufficient to acquire a single entry of the contingency table to determine the joint distribution.

Regardless of the contingencies, the bias in the quantile-based framework can simply be assessed by the absolute or relative quantile difference QD or QD' (Appendix B). Contrary to the approach of FERRO et al. (2005), quantile differences are not transformed here, but reveal differences of rainfall distributions both in location and scale. For instance, an underforecasting of light rainfall can be easily distinguished from an overforecasting of heavy rainfall. The overall offset, covering the entire distribution, can be reproduced by a weighted integral over the absolute value of relative quantile differences ($\overline{QD'}$, Appendix B, Eq. 2.3). In this way, amplitude errors are summarized and dissimilarities of the entire distributions are quantified.

As noted earlier, the PSS is able to measure skill without being perturbed by the base rate (e.g. WOODCOCK, 1976; MASON, 1989). The sensitivity to hedging increases with the rarity of events (STEPHENSON, 2000), as the PSS converges towards the POD. However, quantile-based contingency tables overcome possible biases. If the number of predicted events is modified in the model output or later in the issued forecast, the definition of a quantile requires that it is compared

to the same number of observed events. Consequently, quantile-based scores circumvent the problem of hedging (Appendix C). In respect thereof, the quantile-based PSS is ideally suited to display the potential accuracy, i.e. the potential skill, for a calibrated forecast. In other words, the debiased PSS measures the pixel overlap, i.e. the matching or the shift of pixels, in an amplitude balanced setup (Appendix C). Due to equal marginal totals (Appendix A, Eq. 1.1), the debiased PSS merges into a pure ranking of misses (Appendix C, Eq. 3.2). If the proportion of misses with respect to their random expectation remains fixed, all forecasts therefore receive the same constant rating. Given that a certain ratio fulfills the definition of a constant forecast, it is permissible to compare PSS values for different quantile probabilities or base rates directly. In spite of this ability, it has to be kept in mind that the PSS measures skill in reference to random chance which might not always be appropriate for extreme events (compare with STEPHENSON et al., 2008). Note that the advantages over the conventional formulation in an uncalibrated setting can be illustrated by means of a simple forecast example (Appendix D and Fig. 11). The debiased PSS focuses exclusively on a shift error and therefore is easier to interpret than its conventional counterpart. Similarly to the QD, the weighted integral of the debiased PSS over the whole range of quantiles (\overline{PSS} , Appendix C, Eq. 3.5) summarizes shift errors

over the entire range of intensities and the overall matching characteristics can be condensed in a single number.

Following the above derivation, we are left with two complementary measures QD (or QD') and debiased PSS which concertedly characterize the overall forecast error, namely bias and debiased overlap/matching error. Accordingly, the sample uncertainty also subdivides into the two parts (Appendix E). From the course of the confidence intervals in Fig. 12, it can be gleaned that the quantile-based PSS remains much better defined for rare events than the conventional amplitude-based PSS. Quantile probabilities are not affected by amplitude uncertainties, but their transformation to precipitation amounts is. We can achieve a high confidence for the PSS value of a certain quantile, but still hold a low confidence for the precipitation estimation of the quantile. Some applications require a link to rainfall amounts whereas others only require a link to frequencies, i.e. return periods, which are automatically provided by quantiles.

4 Seasonal error climatology

In the following, we apply the proposed verification framework to the 6.5-year forecast climatology. Initially, the focus is on seasonal error discrepancies. Quantile differences and debiased PSS values are computed for all seasons independently. First of all, error charts for the 90 % quantiles are discussed revealing gridbox-scale error variations. The corresponding thresholds are exceeded on every tenth day and expected to reflect universal QPF errors of reasonably strong events. Then, the grid-point based scores are aggregated, i.e. applied to regions as a whole. Verification measures are computed for the compound data (grid points \times days) of our six predefined domains (Fig.1) and discussed with respect to the whole range of intensities. Finally, the regional model performance is reviewed with the aid of integral error values summarizing the overall performance.

4.1 Grid-point based verification of 90 % quantiles

Throughout all seasons, a strong wet bias is evident over the Northern Alps (Fig.2). During the course of the year, the observed 90 % quantiles roughly account for values between 10 mm/day (winter) and 25 mm/day (summer). However, they are more than 20 mm/day or locally 30 mm/day higher in the model forecasts. This bias is confined to few grid points during spring and fall but affects much larger areas during winter and summer. In addition to the Northern Alps, some parts of the Valais and the Grisons are also significantly overforecast. Even though the observed 90 % quantiles mostly remain below 10 mm/day during winter, spring and fall, the overestimation in the model partly amounts up to +20 mm/day at these times of the year. Thus, the relative amplitude

errors are largest by far in these areas. In the Ticino, the bias is similar to the Northern Alps during spring and summer, but it almost vanishes during winter and only is noticeable at few grid points during fall. A closer look reveals that the bias is linked directly to the hillside of the topography. In the northern and central parts of the Alps, northwest aligned slopes are strongly overforecast whereas southeast aligned slopes are weakly or moderately underforecast. In the Ticino, the strongest overprediction resides over southerly oriented slopes, but only appears clearly during spring and summer (Figs. 2b, c) and at specific grid points during fall (Fig. 2d). The deep valleys in the interior of the Alps (mainly the upper Rhône and Rhine Valleys) and their direct surroundings are clearly underforecast. The strongest dry bias is found directly to the south of the northern Alpine crest. It is most pronounced during summer with peak values around -20 mm/day (Fig. 2c).

Regardless of the bias, seasonal matching characteristics vary substantially within Switzerland (Fig. 3). Generally, the regional PSS pattern is noisiest during summer and smoothest during fall. During winter, the matching clearly is linked to the topography. Western slopes display much higher PSS values than their eastern counterparts at this time. In particular, there is a distinct spatial PSS minimum (Fig. 3a) to the southeast of the highest point in the northern Alpine crest (marked in Fig. 1). Noteworthy, the best pixel overlap is found in the Ticino throughout the year. Both during winter and fall, the debiased PSS locally reaches values over 0.8 in the Ticino corresponding to a debiased POD over 0.82 (see Appendix C, Eq. 3.4 for this calculation). All over Switzerland, the poorest matching is detected during summer (Fig. 3c). The PSS values barely pass 0.6 in the Ticino and some places elsewhere and only vary between 0.2 and 0.4 in the Middleland, the Valais and the Grisons. Note that a debiased PSS of 0.2 only implies a debiased POD of 0.28 (Appendix C, Eq. 3.4). During spring and fall, the matching is different on both sides of the Alps. The pixel overlap in the Ticino is poorer during spring compared to fall, even though it is still superior to other regions. In contrast, it is slightly improved in the Valais and the Grisons during spring compared to fall. On the Alpine north side, there are distinct sectors with PSS values around 0.35 and 0.55 during spring, whereas the matching is surprisingly uniform with values around 0.5 during fall. The clearest and most consistent regional separation is found late in the year. During fall, the pixel overlap is superior in the Ticino, inferior in the Valais and in the Grisons and middle-rate further to the north. Simultaneously, a prominent PSS gradient is present to the north of the Ticino and over the northern Alpine crest.

4.2 Regional verification of quantile courses

To detect possible error variations for different intensities, it is necessary to consider the whole range of quan-

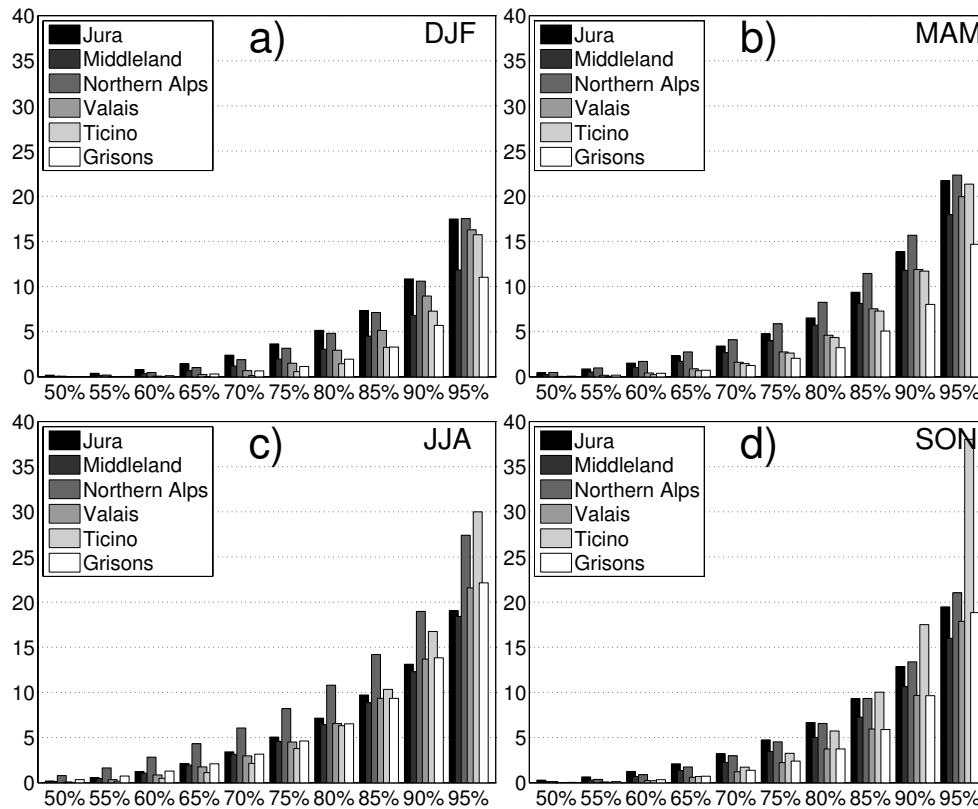


Figure 4: Observed quantile values [mm/day] for discrete quantile probabilities: a) winter, b) spring, c) summer and d) fall.

Table 2: Weighted integral \overline{QD} values for different domains and seasons (winter, spring, summer, fall). The scoring shows the relative amplitude deviation and is explained in Appendix B [$0.1 \hat{=} 10.5\%$ deviation, $0.2 \hat{=} 22.2\%$ deviation, $0.3 \hat{=} 35.3\%$ deviation]. The darker the shading the higher are overall observed rainfall amounts. CH stands for the whole of Switzerland.

	Jura	Middlel.	N. Alps	Valais	Ticino	Grisons	CH
DJF	0.27	0.13	0.45	0.18	0.15	0.33	0.24
MAM	0.17	0.08	0.25	0.16	0.36	0.22	0.13
JJA	0.10	0.13	0.23	0.24	0.30	0.04	0.11
SON	0.22	0.13	0.19	0.12	0.11	0.22	0.06

tiles. Therefore, our verification measures are applied to entire regions. The respective quantile values of the observations are given in Fig. 4 as reference. Since the 50 % quantiles mostly fall below 0.5 mm/day, dry quantiles beneath are omitted in the graphs. Most notably, the Jura and the Northern Alps consistently display highest quantile values during winter and spring, meaning that rainfall amounts are highest here. Later in the year, the Ticino clearly entails highest quantile values for quantile probabilities above 90 %.

