

Multivariate coefficient of variation for functional data

Mirosław Krzyśko¹ and **Łukasz Smaga²**

¹Interfaculty Institute of Mathematics and Statistics
The President Stanisław Wojciechowski State University of Applied Sciences in Kalisz

²Faculty of Mathematics and Computer Science
Adam Mickiewicz University

XLVIII International Biometrical Colloquium
September 9-13, 2018, Szamotuły

- The coefficient of variation (CV) being the ratio of the standard deviation to the population mean is widely used relative variation measure.
- The CV is a dimensionless quantity, which can be expressed in percent.
- This quantity is usually used to compare the variability of several populations, even when they are characterized by variables expressed in different units as well as have really different means. In particular, the CV is often used to assess the performance or reproducibility of measurement techniques or equipments.

- For multivariate data, computing the CV for each variable is a common practice, although this ignores the correlation between them, and this does not summarize the variability of the multivariate data into a single index.
- The known multivariate extensions of the CV are less considered in the literature. Perhaps it is due to the fact that generalizing the univariate CV to the multivariate setting is not straightforward, and the multivariate CV's do not generally measure the same quantity, when the number of variables is greater than one. Nevertheless, they all reduce to the CV in the univariate case.

- Let $\mathbf{X} = (X_1, \dots, X_p)^\top$ be p -dimensional random vector with mean vector $\mathbf{u} \neq \mathbf{0}_p$ and covariance matrix $\mathbf{\Sigma}$.
- The multivariate coefficients of variation (MCV) by Reyment (1960), Van Valen (1974), Voinov and Nikulin (1996, p. 68) and Albert and Zhang (2010) are:

$$\text{MCV}_R = \sqrt{\frac{(\det \mathbf{\Sigma})^{1/p}}{\mathbf{u}^\top \mathbf{u}}}, \text{MCV}_{VV} = \sqrt{\frac{\text{tr} \mathbf{\Sigma}}{\mathbf{u}^\top \mathbf{u}}}, \text{MCV}_{VN} = \sqrt{\frac{1}{\mathbf{u}^\top \mathbf{\Sigma}^{-1} \mathbf{u}}}, \text{MCV}_{AZ} = \sqrt{\frac{\mathbf{u}^\top \mathbf{\Sigma} \mathbf{u}}{(\mathbf{u}^\top \mathbf{u})^2}},$$

respectively.

- The MCV_R and MCV_{VV} are based on the generalized variance $\det \mathbf{\Sigma}$ and the total variance $\text{tr} \mathbf{\Sigma}$, respectively. In the MCV_{VN} , the Mahalanobis distance $\mathbf{u}^\top \mathbf{\Sigma}^{-1} \mathbf{u}$ appears to be a natural extension of the CV. Finally, the MCV_{AZ} is derived based on a matrix generalizing the square of the CV.

- We have:

$$\text{MCV}_{AZ} = \frac{\sqrt{\text{Var}(\mathbf{u}_*^\top \mathbf{X})}}{\|\mathbf{u}\|}, \quad (1)$$

where $\mathbf{u}_* = \mathbf{u}/\|\mathbf{u}\|$, i.e., MCV_{AZ} is the univariate coefficient of variation for $\mathbf{u}_*^\top \mathbf{X}$.

- Let $\mathbf{X}(t) = (X_1(t), \dots, X_p(t))^\top$, $t \in [a, b]$, $a, b \in \mathbb{R}$ be p -dimensional random process with mean function $\boldsymbol{\mu}(t) = (\mu_1(t), \dots, \mu_p(t))^\top \neq \mathbf{0}_p$. We also assume that $\mathbf{X}(t)$, $t \in [a, b]$ belongs to the Hilbert space $L_2^p[a, b]$ of p -dimensional vectors of square integrable functions on $[a, b]$. Let $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ denote the inner product and the norm in $L_2^p[a, b]$.

