



House of
Energy Markets
& Finance



Marginal-Copula-Scores for Multivariate Forecasting Evaluation

Florian Ziel

University of Duisburg-Essen

February 28, 2019

UNIVERSITÄT
DUISBURG
ESSEN

Open-Minded

Motivation: forecasting and evaluation

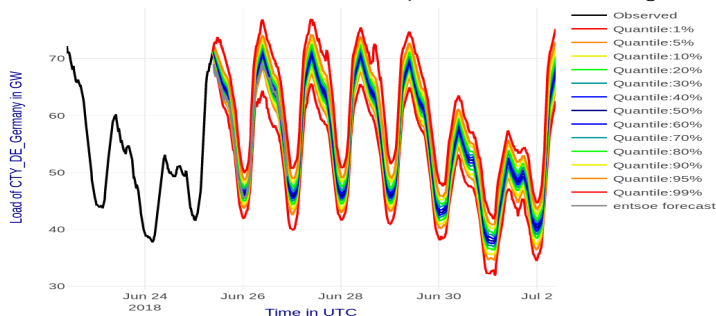
Standard setting:

- ▶ given historic data (and a statistical model) creating a forecast for a target of interest
- ▶ target of interest $\mathbf{Y} \sim \mathbf{F}_{\mathbf{Y}}$, H -dimensional random variable
e.g. maximum load of tomorrow and day after tomorrow (2-dim.)
- ▶ in practice we never know the true $\mathbf{F}_{\mathbf{Y}}$ we just observe \mathbf{y}
- ▶ if we have a *forecast* we can only compare the performance by comparing it with \mathbf{y}
- ▶ evaluation relies on some repeatability of the forecasting experiment
- ▶ ways to report forecast
 - **point forecasting:**
 - $\hat{\mathbf{X}}$ estimator for e.g. $\mathbb{E}(\mathbf{Y})$ or $\text{Med}(\mathbf{Y})$
 - evaluation based on forecasting error $\mathbf{Y} - \hat{\mathbf{X}}$, resp. $\mathbf{y} - \hat{\mathbf{X}}$
 - \mathbb{E} can be strictly proper evaluated using MSE (mean square error)
 - Med can be strictly proper evaluated using MAE (mean absolute error)

Forecasting evaluation

probabilistic forecasting: (everything beyond point forecasting)

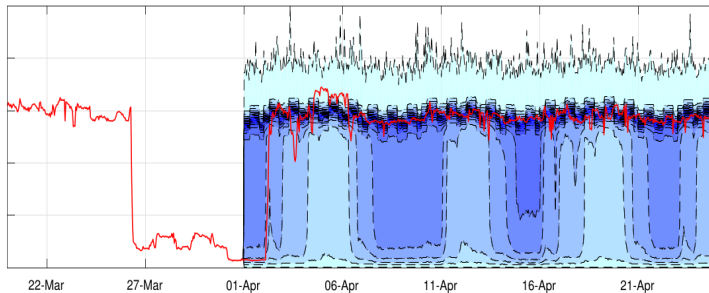
- ▶ characterises the uncertainty in the forecast, e.g.:



- ▶ usually forecast of marginal distributions F_{Y_h} of $\mathbf{Y} = (Y_1, \dots, Y_H)'$, marginal densities f_{Y_h} , or quantiles $q_\alpha(Y_h)$ for $\alpha \in \mathcal{A}$
- ▶ strictly proper evaluation methods available, e.g. continuous rank probability score (**CRPS**) for F_{Y_h}
- ▶ used in e.g. Global Energy Forecasting Competitions (99%-tiles)

Probabilistic forecasting evaluation

- ▶ Problem with standard probabilistic methods:
 - forecasting only the marginals distributions
 - ignoring the dependency structure



(source: Berk, Hoffmann, Müller (2017) *International Journal of Forecasting*)

- ▶ require full forecast F_X for F_Y and strictly proper evaluation method

Evaluation measures for multivariate distributions

some measures available

► Energy score

$$\text{ES}_\beta(\mathbf{F}_X, \mathbf{y}) = \mathbb{E} \left(\|\mathbf{X} - \mathbf{y}\|_2^\beta \right) - \frac{1}{2} \mathbb{E} \left(\|\mathbf{X} - \widetilde{\mathbf{X}}\|_2^\beta \right) \quad (1)$$

- $\beta > 0$, $\mathbf{X}, \widetilde{\mathbf{X}} \stackrel{iid}{\sim} \mathbf{F}_X$
- if $H = 1$ and $\beta = 1 \rightsquigarrow$ CRPS
- strictly proper

► Variogram score

$$\text{VS}_p(\mathbf{F}_X, \mathbf{y}; \mathbf{W}) = \sum_{i=1}^H \sum_{j=1}^H w_{i,j} (|y_i - y_j|^p - \mathbb{E}|X_i - X_j|^p)^2$$

- with $p > 0$ and weight matrix $\mathbf{W} = (w_{i,j})_{i,j}$ (usually $w_{i,j} = c$)
- not strictly proper (*forecasts with shifted mean have same score*)

Evaluation measures for multivariate distributions

► Log-score

$$\text{LogS}(\mathbf{F}_X, \mathbf{y}) = \log(\mathbf{f}_X(\mathbf{y})).$$

- where \mathbf{f}_X is density of \mathbf{F}_X
- strictly proper
- density forecast for \mathbf{X} often not available (even if \mathbf{X} is continuous)

► Dawid-Sebastiani score

$$\text{DSS}(\mathbf{F}_X, \mathbf{y}) = \log(|\Sigma_X|) + (\mathbf{y} - \boldsymbol{\mu}_X)' \Sigma_X^{-1} (\mathbf{y} - \boldsymbol{\mu}_X)$$

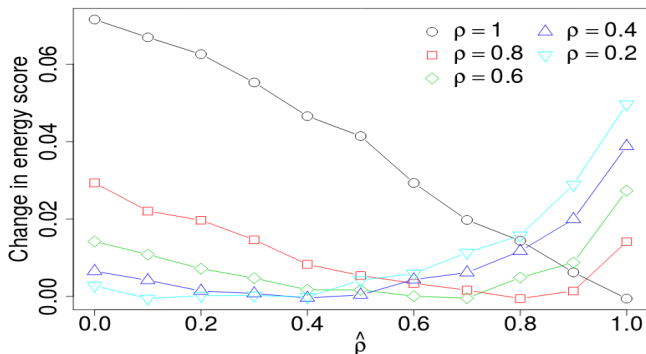
- with $\boldsymbol{\mu}_X$ and Σ_X as mean and covariance matrix of \mathbf{X} .
- optimal if \mathbf{Y} is normally distributed
- not strictly proper

► Summary:

- only energy score and log-score strictly proper
- log-score not useful for practice as density forecast is required

Energy Score

- ▶ Energy score seems to be preferable, still it is hardly applied
- ▶ Pinson, Tatsu(2013) state that **energy score is not sensitive in changes in dependency structure**, based on simulation results from a bivariate normal distribution



(source: Pinson, Tatsu (2013))

- ▶ Conclusion derived by looking at relative change in scores with respect to the true distribution

Marginal-Copula Scores

Idea:

- ▶ instead of full distribution F_Y evaluate
 - marginal distributions F_{Y_h}
 - copula C_Y of Y
 - apply copula theory (Sklar's theorem)
 - hope: control somehow marginal and dependency measures
- ▶ **marginal score:** MS

$$MS(\mathbf{a}) = \mathbf{a}'\mathbf{MS} = \sum_{h=1}^H a_h MS_h.$$

- where MS_h is a univariate scoring rule for Y_h (e.g. CRPS)
 - $\mathbf{a} = (a_1, \dots, a_H)'$ a weight vector (usually $a_h = 1$)
- ▶ **copula score:** CS
 - CS is a multivariate score for the copula C_X of the copula X

Marginal-Copula Scores

Problem:

- ▶ Combine marginal score MS and copula score CS to one score
- ▶ looking for $g : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ defines $CES = g(\text{MS}, \text{CS})$ such that strictly proper scoring rules can be achieved
- ▶ g must be strictly isotonic:

$$g(x_1, y_1) - g(x_1, y_2) - g(x_2, y_1) + g(x_2, y_2) > 0$$

for $x_1, x_2 \in \text{supp}(\text{MS})$ and $y_1, y_2 \in \text{supp}(\text{CS})$ with $x_1 < x_2, y_1 < y_2$

