

# statistical forecasts of sunshine duration

li zhihong and seijo kruizinga

scientific reports; WR 89-05

wetenschappelijke rapporten; WR 89-05

---

de bilt 1989      publicationnumber: Scientific reports = wetenschappelijke  
rapporten; WR 89-05 (DM)

p.o. box 201  
3730 AE de bilt  
wilhelminalaan 10  
tel.+31 30 206911  
telex 47096

Dynamical Meteorology Department

U.D.C.: 551.509.314  
551.509.32  
(492)

ISSN: 0169-1651

© KNMI, De Bilt. All rights reserved. No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher.

# statistical forecasts of sunshine duration

li zhihong and seijo kruizinga

scientific reports; WR 89-05

wetenschappelijke rapporten; WR 89-05

## 1. Introduction

The total amount of sunshine is a useful parameter in describing the weather conditions in the Netherlands on a given day. Usually this parameter is given in relative units ( percent ) indicating the actual duration relative to the maximum possible duration. We will refer to this as the relative sunshine duration RSD. The MOS guidance for the Netherlands introduced in December 1983 (Lemcke and Kruizinga,1988) contains an equation for the prediction of the RSD. This guidance, covering day 1 to day 5 inclusive is based on the products issued by the European Centre for Medium range Weather Forecasting (ECMWF). Because only a limited amount of development data was available at that time a single equation for the whole year was developed. However, the frequency distribution of the RSD is strongly seasonally dependent in the Netherlands. Since there are seven years of data available now, it was decided to develop seasonal equations for this parameter.

Usually linear regression is used to develop equations for parameters like temperature, wind speed and RSD. However, the RSD is not very suitable for linear regression due to its frequency distribution. The RSD is bounded between 0% and 100% and especially in winter there are many days without sunshine. This problem was solved in the following way. First we derived three equations ( per season ) giving the probability that the RSD exceeds 0%, 29% and 59% respectively (RSD values are rounded to the nearest integer). From these probabilities<sup>1</sup> the median of the forecast distribution is inferred through piecewise linear interpolation. This median value is used as a point forecast for the RSD. ( Jensenius, 1988, used a similar procedure). Furthermore it proved to be possible to compute credibility intervals as well using these probability forecasts.

The guidance forecast which is now in operational use contains forecast values verifying at station de Bilt but these

<sup>1</sup>In this report probabilities will be given on a scale ranging from 0.0 to 1.0 to avoid confusion with the RSD which is given in percent.

values are assumed to be valid for the whole area of the Netherlands. In the final section of this paper the possibility of regionalization is examined for the first days of the forecast.

## 2. Potential predictors

All potential predictors that will be examined are derived from the model output produced by the ECMWF. They are listed in Table 1. The first twelve potential predictors are derived directly from the geopotential height field forecasts for lead times +24,+48 out to +144 hours (forecasts valid at 12 GMT). First these forecasts are interpolated to the grid given in Fig. 1. Then the geostrophic winds and the geostrophic vorticity are computed from differentials on this grid. The height, the u and v components of the geostrophic wind and the geostrophic vorticity at the central gridpoint (located at de Bilt) at the levels 500,850 and 1000 hPa are used as potential predictors.

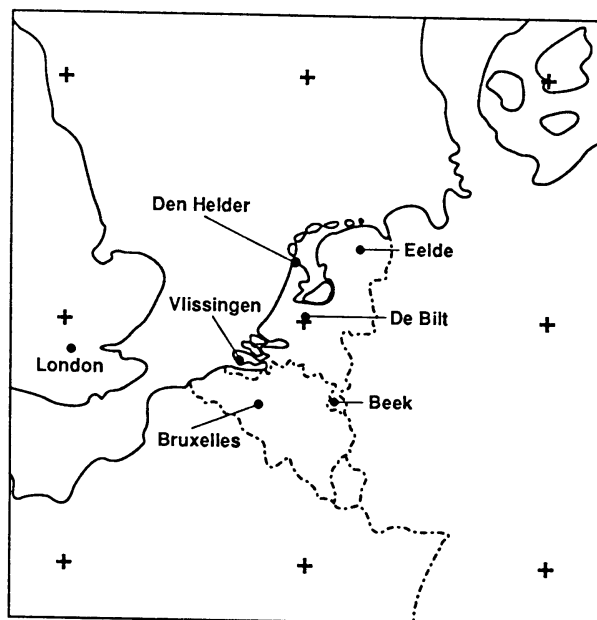


Figure 1. Location of the Dutch stations and the grid used for the ECMWF fields.

Table 1. Potential predictors.

No	Description	No	Description
1	height of 500 hPa	11	u-comp at 1000 hPa
2	v-comp at 500 hPa	12	vorticity 1000 hPa
3	u-comp at 500 hPa	13	cloudiness
4	vorticity 500 hPa	14	average cloudiness
5	height of 850 hPa	15	fifteenth analogue
6	v-comp at 850 hPa	16	analogues N(RSD>0%)
7	u-comp at 850 hPa	17	analogues N(RSD>29%)
8	vorticity 850 hPa	18	analogues N(RSD>59%)
9	height of 1000 hPa	19	height anomaly 500 hPa
10	v-comp at 1000 hPa	20	height anomaly 850 hPa

Furthermore the cloudiness forecasts of the model are used. After interpolation of these forecasts to the grid shown in Fig. 1, the value at the central gridpoint is used. We used the +24,+48 and +72 hour forecasts as well as the time averaged forecasts (+12,+24,+36), (+36,+48,+60) and (+60,+72) as potential predictors. The next four predictors, 15-18, are derived from the 500 hPa forecasts using the analogue technique (Kruizinga and Murphy, 1983). For a given forecast a historical database covering 1949 to 1985 is scanned for similar (or analogue) fields. Only historical fields in the same part (less than twenty days difference in date) of the year are allowed. The sunshine observations of the thirty best analogues are ordered and the 15th value is used as potential predictor. Furthermore we used as predictors the number of analogues with a RSD greater than 0%, 29% or 59% respectively.

Finally the height anomalies of 850 hPa and 500 hPa are offered as predictors. The climatological values were derived from analysed fields covering the period 1972 upto 1979 inclusive.

### 3. Selection of the predictors and development of the equations

As was said before we derived forecast equations for the probability that the RSD exceeds a given level. For each of the three levels 0, 29 and 59% a separate equation was developed for

each of the seasons, winter: December-February, spring: March-May, and so on. For the development of the probabilistic equations the logit model (Brelford and Jones, 1967) was used. For the selection of the predictors from the set of potential predictors we used the stepwise multiple regression technique. Therefore we introduced three 0/1 valued predictands which assume the value 1 if its corresponding level is exceeded. The datasets for each season and each lead time were offered to a forward stepwise regression scheme. This resulted in three lists of predictors (per season and lead time), one for each of the predictands. In order to obtain consistent probabilities these lists were combined manually into one overall list of predictors to be used in the three probabilistic equations used for a given lead time and season. In Table 2 the selected predictors are given in numerical order. As can be seen the model cloudiness ( 13 or 14 ) is selected each time it was available. The analogue predictors proved to be important as well. This is in agreement with the experience obtained with other predictands (Lemcke and Kruizinga, 1988).