Definition 1

The functional multivariate coefficient of variation for $\mathbf{X}(t)$, $t \in [a, b]$ is defined as follows ($\boldsymbol{\mu}_(t) = \boldsymbol{\mu}(t)/\|\boldsymbol{\mu}\|$, $t \in [a, b]$):*

$$\text{FMCV} = \frac{\sqrt{\text{Var}(\langle \boldsymbol{\mu}_*, \mathbf{X} \rangle)}}{\|\boldsymbol{\mu}\|}. \quad (2)$$

Theorem 1

If $X_i(t)$, $t \in [a, b]$, $i = 1, \dots, p$, are square integrable, i.e.,

$$E\|X_i\|^2 = E \int_a^b X_i^2(t) dt < \infty,$$

then

$$\text{Var}(\langle \mu_*, \mathbf{X} \rangle)$$

exists. Furthermore, the FMCV defined in (2) is the CV of the random variable

$$\langle \mu_*, \mathbf{X} \rangle.$$

- Let $\mathbf{X}(t)$ belong to a finite dimensional subspace $\mathcal{L}_2^p[a, b]$ of $L_2^p[a, b]$, where the components of $\mathbf{X}(t)$ can be represented by a finite number of basis functions, i.e.,

$$X_k(t) = \sum_{l=1}^{B_k} \alpha_{kl} \varphi_{kl}(t), \quad (3)$$

where $k = 1, \dots, p$, $t \in [a, b]$, $B_k \in \mathbb{N}$, α_{kl} are random variables with finite variance and $\{\varphi_{kl}\}_{l=1}^{\infty}$, $k = 1, \dots, p$ are bases in the space $L_2^1[a, b]$. The equations (3) can be expressed in the following matrix notation:

$$\mathbf{X}(t) = \mathbf{\Phi}(t)\boldsymbol{\alpha}, \quad (4)$$

where

$$\mathbf{\Phi}(t) = \text{diag} \left(\boldsymbol{\varphi}_1^\top(t), \dots, \boldsymbol{\varphi}_p^\top(t) \right)$$

is the block diagonal matrix of $\boldsymbol{\varphi}_k^\top(t) = (\varphi_{k1}(t), \dots, \varphi_{kB_k}(t))$, $k = 1, \dots, p$ and $\boldsymbol{\alpha} = (\alpha_{11}, \dots, \alpha_{1B_1}, \dots, \alpha_{p1}, \dots, \alpha_{pB_p})^\top$.

$$\text{FMCV} = \frac{\sqrt{\text{Var} \left(\frac{\mathbf{a}^\top \mathbf{J}_\Phi \boldsymbol{\alpha}}{\|\mathbf{J}_\Phi^{1/2} \mathbf{a}\|} \right)}}{\|\mathbf{J}_\Phi^{1/2} \mathbf{a}\|} = \sqrt{\frac{\mathbf{a}^\top \mathbf{J}_\Phi \boldsymbol{\Sigma}_\alpha \mathbf{J}_\Phi \mathbf{a}}{(\mathbf{a}^\top \mathbf{J}_\Phi \mathbf{a})^2}}, \quad (5)$$

where $\mathbf{a} = \mathbf{E}(\boldsymbol{\alpha})$, $\boldsymbol{\Sigma}_\alpha = \text{Cov}(\boldsymbol{\alpha})$ and $\mathbf{J}_\Phi = \text{diag}(\mathbf{J}_{\varphi_1}, \dots, \mathbf{J}_{\varphi_p})$ and $\mathbf{J}_{\varphi_k} = \int_a^b \varphi_k(t) \varphi_k^\top(t) dt$ is the $B_k \times B_k$ cross product matrix corresponding to the basis $\{\varphi_{kl}\}_{l=1}^\infty$, $k = 1, \dots, p$.

Theorem 2

Under the above assumptions and notation, the functional multivariate coefficient of variation for the random process $\mathbf{X}(t)$, $t \in [a, b]$, is the multivariate coefficient of variation of Albert-Zhang type for random vector $\mathbf{J}_\Phi^{1/2} \boldsymbol{\alpha}$, if the matrix $\mathbf{J}_\Phi^{1/2}$ exists.