- ▶ possible options for g :
 - a) $g(x, y) = w_1x + w_2y$ for weights $w_i > 0$ works
 - b) $g(x, y) = wxy$ for weights $w > 0$ works on e.g. $(0, \infty) \times (0, \infty)$
- ▶ option a) is not intuitive (due to scaling/ units)

Consider forecast for \mathbf{Y} and $c\mathbf{Y}$ with $c > 0$. For most marginal scores it follows that $\text{MS}(F_{\mathbf{Y}}, \mathbf{Y}; \mathbf{a}) \neq \text{MS}(F_{c\mathbf{Y}}, c\mathbf{Y}; \mathbf{a})$. For the popular CRPS we even have $\text{MS}(F_{\mathbf{Y}}, \mathbf{Y}; \mathbf{a}) = \frac{1}{c} \text{MS}(F_{c\mathbf{Y}}, c\mathbf{Y}; \mathbf{a})$, but for the copula it holds $\text{CS}(\mathbf{C}_{\mathbf{Y}}, \mathbf{U}_{\mathbf{Y}}) = \text{CS}(\mathbf{C}_{c\mathbf{Y}}, \mathbf{U}_{c\mathbf{Y}})$ with $\mathbf{U}_{\mathbf{Y}} \sim \mathbf{C}_{\mathbf{Y}}$ and $\mathbf{U}_{c\mathbf{Y}} \sim \mathbf{C}_{c\mathbf{Y}}$.

Theorem

If the MS_h is a strictly proper score for Y_h and CS is a strictly proper score for the copula \mathbf{C}_Y of $\mathbf{Y} = (Y_1, \dots, Y_H)'$ then the marginal-copula score

$$\begin{aligned} MCS(\mathbf{F}_X, \mathbf{y}; \mathbf{a}) &= MS((F_{X_1}, \dots, F_{X_H})', \mathbf{y}; \mathbf{a}) CS(\mathbf{C}_X, \mathbf{u}_y) \\ &= CS(\mathbf{C}_X, \mathbf{u}_y) \sum_{h=1}^H a_h MS_h(F_{X_h}, y_h) \end{aligned}$$

with \mathbf{F}_X as cumulative distribution function, with continuous marginals F_{X_1}, \dots, F_{X_H} and copula \mathbf{C}_X , of $\mathbf{X} = (X_1, \dots, X_H)'$ which forecasts \mathbf{Y} , observation vector \mathbf{y} and copula observations $\mathbf{u}_Y = (u_{Y,1}, \dots, u_{Y,H})' = (F_{Y_1}(y_1), \dots, F_{Y_H}(y_H))'$ and $\mathbf{a} = (a_1, \dots, a_H)'$ with $a_h > 0$ is a strictly proper scoring rule.

Marginal-Copula Scores

- ▶ Possible choices for marginal scores F_{Y_i} of $\mathbf{Y} = (Y_1, \dots, Y_H)'$
 - **CRPS** (univariate energy score)
 - Log score
 - Dawid-Sebastiani score
 - pinball score / quantile loss on a dense grid on $(0, 1)$

many properties known

- ▶ Possible choice of the copula score for copula $\mathbf{C}_\mathbf{Y}$ of \mathbf{Y}
 - **Energy score**
 - **Variogram score**
 - Log score
 - **Dawid-Sebastiani score**
- ▶ (originally) proposed score:
 - MS: CRPS
 - CS: Energy score
- ▶ Notation:
 - $\mathbf{U}_\mathbf{Y} = (U_{\mathbf{Y},1}, \dots, U_{\mathbf{Y},H})' = (F_{Y_1}(Y_1), \dots, F_{Y_H}(Y_H))' = \mathbf{F}_\mathbf{Y}(\mathbf{Y})$
 - $\mathbf{u}_\mathbf{Y} = (u_{\mathbf{Y},1}, \dots, u_{\mathbf{Y},H})' = (F_{Y_1}(y_1), \dots, F_{Y_H}(y_H))' = \mathbf{F}_\mathbf{Y}(\mathbf{y})$

Copula Energy Scores (CES)

CES: energy score of the copula minus its lower bound scaled by $H^{-\frac{1}{2}}$

$$\begin{aligned}\text{CES}(\mathbf{C}_X, \mathbf{u}_Y) &= \frac{1}{\sqrt{H}} (\text{ES}(\mathbf{C}_X, \mathbf{u}_Y) - \text{lb}_{\text{CES}}) \\ &= \frac{1}{\sqrt{H}} \left(\mathbb{E}(\|\mathbf{U}_X - \mathbf{u}_Y\|_2) - \frac{1}{2} \mathbb{E}(\|\mathbf{U}_X - \tilde{\mathbf{U}}_X\|_2) - \text{lb}_{\text{CES}} \right)\end{aligned}$$

where $\text{lb}_{\text{CES}} = \frac{1}{4} - \frac{1}{2} \frac{1}{\sqrt{6}}$ due to

Lemma

$$\frac{\sqrt{H}}{4} \leq \mathbb{E} \|\mathbf{U}_X - \mathbf{u}_Y\|_2 \leq \frac{\sqrt{H}}{3} \text{ and } \frac{\sqrt{H}}{3} \leq \mathbb{E} (\|\mathbf{U}_X - \tilde{\mathbf{U}}_X\|_2) \leq \frac{\sqrt{H}}{\sqrt{6}}$$

► not a strict bound (ongoing research)

Copula Variogram Score (CVS)

$$\text{CVS}_p(\mathbf{C}_X, \mathbf{u}_Y; \mathbf{W}) = \frac{1}{\mathbf{1}'\mathbf{W}\mathbf{1}} \text{CS}_p(\mathbf{C}_X, \mathbf{u}_Y; \mathbf{W}) \quad (2)$$

$$= \sum_{i=1}^H \sum_{j=1}^H w_{i,j} (|u_{Y,i} - u_{Y,j}|^p - \mathbb{E}|U_{X,i} - U_{X,j}|^p)^2$$

- ▶ with $\mathbf{U}_X = (U_{X,1}, \dots, U_{X,H})' \sim \mathbf{C}_X$, $p > 0$ and weight matrix $\mathbf{W} = (w_{i,j})_{i,j}$.
- ▶ upper bound:

$$\text{VS}_p(\mathbf{C}_X, \mathbf{u}_Y; \mathbf{W}) = \sum_{i=1}^H \sum_{j=1}^H w_{i,j} (|u_{Y,i} - u_{Y,j}|^p - \mathbb{E}|U_{X,i} - U_{X,j}|^p)^2 \quad (3)$$

$$\leq \sum_{i=1}^H \sum_{j=1}^H w_{i,j} (1^p - 0)^2 = \mathbf{1}'\mathbf{W}\mathbf{1}, \quad (4)$$

which justifies the scaling constant.

- ▶ lower bound is zero, as for is holds $\text{VS}_p(\mathbf{M}_H, \mathbf{U}_Y; \mathbf{W}) = 0$ if $\mathbf{U}_Y \sim \mathbf{M}_H$.