Figure 5 displays the quantile-based measures QD and PSS for different quantiles. Concerning the bias, there is the general tendency that weak intensities are slightly overforecast, medium intensities are slightly under- or overforecast and strong intensities are either strongly overforecast or moderately underforecast. However, distinct regions are exceptions from this archetype. On the

one hand, the following domains are severely underforecast for almost all quantiles: the Jura during winter, spring and fall and the Valais during summer. On the other hand, the following domains are greatly overforecast for all quantiles: the Northern Alps all over the year and the Ticino during spring and summer. Interestingly, heavy precipitation in the Ticino is underforecast during winter and fall. Obviously, the bias on the Alpine south side substantially changes between the two halves of the year.

Concerning the matching, the Ticino mostly holds the highest PSS values compared to the rest of Switzerland. As an indication of good performance, the debiased PSS in the Ticino is almost independent from the quantile. Only above the 90 % quantile, a sharp drop occurs which is most prominent during winter and fall. Low quantiles (50 %–65 %) usually entail a better matching in north-

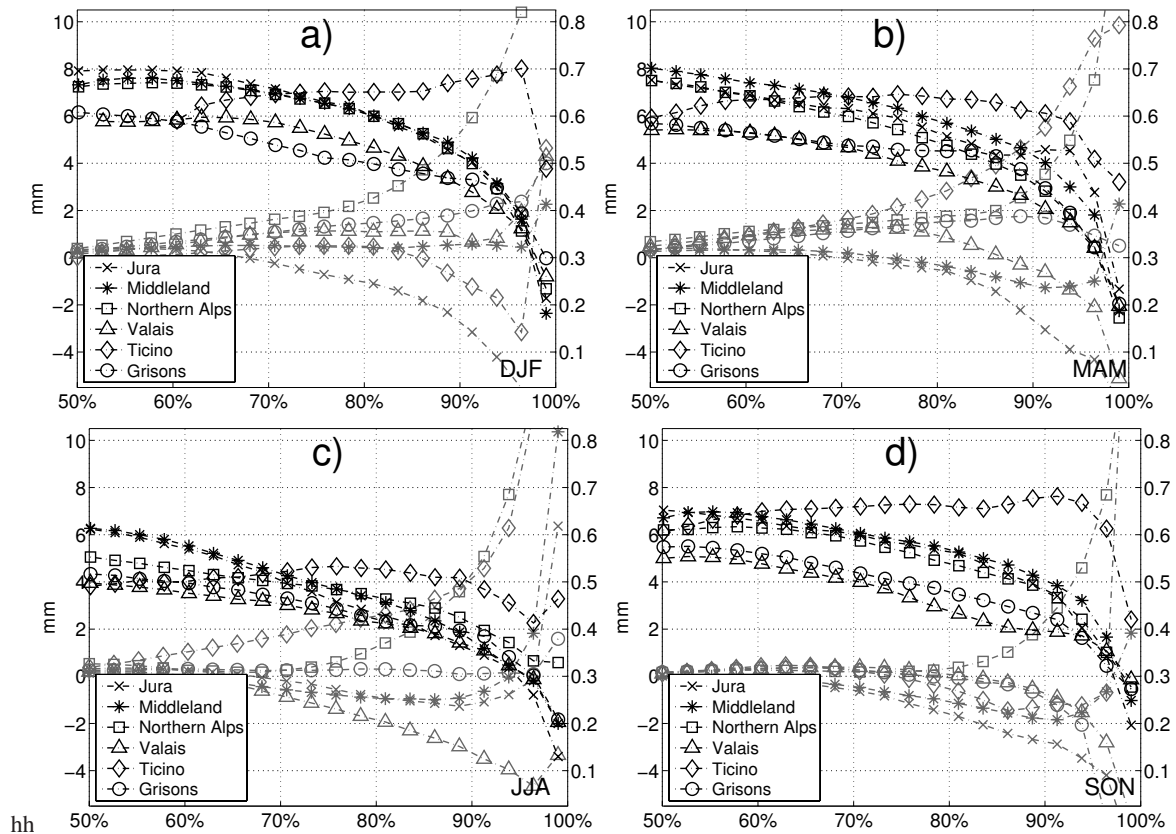


Figure 5: Area-average of the quantile difference (gray, left scale) [mm/day] and the Peirce skill score (black, right scale) for quantiles between 50 % and 99 %: a) winter, b) spring, c) summer and d) fall.

Table 3: Weighted integral $\overline{\text{PSS}}$ values for different domains and seasons (winter, spring, summer, fall). The higher the value the better is the overall matching performance. The darker the shading the higher are overall observed rainfall amounts. CH stands for the whole of Switzerland.

	Jura	Middlel.	N. Alps	Valais	Ticino	Grisons	CH
DJF	0.48	0.46	0.46	0.41	0.63	0.43	0.46
MAM	0.49	0.47	0.43	0.38	0.56	0.41	0.45
JJA	0.34	0.36	0.41	0.34	0.47	0.36	0.39
SON	0.43	0.46	0.43	0.38	0.60	0.38	0.46

ern regions, meaning that light rainfall events are less accurately predicted in the Ticino than elsewhere. As opposed to other regions, the PSS curve is very consistent in the north where the error distributions of the Jura, the Middleland and the Northern Alps are very similar. Only during summer and for low quantiles, PSS values over the Northern Alps diverge slightly from the regions further to the west. In comparison to northern territories, the PSS is up to 0.1 lower/worse in the central Alps. Except for some extreme quantile values, the pixel overlap always is most critical either in the Valais or in the Grisons.

4.3 Rating of integral values

To survey the overall model performance with respect to individual seasons and domains, we discuss weighted error integrals (Tabs. 2 and 3) as a condensed view of Fig. 5 (see Appendix B, Eq. 2.3 and Appendix C, Eq. 3.5 for further details). The integral bias ($\overline{\text{QD}'}$) is highest during winter. $\overline{\text{QD}'}$ constitutes 0.45 over the Northern Alps which is equivalent to an amplitude deviation of more than 50 %. This means in this case that all intensities are overestimated on average by half of their rainfall value. At the same time, $\overline{\text{QD}'}$ amounts to 0.27 and 0.33 over the Jura and over the Grisons, respectively. This is equivalent to an amplitude deviation of more than 30 %. On the Alpine south side, the rela-

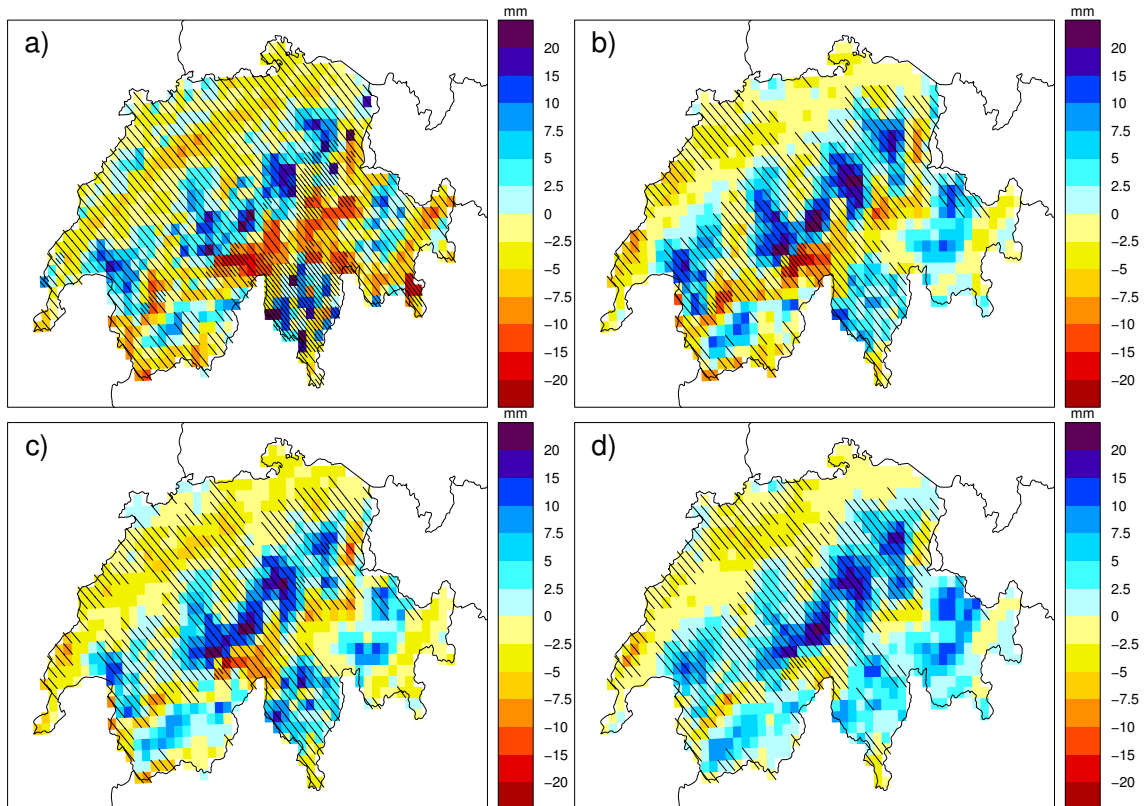


Figure 6: Difference [mm/day] of 90 % quantile values (COSMO-observations) for different model generations: a) IGD, b) NGD, c) NID and d) NIP. Lightly hatched grid points indicate observed 90 % quantiles over 10 mm/day and densely hatched grid points indicate observed 90 % quantiles over 20 mm/day.

tive bias is highest during spring. The Ticino entails a $\overline{QD'}$ of 0.36 between March and May which implies a deviation of more than 40 %. Lowest regional model biases are observed over the Middleland during spring and over the Grisons during summer. Interestingly, $\overline{QD'}$ only accounts for 0.04 over the Grisons during summer, meaning that the rainfall distributions only differ by 4 % in this region. Throughout Switzerland, $\overline{QD'}$ is lowest during fall. However, regional values are partly higher than for the rest of the year, meaning that regional biases cancel each other out in the overall average. The integral matching score (\overline{PSS}) generally is lowest during summer. Throughout Switzerland, the integral value is 0.06 lower between June and August than for the rest of the year. In all seasons, the PSS is highest in the Ticino. Especially during fall, the integral value is almost 0.2 higher here than for the rest of Switzerland. Other outstanding regions are the Jura and the Middleland during spring and to a lesser extent the Northern Alps during winter. These areas exhibit a good matching and rainfall shifts are comparatively small here. In contrast, the overall worst matching is seen in the Jura, the Middleland, the Valais and the Grisons during summer as well as to a slightly smaller extent in the Valais and the Grisons during fall. The lowest value of 0.34 for the Jura and the Valais during summer corresponds to roughly 75 % more mismatching grid points (misses or false alarms)

than for the largest value of 0.63 for the Ticino during winter.