- Although the FMCV is defined for univariate and multivariate functional data, we note that even when $p = 1$, the FMCV reduces to the MCV_{AZ} (no to the CV), since $B_1 > 1$ usually.

- In practice, we have to estimate the unknown vector α in (4) as well as its parameters \mathbf{a} and Σ_{α} appearing in the FMCV given in (5).
- Let $\mathbf{x}_1(t), \dots, \mathbf{x}_n(t)$, $t \in [a, b]$ be a random sample containing realizations of the process $\mathbf{X}(t)$. These observations are represented similarly as in (4), i.e.,

$$\mathbf{x}_i(t) = \Phi(t)\alpha_i,$$

where $t \in [a, b]$ and $i = 1, \dots, n$.

- Then, the vectors α_i , $i = 1, \dots, n$ can be estimated by the least squares method or the roughness penalty approach.
- The expansion lengths B_k in (3) can be selected deterministically or by using information criteria as the Akaike and Bayesian information criteria.

- Using the estimators of α_i , say $\hat{\alpha}_i$, $i = 1, \dots, n$, we can estimate the mean vector \mathbf{a} and the covariance matrix Σ_{α} .
- The classical estimators are the sample mean and the sample covariance matrix, i.e.,

$$\hat{\mathbf{a}} = \frac{1}{n} \sum_{i=1}^n \hat{\alpha}_i, \quad \hat{\Sigma}_{\alpha} = \frac{1}{n} \sum_{i=1}^n (\hat{\alpha}_i - \hat{\mathbf{a}})(\hat{\alpha}_i - \hat{\mathbf{a}})^{\top}. \quad (6)$$

- However, these estimators may break down, when the data contain outliers. Thus, many authors recommend the use of robust estimators of location and scatter in the presence of outlying observations.
- Similarly to Aerts et al. (2015), we will mainly use the two commonly used ones, i.e., the minimum covariance determinant (MCD) estimator and the S-estimator.

- For a given breakdown point α , the MCD estimator is based on a subset of $\{\hat{\alpha}_1, \dots, \hat{\alpha}_n\}$ of size $h = \lfloor n(1 - \alpha) \rfloor$ minimizing the generalized variance (i.e., the determinant of covariance matrix) among all possible subsets of size h . Then, the MCD estimators of \mathbf{a} and Σ_α are the sample mean and the sample covariance matrix (multiplied by a consistency factor) computed from this subset.
- The location and scatter S-estimators are the vector \mathbf{a}_n and the positive definite symmetric matrix Σ_n which minimizes $\det(\Sigma_n)$ subject to

$$\frac{1}{n} \sum_{i=1}^n \rho \left(\sqrt{(\hat{\alpha}_i - \mathbf{a}_n)^\top \Sigma_n^{-1} (\hat{\alpha}_i - \mathbf{a}_n)} \right) = b_0,$$

where $\rho : \mathbb{R} \rightarrow [0, \infty)$ is a given non-decreasing and symmetric function (e.g., Tukey's biweight) and b_0 a constant needed to ensure consistency of the estimator.