Copula Dawid-Sebastiani Score (CDSS)

$$\begin{aligned}\text{CDSS}(\mathbf{C}_X, \mathbf{u}_y) &= \text{DSS}(\mathbf{C}_X, \mathbf{u}_y) \\ &= \log(\det(\boldsymbol{\Sigma}_{U_X})) + (\mathbf{u}_y - \boldsymbol{\mu}_{U_X})' \boldsymbol{\Sigma}_{U_X}^{-1} (\mathbf{u}_y - \boldsymbol{\mu}_{U_X}) \quad (5)\end{aligned}$$

- ▶ $\boldsymbol{\mu}_{U_X}$ and $\boldsymbol{\Sigma}_{U_X}$ as mean and covariance matrix of $U_X \sim \mathbf{C}_X$.
- ▶ as $\boldsymbol{\mu}_{U_X} = \frac{1}{2}\mathbf{1}$ and $\boldsymbol{\Sigma}_{U_X} = \mathbf{S}\mathbf{R}_{U_X}\mathbf{S}$ where $\mathbf{S} = \frac{1}{\sqrt{12}}\mathbf{I}$ and correlation matrix \mathbf{R}_{U_X} it holds

$$\text{CDSS}(\mathbf{C}_X, \mathbf{u}_y) = -H \log(12 \det(\mathbf{R}_{U_X})) + \frac{1}{12} \left(\mathbf{u}_y - \frac{1}{2}\mathbf{1} \right)' \mathbf{R}_{U_X}^{-1} \left(\mathbf{u}_y - \frac{1}{2}\mathbf{1} \right)$$

- ▶ \Rightarrow only measures dependency in correlation of the copula
- ▶ CDSS is unbounded:
for $\mathbf{R}_{U_X}(\delta) = (1 - \delta)\mathbf{I} + \delta\mathbf{1}\mathbf{1}'$ then it holds $\lim_{\delta \rightarrow 1} \det(\mathbf{R}_{U_X}(\delta)) = 0$.
 \Rightarrow not applicable for multiplicative marginal-copula score

Reporting multivariate forecasts

- ▶ for sophisticated problems forecast distribution F_X (or density f_X , or characteristic function φ_X) is not explicitly available.
- ▶ reporting forecast as a large ensemble $X^{(1)}, \dots, X^{(M)}$ for forecasting Y
- ▶ repeat N (similar) forecasting experiments in a rolling window forecasting study: forecasts X_1, \dots, X_N for Y_1, \dots, Y_N
- ▶ realised ensemble forecasts $\mathcal{X}_i = (x_i^{(1)}, \dots, x_i^{(M)})'$ of the forecasting distribution X_i for Y_i

Illustration rolling window forecasting study

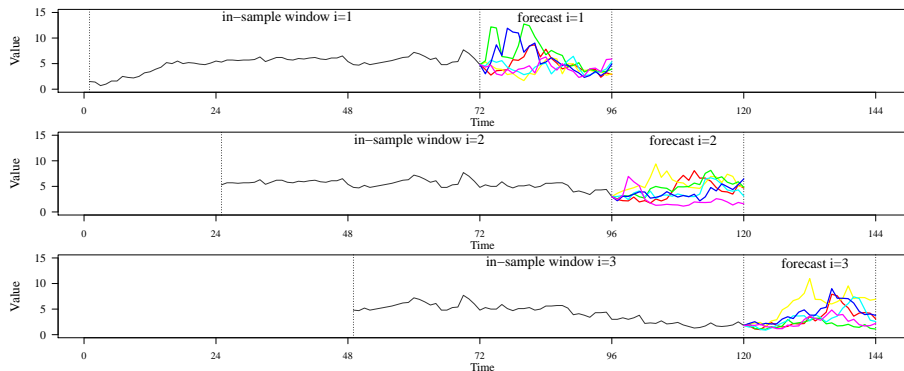


Figure: Illustration of a rolling window forecasting study with non-overlapping windows ($s_i = H(i - 1)$) for $i = 1, \dots, 3$ windows and $M = 6$ forecast samples $\mathbf{x}_{T,i}^{(1)}, \dots, \mathbf{x}_{T,i}^{(M)}$ for each window i .

Estimating the scores

- ▶ for standard (multivariate) scores estimation straight forward, e.g.

$$\text{ES}_{i,\beta}(\mathbf{F}_{\mathbf{X}_i}, \mathbf{y}_i) = \mathbb{E} \left(\|\mathbf{X}_i - \mathbf{y}_i\|_2^\beta \right) - \frac{1}{2} \mathbb{E} \left(\|\mathbf{X}_i - \widetilde{\mathbf{X}}_i\|_2^\beta \right) \quad (6)$$

$$= \text{ED}_{\beta,i}(\mathbf{X}_i, \mathbf{y}_i) - \frac{1}{2} \text{El}_{\beta,i}(\mathbf{X}_i, \mathbf{y}_i) \quad (7)$$

- ▶ estimated by

$$\widehat{\text{ED}}_{i,\beta} = \frac{1}{M} \sum_{j=1}^M \left\| \mathbf{X}_i^{(j)} - \mathbf{y}_i \right\|_2^\beta.$$

and

$$\widehat{\text{El}}_{i,\beta}^{K\text{band}} = \frac{1}{M} \sum_{j=1}^M \sum_{k=j}^K \left\| \mathbf{X}_i^{(j)} - \mathbf{X}_i^{(j+k)} \right\|_2^\beta$$

- $K = M$ computationally expensive - but optimal
- $K = 1$ fast

Estimating the copula scores

- ▶ $\mathbf{U}_{Y_i} = (U_{Y_{i,1}}, \dots, U_{Y_{i,H}})' = (F_{Y_{i,1}}(Y_{i,1}), \dots, F_{Y_{i,H}}(Y_{i,H}))'$ with copula observations
 $\mathbf{u}_{Y_i} = (u_{y_{i,1}}, \dots, u_{y_{i,H}})' = (F_{Y_{i,1}}(y_{i,1}), \dots, F_{Y_{i,H}}(y_{i,H}))'$ depend on true marginals $F_{Y_{i,h}}$.

- ▶ estimate $F_{Y_{i,h}}$ empirical distribution function (ecdf), e.g.

$$\hat{F}_{Y_{i,h}}(z) = \hat{F}_{Y_{i,h}}(z; \mathcal{X}_i) = \frac{1}{M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\},$$

(or the mid-point rule $\hat{F}_{Y_{i,h}}^{\text{mid}}(z) = \frac{1}{2M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\} + \mathbb{1}\{\mathbf{x}_i^{(j)} < z\}$)

$$\hat{\mathbf{u}}_{Y_i} = \hat{\mathbf{u}}_{Y_i}(\mathcal{X}_i) = (\hat{F}_{Y_{i,1}}(y_{i,1}; \mathcal{X}_i), \dots, \hat{F}_{Y_{i,H}}(y_{i,H}; \mathcal{X}_i))'.$$

Estimating the copula scores

- ▶ $\mathbf{U}_{Y_i} = (U_{Y_{i,1}}, \dots, U_{Y_{i,H}})' = (F_{Y_{i,1}}(Y_{i,1}), \dots, F_{Y_{i,H}}(Y_{i,H}))'$ with copula observations
 $\mathbf{u}_{Y_i} = (u_{y_{i,1}}, \dots, u_{y_{i,H}})' = (F_{Y_{i,1}}(y_{i,1}), \dots, F_{Y_{i,H}}(y_{i,H}))'$ depend on true marginals $F_{Y_{i,h}}$.
- ▶ estimate $F_{Y_{i,h}}$ empirical distribution function (ecdf), e.g.

$$\hat{F}_{Y_{i,h}}(z) = \hat{F}_{Y_{i,h}}(z; \mathcal{X}_i) = \frac{1}{M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\},$$
 (or the mid-point rule $\hat{F}_{Y_{i,h}}^{\text{mid}}(z) = \frac{1}{2M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\} + \mathbb{1}\{\mathbf{x}_i^{(j)} < z\}$)

$$\hat{\mathbf{u}}_{Y_i} = \hat{\mathbf{u}}_{Y_i}(\mathcal{X}_i) = (\hat{F}_{Y_{i,1}}(y_{i,1}; \mathcal{X}_i), \dots, \hat{F}_{Y_{i,H}}(y_{i,H}; \mathcal{X}_i))'.$$
- ▶ Problem:
 misspecified marginal can lead to estimated copula lower scores than the true model

Estimating the copula scores

- ▶ $U_{Y_i} = (U_{Y_{i,1}}, \dots, U_{Y_{i,H}})' = (F_{Y_{i,1}}(Y_{i,1}), \dots, F_{Y_{i,H}}(Y_{i,H}))'$ with copula observations
 $u_{Y_i} = (u_{y_{i,1}}, \dots, u_{y_{i,H}})' = (F_{Y_{i,1}}(y_{i,1}), \dots, F_{Y_{i,H}}(y_{i,H}))'$ depend on true marginals $F_{Y_{i,h}}$.