5 Differentiation of model designs

There are three decisive model updates within the operational phase of the COSMO model which are expected to affect QPF to a large extent. First of all, a continuous assimilation cycle (STAUFFER and SEAMAN, 1994) replaced a pure interpolation of initial values from the driving model. Wind, pressure, temperature and humidity have been nudged towards surface and upper-air observations since then. Later, there was a changeover of the driving model at the boundaries. The COSMO model was driven by the GME in the early stages, before the IFS has been used in exchange. The switch-over allowed the COSMO model to benefit from high-quality features of the ECMWF model such as the 4D Var analysis. Finally, the prognostic precipitation scheme substituted a diagnostic treatment with the assumption of column equilibrium for precipitation particles. Hence, hydrometeors have been advected by the ambient flow with different sedimentation velocities for snow and rain.

We apply our refined measures in the same manner as for the seasons (Section 4) and quantify the impacts of the model updates on QPF performance. The four consecutive periods which are separated by these three

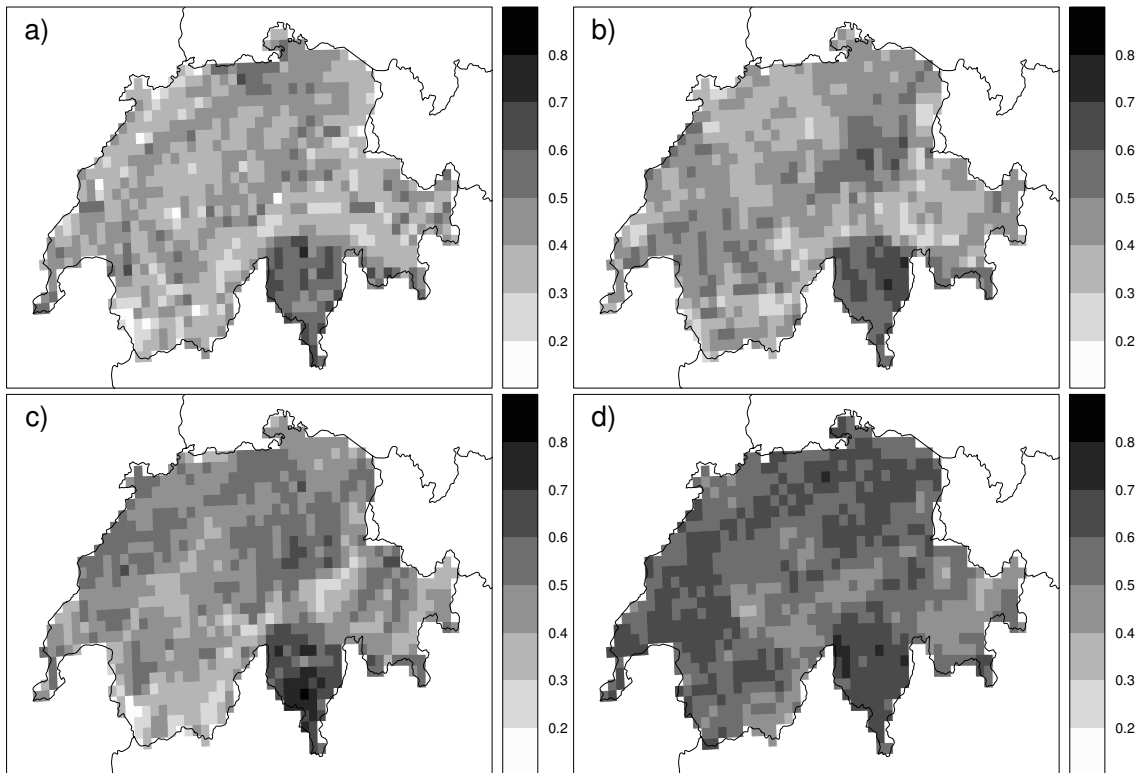


Figure 7: Peirce skill score for 90 % quantiles of daily precipitation sums [1: perfect, 0: random forecast]: a) IGD, b) NGD, c) NID and d) NIP.

updates are evaluated in the following. The first phase (named IGD) comprises all days between 01/07/2000 and 30/10/2001. It features the interpolation for the initialization, the GME at the boundaries and the diagnostic precipitation. The second phase (NGD) ranges from 31/10/2001 until 15/09/2003 and already features the nudging assimilation. The third phase (NID) uses the IFS instead of the GME boundary fields and covers all days between 16/09/2003 and 15/11/2004. The latest model setup (NIP) runs additionally with the prognostic precipitation and covers the period between 16/11/2004 and 31/12/2006. Note that the seasonal composition of the four slices is very similar. Despite the interannual variability of the errors, we consciously consider all seasons together to retain reasonable sample sizes. Strongest climatological anomalies are found during IGD with much higher rainfall amounts than afterwards. In fact, there were some persistent rainfall episodes from October until November 2000 and from March until April 2001. The former were active on the Alpine south side whereas the latter mostly affected the Alpine north side.

5.1 Grid-point based verification of 90 % quantiles

During IGD, error patterns are significantly noisier than afterwards (Figs. 6a and 7a). Reasons are found during the first four months of the preoperational phase,

when the model ran without a filtered orography (also pointed out in the outlook of KAUFMANN et al., 2003). Daily precipitation fields are very spotty at that time and do not succeed in capturing area-wide rainfall. However, principal strengths and weaknesses of the COSMO model are already evident in the primary model setup. The overestimation over the northern Alpine foothills, the underestimation in the interior of the Alps, the good matching in the Ticino as well as the poor matching in deep Alpine valleys are obvious during IGD. The dry bias along the low elevations of the Middleland and especially inside the Alps affects a much larger area than afterwards. More than 40 % of all Alpine grid points over 1500 m MSL are underforecast with an offset worse than -5 mm/day. However, observed 90 % quantiles are roughly 5 mm/day higher than later on, meaning that the relative offset is comparatively low. The pixel overlap is most deficient at the western edge of the Valais and over the Jura mountains. At some grid points, the debiased PSS only constitutes about 0.15 which is equivalent to a debiased POD around 0.24 (Appendix C, Eq. 3.4).

The next two phases NGD and NID are unlike IGD but similar to each other. The percentage of severely underforecast grid points over 1500 m MSL only drops down from 20 % during NGD to 14 % during NID (Figs. 6b, c). The only outstanding difference is given by the diverging matching evolution in the south (Figs. 7b,c). During NGD, the pixel overlap measured by the debiased PSS is average both in the Valais and the Ticino in

comparison to other phases. During NID, the pixel overlap is inferior in the Valais and superior in the Ticino. Reasons for this disparate behavior are not obvious and need to be investigated in further detail. The matching with the IFS at the boundaries leads to a considerably better performance in the Jura and the Middleland than before. In contrast to the Valais and the Grisons, the debiased PSS exclusively remains over 0.3 here.

The change of the precipitation scheme implicates the most conspicuous forecast improvements. Very clearly, the matching improved all over the country (Fig. 7d). During NIP, the debiased PSS varies between 0.3 and 0.8 without exception. Above all, the pixel overlap in the main Alpine valleys improves significantly. Former PSS values around 0.3 rise to values around 0.5 which implicates that the hits H are 1.5 times more frequent than before. The poorest matching still is found throughout the Grisons and over some parts of the northern Alpine ridges. Aggregated over the Valais, the PSS for the 90 % quantiles rises from 0.33 during NID to 0.54 during NIP (Figs. 7c, d; compare with the 90 % quantile location in Figs. 9c, d). Aggregated over Switzerland, it rises from 0.47 to 0.56. In contrast to the pixel overlap, the bias over the Alps displays an ambivalent trend during NIP (Fig. 7d). Precipitation now is advected to areas which have been underforecast before and maxima of the underestimation are less severe. The dry bias rarely falls below -5 mm/day and no longer falls below -10 mm/day. However, maxima of the overestimation simultaneously are broadened and partly extend to the lee inside the Alps. Peak values of the bias only drop down from $+30$ mm/day to $+20$ mm/day. Therefore the areal overestimation strongly worsens over some Alpine slopes and in particular in the center of the Grisons. Aggregated over the Grisons, the overcharge of the 90 % quantiles rises from ~ 0 mm/day during NID to $+2.9$ mm/day during NIP (Figs. 6c, d; compare with the 90 % quantile location in Figs. 9c, d). Aggregated over Switzerland, it changes from -0.1 mm/day to $+1.6$ mm/day.

5.2 Regional verification of quantile courses

Developments of regional error characteristics for a range of rainfall intensities can be gleaned from Fig. 9. The respective quantile values of the observations are given in Fig. 8 as reference. It is obvious that observed rainfall amounts are much higher during IGD than afterwards. Most notably, the quantile values are about 50 % higher in the Ticino and the Grisons than later on. In contrast, least rainfall is observed in the latest phase NIP.

Astonishingly, the regional bias in the first model setup IGD is remarkably small, especially for low intensities (Fig. 9a), even though related quantiles stand for higher rainfall amounts than afterwards (Fig. 8a). However, it has to be kept in mind that the aggregated quantile difference does not display local rainfall shifts

within one region. Thus, a simultaneous overestimation and underestimation within one domain cancel each other out. Given the noisiness of the IGD field, this might explain the small QD values. During NGD, regional amplitude errors increase considerably, above all for lower intensities. Simultaneously, the debiased PSS rises in all regions considerably. The improvements are largest in the northern regions. PSS values in the Jura, the Middleland and the Northern Alps are enhanced by almost 0.1 for small quantiles (~ 70 %). The alteration stands out least in the Ticino where the performance already has been very good. Improvements between NGD and NID are not self-evident. In central Alpine regions, the matching remains insufficient for most quantiles. Above all, the debiased PSS is lowered in the Valais for quantile probabilities between 75 % and 95 %. In contrast, the outstanding pixel overlap in the Ticino even changes for the better for almost all quantiles. In regions to the north, the matching is slightly worse during NID than before for quantile probabilities below 75 %, but it is slightly improved during NID for quantiles above. Considering NIP, improvements for all quantiles are remarkable. The values of the PSS only slightly drop off towards high quantiles. The performance for 60 %, 70 % and 80 % quantiles is very similar now. However, an increase of the overestimation is also obvious during NIP for all quantiles. Only the comparatively flat Middleland remains almost bias-free. The amplitude error in the Ticino remains constant, but the bias is drastically increased to the north. In particular, low and medium intensities are now significantly overforecast in all Alpine areas. For example, the 70 % quantile offset increases in the Northern Alps from $+0.6$ mm/day to $+1.8$ mm/day (compare Figs. 9c, d). Note that the related 70 % quantiles values roughly constitute 3 mm/day in the observations (Fig. 8c, d).