The dataset used for the derivation of the equations covered the period December 1980 to November 1986. The equations were tested on the subsequent year December 1986 to November 1987. The logit model used to construct the probabilistic equations models the probability of the occurrence of an event as:

$$P(\text{RSD} > p\%) = 1 / (1 + e^{fx})$$

$$\text{with } fx = a_0 + a_1 * x_1 + a_2 * x_2 \dots\dots$$

The event is the exceedence of a certain level of the RSD. The predictors  $x_1, x_2, \dots$  are selected in the previous step. The coefficients  $a_1, a_2, \dots$  are computed with an iterative procedure which optimizes the likelihood when this equation is applied to the development dataset.

Table 2. Selected predictors per season and lead time.

Season	+24	+48	+72	+96	+120	+144
Winter	1	8	8	1	1	9
	8	11	14	10	10	15
	13	13	16	12	12	16
	16	16	19	16	15	18
		18		17	20	20
Spring	8	3	10	5	11	11
	10	8	12	12	12	12
	11	10	14	16	16	16
	13	14	16	18	18	18
	16	17	20	20	20	20
Summer	1	4	1	5	1	1
	3	10	10	15	5	4
	10	13	14	17	15	15
	13	18	16	18	17	17
	20	20	18	19		
Autumn	1	8	8	5	1	1
	10	10	10	12	8	2
	13	14	14	16	10	5
	16	16	15	18	15	15
	20	20	20	19	16	16

#### 4. Verification of probability forecasts

For the verification of a probability forecast on a yes/no predictand the Brierscore (Brier, 1957) is commonly used. For a series of  $N$  forecasts this Brierscore can be expressed as:

$$BS = \frac{1}{N} \sum (P_n - O_n)^2$$

where  $P_n$  are the predicted probabilities (ranging from 0 to 1) and  $O_n$  are the observations, either 0 or 1. By definition the Brierscore is always positive and will be zero in the case of a perfect forecast. Usually this score is compared with a reference score computed from a reference forecast. In our case we used a monthly climatology as the reference forecast (BSC). The Brierscore of a skilful forecast should be between zero and BSC. In Tables 3 and 4 some verification results on dependent as well as independent data are given as an example.



These tables show that the Brierscores of the equations are generally lower than the Brierscores of climatology on dependent as well as independent data indicating that the forecast scheme improves over climatology. The Brierscore can be transformed into a skillscore BSS by

$$BSS = 100 * \frac{(BSC - BS)}{BSC}$$

BSS indicates the improvement in percent over climatology. In Figures 2 to 4 the seasonal skillscores are plotted for lead times +24,+48 and +72 hours, for dependent as well as independent data. In Fig 2 we can see that the predictions for the RSD>0% are hardly skilful in summer this is mainly due to the fact that in summer nearly every day has sunshine. Therefore it is not easy to improve over climatology in this case. Figure 4 shows that the skill for RSD>59% reaches its maximum during summer and spring. The annual mean indicates moderate skill for the lead times +24,+48 and +72 hours.

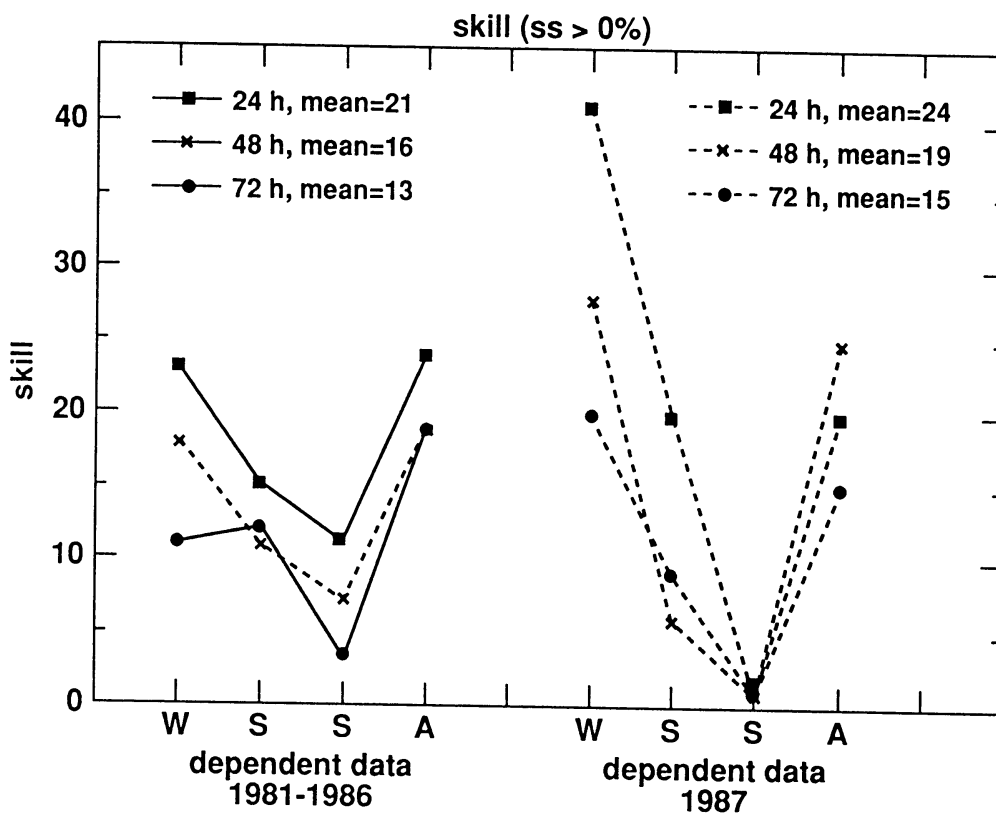


Figure 2. Seasonal Brier Skill Scores for P(RSD>0%) for dependent and independent dataset.