- ▶ estimate $F_{Y_{i,h}}$ empirical distribution function (ecdf), e.g.

$$\hat{F}_{Y_{i,h}}(z) = \hat{F}_{Y_{i,h}}(z; \mathcal{X}_i) = \frac{1}{M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\},$$

(or the mid-point rule $\hat{F}_{Y_{i,h}}^{\text{mid}}(z) = \frac{1}{2M} \sum_{j=1}^M \mathbb{1}\{\mathbf{x}_i^{(j)} \leq z\} + \mathbb{1}\{\mathbf{x}_i^{(j)} < z\}$)

$$\hat{u}_{Y_i} = \hat{u}_{Y_i}(\mathcal{X}_i) = (\hat{F}_{Y_{i,1}}(y_{i,1}; \mathcal{X}_i), \dots, \hat{F}_{Y_{i,H}}(y_{i,H}; \mathcal{X}_i))'.$$

- ▶ Problem:
misspecified marginal can lead to estimated copula lower scores than the true model
- ▶ Solution:
Force the marginals to be uniform, while preserving the dependency structure

Estimating the copula scores

- ▶ $R_{i,h}$ the rank of $\hat{u}_{Y_{i,h}}$ within $\hat{u}_{Y_{1,h}}, \dots, \hat{u}_{Y_{N,h}}$.
- ▶ define the *adjusted estimated copula observations* by

$$\hat{u}_{Y_{i,h}}^* = \frac{2R_{i,h} - 1}{2N}.$$

taking values on $\frac{1}{2N}, \dots, \frac{2N-1}{2N}$

- ▶ resulting ecdf has minimal Komogorov-Smirnov (KS) distance to the uniform distribution. (no other justification)

Estimating the copula scores

- ▶ $R_{i,h}$ the rank of $\hat{u}_{Y_{i,h}}$ within $\hat{u}_{Y_{1,h}}, \dots, \hat{u}_{Y_{N,h}}$.
- ▶ define the *adjusted estimated copula observations* by

$$\hat{u}_{Y_{i,h}}^* = \frac{2R_{i,h} - 1}{2N}.$$

taking values on $\frac{1}{2N}, \dots, \frac{2N-1}{2N}$

- ▶ resulting ecdf has minimal Komogorov-Smirnov (KS) distance to the uniform distribution. (no other justification)
- ▶ estimate C_{X_i} by the empirical copula, we suggest

$$\hat{C}_{X_i}(u_1, \dots, u_H) = \frac{1}{M} \sum_{j=1}^M \mathbb{1}\{\tilde{R}_{i,j,1}/M \leq u_1, \dots, \tilde{R}_{i,j,H}/M \leq u_H\}$$

with the ranks

$$\tilde{R}_{i,j,h} = \frac{1}{2} \sum_{k=1}^M \mathbb{1}\{x_{i,h}^{(k)} \leq x_{i,h}^{(j)}\} + \mathbb{1}\{x_{i,h}^{(k)} < x_{i,h}^{(j)}\}$$

which break ties by the mid-point rule.

Estimating the copula scores

- ▶ Further problem: We know

$$\mathbb{E}[MS \cdot CS] = \mathbb{E}[MS]\mathbb{E}[CS] + \mathbb{Cov}[MS, CS]$$

- ▶ Thus we estimate

$$\widehat{MS(\mathbf{a})-CS}_i = \widehat{MS}_i(\mathbf{a})\widehat{CS}_i - \widehat{\sigma}_{MS,CS}$$

with $\widehat{\sigma}_{MS,CS}$ as estimator for $\mathbb{Cov}[MS, CS]$.

Application in simulation studies

9 scores:

- i) Energy score (ES)
- ii) Variogram score (VS)
- iii) Dawid-Sebastiani score (DSS)
- iv) CRPS-copula energy score
- v) CRPS-copula variogram score
- vi) CRPS
- vii) Copula energy score
- viii) Copula variogram score
- ix) Copula DSS score

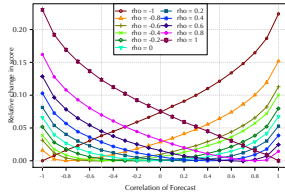
→ evaluation for score SC using two criteria

I) relative change in score with respect to best: $\text{RelCh}(\text{SC}) = \frac{\overline{\text{SC}} - \overline{\text{SC}}^*}{\overline{\text{SC}}^*}$

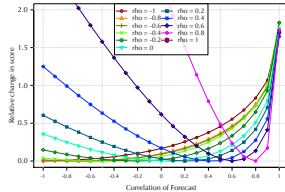
II) DM-test statistics with respect to the best

first study of Pinson, Tatsu (2013) [change in correlation of bivariate normal] using

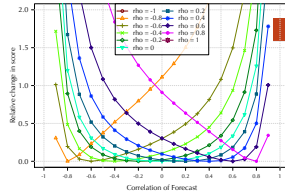
► $N = 2^9 = 512$ (window length), $M = 2^{14} = 16384$ (ensemble size)



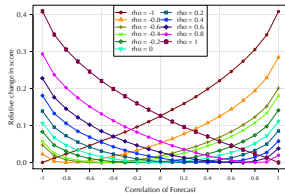
ES



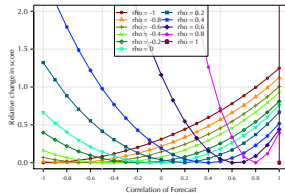
VS



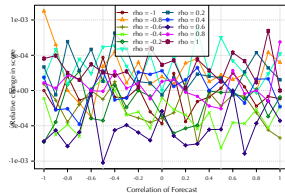
DSS



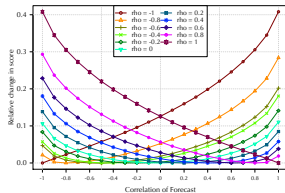
CRPS-CES



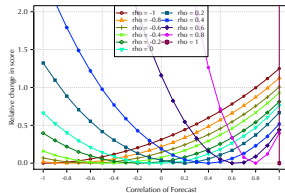
CRPS-CVS



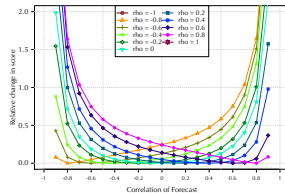
CRPS



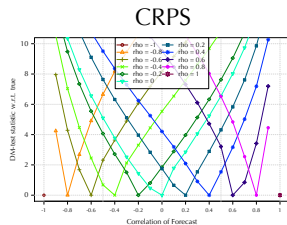
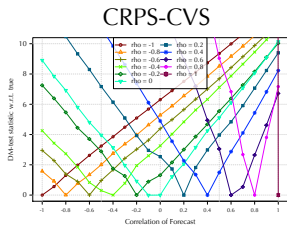
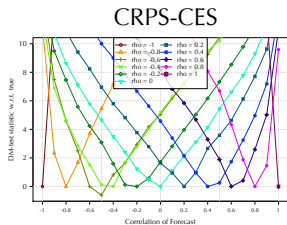
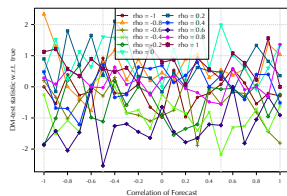
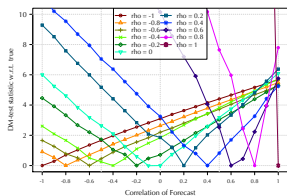
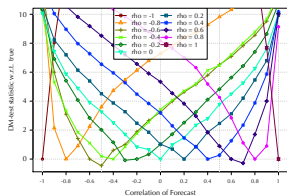
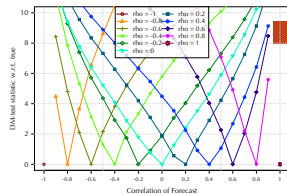
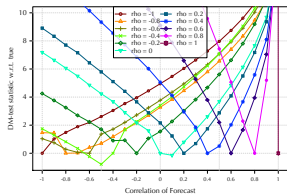
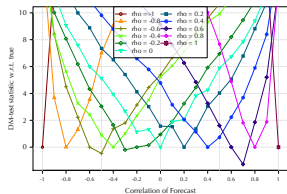
CES



CVS



CDSS



2nd Experiment: on bivariate normal distribution

- With $\boldsymbol{\mu} = (0, 0)'$ and $\boldsymbol{\Sigma}(\rho) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$
- i) (true setting): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\rho = \sqrt{2}/2$
 - ii) (symmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + a_1 \mathbf{1}, \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iii) (asymmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + (a_2, -a_2)', \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iv) (smaller variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_3 \boldsymbol{\Sigma}(\rho))$ with $a_3 < 1$ and $\rho = \sqrt{2}/2$
 - v) (larger variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_4 \boldsymbol{\Sigma}(\rho))$ with $a_4 > 1$ and $\rho = \sqrt{2}/2$
 - vi) (smaller correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_5))$ with $a_5 < \rho$
 - vii) (larger correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_6))$ with $a_6 > \rho$