5.3 Model description

5.4 Rating of integral values

The integrated scores recapitulate the QPF behavior explained above. The integral performance strongly depends on regional characteristics and QPF improvements deeply vary among the domains. The average quantile difference (Tab. 4) clearly increases over the Jura between IGD and NGD. \overline{QD} rises from 0.07 to 0.29 implying that relative deviations increase from 7 % to almost 35 % along the distributions under comparison. At the same time, amplitude errors level off over the Middleland and diversely change over the Alps. However, the most conspicuous worsening of the bias is obvious during NIP. All over Switzerland \overline{QD} rises from 0.09 to 0.22 implying that deviations grow from 9 % to 25 %. None of the Alpine regions has its lowest \overline{QD} value in the last period with the most sophisticated model version. The change is most critical over

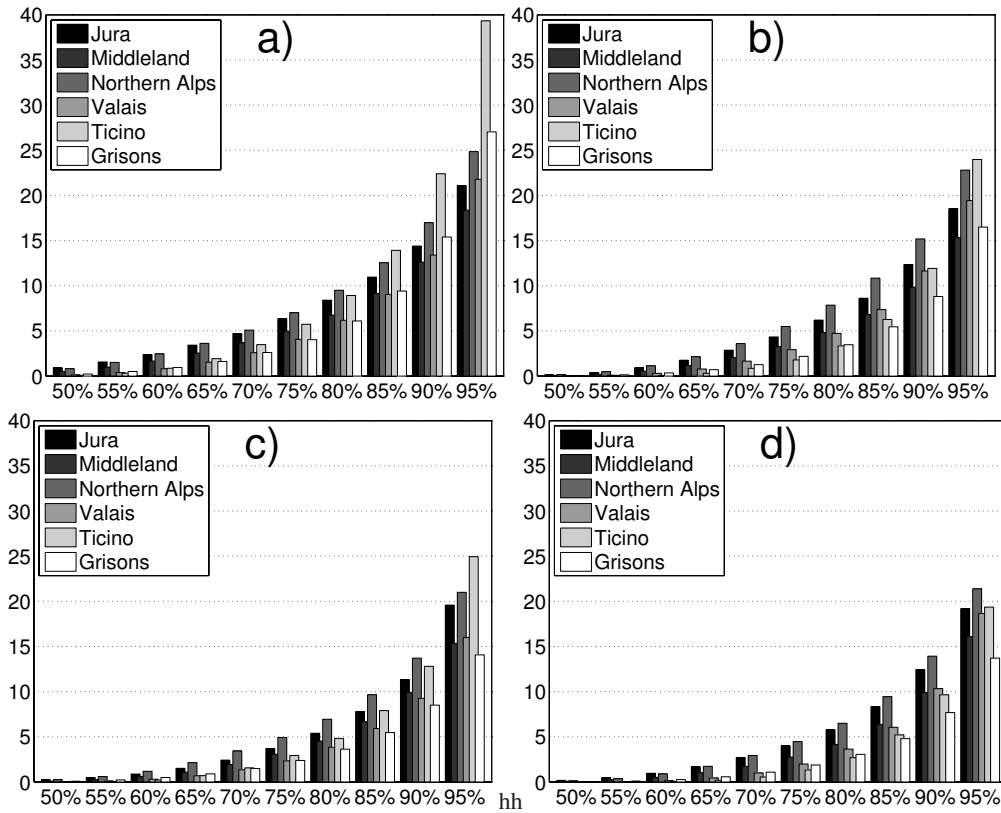


Figure 8: Observed quantile values [mm/day] for discrete quantile probabilities: a) IGD, b) NGD, c) NID and d) NIP.

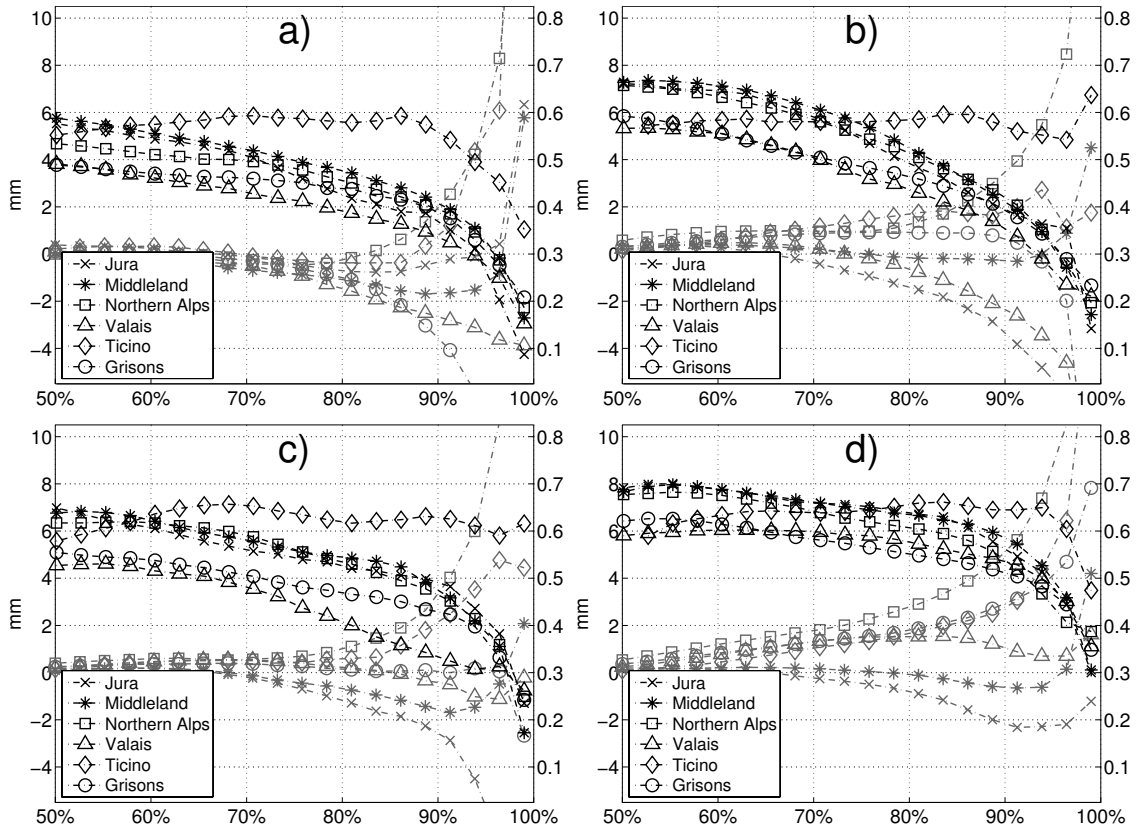


Figure 9: Area-average of the quantile difference (gray, left scale) [mm/day] and the Peirce skill score (black, right scale) for quantiles between 50 % and 99 %: a) IGD, b) NGD, c) NID and d) NIP.

Table 4: Weighted integral $\overline{QD'}$ values for different domains and periods. The scoring shows the relative amplitude deviation and is explained in Appendix B [$0.1 \hat{=} 10.5\%$ deviation, $0.2 \hat{=} 22.2\%$ deviation, $0.3 \hat{=} 35.3\%$ deviation]. The darker the shading the higher are overall observed rainfall amounts. CH stands for the whole of Switzerland.

	Jura	Middlel.	N. Alps	Valais	Ticino	Grisons	CH
IGD	0.07	0.13	0.16	0.18	0.14	0.22	0.07
NGD	0.29	0.07	0.26	0.23	0.15	0.19	0.09
NID	0.28	0.12	0.28	0.08	0.14	0.07	0.09
NIP	0.13	0.07	0.38	0.17	0.30	0.37	0.22

Table 5: Weighted integral \overline{PSS} values for different domains and periods. The higher the value the better is the overall matching performance. The darker the shading the higher are overall observed rainfall amounts. CH stands for the whole of Switzerland.

	Jura	Middlel.	N. Alps	Valais	Ticino	Grisons	CH
IGD	0.34	0.38	0.37	0.30	0.49	0.34	0.39
NGD	0.40	0.38	0.39	0.32	0.58	0.35	0.40
NID	0.44	0.42	0.44	0.34	0.61	0.39	0.44
NIP	0.52	0.52	0.51	0.49	0.59	0.48	0.51

the Ticino and the Grisons where $\overline{QD'}$ rises from 0.14 to 0.30 and from 0.07 to 0.37, respectively. The bias only improves over the Jura and the Middleland between NID and NIP. Interestingly, $\overline{QD'}$ continuously increases within the four phases for the Northern Alps. The value rises from 0.16 during IGD to 0.38 during NIP. Concerning \overline{PSS} (Tab. 5), QPF performance definitely is best for the latest model version. Except for the Ticino, where the pixel overlap already has been superior before, all regions display their highest \overline{PSS} during NIP. Thus, the great improvement of advecting precipitation by the ambient flow is brought out clearly by the summary measure. Considering all phases, the matching is upgraded most in the Jura and the Valais where \overline{PSS} rises by almost 0.2 between IGD and NIP.

6 Interpretation of verification results

Some aspects are worthwhile noting in the interpretation of our results. Firstly, the observational precipitation analysis is not perfect. In particular, it exhibits a negative bias (too low values) due to systematic measurement errors. In the Alpine region, this bias is less than 12 % from spring to fall, but can reach several tens of percent in winter for exposed stations above 1500 m MSL (SEVRUK, 1985; RICHTER, 1995). Although the rain-gauge under-catch may possibly explain some part of the apparent quantile overestimation in winter, the overall characteristics of the model errors remain valid. The magnitude of the model bias is substantially larger than the expected measurement bias. Note that only about 10 % of the rain gauge stations are at elevations above 1500 m MSL even in inner Alpine regions (see e.g. Fig. 6 in FREI and SCHÄR, 1998) and hence the bias in the observational analysis is much smaller than that at single exposed stations. Moreover, our primary focus is on

high quantiles (i.e. intense rainfall or heavy snowfall) for which the measurement bias is considerably smaller than for light events.

Generally, our verification results agree with those found by ELEMENTI et al. (2005) in case studies. First and foremost, the severe overestimation at the Alpine fringe is confirmed. The COSMO model versions under consideration seem to have problems to represent the correct flow and moisture field around orography. The windward side receives too much and the lee side too little precipitation. This behavior is most pronounced with the diagnostic precipitation scheme, but still holds for the prognostic scheme. Further investigations (not shown) have proven that the overestimation mostly stems from the resolved part of the total rainfall (see KAUFMANN et al., 2003, for comparison) which holds especially for very strong events. Recent tests with the COSMO model revealed that a three-step Runge-Kutta time integration scheme (WICKER and SKAMAROCK, 2002) partly rectifies the overestimation on the windward side of a mountain (D. LEUENBERGER (MeteoSwiss), pers. comm.). However, it remains ambiguous, how inconsistencies of the used leapfrog scheme affect the formation of precipitation and cause the offsets on the windward side.