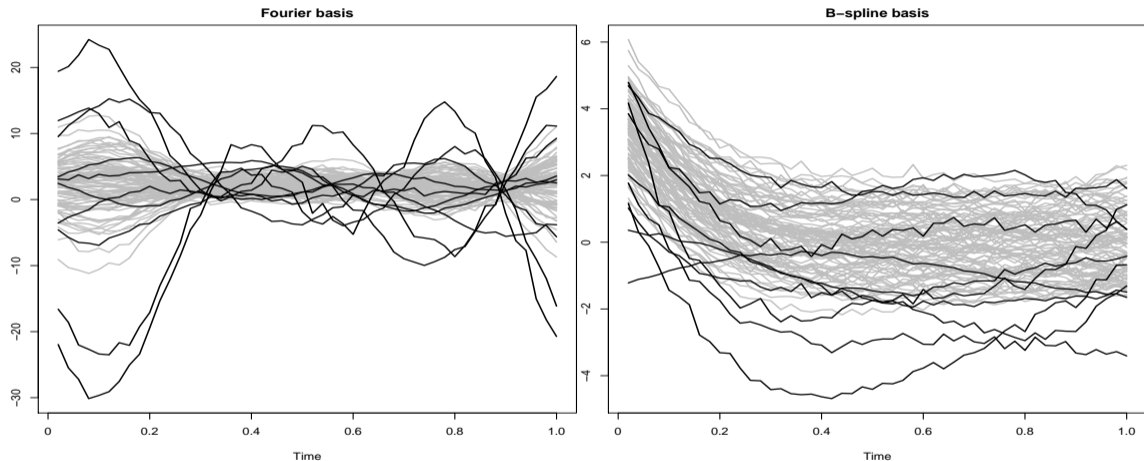
- The (classical and robust) estimators of the FMCV are obtained by substituting the parameters \mathbf{a} and Σ_α in (5) by their estimators $\hat{\mathbf{a}}$ and $\hat{\Sigma}_\alpha$.

- The functional sample $\mathbf{x}_1(t), \dots, \mathbf{x}_n(t)$ of size $n = 100$ contains realizations of $\mathbf{X}(t)$, $t \in [0, 1]$ with $p = 5$. These observations are generated as follows:

$$\mathbf{x}_i(t_j) = \Phi(t_j)\alpha_i + \epsilon_{ij},$$

where $i = 1, \dots, n$, t_j , $j = 1, \dots, 50$ are equally spaced design time points in $[0, 1]$, the matrix $\Phi(t)$ contains basis functions with $B_k = 5$, $k = 1, \dots, p$, α_i are $5p$ -dimensional random vectors, and $\epsilon_{ij} = (\epsilon_{ij1}, \dots, \epsilon_{ijp})^\top$ are the measurement errors such that $\epsilon_{ijk} \sim N(0, 0.025r_{ik})$ and r_{ik} is the range of the k -th row of $(\Phi(t_1)\alpha_i \dots \Phi(t_{50})\alpha_i)$, $k = 1, \dots, p$. We use the Fourier and B-spline bases.

- The vectors α_i , $i = 1, \dots, n$ were generated from multivariate normal or t_5 - distributions with mean \mathbf{a} and covariance matrix Σ_α . Similarly to Aerts et al. (2015), we set $\mathbf{a} = \mathbf{a}_1 := a\mathbf{e}_1$ or $\mathbf{a} = \mathbf{a}_2 := (a/(5p)^{1/2})\mathbf{1}_{5p}$ and $\Sigma_\alpha = (1 - \rho)\mathbf{I}_{5p} + \rho\mathbf{1}_{5p}\mathbf{1}_{5p}^\top$, where a is chosen to get a given value of the FMCV, $\mathbf{e}_1 = (1, 0, \dots, 0)^\top$ and $\rho = 0, 0.5, 0.8$. We set FMCV = 0.1, 0.5, 0.9. Moreover, to obtain uncontaminated and contaminated functional data, $\varepsilon\%$ of the observations are generated with $10\Sigma_\alpha$, where $\varepsilon = 0, 10, 20, 30, 40, 50$.



Exemplary realizations of the first functional variables of simulated data with 10% of outlying observations. The uncontaminated (resp. contaminated) data are depicted in gray (resp. black).