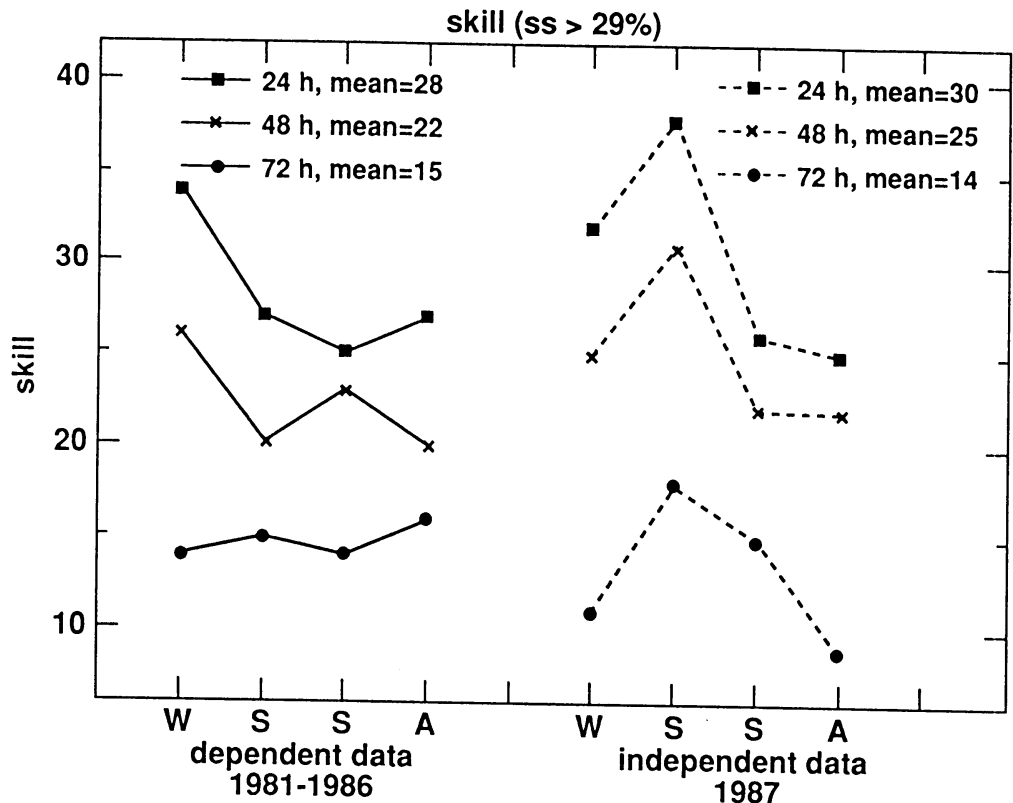


Figure 3. Seasonal Brier Skill Scores for  $P(RSD > 29\%)$  for dependent and independent dataset.

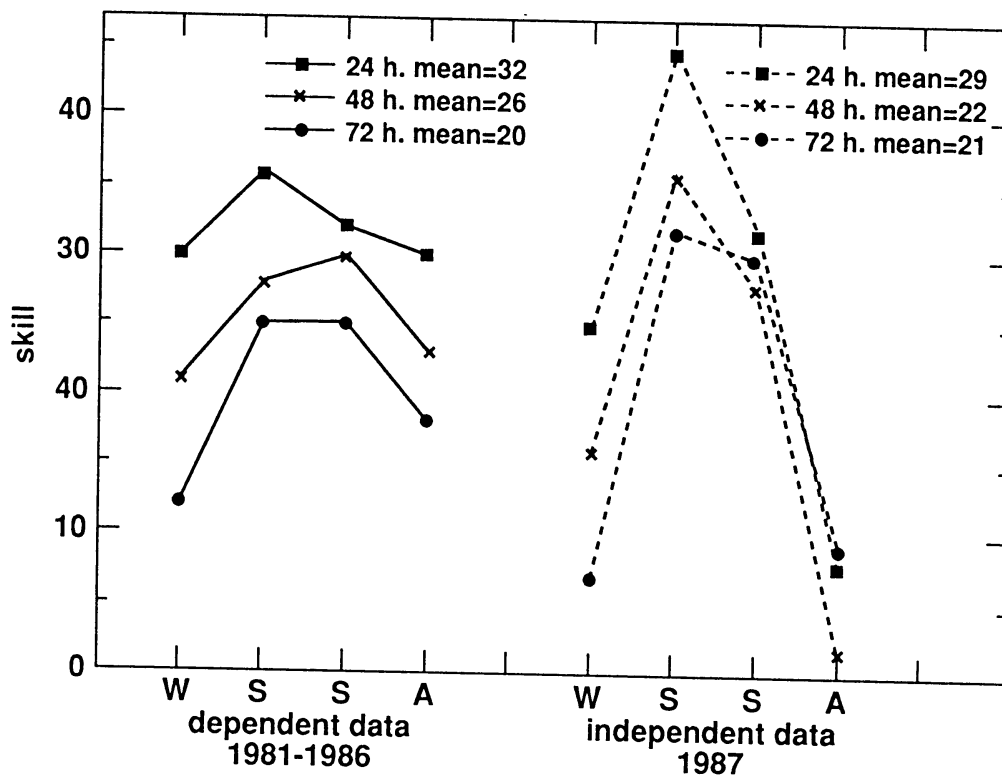


Figure 4. Seasonal Brier Skill Scores for  $P(RSD > 59\%)$  for dependent and independent dataset.

Table 3: Verification results for +48 hour lead time.

Dependent data December 1980 to November 1986 inclusive.

		WINTER	SPRING	SUMMER	AUTUMN	MEAN
RSD>0%	BS	0.198	0.124	0.051	0.150	0.131
	BSC	0.242	0.139	0.055	0.185	0.155
RSD>29%	BS	0.171	0.199	0.185	0.194	0.187
	BSC	0.230	0.250	0.240	0.242	0.241
RSD>59%	BS	0.119	0.119	0.139	0.119	0.124
	BSC	0.152	0.165	0.200	0.154	0.168

Independent data December 1986 to November 1987 inclusive.

		WINTER	SPRING	SUMMER	AUTUMN	MEAN
RSD>0%	BS	0.175	0.126	0.092	0.129	0.130
	BSC	0.244	0.134	0.093	0.172	0.161
RSD>29%	BS	0.154	0.172	0.213	0.182	0.180
	BSC	0.206	0.250	0.274	0.232	0.240
RSD>59%	BS	0.089	0.138	0.095	0.155	0.119
	BSC	0.107	0.215	0.131	0.159	0.153

Apart from being skilful probability forecasts need to be reliable. This means that a forecast of probability P should be followed on average by an observation of the event in a fraction P of all cases. The reliability can be tested with the so-called reliability diagrams. To construct such a diagram the forecast probabilities are grouped into classes 0-0.1,0.1-0.2,0.2-0.3 etc. After this the relative frequency of occurrence of the event is plotted versus the midpoint of the classes. For reliable forecasts the plots should be near to the diagonal (0,0) to (1,1). We used dependent and independent data of all seasons together because a large set of data is needed to construct such reliability diagrams. In the figures 5 and 6 results for lead times +24 and +72 hours are shown. The relative frequencies of forecasts within the probability classes are also given in these figures.