A closer look reveals that the overestimation is most pronounced over higher elevations during winter. Even though the results are slightly exaggerated by wintertime measurement errors, snowfall seems to be more overstated than rainfall. This can be attributed to problems with the cloud ice scheme for the winters 2003/2004 and 2004/2005 (F. SCHUBIGER (MeteoSwiss), pers. comm.).

Interestingly, the common pattern with a wet bias at the foot of mountains does not always hold for the Alpine south side. In opposition to the findings discussed above, there is a slight dry bias over the Ticino

during winter and fall. Obviously, a drying of the water cycle is superimposed to the usual overestimation of precipitation at the Alpine fringe. Comparisons of COSMO forecasts and COSMO analyses revealed that forecasts over the Po valley tend to be significantly dryer than in the corresponding assimilation cycle. Apparently, the water balance in the COSMO model is predisposed to leak over the northwest of Italy. However, the cause of this is not clear and requires further investigation.

A positive correlation between the quantile difference and the debiased PSS can be confirmed statistically for the diagnostic precipitation scheme. In other words, the debiased PSS usually is higher for positive quantile differences than for negative ones. The reason is the absent horizontal transport of hydrometeors which produces both an underestimation and a poor matching on the lee side of mountains. By introducing the prognostic precipitation scheme, this correlation vanishes. The bias no longer interferes significantly with the matching and a main error source which affected both verification components seems to be removed. Note that this behavior only is proven so clearly, because all intensities are related to quantiles/frequencies and not to amplitudes.

At the same time, a significant worsening of the wet bias is evident during the latest model phase. In particular over the Ticino and the Grisons, the relative offset of the distributions under comparison increased drastically. Most notably, southerly flow now entails a much higher wet bias than before. Possible explanations relate to unmentioned model updates like a change of the cloud ice scheme, the introduction of prognostic turbulent kinetic energy or the switch of the IFS boundary fields to higher resolutions. Overall, it is most suspect that too much moisture is advected towards the Alps by the driving model. Again, implications of verification results do not transfer directly to model diagnosis and the original error source remains to be investigated.

One of the most outstanding findings is the high quality of the pixel overlap on the Alpine south side. As shown in Section 4, the pixel overlap in the Ticino is much better than elsewhere and is present throughout the year and throughout all model versions. It is most manifest for strong intensities and most pronounced during winter and fall (Figs. 5a,d). Explanations are found in the special geographical position of southern Switzerland in connection to the prevailing synoptic flow patterns. Meridionally aligned stratospheric intrusions determine the large-scale predictability of heavy precipitation (e.g. FEHLMANN and QUADRI, 2000; MARTIUS et al., 2006) and primarily support the enhanced pixel overlap during fall. The relief of the Ticino is uniformly aligned to the south and lee effects play a minor role for a southerly flow. The high predictability also is supported by idealized model simulations of GHEUSI and DAVIES (2004). They found that orographically induced precipitation enhancements in southern Ticino are comparatively insensitive to the changes in direction and

speed of the incident flow from south to southwest or from 10 to 30 m/s. In addition to its favorable exposure to the south, the Ticino is shielded by the Alpine crest for northerly flow directions. Embedded shallow fronts provide a potential for misforecasts on the windward side, but are obstructed by the Alpine crest and do not affect the Ticino (E. ZALA (MeteoSwiss), pers. comm.). Their frequency maximum on the windward side during the cold season (JENKNER et al., in press) supports this explanation for an enhanced wintertime pixel overlap in the Ticino.

Concerning the seasonal cycle, it is obvious that the pixel overlap is highest during winter (in some regions also during fall) and lowest during summer. In this regard, model performance strongly anticorrelates with the amount of convective precipitation or likewise with the boundary layer height. The dependency emerges clearly in our results and is worthwhile mentioning, even though convection schemes are already well-known to present difficulties to QPF (e.g. ELEMENTI et al., 2005). The interaction of the boundary layer and the free atmosphere is well-developed during the convective season and challenges the interplay of parameterized and resolved processes in current NWP models. The final outcome is a degradation of the local skill in convective situations.

7 Summary

In the present study a novel treatment of the traditional categorical verification has been presented. By using frequency thresholds instead of amplitude thresholds, deterministic verification applications clearly benefit from the quantile-based formulation of the verification problem. We propose to use the absolute or relative quantile difference to describe the bias and the debiased/calibrated Peirce skill score to describe the potential pixel overlap. In this way, the total error is split up into an amplitude part and a matching part. If multiple quantiles are taken into account, spectral performance can be assessed. This setup allows for a meaningful juxtaposition of different value ranges and renders possible a meaningful evaluation of individual intensities. In this context, our distribution-oriented approach makes a decisive step towards equitable scores, as defined by GANDIN and MURPHY (1992). In our opinion, it is more meaningful for model developers to relate verification results to characteristic numbers of the model output (e.g. quantiles) than to a priori fixed limits such as amplitudes. The main advantages and challenges of the refined approach are itemized here for recapitulation:

- The degrees of freedom within the contingency table reduce to one. The information content of a calibrated forecast can be displayed by a single score.

- Owing to the use of quantiles, the conducted debiasing has a physical validity and only relates to the characteristics of the distributions under comparison. Such a straightforward calibration cannot be achieved so easily by other approaches.
- The whole range of precipitation intensities can be assessed by looping over quantiles. Different value ranges are related to each other in a consistent and comprehensive way.
- The Peirce skill score strongly simplifies with the use of quantiles. The debiased PSS quantifies the potential skill in a calibrated forecast and is no longer susceptible to hedging. In the presented formulation, the PSS offers the possibility to compare results for different base rates without being affected by the marginal distributions.
- Synoptic biases or gridding errors do not influence the matching error score, if they are applied to homogenous subsets (in our case orographically distinct areas) with a consistent bias behavior.
- Scores can be integrated over all quantiles while weighting them accordingly. In such a way, verification results do not only refer to a single threshold, but take into account spectral performance with a tunable resolution. Thus, the overall performance is condensed meaningfully and allows for a quick QPF overlook, for example for administrative purposes.
- As a slight drawback, quantiles universally are less intuitive than fixed thresholds. A verification which is issued to the public requires values to be linked to proper rainfall amounts due to convenience. It can be speculated, whether model developers, forecasters and end-users become more familiar with the reference to quantiles in the future.
- The interpretation of the debiased PSS might remain slightly ambivalent for rare events. It is not clear, if random chance always provides a meaningful reference for the definition of forecast skill. Especially in the tail of rainfall distributions, special reference distributions might be beneficial (STEPHENSON et al., 2008).
- Error sources in NWP models can be tackled in a straightforward way. If a model always exhibits a distinct bias, the forecast can never be perfect. Thus, it is meaningful to rate the bias beforehand and evaluate the residual error afterwards.
- Regional QPF performance is strongly determined by local orography in connection with the impinging flow direction. Distinct areas with a good and poor pixel overlap can be identified. The fine-scale error structure emerges most clearly during winter.
- The COSMO model exhibits a persistent overestimation over the Alpine foreland and a transient underestimation over interior valleys (primarily during winter and summer and with the diagnostic precipitation scheme). The Alpine south side partly is an exception with comparatively low seasonal amplitude errors during winter and fall.
- Overall matching characteristics are worst during summer. All over Switzerland, the pixel overlap is roughly 6–7 % worse during summer than during the rest of the year. Seasonal disparities are even higher over the Jura and the Ticino. Altogether, the matching is worst over the Jura and the Valais during summer.
- The potential skill of the calibrated forecast is much higher on the Alpine south side than on the north side. On average, the pixel overlap is roughly 15 % better over the Ticino than elsewhere.
- The matching continuously improves within the validated period and can be clearly attributed to updates of the COSMO model. All over Switzerland, the pixel overlap is 12 % better in the latest considered phase than in the first one.
- The overall bias is worst by far in the latest considered phase. The percentaged amplitude error constitutes about 9 % between 2000 and 2004, but it rises to 25 % between 2005 and 2006. Note that the better matching characteristics with the prognostic rainfall scheme in the latest phase could not be visible so clearly with the conventional PSS, because the better rainfall overlap is masked by the large bias in an uncalibrated forecast.

The steep orography in Switzerland imposes a severe constraint on QPF performance. The discussed QPF errors above exhibit a strong dependency on weather patterns with specific flow features in conjunction with the complex terrain. Thus, it is instructive to study model performance with respect to different synoptic situations. Such a study will be conducted with a consistent version of the COSMO model in a follow-up paper.

Acknowledgement

The authors cordially thank Daniel LEUENBERGER for providing excellent technical support and together with Emanuele ZALA and Francis SCHUBIGER for giving helpful comments concerning the interpretation of results. Additional comments provided by Pertti NURMI

The methodology has been applied to 6.5 years of pre-operational and operational forecasts from the Swiss implementation of the COSMO model. Seasons and model versions have been evaluated separately. 90 % quantiles have been investigated in detail and quantile courses have been discussed for six predefined regions. The most important results are recapitulated in the following:

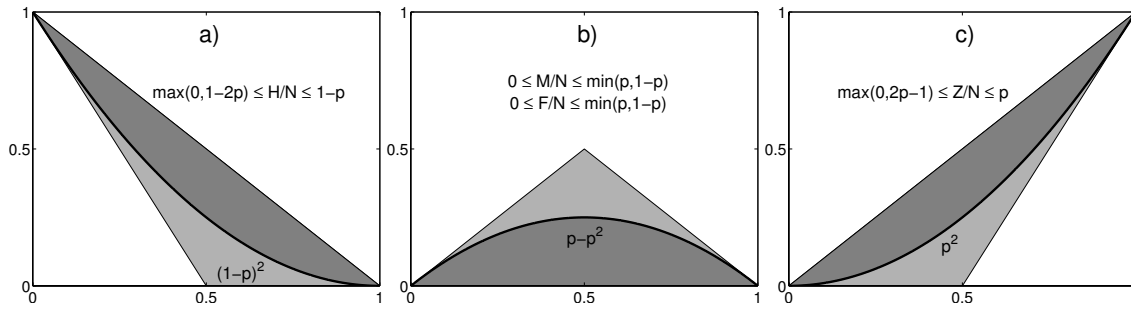


Figure 10: Possible entries (gray shaded) and random expectation values (thick black lines) of the four entries in the quantile-based contingency table: a) hits $H(p)$, b) misses $M(p)$ or false alarms $F(p)$, c) correct negatives $Z(p)$. Light gray indicates the area without skill whereas dark gray indicates the area with skill. The thick dividing line depicts the fractions for a random forecast. The x-axis displays the range of quantile probabilities p whereas the y-axis shows the proportion of the sample size H/N , M/N or F/N and Z/N .

are acknowledged as well. This study was funded by the Swiss National Science Foundation (SNF) under grant number 200021-105591.

