

An aerial photograph of a town, likely in the Alps, is shown from a high angle. The town is surrounded by green hills and is partially obscured by thick, white clouds. Overlaid on the bottom left of the image is a weather map with white contour lines and arrows indicating wind direction and speed. The contour lines are labeled with values such as 1010, 1015, 1020, 1025, 1030, 1035, 1040, 1045, 1050, 1055, 1060, 1065, 1070, 1075, 1080, 1085, 1090, 1095, 1100, 1105, 1110, 1115, 1120, 1125, 1130, 1135, 1140, 1145, 1150, 1155, 1160, 1165, 1170, 1175, 1180, 1185, 1190, 1195, 1200, 1205, 1210, 1215, 1220, 1225, 1230, 1235, 1240, 1245, 1250, 1255, 1260, 1265, 1270, 1275, 1280, 1285, 1290, 1295, 1300, 1305, 1310, 1315, 1320, 1325, 1330, 1335, 1340, 1345, 1350, 1355, 1360, 1365, 1370, 1375, 1380, 1385, 1390, 1395, 1400, 1405, 1410, 1415, 1420, 1425, 1430, 1435, 1440, 1445, 1450, 1455, 1460, 1465, 1470, 1475, 1480, 1485, 1490, 1495, 1500. The arrows are labeled with values such as 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150. The background of the slide is a dark blue gradient with a white wave-like pattern at the top left and bottom right.

# Local forecast of probability distributions by statistical adaptation

S. FARGES ([serge.farges@meteo.fr](mailto:serge.farges@meteo.fr))

DPREVI/GCRI

April 05 2012



**METEO FRANCE**  
Toujours un temps d'avance

## Introduction

In an atmospheric model, coasts, orography, soil type and vegetation status are not perfectly well represented. In addition for computation purposes that space is discretized ; by a way of consequence **a statistical post-processing is necessary to account for the local climatology of the site (station).**

This is called **statistical adaptation** of numerical model outputs (SA). Meteo-France currently produces approximately **3.5 billion forecasts using SA per day** (for a station, a parameter and time range given).

Uncertainty is present since the start of the forecast covered time range. **We want to increase the use of a probabilistic approach**, including, if possible, by forecasting the distributions of the parameters. **We have several methods available to us**, as we shall see...



## Forecast of a probabilistic distribution : different ways

### 1. Dynamic approach : the ensemble forecast

- ▶ The ensemblist systems (varEPS, NCEP, PEARP...) :
  - ▶ Advantages : realizations of trajectories and maps, extreme phenomena forecasting (if enough members).
  - ▶ Disadvantage : cost ( $\Rightarrow$  limited resolution).
- ▶ The multi-models systems (IFS + ARPEGE + ...) :
  - ▶ Advantages : realizations of trajectories and maps, quality of the deterministic forecast in the short time range.
  - ▶ Disadvantages : more time to collect the data from the various producers and limited number of ensemble members.
- ▶ The multi-ensemblists systems (TIGGE) :
  - ▶ Advantages and disadvantages of previously described approaches.

**Often require post-processing statistics** (because the produced probabilities are often unreliable) : Ensemble Dressing, Bayesian Model Averaging (BMA), Nonhomogenous Gaussian Regression (NGR) or Ensemble Regression.



### 2. Statistical approach : probabilistic statistical adaptation of an atmospheric model

- ▶ Discrimination models (LDA, logistic regression, neural networks...) :
  - ▶ Advantages : cost and robustness of linear methods.
  - ▶ Disadvantage : production of occurrence probabilities but no production of distribution.
- ▶ Generalized Linear Models :
  - ▶ Advantages : cost and possibility to use the underlying probability distributions.
  - ▶ Disadvantages : limited number of supported distributions (exponential family) not always perfectly calibrated after fitting.



### **Proposed new approach** : *the generalized linear regression by Gaussian anamorphosis*

- ▶ Advantages : cost, robustness (linear model), no constraint on probabilistic distributions and strongly calibrated (in theory). Can be used for the postprocessing of an ensemble forecast coupled with a BMA.
- ▶ Disadvantage : generally no direct computable formulation of deterministic forecast (expectation).

**Remark** : statistical methods have the disadvantage of not being able to easily produce realistic realizations of trajectories or map.