Simulation studies - mean squared error

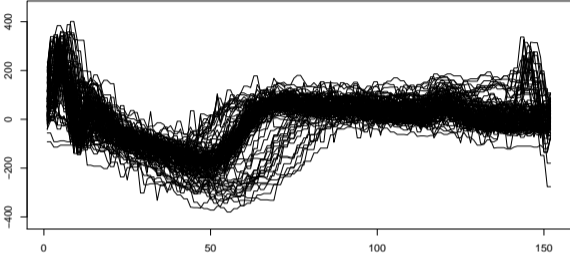
	class.	MCD	S	class.	MCD	S	class.	MCD	S	
	FMCV = 0.1			FMCV = 0.5			FMCV = 0.9			
ε										
F	0	0.0001	0.0002	0.0001	0.0209	0.0314	0.0212	0.1197	0.1654	0.1216
	10	0.0019	0.0002	0.0002	0.1278	0.0347	0.0304	0.6283	0.1763	0.1662
	20	0.0053	0.0002	0.0006	0.2919	0.0368	0.0519	1.1743	0.1909	0.2820
	30	0.0095	0.0003	0.0019	0.4902	0.0412	0.1177	1.7589	0.2138	0.5742
	40	0.0147	0.0011	0.0069	0.7307	0.0936	0.3261	2.3889	0.4412	1.4175
	50	0.0204	0.0065	0.0140	0.9976	0.3683	0.6583	2.9514	1.3361	2.5026
B	0	0.0002	0.0004	0.0002	0.0433	0.0616	0.0441	0.2508	0.3492	0.2533
	10	0.0023	0.0004	0.0003	0.2143	0.0644	0.0582	1.0780	0.3563	0.3268
	20	0.0060	0.0004	0.0007	0.4704	0.0687	0.0882	2.0930	0.3842	0.4986
	30	0.0109	0.0005	0.0022	0.7978	0.0717	0.1672	3.2243	0.3956	0.9228
	40	0.0166	0.0013	0.0074	1.1598	0.1435	0.4648	4.3543	0.7577	2.2992
	50	0.0227	0.0078	0.0153	1.5075	0.6046	0.9543	5.3330	2.3050	4.1944

- MSE's are usually similar for different bases, but greater differences can appear for greater FMCV. Moreover, MSE for the B-spline basis is often greater than for the Fourier basis.

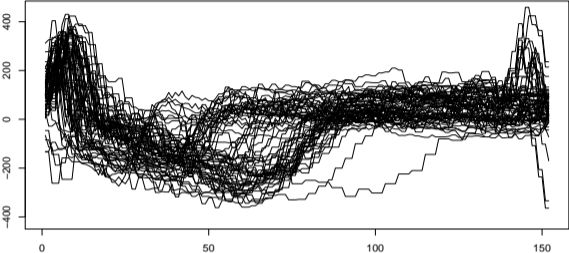
- We consider the ECG data set originated from Olszewski (2001).
- During each of 200 heartbeats, two electrodes were used to measure the ECG. For one heartbeat and one electrode, the ECG was measured in 152 design time points, and the resulting values of ECG form a curve, which can be treated as discrete functional observation.
- We have 200 two-dimensional discrete functional data observed in 152 design time points ($n = 200$, $p = 2$, $m_i = 152$, $i = 1, \dots, n$).
- The heartbeats were assigned to normal or abnormal group.
- Abnormal heartbeats are representative of a cardiac pathology known as supraventricular premature beat.
- The normal and abnormal groups consist of 133 and 67 functional observations, respectively.