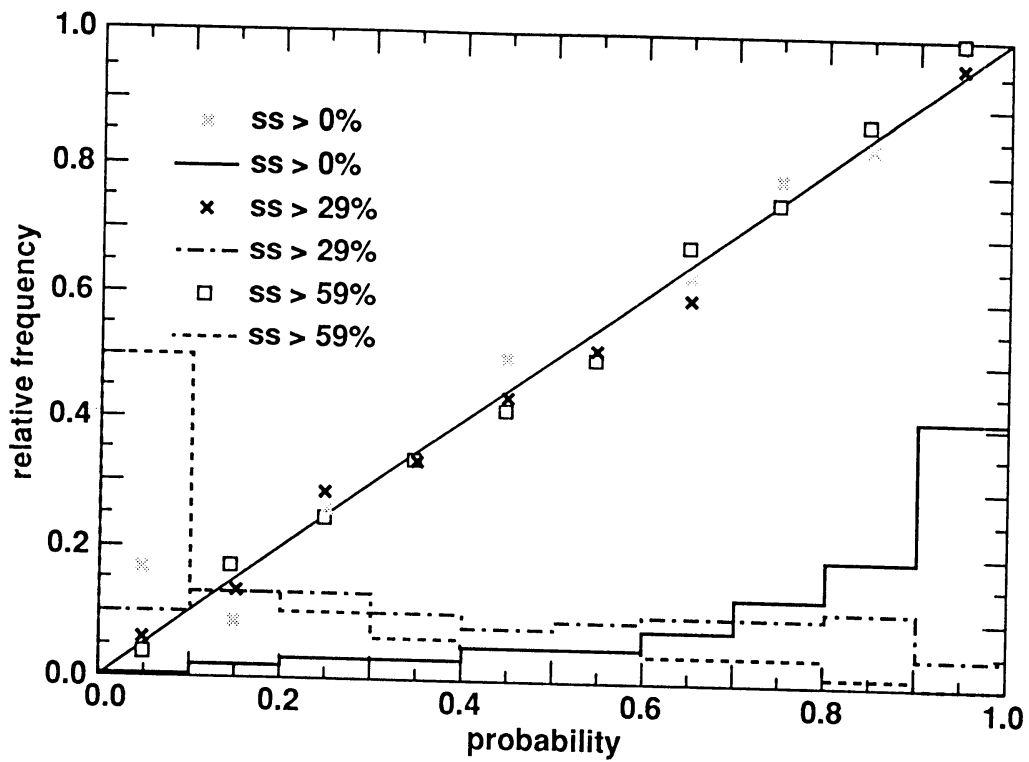


Figure 5. Reliability diagrams based on dependent and independent data, forecast period 24 hours.

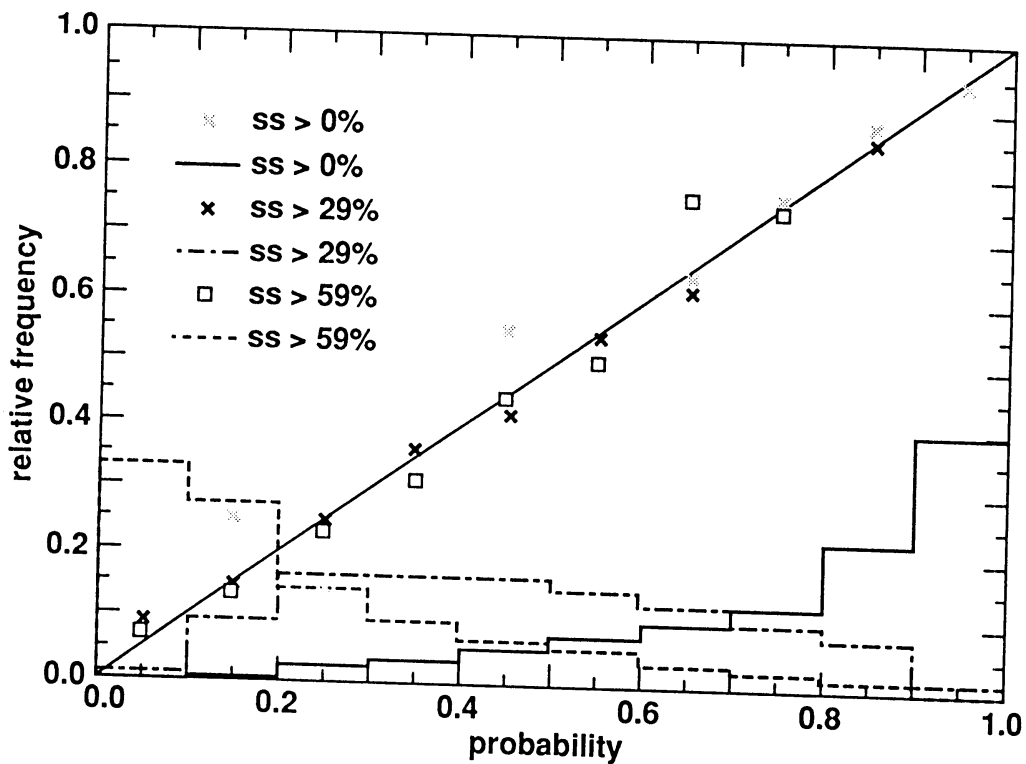


Figure 6. Reliability diagrams based on dependent and independent data, forecast period 72 hours.

Table 4: Verification results for +96 hour lead time.

Dependent data december 1980 to november 1986 inclusive.

		WINTER	SPRING	SUMMER	AUTUMN	MEAN
RSD>0%	BS	0.224	0.130	0.057	0.154	0.141
	BSC	0.244	0.141	0.059	0.180	0.156
RSD>29%	BS	0.213	0.227	0.210	0.218	0.217
	BSC	0.229	0.250	0.239	0.243	0.240
RSD>59%	BS	0.133	0.135	0.175	0.140	0.146
	BSC	0.146	0.170	0.200	0.155	0.168

Independent data december 1986 to november 1987 inclusive.