2nd Experiment: on bivariate normal distribution

- ▶ With $\boldsymbol{\mu} = (0, 0)'$ and $\boldsymbol{\Sigma}(\rho) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$
 - i) (true setting): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\rho = \sqrt{2}/2$
 - ii) (symmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + a_1 \mathbf{1}, \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iii) (asymmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + (a_2, -a_2)', \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iv) (smaller variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_3 \boldsymbol{\Sigma}(\rho))$ with $a_3 < 1$ and $\rho = \sqrt{2}/2$
 - v) (larger variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_4 \boldsymbol{\Sigma}(\rho))$ with $a_4 > 1$ and $\rho = \sqrt{2}/2$
 - vi) (smaller correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_5))$ with $a_5 < \rho$
 - vii) (larger correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_6))$ with $a_6 > \rho$
- ▶ change in biased model so that change theoretical likelihood is the same for all settings ii) - vii)

2nd Experiment: on bivariate normal distribution

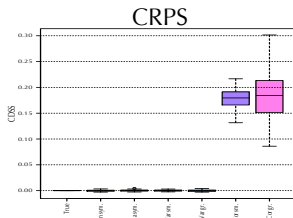
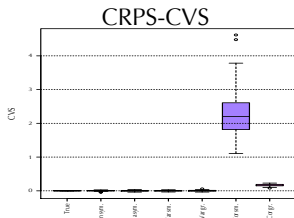
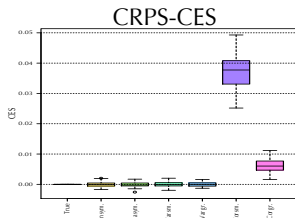
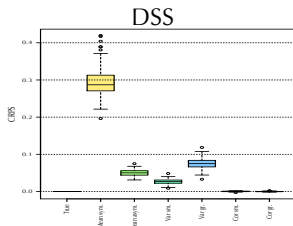
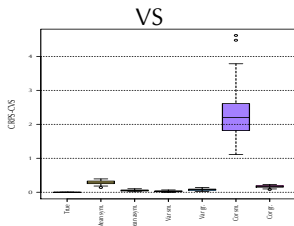
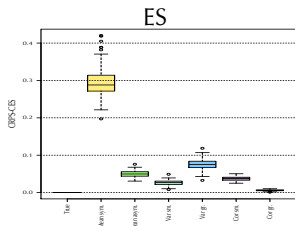
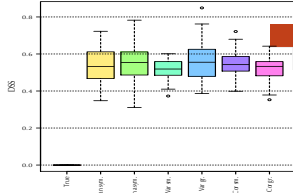
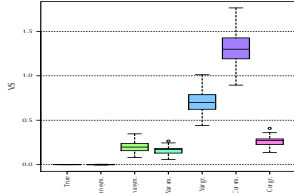
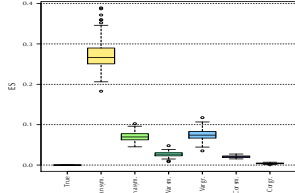
- ▶ With $\boldsymbol{\mu} = (0, 0)'$ and $\boldsymbol{\Sigma}(\rho) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$
 - i) (true setting): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\rho = \sqrt{2}/2$
 - ii) (symmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + a_1 \mathbf{1}, \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iii) (asymmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + (a_2, -a_2)', \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iv) (smaller variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_3 \boldsymbol{\Sigma}(\rho))$ with $a_3 < 1$ and $\rho = \sqrt{2}/2$
 - v) (larger variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_4 \boldsymbol{\Sigma}(\rho))$ with $a_4 > 1$ and $\rho = \sqrt{2}/2$
 - vi) (smaller correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_5))$ with $a_5 < \rho$
 - vii) (larger correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_6))$ with $a_6 > \rho$
- ▶ change in biased model so that change theoretical likelihood is the same for all settings ii) - vii)
- ▶ consider $\rho = \sqrt{2}/2 \approx 0.707$, $a_5 = 0$ (1 degree of freedom) \Rightarrow likelihood reduction $\delta = \frac{1}{2} \log(2)$

$$(a_1 = \sqrt{\frac{\delta}{2-\sqrt{2}}}, a_2 = \sqrt{\frac{\delta}{2+\sqrt{2}}}, a_3 \approx 0.48124, a_4 \approx 2.62729 \text{ and } a_6 \approx 0.89032)$$

2nd Experiment: on bivariate normal distribution

- ▶ With $\boldsymbol{\mu} = (0, 0)'$ and $\boldsymbol{\Sigma}(\rho) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$
 - i) (true setting): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\rho = \sqrt{2}/2$
 - ii) (symmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + a_1 \mathbf{1}, \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iii) (asymmetric mean bias): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu} + (a_2, -a_2)', \boldsymbol{\Sigma}(\rho))$ with $\rho = \sqrt{2}/2$
 - iv) (smaller variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_3 \boldsymbol{\Sigma}(\rho))$ with $a_3 < 1$ and $\rho = \sqrt{2}/2$
 - v) (larger variance): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, a_4 \boldsymbol{\Sigma}(\rho))$ with $a_4 > 1$ and $\rho = \sqrt{2}/2$
 - vi) (smaller correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_5))$ with $a_5 < \rho$
 - vii) (larger correlation): $\mathbf{X} \sim \mathcal{N}_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}(a_6))$ with $a_6 > \rho$
- ▶ change in biased model so that change theoretical likelihood is the same for all settings ii) - vii)
- ▶ consider $\rho = \sqrt{2}/2 \approx 0.707$, $a_5 = 0$ (1 degree of freedom) \Rightarrow likelihood reduction $\delta = \frac{1}{2} \log(2)$

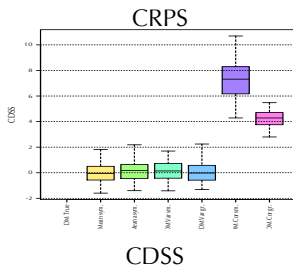
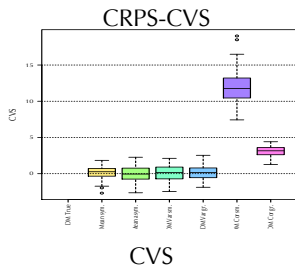
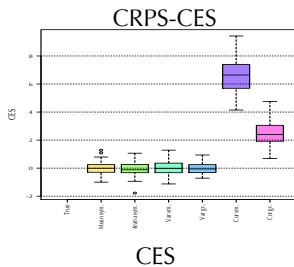
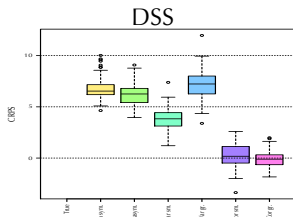
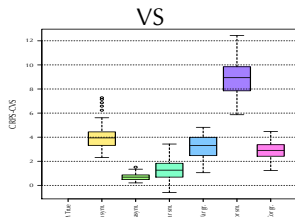
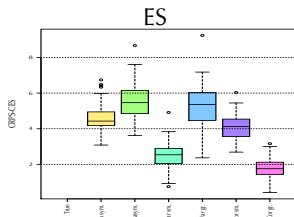
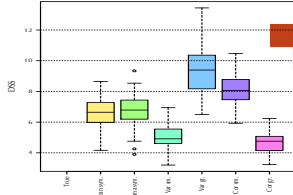
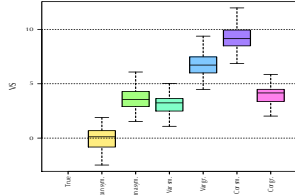
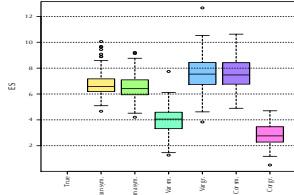
$(a_1 = \sqrt{\frac{\delta}{2-\sqrt{2}}}, a_2 = \sqrt{\frac{\delta}{2+\sqrt{2}}}, a_3 \approx 0.48124, a_4 \approx 2.62729 \text{ and } a_6 \approx 0.89032)$
- ▶ $M = 2^{13} = 8192$ (ensemble sample size), $N = 2^8 = 256$ (rolling window length), $L = 2^6 = 64$ (replications)