A Quantile-oriented binary contingency tables

The entries of the standard 2x2 contingency table (Tab. 1) are defined by means of sample quantiles instead of a preassigned threshold value. Thus, the marginal distributions ($H+F$, $Z+M$, etc.) are fixed automatically. The forecast data are calibrated in such a way that the decision criterion is changed while the underlying distributional overlap is maintained (compare with signal detection theory, SWETS, 1988). If the same quantile probability is chosen for the observations and the forecast ($p \equiv p_{obs} = p_{mod}$), the event frequency is the same in both datasets. Two additional interrelations are then incorporated into the conceptual formulation:

$$M = F \quad M + Z = pN \quad (1.1)$$

Firstly, the misses M equal the false alarms F and the bias automatically is removed from subsequent considerations. Secondly, the quantile probability p by definition sets the base rate to $1 - p$ (note that the definition here is not p as in most of the literature) and divides the sample into $(1 - p)N$ events and pN non-events. Thus, all four entries H , M , F and Z additionally are linked to p itself. Since the sample size $N=H+M+F+Z$ is given, the four numbers are connected by a total of three constraints. Only one degree of freedom is left and uniquely describes the joint distribution.

The determination of p imposes stringent restrictions on the valid range of the four entries in the contingency table. Depending on p , the four counts only vary within a limited span (Fig. 10). If the quantile probability is below 0.5, there are always some hits H by definition, as exceeding events cover more than half of the sample. If the quantile probability is above 0.5, there are always some correct negatives Z by definition, as exceeding events cover less than half of the sample. The misses

M and false alarms F consistently are limited at the top. The maximum numbers M_{max} and F_{max} are equal and restricted either by the number of non-events ($p \leq 0.5$) or by the number of events ($p \geq 0.5$). Just in case of the median ($p = 0.5$), all four counts hold the same range of values. Then the setting is balanced and the number of hits H by definition equals the number of correct negatives Z .

Due to rules of combinatorics, the probability for the worst possible forecast with the maximum number of non-matching pixels is a function of both p and N and can be written as:

$$P\langle M = M_{max} \rangle = \begin{cases} \frac{\binom{(1-p)N}{pN}}{\binom{N}{pN}} & \text{for } p \leq 0.5 \\ \frac{\binom{pN}{(1-p)N}}{\binom{N}{(1-p)N}} & \text{for } p \geq 0.5 \end{cases} \quad (1.2)$$

Note that the binomial coefficients only can be computed, if pN and $(1 - p)N$ denote integers. Owing to Eq. 1.2, it is hardest to miss all events or all non-events in case of the median ($p = 0.5$) and easiest in case of extreme quantile probabilities.

A random forecast divides the joint distribution into a region with skill and one without skill. The borderline is $(1 - p)^2N$ for the hits, $(p - p^2)N$ for the misses or false alarms and p^2N for the correct negatives (Fig. 10). Thus, a valuable forecast always exhibits less than $(p - p^2)N$ misses/false alarms. In case of the median ($p = 0.5$), the border consistently resides at $N/4$ which is in the center of valid ranges. It gradually approaches a margin for p converging towards 0 or 1.

B Quantile differences

Certain sample quantiles are useful in an exploratory rainfall verification. The quantile-quantile plot is useful in small datasets, but it is beneficial to map quantile differences in gridded samples (FERRO et al., 2005). In our context, the quantile difference directly exhibits a distinctive amplitude error, i.e. bias between the datasets:

$$QD(p) = q_{mod}(p) - q_{obs}(p) \quad (2.1)$$

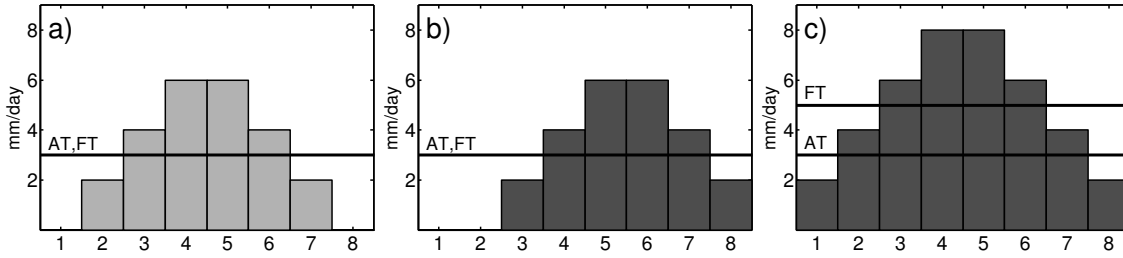


Figure 11: Forecast example of predicting daily rainfall amounts for 8 days: a) observations, b) first forecast and c) second forecast. The selected amplitude threshold (AT) constitutes 3 mm/day. The selected frequency threshold (FT) corresponds to the 50 % quantile. The conventional PSS based on the AT accounts for 0.5 in both forecasts. The debiased PSS based on the FT accounts for 0.5 and 1 for the first and second forecast, respectively.

As an option, QD can also be standardized by individual quantiles. Similar to the definition of the amplitude component of the SAL measure (WERNLI et al., 2008), the scaling of the error preferably is done by the arithmetic mean of the rainfall amounts under comparison:

$$\text{QD}'(p) = \frac{2 \text{QD}(p)}{q_{\text{obs}}(p) + q_{\text{mod}}(p)} \quad (2.2)$$

As a matter of fact, QD describes the absolute quantile difference whereas QD' describes the relative quantile deviation. QD is measured in physical units (e.g. mm). QD' is dimensionless and varies between -2 and $+2$. A positive QD or QD' corresponds to a wet bias, i.e. an overestimation of precipitation. A negative QD or QD' stands for a dry bias, i.e. an underestimation of precipitation. Note that the percentaged bias constitutes $\frac{\text{QD}(p)}{q_{\text{obs}}} = \frac{2 \text{QD}'(p)}{2 - \text{QD}'(p)}$, but remains undefined in case of dry observations ($q_{\text{obs}} = 0$ mm). As the case may be, it is more convenient to display QD or QD'.

It is advisable to weight the relative quantile deviation while expressing the overall performance over precipitation magnitudes. The weighted integral over quantiles preferably is computed with the absolute value, because deviations with a varying sign cancel each other out in the average. Since the quantile difference originally is an additive quantity, we use the arithmetic mean of observed and modeled quantiles as weighting function:

$$\overline{\text{QD}'} = \frac{1}{\int w(p) dp} \int_0^1 \underbrace{w(p) |\text{QD}'(p)|}_{|\text{QD}(p)|} dp \quad (2.3)$$

$$w(p) = \frac{q_{\text{obs}}(p) + q_{\text{mod}}(p)}{2}$$

Note that $\overline{\text{QD}'}$ only varies between 0 and $+2$ and does not scale linearly with the bias. Amplitude deviations of 5 %, 10 %, 20 %, 40 % and 80 % lead to $\overline{\text{QD}'}$ values of 0.049, 0.095, 0.182, 0.333 and 0.571, respectively. The lower the value of $\overline{\text{QD}'}$ the more alike are the compared distributions in terms of intensities and the smaller is the overall amplitude error.

C The Peirce skill score with quantiles

The conventional Peirce skill score (PEIRCE, 1884) can be reformulated with an offset for selected quantile probabilities $dp = p_{\text{mod}} - p_{\text{obs}}$ with p_{mod} and p_{obs} denoting the quantile probabilities of the considered amplitude threshold:

$$\text{PSS}(p) = \frac{H}{H+M} - \frac{F}{F+Z} = 1 - \frac{M}{\underbrace{(p_{\text{obs}} - p_{\text{obs}}^2)N}_{c_{\text{match}}}} - \frac{dp}{\underbrace{p_{\text{obs}}}_{c_{\text{bias}}}} \quad (3.1)$$

Hence, forecast deficiencies which affect the PSS are clearly separated into a contribution resulting from an insufficient matching (c_{match}) and a contribution resulting from a bias (c_{bias}). Note that the difference of forecast frequencies (dp) changes the PSS more intensely for large base rates ($1 - p_{\text{obs}}$) than for small ones. The problem of hedging results from the fact that people may change p_{mod} while p_{obs} remains fixed. However, it is eluded, if quantile probabilities coincide to each other ($p \equiv p_{\text{obs}} = p_{\text{mod}}$, see Appendix A for details) and a change of p_{mod} always implicates a change of p_{obs} . Several skill scores merge into each other in that situation. In particular, the Peirce skill score (PEIRCE, 1884) gets equal to the Heidke skill score (HEIDKE, 1926) and to the Clayton's skill score (CLAYTON, 1934). In the debiased formulation, the formula of the PSS further simplifies to:

$$\text{PSS}(p) = 1 - \frac{M}{M_{\text{rand}}} \quad M_{\text{rand}} = (p - p^2)N \quad (3.2)$$

Now, the PSS exclusively focuses on the proportion of mismatching pixels, i.e. the magnitude of the overlap of rain pixels. The symmetry of the random misses M_{rand} (see Fig. 10b for visualization) guarantees that the debiased PSS is invariant with respect to taking the complement of the events (see STEPHENSON, 2000, for details). Therefore it is equivalent to define the quantiles going upward or downward along the distribution. In addition, a swapping of forecasts and observations is allowed in our setup and the debiased PSS fulfills the requirements

for transpose symmetry (see STEPHENSON, 2000, for definition).

Positive PSS values indicate skill compared to a shuffled sample, since a random forecast results in $PSS = 0$ at all times. Owing to the limited range of the four entries in the contingency table (see Appendix A for details), the debiased PSS only varies between $1 - 1/(1-p)$ and 1 for $p \leq 0.5$ or accordingly between $1 - 1/p$ and 1 for $p \geq 0.5$. Different probabilities for a complete miss of all events or all non-events (Appendix A, Eq. 1.2) are included in the concept and it is not possible to score much worse than random chance in the case of rare events or rare non-events.

Using the definitions introduced by GANDIN and MURPHY (1992) the PSS can be computed in different ways corresponding to the four entries of the matrix product between the symmetric scoring matrix \mathbf{S} and the symmetric performance matrix \mathbf{P} :

$$\begin{aligned} \mathbf{SP} &= \begin{pmatrix} 1/p & -1/p \\ -1/(1-p) & 1/(1-p) \end{pmatrix} PSS(p) \\ \mathbf{S} &= \begin{pmatrix} p/(1-p) & -1 \\ -1 & (1-p)/p \end{pmatrix} \quad (3.3) \\ \mathbf{P} &= \frac{1}{N} \begin{pmatrix} H & M \\ M & Z \end{pmatrix} \end{aligned}$$

Note that the scoring matrix \mathbf{S} determines the nature of the PSS, i.e. the way how the PSS is defined using the performance matrix \mathbf{P} .