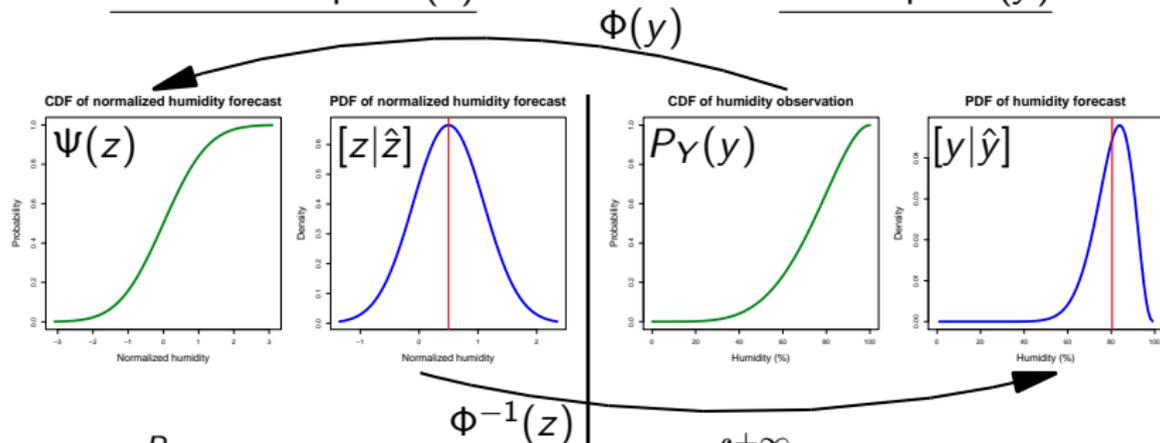


# Generalized linear regression by Gaussian anamorphosis

## 1. Methodology (beta distribution for humidity)

Normalized space (z)

Initial space (y)



$$\hat{z} = \sum_{p=1}^P a_p x_p \quad (\epsilon_{\hat{z}} = \hat{z} - z)$$

$$z|\hat{z} \sim \mathcal{N}(\hat{z}, \sigma_{\epsilon_{\hat{z}}})$$

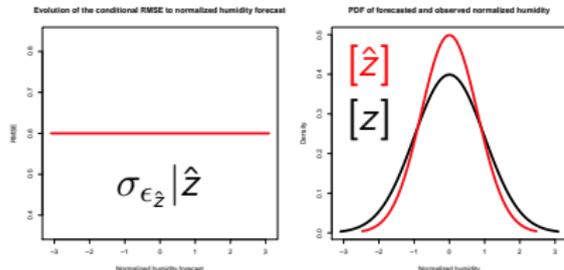
$$\hat{y} = \int_{-\infty}^{+\infty} \Phi^{-1}(z) \varphi(z, \hat{z}, \sigma_{\epsilon_{\hat{z}}}) dz$$

$$[y|\hat{y}] = [y] \frac{\varphi(\Phi(y), \hat{z}, \sigma_{\epsilon_{\hat{z}}})}{\varphi(\Phi(y))}$$

$$\Phi(y) = \Psi^{-1} \circ P_Y(y)$$

# Generalized linear regression by Gaussian anamorphosis

## Normalized space ( $z$ )

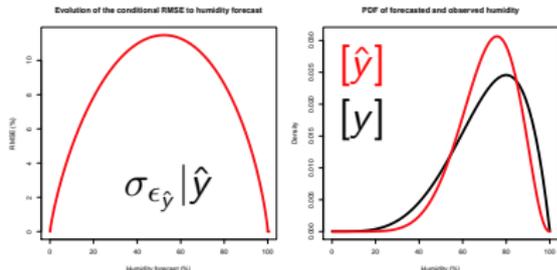


$$\sigma_z^2 = \sigma_{\hat{z}}^2 + \sigma_{\epsilon_{\hat{z}}}^2$$

$$\text{corr}(\hat{z}, \epsilon_{\hat{z}}) = 0 \text{ et } E(\epsilon_{\hat{z}}|\hat{z}) = 0$$

- ▶ Independance between  $\hat{z}$  and  $\epsilon_{\hat{z}}$
- ▶ Homoscedasticity
- ▶ Probabilistic and marginal calibration of  $[z|\hat{z}]$

## Initial space ( $y$ )



$$\sigma_y^2 = \sigma_{\hat{y}}^2 + \sigma_{\epsilon_{\hat{y}}}^2$$

$$\text{corr}(\hat{y}, \epsilon_{\hat{y}}) = 0 \text{ et } E(\epsilon_{\hat{y}}|\hat{y}) = 0$$