		WINTER	SPRING	SUMMER	AUTUMN	MEAN
RSD>0%	BS	0.211	0.126	0.083	0.110	0.132
	BSC	0.246	0.127	0.088	0.160	0.155
RSD>29%	BS	0.186	0.215	0.264	0.212	0.219
	BSC	0.204	0.251	0.273	0.233	0.240
RSD>59%	BS	0.071	0.163	0.100	0.166	0.125
	BSC	0.090	0.216	0.135	0.166	0.152

## 5. Transformation of the probabilities to RSD forecasts

The forecast probabilities for a given day can be interpreted as points of the complete conditional probability distribution of the RSD. Forecasts of the RSD itself can be found by estimating for instance the mean of this distribution or the median of this distribution. We preferred to use the median because of the special properties of the distribution of the RSD. The value of the median is estimated through linear interpolation of the forecast probabilities. This interpolation is based on a piecewise linear approximation of forecast probability distribution, (see Figure 7):

1. Define the probabilities:

$$P_0 = P(\text{RSD} > 0\%) \quad P_1 = P(\text{RSD} \geq 30\%) \quad P_2 = P(\text{RSD} \geq 60\%)$$

2. Now we can discriminate four cases with four rules for the computation of the RSD forecast:
1.  $P_0 < 0.5$  then RSD=0%
  2.  $P_0 \geq 0.5 \wedge P_1 < 0.5$  then  $RSD = 30 * (P_0 - 0.5) / (P_0 - P_1)$
  3.  $P_1 \geq 0.5 \wedge P_2 < 0.5$  then  $RSD = 30 * (P_1 - 0.5) / (P_1 - P_2) + 30$
  4.  $P_2 \geq 0.5$  then  $RSD = 40 * (P_2 - 0.5) / P_2 + 60$

The maximum RSD predicted by these rules is 80 percent. However in the Netherlands this is not a serious problem because RSD values greater than 80% are observed on about 22 days in a year only.

The same type of computational procedure was used to estimate the 0.10 and 0.90 percentiles of the conditional distribution. These percentiles were used as boundaries of a credibility interval associated with the RSD forecast.

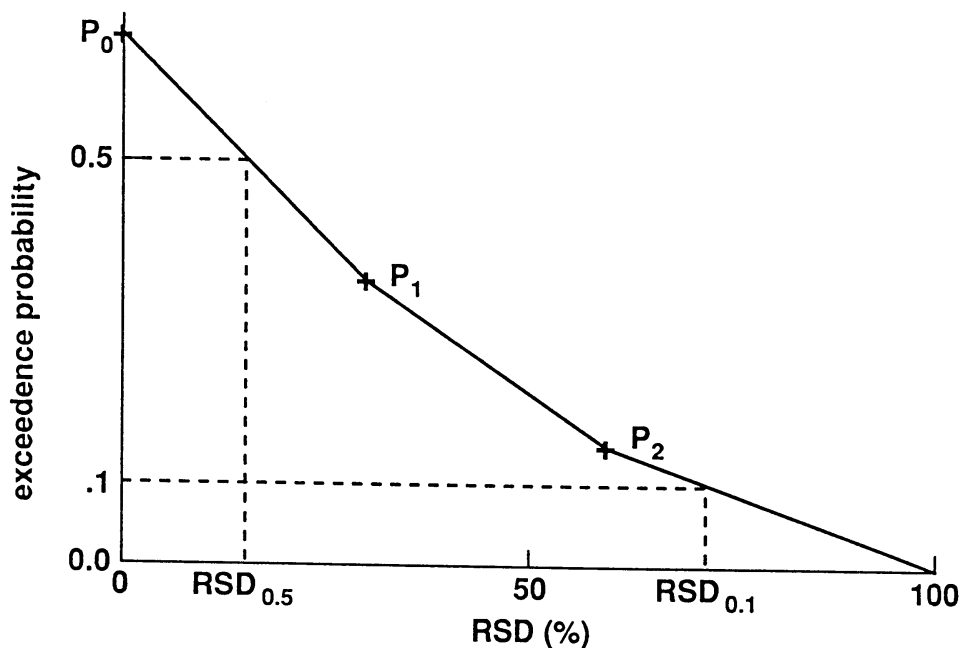


Figure 7. Piecewise linear approximation of the forecast probability distribution and the estimation of the 0.90, 0.50 and 0.10 percentiles of this distribution (see text)

## 6. Verification of RSD forecasts and credibility intervals

The RSD forecasts are verified in terms of mean absolute error (MAE) and correlation coefficient for each lead time and each season. In Figure 8 the seasonal correlation coefficients between forecast and observation are plotted for both the dependent and the independent data. As can be seen the correlation is high during spring and summer. The mean absolute error (Fig. 9) is also high in summer and spring. This is mainly due to the probability distribution of the RSD in the different seasons. The MAE ranges from a maximum of 24% for +144 forecast in spring to a minimum of 15% for the +24 forecast in winter. In Table 4 the yearly averaged mean absolute errors are given for the independent data (December 1986 to November 1987). In this table the verification results of the operational guidance forecast, the final forecast of the forecaster and climatology are also given. Due to the long cutoff time of the ECMWF the forecaster has to use +48 hour forecast for the day 1 forecast and so on. Therefore the the day 1 verification results of the forecaster have been shifted to the +48 hour column and so on. As can be seen the equations improve substantially over the old equations, especially at the short forecast ranges. In Figures 10 to 13 the frequency distributions of the 24 hour forecast errors are plotted for each season separately.

Table 4. Yearly averaged MAE for the new equations (NEW), the operational equations (OLD), the forecaster (FRC) and climatology (CLIM), for the independent dataset (December 1986 to November 1987).

MAE	Lead Time in hours					
	+24	+48	+72	+96	+120	+144
NEW	16.5	17.9	19.4	19.8	20.0	20.2
OLD	18.9	20.2	21.0	20.9	20.4	21.0
FRC		17.0	20.0	20.0	21.0	23.0
CLIM	23.7	23.8	23.7	23.8	23.8	23.8

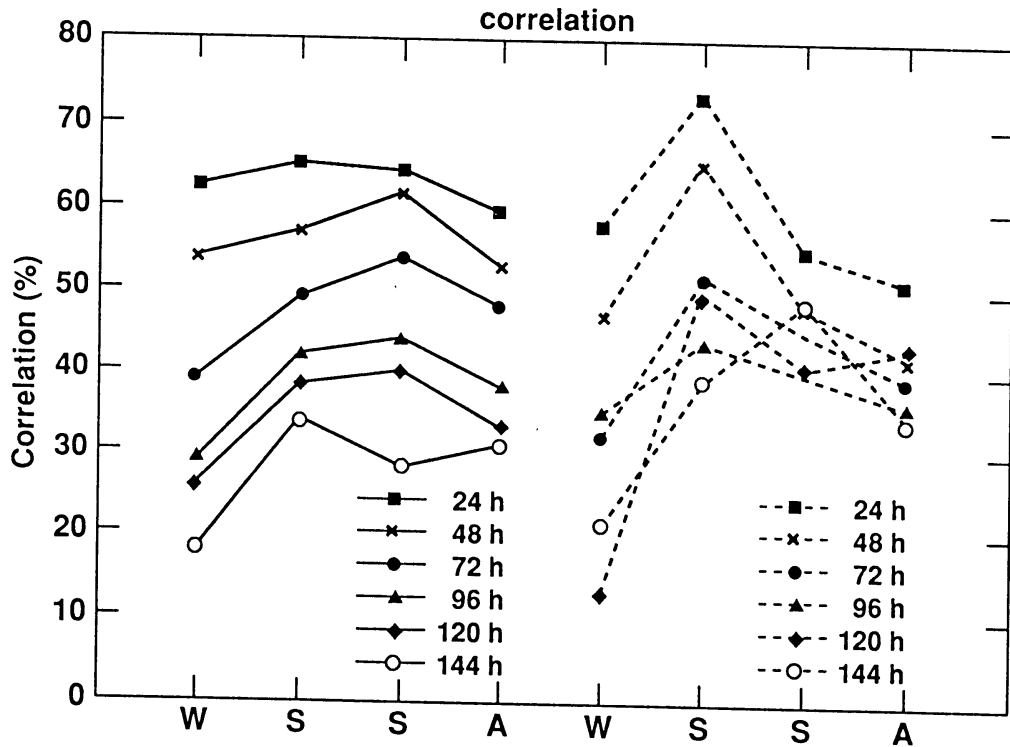


Figure 8. Seasonal correlation coefficient between forecast and observation for dependent ( right ) and independent ( left ) data.