Application to ECG data

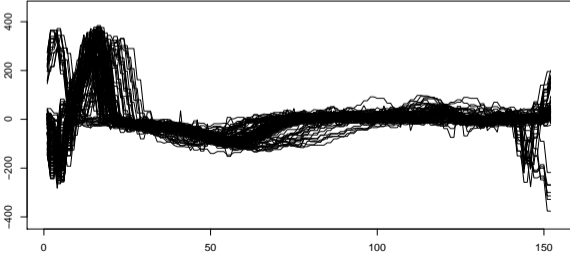
First variable – normal



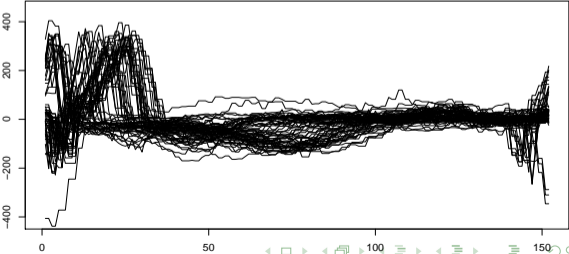
First variable – abnormal



Second variable – normal

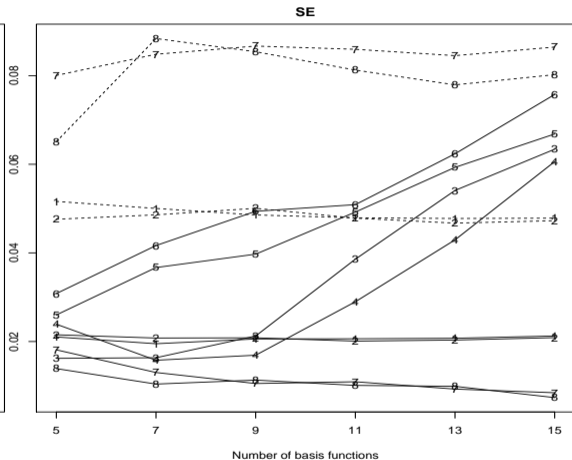
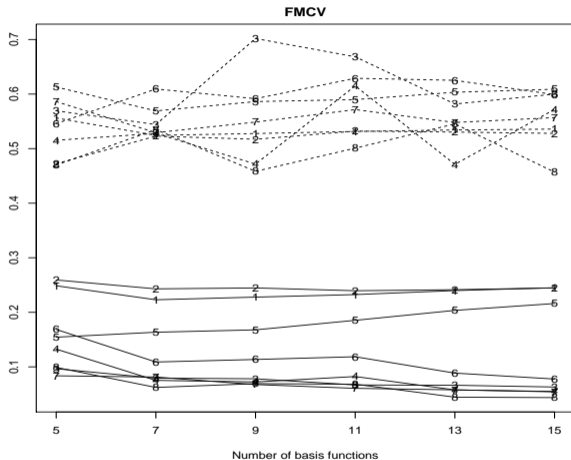


Second variable – abnormal



- The ECG database was used to discriminate between normal and abnormal heartbeats (Olszewski, 2001).
- For illustrative purposes, we show that this has also sense from variability point of view. To do this, we compute the FMCV for both normal and abnormal heartbeats separately.
- The basis functions representation of the data was obtained by using the Fourier and B-spline bases and $B_1 = B_2 = 5, 7, 9, 11, 13, 15$, if it was possible.
- To estimate the FMCV, we used the same estimators as in simulation experiments, i.e., the classical, MCD and S estimators, as well as the pairwise estimator (OGK).
- The standard errors (SE) were obtained by the bootstrap method, based on 1000 bootstrap samples.

Application to ECG data



normal - solid line, abnormal - dashed line, 1 - classical, Fourier, 2 - classical, B-spline, 3 - MCD, Fourier, 4 - MCD, B-spline, 5 - S, Fourier, 6 - S, B-spline, 7 - OGK, Fourier, 8 - OGK, B-spline

- ① Aerts, S., Haesbroeck, G., Ruwet, C. (2015). Multivariate coefficients of variation: comparison and influence functions. *Journal of Multivariate Analysis* 142, 183–198.
- ② Albert, A., Zhang, L. (2010). A novel definition of the multivariate coefficient of variation. *Biometrical Journal* 52, 667–675.
- ③ Olszewski, R.T. (2001). Generalized Feature Extraction for Structural Pattern Recognition in Time-Series Data. Ph.D. Thesis, Carnegie Mellon University, Pittsburgh, PA.
- ④ Reyment, R.A. (1960). Studies on Nigerian Upper Cretaceous and Lower Tertiary Ostracoda: part 1. Senonian and Maastrichtian Ostracoda. *Stockholm Contributions in Geology* 7, 1–238.
- ⑤ Van Valen, L. (1974). Multivariate structural statistics in natural history. *Journal of Theoretical Biology* 45, 235–247.
- ⑥ Voinov, V.G., Nikulin, M.S. (1996). *Unbiased Estimators and Their Applications, Vol. 2, Multivariate Case*. Kluwer, Dordrecht.

Thank you for your attention