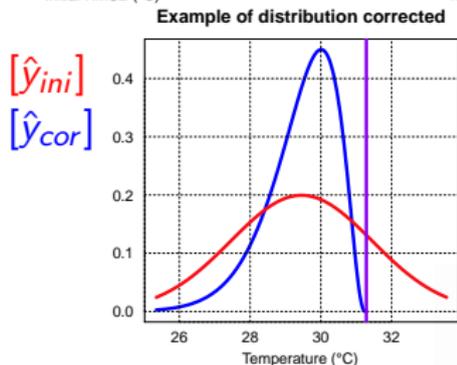
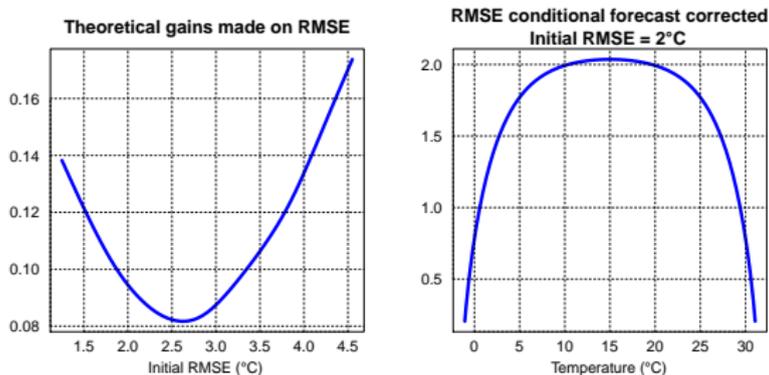
- ▶ Only linear independance between  $\hat{y}$  and  $\epsilon_{\hat{y}}$
- ▶ No homoscedasticity
- ▶ Probabilistic and marginal calibration of  $[y|\hat{y}]$



# Generalized linear regression by Gaussian anamorphosis

## 3. Application examples

Using a truncated normal distribution for a better forecasting of extreme temperatures (theoretical simulation)



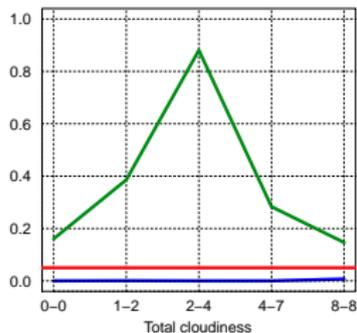
— Record value

# Generalized linear regression by Gaussian anamorphosis

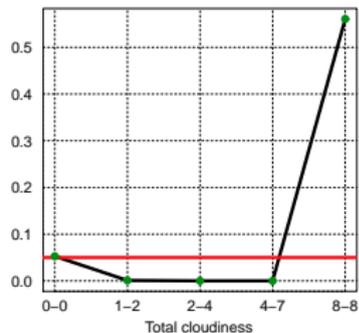
Forecast of the cloudiness : comparison of the qualities of the model against those of logistic regressions (for probabilities forecast) and linear regression (for deterministic forecast)

- New approach
- Classical approach

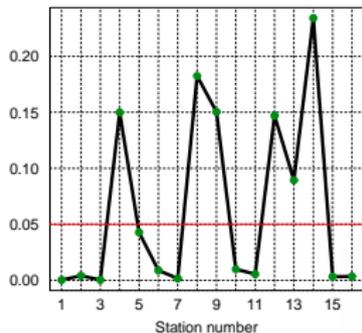
Probability of the hypothesis reliability



Probability that the ROC AUC are different



Probability that the RMSE are different



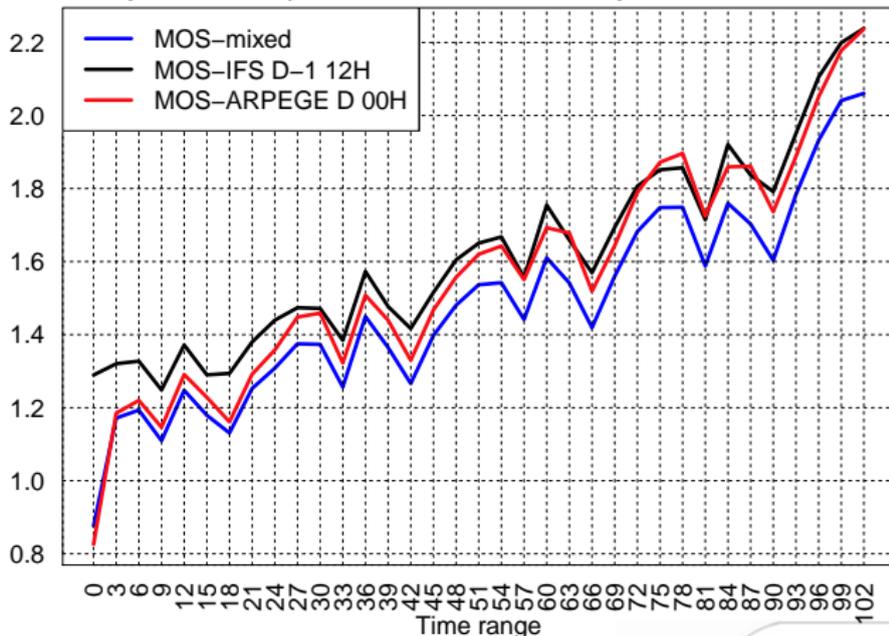
- New approach is better
- Classical approach is better