CES

CVS

CDSS



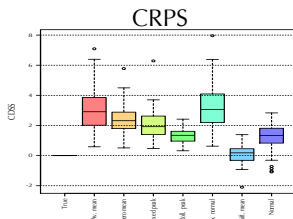
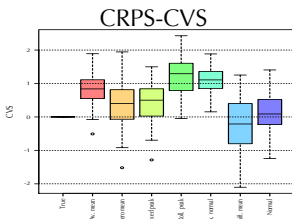
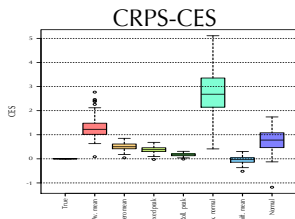
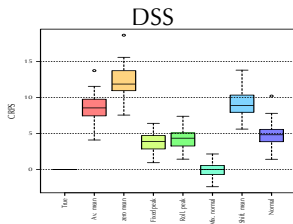
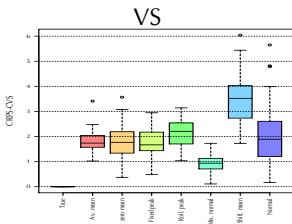
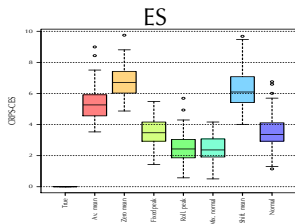
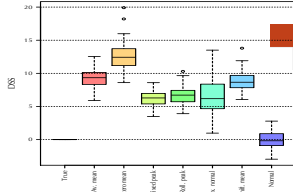
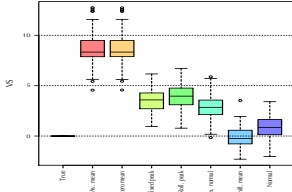
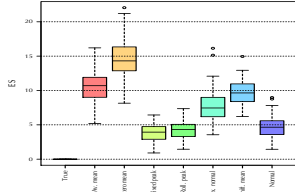
3rd Experiment: Random peak study

- ▶ i) (true) $\mathbf{X}_i = \mathbf{Y}_i + Q\mathbf{Z}_i$ with $\mathbf{Y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mathbf{0}, \mathbf{I})$ and $\mathbf{Z}_i \stackrel{\text{iid}}{\sim} \mathcal{U}(\{\mathbf{e}_1, \dots, \mathbf{e}_H\})$
- ii) (average mean) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mu \mathbf{1}, \mathbf{I})$ with $\mu = \frac{Q}{H}$
- iii) (zero mean) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mathbf{0}, \mathbf{I})$
- iv) (fixed peak) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}, \mathbf{I})$ with $\boldsymbol{\mu} = (Q, 0, \dots, 0)'$
- v) (rolling peak) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}_i, \mathbf{I})$ with $\boldsymbol{\mu}_i = \mathbf{e}_{1+(i-1) \bmod H}$ with as unit vector for the i 'th coordinate.
- vi) (mixture normal with same marginals) $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,H})'$ with

$$X_{i,j} \stackrel{\text{iid}}{\sim} \begin{cases} \mathcal{N}_1(0, 1) & U_j \leq \frac{H-1}{H} \\ \mathcal{N}_1(Q, 1) & U_j > \frac{H-1}{H} \end{cases} \text{ where } U_j \stackrel{\text{iid}}{\sim} \mathcal{U}([0, 1]).$$
- vii) (shifted mean) $\mathbf{X}_i = \mathbf{Y}_i + Q\mathbf{Z}_i$ with $\mathbf{Y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(Q/H \mathbf{1}, \mathbf{I})$ and $\mathbf{Z}_i \stackrel{\text{iid}}{\sim} \mathcal{U}(\{\mathbf{e}_1, \dots, \mathbf{e}_H\})$
- viii) (normal with true mean and covariance) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\mu} = \frac{Q}{H} \mathbf{1}$ and $\boldsymbol{\Sigma} = (H + Q^2)/H \mathbf{I} - Q^2/H^2 \mathbf{1}\mathbf{1}'$

3rd Experiment: Random peak study

- ▶ i) (true) $\mathbf{X}_i = \mathbf{Y}_i + Q\mathbf{Z}_i$ with $\mathbf{Y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mathbf{0}, \mathbf{I})$ and $\mathbf{Z}_i \stackrel{\text{iid}}{\sim} \mathcal{U}(\{\mathbf{e}_1, \dots, \mathbf{e}_H\})$
- ii) (average mean) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mu \mathbf{1}, \mathbf{I})$ with $\mu = \frac{Q}{H}$
- iii) (zero mean) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\mathbf{0}, \mathbf{I})$
- iv) (fixed peak) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}, \mathbf{I})$ with $\boldsymbol{\mu} = (Q, 0, \dots, 0)'$
- v) (rolling peak) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}_i, \mathbf{I})$ with $\boldsymbol{\mu}_i = \mathbf{e}_{1+(i-1) \bmod H}$ with as unit vector for the i 'th coordinate.
- vi) (mixture normal with same marginals) $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,H})'$ with
$$X_{i,j} \stackrel{\text{iid}}{\sim} \begin{cases} \mathcal{N}_1(0, 1) & U_j \leq \frac{H-1}{H} \\ \mathcal{N}_1(Q, 1) & U_j > \frac{H-1}{H} \end{cases} \text{ where } U_j \stackrel{\text{iid}}{\sim} \mathcal{U}([0, 1]).$$
- vii) (shifted mean) $\mathbf{X}_i = \mathbf{Y}_i + Q\mathbf{Z}_i$ with $\mathbf{Y}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(Q/H \mathbf{1}, \mathbf{I})$ and $\mathbf{Z}_i \stackrel{\text{iid}}{\sim} \mathcal{U}(\{\mathbf{e}_1, \dots, \mathbf{e}_H\})$
- viii) (normal with true mean and covariance) $\mathbf{X}_i \stackrel{\text{iid}}{\sim} \mathcal{N}_H(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\mu} = \frac{Q}{H} \mathbf{1}$ and $\boldsymbol{\Sigma} = (H + Q^2)/H \mathbf{I} - Q^2/H^2 \mathbf{1}\mathbf{1}'$
- ▶ First: $H = 3$ -dim. case with a peak size of $Q = 5$
- ▶ $M = 2^{14} = 16384$ (ensemble sample size), $N = 2^5 = 32$ (rolling window length), $L = 2^6 = 64$ (replications), only DM-test



Variants of 3rd Experiment: Random peak study

- a) Effect of ensemble sample size M

$H = 3$ -dim. case with a peak size of $Q = 5$

$$M \in \mathcal{M} = \{2^i | i \in \{4, \dots, 14\}\} = \{2^4, 2^5, \dots, 2^{14}\} = \{16, 32, \dots, 16384\}$$

Variants of 3rd Experiment: Random peak study

- a) Effect of ensemble sample size M

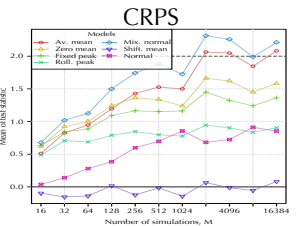
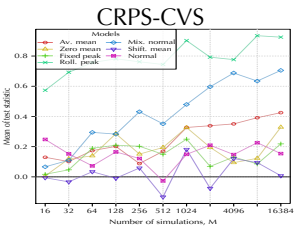
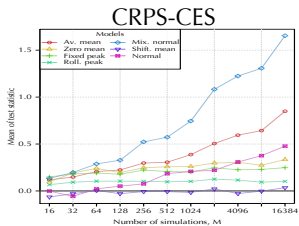
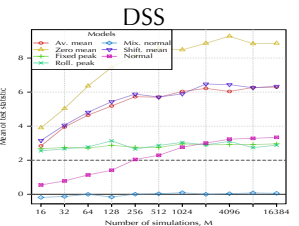
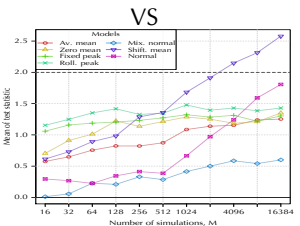
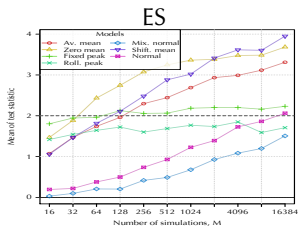
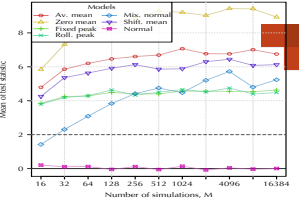
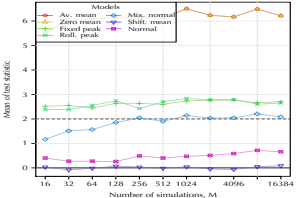
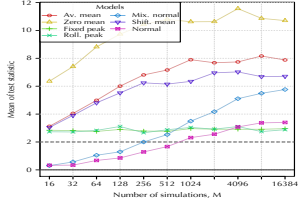
$H = 3$ -dim. case with a peak size of $Q = 5$