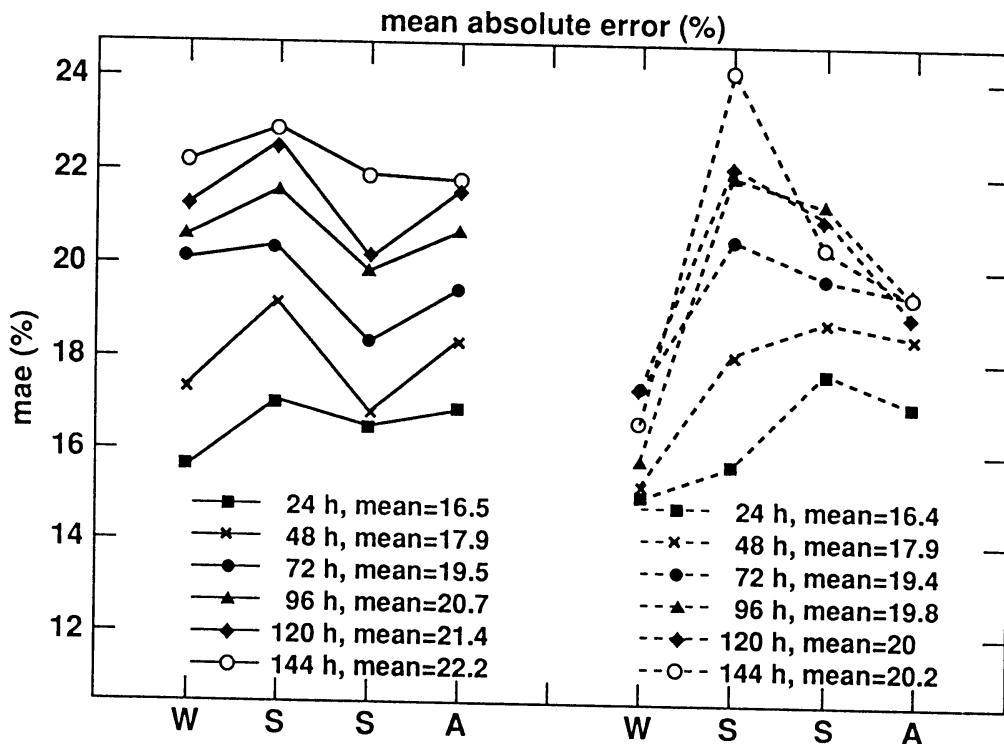


Figure 9. Seasonal Mean Absolute Error for dependent (right) and independent (left) data.



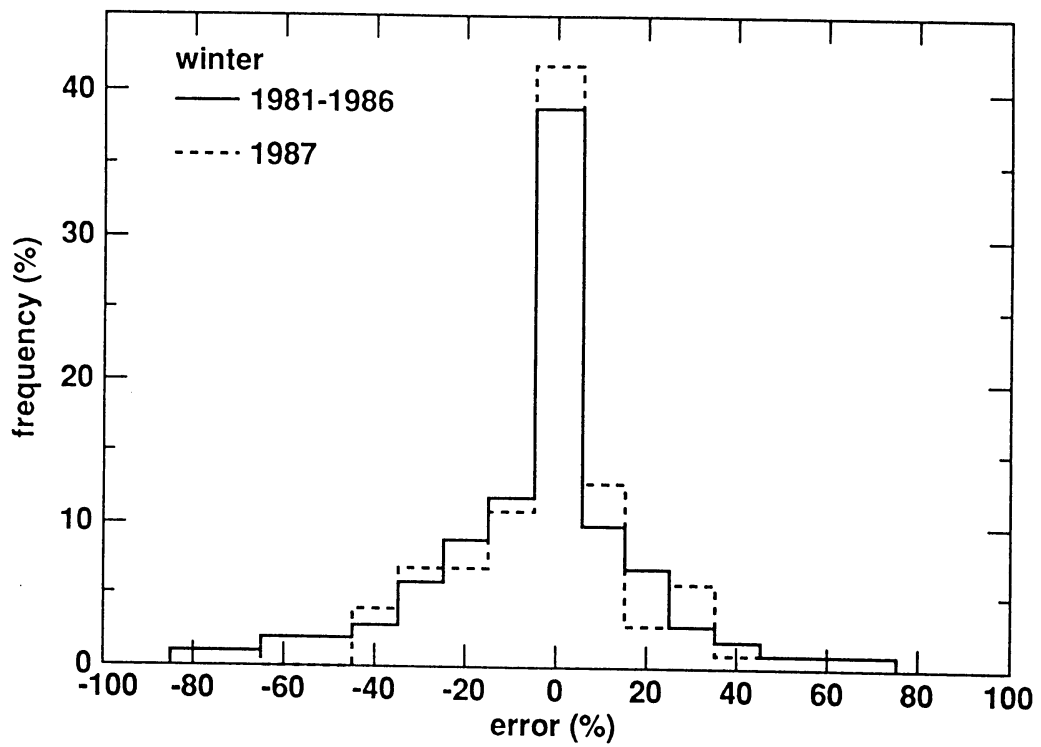


Figure 10. Frequency distribution of the forecast errors of RSD in the winter, forecast period 24 hours.