To analyze the debiased PSS, it can be additionally transformed into the debiased POD. However, the convenient characteristics of the PSS are then lost:

$$POD(p) = 1 - p(1 - PSS(p)) \quad (3.4)$$

Individual PSS values can be integrated to express the overall matching performance. The weighted average over all quantiles accounts for the collective skill. Note that the PSS depends on the definition of the contingency table, if any one of the quantile values q_{obs} or q_{mod} vanishes, because rainfall distributions usually exhibit a large amount of ties at 0 mm. To avoid problems at 0 mm, the geometric mean of observed and modeled quantiles is used as weighting function:

$$\begin{aligned} \overline{PSS} &= \frac{1}{\int w(p) dp} \int_0^1 w(p) PSS(p) dp \quad (3.5) \\ w(p) &= \sqrt{q_{obs}(p)q_{mod}(p)} \end{aligned}$$

The higher the value of \overline{PSS} the larger is the pixel overlap and the smaller is the overall shift error.

D The interpretation of the Peirce skill score

To clarify the peculiarity of the PSS with quantiles, a simple forecast situation can be considered covering

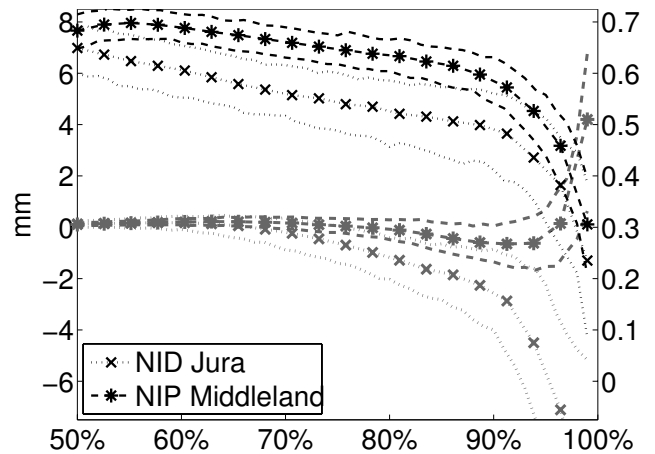


Figure 12: Examples for 95 % confidence intervals (lines without markers) obtained from a bootstrap sample with 500 repetitions for the QD (gray, left scale) [mm/day] and the debiased PSS (black, right scale): The Jura during NID (16/09/2003 until 15/11/2004, 427 days) and the Middleland during NIP (16/11/2004 until 31/12/2006, 776 days) are displayed.

daily rainfall amounts for 8 days. We assume that the observations peak on day four and five and have a distribution as in Fig. 11a. Then two constructed forecasts can be evaluated with distributions as in Figs. 11b, c. The first forecast exhibits a pure shift error of 1 day whereas the second forecast exhibits a pure bias with an overestimation of $QD = 2$ mm/day. Initially, we compute the contingency table based on an amplitude threshold with 3 mm/day and therefore verify the forecasts without a calibration. The two error types produce the same skill, as the PSS constitutes 0.5 for both forecasts. Thus, the origin of the error (a temporal shift or an amplitude overestimation) is not detected by the PSS in the conventional formulation. Then we compute the contingency table based on the median of daily rainfall amounts. Note that the effective threshold remains 3 mm/day in the observations and the first forecast, but changes to 5 mm/day in the second forecast. Therefore we now verify the forecasts after being calibrated. The PSS automatically is debiased and only detects the shift error. The resulting PSS values now are 0.5 for the first forecast and 1 for the second forecast. The displayed skill is no longer influenced by the bias and differs considerably in the two predictions. However, the skill might be regarded as potential, as the calibrated forecast is not issued, but only is used to reformulate the verification problem. This simple example illustrates that the two error types are no longer mixed in the PSS. If the bias is rated as well, it is straightforward to interpret the verification outcome in terms of shift and amplitude errors.

E The uncertainty of results

Essentially, the estimation of the sample quantiles is sensitive to the available dataset (e.g. CONOVER, 1999).

Since rainfall distributions usually are highly skewed, only few nonparametric methods exist to quantify the uncertainty of the computed quantiles. A convenient approach is described by CONOVER (1999). If a dataset is considered independent and identically distributed, the binomial distribution describes the probability that a single data point constitutes the targeted quantile. The confidence intervals are determined using the order statistics and the resulting $[r, s]$ interval is then converted to quantiles. The positions r and s in the order statistics can be obtained as:

$$\begin{aligned} r &= Np + y_{\alpha/2} \sqrt{Np(1-p)} \\ s &= Np + y_{1-\alpha/2} \sqrt{Np(1-p)} \end{aligned} \quad (5.1)$$

Intrinsically, $y_{\alpha/2}$ and $y_{1-\alpha/2}$ denote quantiles of the binomial distribution and $Np(1-p)$ represents its variance. If N is sufficiently large, the central limit theorem can be applied. Then $y_{\alpha/2}$ and $y_{1-\alpha/2}$ can be taken from the standard normal distribution.

Let us exemplify the proceeding of CONOVER (1999) by means of the highest evaluated quantile in the smallest and the largest regional sample. The Jura only encompasses $N = 23485$ data items during the period NID (427 days \times 55 grid points, see Section 5 for details). Observed and modeled 99 % quantiles correspond to 38.4 mm and 27.2 mm respectively. Following the explained method, the confidence interval (95 % significance level) ranges from $r = 23220$ to $s = 23280$ which implies [37.4 mm, 39.5 mm] and [25.5 mm, 28.7 mm]. In contrast, the Middleland encompasses $N = 210296$ data items during the period NIP (776 days \times 271 grid points, see Section 5 for details). Observed and modeled 99 % quantiles correspond to 29.4 mm and 33.6 mm respectively. The confidence interval now ranges from $r = 208104$ to $s = 208283$ which implies [29.1 mm, 29.8 mm] and [33.1 mm, 34.2 mm]. On account of the larger sample, the uncertainty is much smaller in the second case.

For the error itself, the overall sample uncertainty subdivides into two different parts. One portion relates to the bias and another one relates to the pixel overlap. On the one hand, the estimation of quantiles directly affects the uncertainty of the quantile difference. Provided that observed and forecast values are independent, variances add up to that of the difference. On the other hand, the determination of matching and non-matching pixels affects the uncertainty of the Peirce skill score. HANSEN and KUIPERS (1965) derived the variance of the PSS by means of parametric assumptions. In our notation, it can be written as:

$$\text{var}(\text{PSS}(p)) = \frac{1}{N} \left(\frac{1}{4p(1-p)} - \text{PSS}(p)^2 \right) \quad (5.2)$$

Following equation 5.2, the uncertainty of the PSS is highest for extreme quantiles and PSS values close to zero. In contrast, it is lowest for quantiles around the

median and PSS values close to 1.

Strictly speaking, the spatial correlation among individual grid points is crucial and must be maintained while estimating uncertainties of regional scores. In contrast, individual days are approximately independent from each other. Thus, meaningful confidence intervals are obtained by fixing the spatial configuration and varying the temporal composition. Proper resampling methods are explained by FERRO et al. (2005). We apply a bootstrap with 500 repetitions and computed quantile-based confidence intervals. In other words, we resample available days with replacement and use ordinary quantiles to determine confidence limits. Once again, we take the Jura during NID and the Middleland during NIP for illustration purposes (Fig. 12). The uncertainty generally is much smaller in the larger sample. Confidence intervals of the quantile difference tend to spread when moving to higher quantiles, because they are not independent from the amplitudes. The uncertainty of the debiased PSS only slightly increases for rare events which is in contrast to a conventional computation based on amplitude thresholds.

References

- ACCADIA, C., S. MARIANI, M. CASAIOLI, A. LAVAGNINI, A. SPERANZA, 2005: Verification of precipitation forecasts from two limited-area models over Italy and comparison with ECMWF forecasts using a resampling technique. – *Wea. Forecast.* **20**, 276–300.
- ASSELIN, R., 1972: Frequency filter for time integrations. – *Mon. Wea. Rev.* **100**, 487–490.
- BENOIT, R., C. SCHÄR, P. BINDER, S. CHAMBERLAND, H.C. DAVIES, M. DESGAGNE, C. GIRARD, C. KEIL, N. KOUWEN, D. LÜTHI, D. MARIC, E. MÜLLER, P. PELLERIN, J. SCHMIDLI, F. SCHUBIGER, C. SCHWIERZ, M. SPRENGER, A. WALSER, S. WILLEMSE, W. YU, E. ZALA, 2002: The real-time ultrafinescale forecast support during the special observing period of the MAP. – *Bull. Amer. Meteor. Soc.* **83**, 85–109.
- BUZZI, A., S. DAVOLIO, M. D'ISIDORO, P. MALGUZZI, 2004: The impact of resolution and of MAP reanalysis on the simulations of heavy precipitation during MAP cases. – *Meteorol. Z.* **13**, 91–97.
- CASATI, B., G. ROSS, D. STEPHENSON, 2004: A new intensity-scale approach for the verification of spatial precipitation forecasts. – *Meteor. Apps.* **11**, 141–154.
- CLAYTON, H., 1934: Rating weather forecasts. – *Bull. Amer. Meteor. Soc.* **15**, 279–283.
- CONOVER, W., 1999: Practical Nonparametric Statistics, chapter 3, pages 123–178 Wiley and Sons.
- DAVIES, H. C., 1976: Lateral boundary formulation for multilevel prediction models. – *Quart. J. Roy. Meteor. Soc.* **102**, 405–418.
- ELEMENTI, M., C. MARSIGLI, T. PACCAGNELLA, 2005: High resolution forecast of heavy precipitation with Lokal Modell: Analysis of two case studies in the Alpine area. – *Nat. Haz. Earth Sys. Sci.* **5**, 593–602.
- FEHLMANN, R., C. QUADRI, 2000: Predictability issues of heavy alpine south-side precipitation. – *Meteor. Atmos. Phys.* **72**, 223–231.