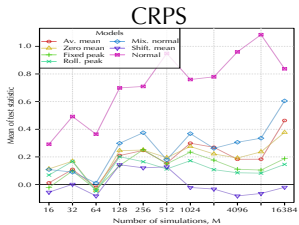
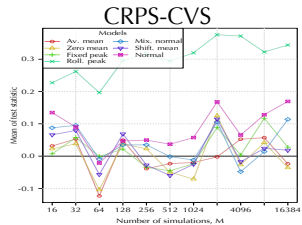
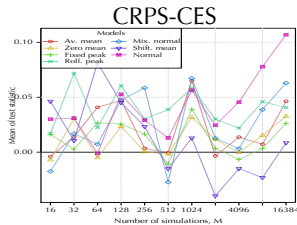
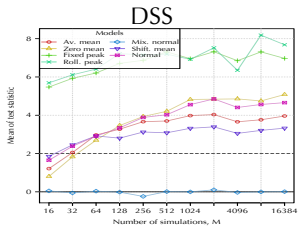
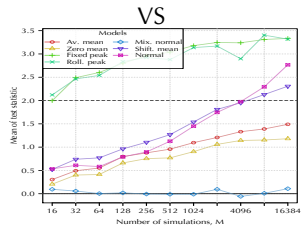
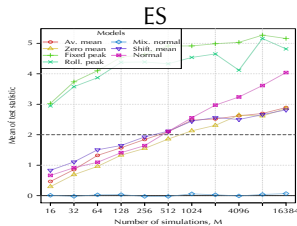
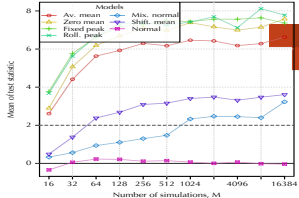
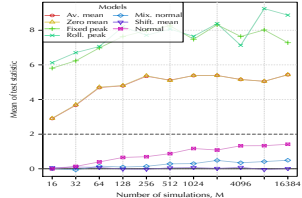
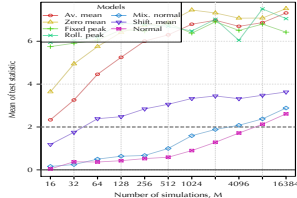
$$M \in \mathcal{M} = \{2^i | i \in \{4, \dots, 14\}\} = \{2^4, 2^5, \dots, 2^{14}\} = \{16, 32, \dots, 16384\}$$

- b) Effect of ensemble sample dimension H

$H = 9$ -dim. case with a peak size of $Q = 5$

$$M \in \mathcal{M} = \{2^i | i \in \{4, \dots, 14\}\} = \{2^4, 2^5, \dots, 2^{14}\} = \{16, 32, \dots, 16384\}$$





CES

CVS

CDSS

Conclusions from the simulation studies

- ▶ Energy score is the only considered score which seems to be suitable
- ▶ In some cases other scores are slightly better in identifying special features

But:

- ? Why is the energy score so powerful?

Conclusions from the simulation studies

- ▶ Energy score is the only considered score which seems to be suitable
- ▶ In some cases other scores are slightly better in identifying special features

But:

- ? Why is the energy score so powerful?
- ▶ Reason: Structure of the energy distance

$$d_E(\mathbf{X}, \mathbf{Y}) = \mathbb{E} \|\mathbf{X} - \mathbf{Y}\|_2^\beta - \frac{1}{2} \mathbb{E} \|\mathbf{X} - \widetilde{\mathbf{X}}\|_2^\beta - \frac{1}{2} \mathbb{E} \|\mathbf{Y} - \widetilde{\mathbf{Y}}\|_2^\beta$$

- ▶ d_E yields energy score for observed $\mathbf{Y}(\omega) = \mathbf{y}$

Properties of the energy distance

- ▶ d_E is zero if and only if $\mathbf{X} \stackrel{d}{=} \mathbf{Y}$
- ▶ special weighted L^2 -distance between characteristic functions:

$$d_E(\mathbf{X}, \mathbf{Y}) = \frac{\pi^{\frac{H+1}{2}}}{\Gamma(\frac{H+1}{2})} \int_{\mathbb{R}^H} \frac{|\varphi_{\mathbf{X}}(z) - \varphi_{\mathbf{Y}}(z)|^2}{\|z\|_2^{H+\beta}} dz \quad (8)$$

for characteristic functions $\varphi_{\mathbf{X}}(z) = \mathbb{E}(e^{iz'\mathbf{X}})$ and $\varphi_{\mathbf{Y}}(z) = \mathbb{E}(e^{iz'\mathbf{Y}})$.

- ▶ If considering the weighted L^2 -distance between $\varphi_{\mathbf{X}}$ and $\varphi_{\mathbf{Y}}$:

$$C \int_{\mathbb{R}^d} \xi(z) |\varphi_{\mathbf{X}}(z) - \varphi_{\mathbf{Y}}(z)|^2 dz$$

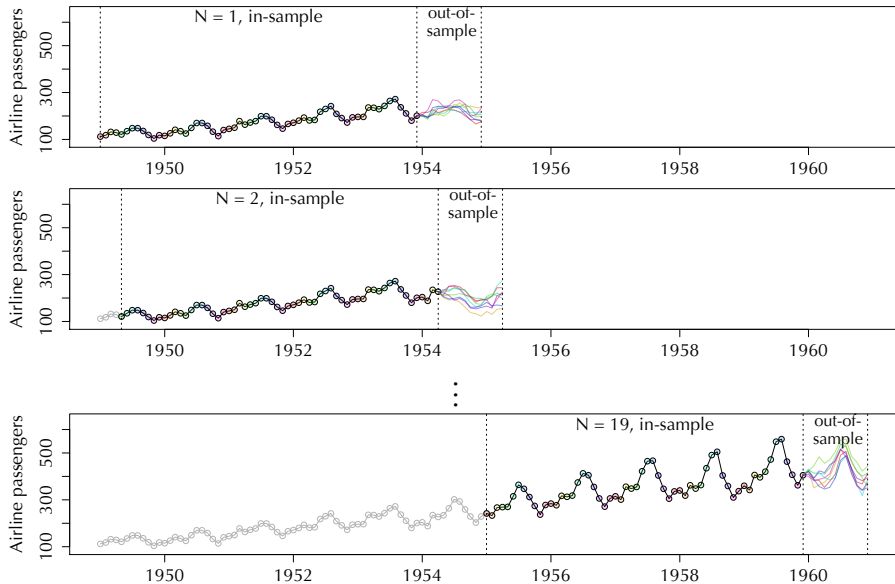
then $\xi(z) = \|z\|_2^{H+\beta}$ is the only choice such that the distance is scale equivariant and rotationally invariant

- ▶ d_E measures distances between \mathbf{X} and \mathbf{Y} in $\mathbb{C} \cong \mathbb{R}^2$ not \mathbb{R}^H
 \Rightarrow should be efficient for $H > 2$.