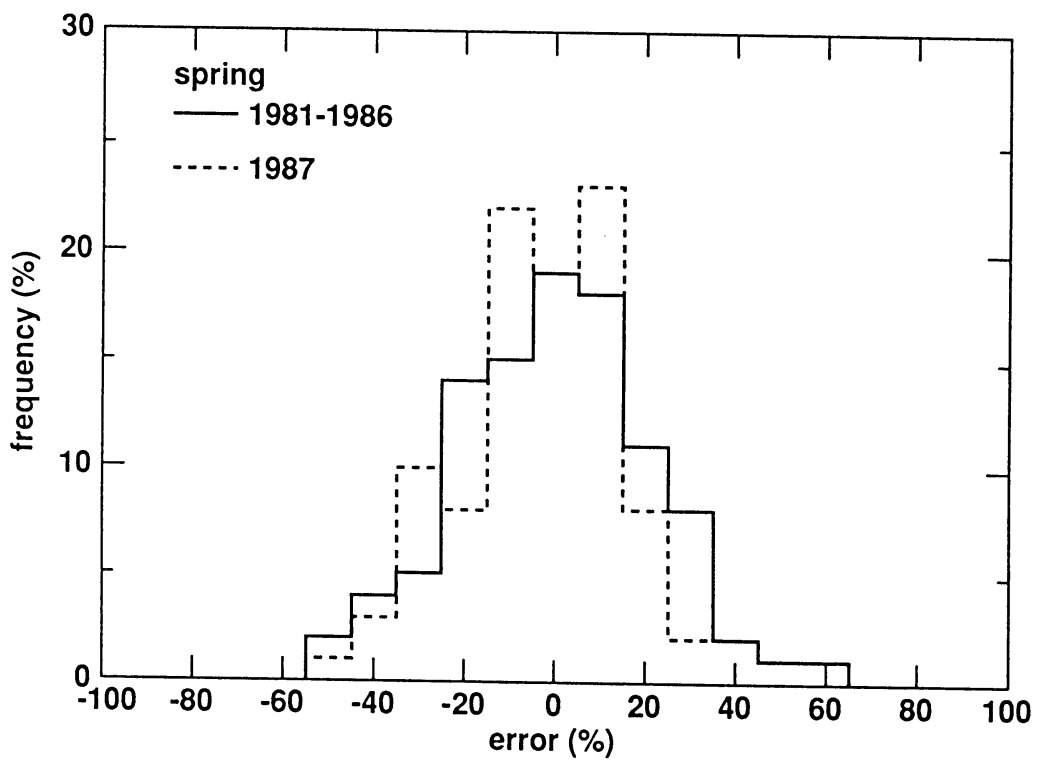


Figure 11. Frequency distribution of the forecast errors of RSD in the spring, forecast period 24 hours.

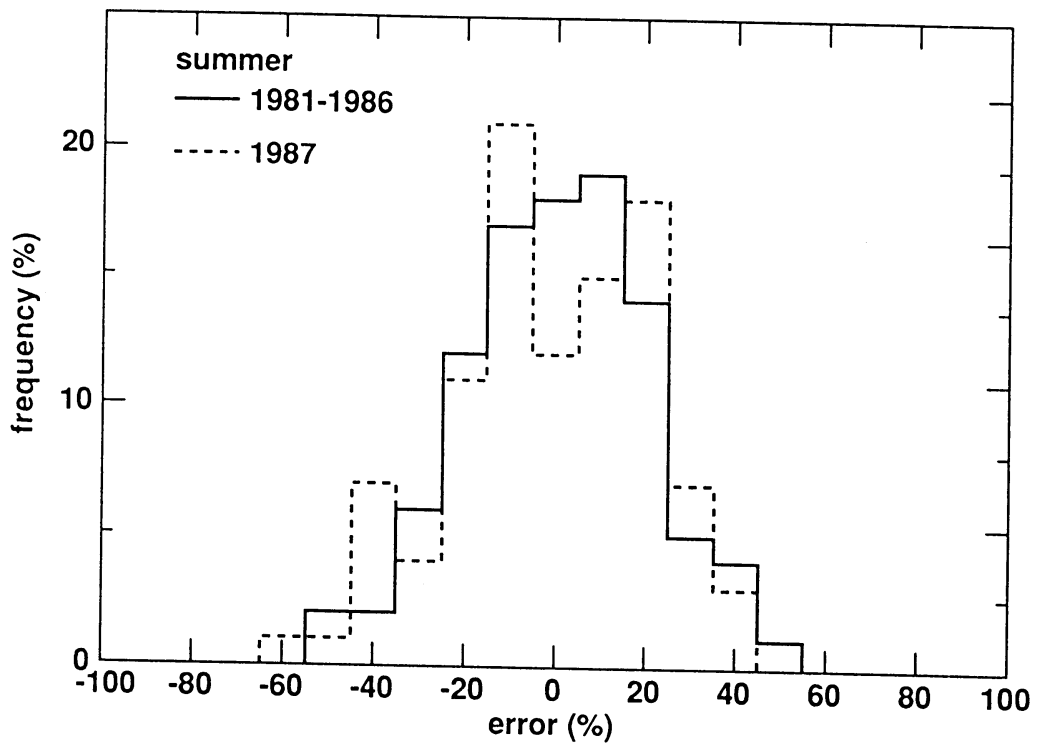


Figure 12. Frequency distribution of the forecast errors of RSD in the summer, forecast period 24 hours.

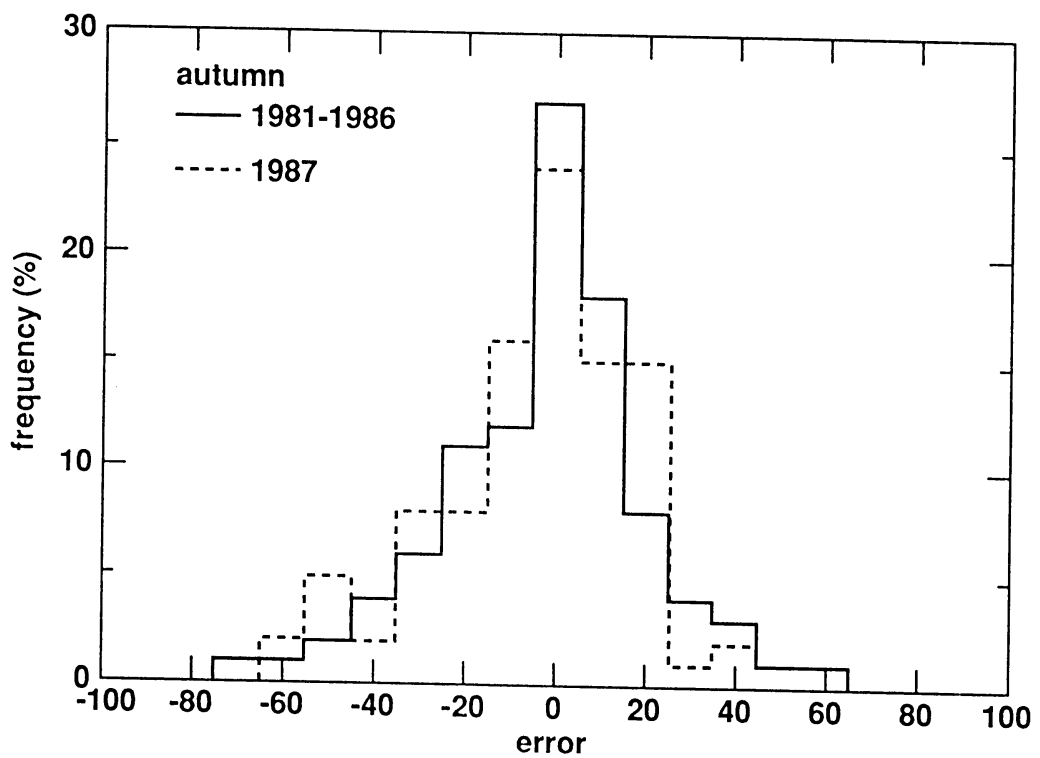


Figure 13. Frequency distribution of the forecast errors of RSD in the autumn, forecast period 24 hours.