## Probabilistic mixed statistical adaptation

This is a statistically post-processed multi-model system.  
Advantage for the deterministic forecast for the short time range :

**RMSE of temperature forecasts mixed versus its components (average of 25 european stations). Calculated results for the period 09/2011 to 01/2012.**



## Probabilistic mixed statistical adaptation

The goal of the probabilistic mixed SA is to fit a linear regression model to each model of an ensemble set, using a BMA.

We have  $K$  models  $M_k$ .  $z$  distribution forms (in the normalized space) given by the BMA :

$$p(z | (M_k)_{k=1, \dots, K}, \Theta, z^T) = \sum_{k=1}^K \alpha_k p(z | M_k, \theta_k, z^T)$$
$$z | M_k, \theta_k, z^T \sim \mathcal{N}\left(\sum_{i=1}^q a_i^k x_i^k, \sigma_k\right)$$

$z^T = (z_1, \dots, z_n)$ ,  $\theta_k = ((a_i^k)_{i=1, \dots, q}, \sigma_k)$  et  $\Theta = ((\alpha_k, \theta_k)_{k=1, \dots, K})$

We estimate the  $\Theta$  parameters using the maximum likelihood with the EM algorithm.

**Step E (expectation)** : we estimate the probabilities  $p(M_k | \Theta, z^T)$

$$p(M_k | \Theta^g, z^T) = \frac{\alpha_k^{g-1} p(z^T | M_k, \theta_k^{g-1})}{\sum_{l=1}^K \alpha_l^{g-1} p(z^T | M_l, \theta_l^{g-1})}$$



## Probabilistic mixed statistical adaptation

Step M (maximisation) : parameters  $\Theta$  are iteratively estimated

$$\alpha_k^g = \frac{1}{n} \sum_{j=1}^n p(M_k | \Theta^g, z_j)$$

$$(a_i^k)_{i=1, \dots, q}^g = (X_k' P_k^g X_k)^{-1} X_k' P_k^g z^T$$

$X_k$  matrix of model predictors  $k$

$$P_k^g = \begin{pmatrix} p(M_k | \Theta^g, z_1) & 0 & \dots & 0 \\ 0 & p(M_k | \Theta^g, z_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & p(M_k | \Theta^g, z_n) \end{pmatrix}$$

$$\sigma_k^g = \sqrt{\frac{1}{n\alpha_k^g} \sum_{j=1}^n p(M_k | \Theta^g, z_j) (\hat{z}_{kj}^g - z_j)^2}$$



## Conclusion

- ▶ Among the various ways to produce probability distributions, the *generalized linear regression by Gaussian anamorphosis* has the advantage of being **inexpensive, robust, to overcome discrimination techniques** requiring a statistical model by class, and **to issue reliable probabilistic forecasts** (if the observations distribution is adjusted properly and the number of predictors sufficiently high).
- ▶ The deterministic forecast produced by this model **has good properties** and may be better than the one obtained by a linear regression.
- ▶ Coupled with BMA, the Gaussian anamorphosis can be used to improve the quality of an ensemble forecast. This is a statistical-dynamical system **particularly promising**.



## Prospects : probabilistic approach and spatialization

- ▶ Multi-parameters probabilistic mixed SA (**course in progress**).
- ▶ **Probabilistic spatialization of SA on a regular grid.** It can be facilitated in a normalized space by anamorphosis. Idea of a procedure :
  1. Normalization of all the explanatory variables before regression.
  2. Report of the regressions based on a spatial classification of the parameter analyzed by AROME model.
  3. Evaluation of the observations distributions on the grid based on the observed data of the stations and analyzed by AROME (suggestion : interpolation of the Hermite's polynomials coefficients with ajustement cubic splines).
- ▶ **Exploiting the forecaster's expertise in probabilistic forecasting**, especially in the case of events strongly bi-modal (effect of the presence or not of low clouds on the temperature).





The end



**METEO FRANCE**  
Toujours un temps d'avance