Properties of the energy distance

- ▶ Allows arbitrary 2-sample test
- ▶ Test for multivariate normality (more powerful than standard tests)
- ▶ Allows construction of ' β -distance-covariance':
$$\text{dstCov}_\beta(\mathbf{X}, \mathbf{Y}) = \frac{1}{C} \int_{\mathbb{R}^H} \int_{\mathbb{R}^H} \frac{|\varphi_{\mathbf{X}, \mathbf{Y}}(z, v) - \varphi_{\mathbf{X}}(z) \varphi_{\mathbf{Y}}(v)|^2}{\|z\|_2^{\beta+H} \|v\|_2^{\beta+H}} d\mathbf{v} d\mathbf{z}$$
- ▶ β -distance-covariance allows tests for (multivariate) independence(!)

Real data example



Real data example: 9 models

- 1a) AR(12): $Y_t = \phi_0 + \sum_{k=1}^{12} \phi_k Y_{t-k} + \varepsilon_t$ with ε_t iid and $\mathbb{E}(\varepsilon_t) = 0$.
- 2a) AR(13): $Y_t = \phi_0 + \sum_{k=1}^{13} \phi_k Y_{t-k} + \varepsilon_t$ with ε_t iid and $\mathbb{E}(\varepsilon_t) = 0$.
- 3a) AR(p): $Y_t = \phi_0 + \sum_{k=1}^p \phi_k Y_{t-k} + \varepsilon_t$ with ε_t iid, $\mathbb{E}(\varepsilon_t) = 0$ and $p \in \{1, \dots, T/2\}$ such that the corresponding Akaike information criterion (AIC) is minimized.
- 1b) AR(12) as in 1a) but with comonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{M}_2)
- 2b) AR(13) as in 2a) but with comonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{M}_2)
- 3b) AR(p) as in 3a) but with comonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{M}_2)
- 1c) AR(12) as in 1a) but with countermonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{W}_2)
- 2c) AR(13) as in 2a) but with countermonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{W}_2)
- 3c) AR(p) as in 3a) but with countermonotone residuals (i.e. $(\hat{\varepsilon}_t, \hat{\varepsilon}_{t+1})$ have the copula \mathbf{W}_2)

Real data example

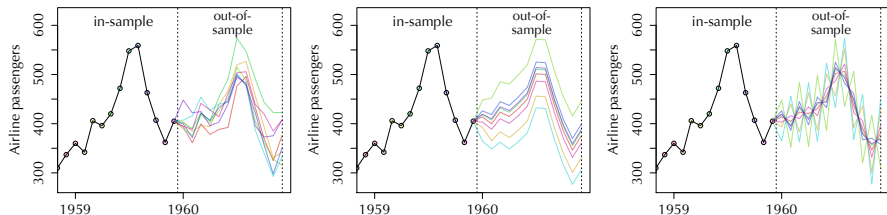


Figure: Illustration of standard (left), comonotone (center) and countermonotone (right) model simulations for the AR(13) with $M = 8$ paths for the last experiment ($N = 19$).

► ensemble size: $M = 2^{16} = 65536$

Real data example: results

Score\Model	AR(12)	AR(13)	AR(p)	AR(12)-M	AR(13)-M	AR(p)-M	AR(12)-W	AR(13)-W	AR(p)-W
ES	120.8	137.5	134.9	197.3	196.2	190.9	203.2	206.6	201.0
VS	158011	120657	111823	205513	133903	122607	175112	123631	114357
DSS	95.92	90.99	91.78	683635	372563	364726	995252	838370	762531
CRPS-CES	3.228	3.887	3.828	8.617	10.240	10.045	8.818	10.537	10.365
CRPS-CVS	0.3346	0.4638	0.4166	0.6357	0.7712	0.7248	0.7989	0.9954	0.9464
CRPS	28.28	33.94	33.36	28.32	33.90	33.36	28.29	33.92	33.35
CES	0.1142	0.1144	0.1146	0.3044	0.3020	0.3010	0.3118	0.3105	0.3106
CVS	0.01184	0.01367	0.01247	0.02247	0.02276	0.02171	0.02826	0.02935	0.02836
CDSS	-26.06	-24.54	-23.78	-	-	-	-	-	-

Table: Score averages \overline{SC} across the $N = 19$ out-of-sample windows for the considered scores and models. -M indicates models with comonotone residuals, -W for countermonotone residuals.

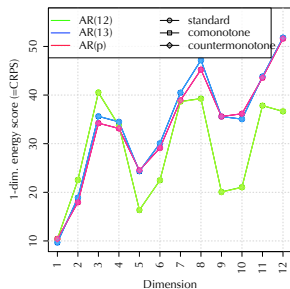
- AR(12) seems to be best concerning, but not uniformly \Rightarrow improvements possible

Real data example: results

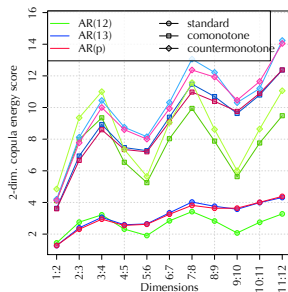
	AR(12)	AR(13)	AR(p)	AR(12)-M	AR(13)-M	AR(p)-M	AR(12)-W	AR(13)-W	AR(p)-W
AR(12)		-4.04 [<0.001]	-2.57 [0.005]	-26.53 [<0.001]	-19.57 [<0.001]	-13.39 [<0.001]	-30.02 [<0.001]	-17.90 [<0.001]	-13.02 [<0.001]
AR(13)	4.04 [>0.999]		0.73 [0.766]	-11.15 [<0.001]	-31.58 [<0.001]	-13.67 [<0.001]	-13.57 [<0.001]	-35.14 [<0.001]	-14.51 [<0.001]
AR(p)	2.57 [0.995]	-0.73 [0.234]		-9.12 [<0.001]	-14.06 [<0.001]	-29.95 [<0.001]	-10.74 [<0.001]	-16.87 [<0.001]	-31.22 [<0.001]
AR(12)-M	26.53 [>0.999]	11.15 [>0.999]	9.12 [>0.999]		0.26 [0.602]	1.08 [0.860]	-7.18 [<0.001]	-1.86 [0.031]	-0.55 [0.292]
AR(13)-M	19.57 [>0.999]	31.58 [>0.999]	14.06 [>0.999]	-0.26 [0.398]		1.39 [0.918]	-1.98 [0.024]	-9.32 [<0.001]	-1.09 [0.138]
AR(p)-M	13.39 [>0.999]	13.67 [>0.999]	29.95 [>0.999]	-1.08 [0.140]	-1.39 [0.082]		-2.28 [0.011]	-4.18 [<0.001]	-9.51 [<0.001]
AR(12)-W	30.02 [>0.999]	13.57 [>0.999]	10.74 [>0.999]	7.18 [>0.999]	1.98 [0.976]	2.28 [0.989]		-0.78 [0.219]	0.36 [0.642]
AR(13)-W	17.90 [>0.999]	35.14 [>0.999]	16.87 [>0.999]	1.86 [0.969]	9.32 [>0.999]	4.18 [>0.999]	0.78 [0.781]		1.39 [0.918]
AR(p)-W	13.02 [>0.999]	14.51 [>0.999]	31.22 [>0.999]	0.55 [0.708]	1.09 [0.862]	9.51 [>0.999]	-0.36 [0.358]	-1.39 [0.082]	

Table: DM-test statistics with corresponding p-value given in squared brackets for the energy score (ES).

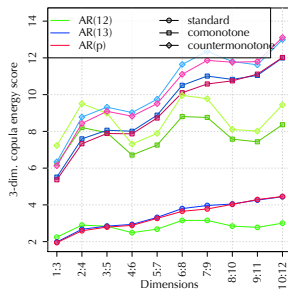
Real data example: results



CRPS



2-dim. CES



3-dim. CES

- CRPS and 2-dim. CES across horizon very useful to detect failures in performance

Summary

- ▶ Energy score is suitable distance for multivariate evaluation (in combination with significance tests)
- ▶ Ensemble sample size should be as large as computationally feasible
- ▶ Additionally consider
 - CRPS for checking individual marginals across the forecasting horizon
 - copula energy score for evaluation 2-way dependencies across the forecasting horizon

Summary

- ▶ Energy score is suitable distance for multivariate evaluation (in combination with significance tests)
- ▶ Ensemble sample size should be as large as computationally feasible
- ▶ Additionally consider
 - CRPS for checking individual marginals across the forecasting horizon
 - copula energy score for evaluation 2-way dependencies across the forecasting horizon

Thank you for your attention.