- FERRO, C., A. HANNACHI, D. STEPHENSON, 2005: Simple nonparametric techniques for exploring changing probability distributions of weather. – *J. Climate* **18**, 4344–4354.
- FREI, C., C. SCHÄR, 1998: A precipitation climatology of the Alps from high-resolution rain-gauge observations. – *Int. J. Climatol.* **18**(8), 873–900.
- FREI, C., R. SCHOLL, S. FUKUTOME, R. SCHMIDLI, P. VIDALED, 2006: Future change of precipitation extremes in Europe: Intercomparison of scenarios from regional climate models. – *J. Geophys. Res.* **111**, D06105.
- GANDIN, L., A. MURPHY, 1992: Equitable skill scores for categorical forecasts. – *Mon. Wea. Rev.* **120**, 361–370.
- GHEUSI, F., H. C. DAVIES, 2004: Autumnal precipitation distribution on the southern flank of the Alps: A numerical-model study of the mechanisms. – *Quart. J. Roy. Meteor. Soc.* **130**, 2125–2152.
- HAMILL, T., J. JURAS, 2006: Measuring forecast skill: Is it real skill or is it the varying climatology?. – *Quart. J. Roy. Meteor. Soc.* **132**, 2905–2923.
- HANSEN, A., W. KUIPERS, 1965: On the relationship between the frequency of rain and various meteorological parameters. – *Meded. Verh.* **81**, 2–15.
- HEIDKE, P., 1926: Berechnung des Erfolges und der Güte der Windstärkevorhersagen im Sturmwarnungsdienst. – *Geogr. Ann.* **8**, 301–349.
- HILLIKER, J., 2004: The sensitivity of the number of correctly forecasted events to the threat score: A practical application. – *Wea. Forecast.* **19**, 646–650.
- HOHENEGGER, C., C. SCHÄR, 2007: Predictability and error growth dynamics in cloud-resolving models. – *J. Atmos. Sci.* **64**, 4467–4478.
- HOHENEGGER, C., D. LÜTHI, C. SCHÄR, 2006: Predictability mysteries in cloud-resolving models. – *Mon. Wea. Rev.* **134**, 2095–2107.
- HOUZE, R., S. MEDINA, 2005: Turbulence as a mechanism for orographic precipitation enhancement. – *Quart. J. Roy. Meteor. Soc.* **62**, 3599–3623.
- JENKNER, J., M. SPRENGER, I. SCHWENK, C. SCHWIERZ, S. DIERER, D. LEUENBERGER, in press: Detection and climatology of fronts in a high-resolution model reanalysis over the Alps. – *Meteor. Apps.*
- KÄLLBERG, P., A. MONTANI, 2006: A case study carried out with two different NWP systems. – *Nat. Haz. Earth Sys. Sci.* **6**, 755–760.
- KAUFMANN, P., F. SCHUBIGER, P. BINDER, 2003: Precipitation forecasting by a mesoscale numerical weather prediction (NWP) model: Eight years of experience. – *Hydrol. Earth Syst. Sci.* **7**, 812–832.
- KESSLER, E., 1969: On the distribution and continuity of water substance in atmospheric circulations. – *Meteor. Monogr.* **10**, 84.
- KONZELMANN, T., R. WEINGARTNER, 2007: Niederschlagsmessnetze. – In: *Hydrological Atlas of Switzerland, Landeshydrol. und Geol., Bern*, plate 2.1.
- LORD, S., H. WILLOUGHBY, J. PIOTROWICZ, 1984: Role of a parameterized ice-phase microphysics in an axisymmetric, nonhydrostatic tropical cyclone model. – *J. Atmos. Sci.* **41**, 2836–2848.
- MAHONEY, K., G. LACKMANN, 2007: The effect of upstream convection on downstream precipitation. – *Wea. Forecast.* **22**, 255–277.
- MANZATO, A., 2005: An odds ratio parameterization for ROC diagram and skill score indices. – *Wea. Forecast.* **20**, 918–930.
- MARTIUS, O., E. ZENKLUSEN, C. SCHWIERZ, H. C. DAVIES, 2006: Episodes of Alpine heavy precipitation with an overlying elongated stratospheric intrusion: A climatology. – *Int. J. Climatol.* **26**, 1149–1164.
- MARZBAN, C., 1998: Scalar measures of performance in rare-event situations. – *Wea. Forecast.* **13**, 753–763.
- MARZBAN, C., V. LAKSHMANAN, 1999: On the uniqueness of Gandin and Murphy's equitable performance measures. – *Mon. Wea. Rev.* **127**, 1134–1136.
- MASON, I., 1989: Dependence of the critical success index on sample climate and threshold probability. – *Aust. Meteor. Mag.* **37**, 75–81.
- , 2003: Binary Events. – In: JOLLIFFE, I., D. STEPHENSON (Ed.): *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, Wiley and Sons, 37–76.
- MATTHEWS, R., 1996: Base-rate errors and rain forecasts. – *Nature* **382**, 766.
- MCBRIDE, J., E. EBERT, 1999: Verification of quantitative precipitation forecasts from operational numerical weather prediction models over Australia. – *Wea. Forecast.* **15**, 103–121.
- MESINGER, F., 2008: Bias adjusted precipitation threat scores. – *Adv. Geosci.* **16**, 137–142.
- MITTERMAIER, M., 2006: Using an intensity-scale technique to assess the added benefit of high-resolution model precipitation forecasts. – *Atmos. Sci. Lett.* **7**, 35–42.
- MURPHY, A., 1993: What is a good forecast – an essay on the nature of goodness in weather forecasting. – *Wea. Forecast.* **8**, 281–293.
- MURPHY, A., E. EPSTEIN, 1967: A note on probability forecasts and “hedging”. – *J. Appl. Meteor.* **6**, 1002–1004.
- MURPHY, A., R. WINKLER, 1987: A general framework for forecast verification. – *Mon. Wea. Rev.* **115**, 1330–1338.
- PAULAT, M., C. FREI, M. HAGEN, H. WERNLI, 2008: A gridded dataset of hourly precipitation in Germany: its construction, climatology and application. – *Meteorol. Z.* **17**, 719–732.
- PEIRCE, C., 1884: The numerical measure of the success of predictions. – *Science* **4**, 453–454.
- PUJOL, O., J. GEORGIS, M. CHONG, F. ROUX, 2005: Dynamics and microphysics of orographic precipitation during MAP IOP3. – *Quart. J. Roy. Meteor. Soc.* **131**, 2795–2819.
- RICHARD, E., S. COSMA, R. BENOIT, P. BINDER, A. BUZZI, P. KAUFMANN, 2003: Intercomparison of mesoscale meteorological models for precipitation forecasting. – *Hydrol. Earth Syst. Sci.* **7**, 799–811.
- RICHARD, E., A. BUZZI, G. ZÄNGL, 2007: Quantitative precipitation forecasting in the Alps: The advances achieved by the Mesoscale Alpine Programme. – *Quart. J. Roy. Meteor. Soc.* **133**, 831–846.
- RICHTER, D., 1995: Ergebnisse methodischer Untersuchungen zur Korrektur des systematischen Messfehlers des Hellmann-Niederschlagsmessers. – In: *Bericht des Deutschen Wetterdienstes* **194**, 93 pp.
- ROTUNNO, R., R. HOUZE, 2007: Lessons on orographic precipitation from the Mesoscale Alpine Programme. – *Quart. J. Roy. Meteor. Soc.* **133**, 811–830.
- SCHMIDLI, J., C. SCHMUTZ, C. FREI, H. WANNER, C. SCHÄR, 2002: Mesoscale precipitation variability in the region of the European Alps during the 20th century. – *Int. J. Climatol.* **22**, 1049–1074.

- SCHWARB, M., C. DALY, C. FREI, C. SCHÄR, 2001: Mean annual and seasonal precipitation in the European Alps 1971–1990. – In: Hydrological Atlas of Switzerland, Landeshydrol. und Geol., Bern, plates 2.6, 2.7.
- SEVRUK, B., 1985: Systematischer Niederschlagsmessfehler in der Schweiz. – In SEVRUK, B. (Ed.): Der Niederschlag in der Schweiz, Geol. Schweiz. Hydrol. **31**, 65–75.
- SHEPARD, D., 1984: Computer mapping: The SYMAP interpolation algorithm. – In: GAILE, G., C. WILLMOTT (Eds.): Spatial Statistics and Models, Dordrecht. 133–145.
- SKAMAROCK, W., J. KLEMP, 1992: The stability of time-split numerical-methods for the hydrostatic and the nonhydrostatic elastic equations. – Mon. Wea. Rev. **120**, 2109–2127.
- STAUFFER, D., N. SEAMAN, 1994: Multiscale 4-dimensional data assimilation. – J. Appl. Meteor. **33**, 416–434.
- STEPHENSON, D., 2000: Use of the “odds ratio” for diagnosing forecast skill. – Wea. Forecast. **15**, 221–232.
- STEPHENSON, D., B. CASATI, C. FERRO, C. WILSON, 2008: The extreme dependency score: A non-vanishing measure for forecasts of rare events. – Meteor. Apps. **15**, 41–50.
- STEPELER, J., G. DOMS, U. SCHÄTTLER, H. BITZER, A. GASSMANN, U. DAMRATH, G. GREGORIC, 2003: Meso-gamma scale forecasts using the nonhydrostatic model LM. – Meteor. Atmos. Phys. **82**, 75–96.
- SWETS, J., 1988: Measuring the accuracy of diagnostic systems. – Science **240**, 1285–1293.
- THORNES, J., D. STEPHENSON, 2001: How to judge the quality and value of weather forecast products. – Meteor. Apps. **8**, 307–314.
- TIEDTKE, M., 1989: A comprehensive mass flux scheme for cumulus parameterization in large-scale models. – Mon. Wea. Rev. **117**, 1779–1800.
- WERNLI, H., M. PAULAT, M. HAGEN, C. FREI, 2008: SAL – a novel quality measure for the verification of quantitative precipitation forecasts. – Mon. Wea. Rev., published online <http://ams.allenpress.com/perlserv/?request=get-abstract&doi=10.1175>
- WICKER, L., W. SKAMAROCK, 2002: Time-splitting methods for elastic models using forward time schemes. – Mon. Wea. Rev. **130**, 2088–2097.
- WIDMANN, M., C. BRETHERTON, 2000: Validation of mesoscale precipitation in the NCEP reanalysis using a new gridcell dataset for the northwestern United States. – J. Climate **13**, 1936–1950.
- WILKS, D., 2006: Statistical Methods in the Atmospheric Sciences. – Academic Press, chapter 3, 23–70.
- WILLMOTT, C., C. ROWE, W. PHILPOT, 1985: Small-scale climate maps – a sensitivity analysis of some common assumptions associated with grid-point interpolation and contouring. – American Cartographer **12**, 5–16.
- WOODCOCK, F., 1976: Evaluation of Yes-No forecasts for scientific and administrative purposes. – Mon. Wea. Rev. **104**, 1209–1214.
- ZÄNGL, G., 2007: To what extent does increased model resolution improve simulated precipitation fields? A case study of two north-Alpine heavy-rainfall events. – Meteorol. Z. **16**, 571–580.
- ZENG, Z., S. YUTER, R. HOuze, D. KINGSMILL, 2001: Microphysics of the rapid development of heavy convective precipitation. – Mon. Wea. Rev. **129**, 1882–1904.