The credibility intervals associated with the RSD forecasts were also verified. However in this case verification is difficult because there exists no observed value with which the forecast interval can be compared. So only indirect verification is possible. It is assumed that the average error of the RSD forecast will be small in case of small forecast credibility intervals and larger in other cases. Therefore the RSD forecasts were grouped according to the width of the credibility intervals (CIW). The first group contained the forecasts with a  $CIW < 0.35$ , the next group the forecasts with  $0.35 \leq CIW \leq 0.70$ . The forecasts with  $CIW > 0.70$  were assigned to the third group. For these groups the Mean Absolute Error (MAE) was computed. In Table 5 the MAE's for the different groups for dependent as well as independent data are given.

Table 5. Verification results of RSD forecasts for three groups of credibility interval width (CIW). dependent data, (DEP), December 1980 to November 1986, independent data (IND), December 1986 to November 1987.

Lead time 24 hours.				
Season	Period	$CIW < 0.35$	$0.35 \leq CIW \leq 0.70$	$CIW > 0.70$
Cool	DEP	5.5(226)	15.6(470)	25.9(348)
	IND	3.8( 36)	16.5( 72)	22.4( 54)
Warm	DEP	8.1( 44)	15.2(578)	18.5(414)
	IND	12.2( 13)	15.8(119)	21.7( 49)

Lead time 48 hours.				
Season	Period	$CIW < 0.35$	$0.35 \leq CIW \leq 0.70$	$CIW > 0.70$
Cool	DEP	8.3(240)	18.0(468)	26.6(331)
	IND	4.0( 35)	18.5( 75)	24.9( 48)
Warm	DEP	11.9( 50)	17.3(579)	19.7(402)
	IND	11.5( 14)	16.7(118)	24.8( 45)

This table clearly shows that RSD forecasts with narrow credibility intervals are also accurate forecasts.

## 7. Regionalization of RSD forecasts.

The RSD forecasts studied in the preceding paragraphs are forecasts which verify in the centre of the Netherlands (station de Bilt). These forecasts are assumed to be useful throughout the Netherlands. For the longer ranges this assumption will hold in the sense that forecast errors will be larger than regional differences. In this paragraph it will be examined if regionalization is useful at shorter ranges. Therefore we developed RSD equations for four other stations in the Netherlands namely Den Helder (DH), Eelde (EE), Vlissingen (VL) and Beek (BK), see Fig. 1. The average distance between these stations is about 150 kilometers, therefore the same set of potential predictors at the same gridpoint is used for all stations. In Table 6 the mean absolute errors of the RSD forecasts for each of the stations is given for the independent period, December 1986 to November 1987.

Table 6. MAE for five stations in the period December 1986 to November 1987. Lead times (LT) +24 and +48 hours.

LT	DH	DB	EE	VL	BK
24	16.5	16.5	17.2	17.1	15.6
48	19.0	17.9	18.2	18.4	17.5

This table indicates that the differences in forecast skill for the five stations are rather small. Of course the stations are very close to each other so this result could be expected. However in order to conclude that regionalization is useful it is necessary that there is skill in the forecast of the daily differences between the stations. In order to verify this we computed for each pair of stations the time series correlation between the difference of the RSD forecast and the difference of the observed RSD. The correlations presented in Table 7 are based on 90 observations. So the level of significance is about 0.22. In general the results indicate that for Beek a separate equation will be useful.

Table 7. Correlation coefficient (\*100) between forecast and observed daily differences. Four seasons and two lead times. Period December 1986 to November 1987.

		Lead time 24 hr.				Lead time 48 hr.			
		DB	EE	VL	BE	DB	EE	VL	BE
Winter	DH	-12	23	-16	46	3	19	0	46
	DB		16	-1	53		23	6	44
	EE			38	41			33	32
	VL				53				48
Spring	DH	26	17	30	41	20	11	5	24
	DB		14	20	25		-16	23	13
	EE			25	22			10	3
	VL				31				25
Summer	DH	-12	20	0	29	-5	33	0	33
	DB		-9	6	30		9	-4	26
	EE			33	22			19	22
	VL				40				35
Autumn	DH	-1	20	28	36	0	24	4	20
	DB		35	23	14		34	6	16
	EE			30	25			31	26
	VL				50				37

## 8. Conclusion

In this study MOS forecast equations for the relative sunshine duration RSD were developed. In paragraph 3 seasonal probabilistic equations were developed. These equations give the probability that the RSD observation will exceed the levels 0%, 29% and 59% respectively. In par. 4 it was shown that these equations have skill. Subsequently these probabilistic forecasts were transformed into a point forecast with an associated credibility interval. The point forecasts obtained in this way were substantially better than those obtained with the operational yearly equations, at least at the forecast ranges +24, +48 and +72 hours. This improvement is caused by the introduction of ECMWF cloudiness forecasts as predictor as well the development of seasonal equations. Furthermore it was shown that narrow credibility intervals are associated with more accurate forecasts. In the last paragraph it was shown that it is useful, even within the Netherlands, to develop at least two separate equations.

## References.

- Brelsford, W.M. and R.H. Jones, 1967:  
Estimating probabilities. *Monthly Weather Review*, Vol. 95,  
No. 8.
- Brier, G.E., 1950:  
Verification of forecasts expressed in terms of probability.  
*Monthly Weather Review*, Vol. 78, No. 1, 1-3
- Jensenus J.S. Jr., 1988:  
Objectively Forecasting Sunshine. *Weather and Forecasting*,  
Vol. 3, 5-17.
- Kruizinga, S. and A.H. Murphy, 1983:  
Use of an Analogue Procedure to Formulate Objective  
Probabilistic Temperature Forecasts in the Netherlands.  
*Monthly Weather Review*, Vol. 111, 2244-2254.
- Lemcke, C. and S. Kruizinga, 1988:  
Model Output Statistics Forecasts: Three Years of Operational  
Experience in the Netherlands. *Monthly Weather Review*, Vol.  
116, 1077-1090.