

# Journal of Statistical Software

MMMMMM YYYY, Volume VV, Issue II.

http://www.jstatsoft.org/

## An Excel Tool for Statistical Analysis

Fanghu Dong University of Hong Kong Guosheng Yin
University of Hong Kong

#### Abstract

This article presents a standalone macro-enhanced Excel file for statistical analysis. The tool was motivated by classroom need to demonstrate the inner workings of classical multivariate statistical methods to students of levels from advanced undergraduate to entry graduate who take the Multivariate Analysis course using the SAS software for computation. One initial idea driving the development of the Excel tool was to allow students reproduce on a spreadsheet the output of several SAS procs including glm, reg, princomp, cancorr, factor, discrim. Beyond reproducing results, the Excel tool keeps a large portion of the formula chain. Users can change anywhere on the formula chain and immediately see how the output respond to the changes as the spreadsheet automatically updates downstream from the change point. Users can also see all the key variables at the same time, allowing them quickly identify some close relationships between the results. The Excel tool also includes macros that can make non-trivial plots. It can run on both Windows and Mac versions of Excel.

Keywords: Canonical Correlation Analysis, Factor Analysis, General Linear Model, Hotelling's  $T^2$ , Longitudinal Analysis, Linear Discriminant Analysis, Microsoft Excel, Multivariate ANOVA, Multivariate Normal Random Vector, Multivariate Regression, Multivariate Statistical, Principal Component Analysis, Test of Covariance structures, Test of multivariate mean location, Two-sample multivariate mean comparison, Visual Basic for Applications (VBA).

### 1. Introduction

Students learning core multivariate statistics need a little bit different kind of statistical software than the one needed by professionals for production purposes. It is by good software engineering principle to hide the implementation details away from the user. But some proper amount of detail is exactly what is needed for understanding the methods. Arguably, if non-present time pressure, the ideal way to learn the subject is to code one's own implementation of a method like Multivariate Regression with the explicit goal of lining the outputs (coefficients

estimates and sd, MANOVA statistics, etc.) up with those of an established software. One good place to carry out such line-up is on a spreadsheet, where one can have a view of the entire "memory" layout and its dynamic updating monitored by an event system and orchestrated by the functional evaluator. Microsoft Excel is a very popular implementation of the spreadsheet. It is fully integrated with the highly productive Visual Basic for Applications (VBA) language. VBA complements the sheet-level functional environment with procedural programming (e.g. loops, state variables, classes) and integrates with Excel so closely that it can automate literally everything that one does manually on Excel. And even the automation itself is automated. Excel carries a macro recording utility to automate the coding of manual operations. Excel also implements a set of data visualizations utilities that produces several types of sophisticated plots that can be made with a few selections and clicks.

Short-comings of Excel may include reduced speed, a hard cap of data size when facing big datasets, and being not free (monetary sense). For Excel add-in development beyond VBA, extensions are commonly written as COM dll using C++ and/or on the .NET platform using C# or VB.NET through the Visual Studio Tools for Office (VSTO) and therefore is currently hinged to the Windows platform. There could be other short-comings perceived with individual developer's experience. Despite of these, Excel is still a popular numerical environment for mathematical modeling. It supports basic matrix mathematics (multiplication, inversion, and determinant), includes many of the building-block functions in mathematics and statistics, and has the basic functions and utilities for text processing. Finally, for Windows users, Excel has access to unlimited number of dll files that exposes functions and objects to COM.

The software has been used together with SAS (9.2) during a course in Multivariate Statistics at advanced undergraduate to entry graduate level. The course covers both theoretical and computational aspects of classical methods for normality-based multivariate methods including theoretical frameworks such as the General Linear Model and several classical works of Hotelling, Fisher, Pearson, and others. These include multivariate extensions of the univariate two-sample t test, canonical correlation (as a measure of association between two sets of variates), linear discrimination (as a supervised classification algorithm), principal component (as a data orthogonalization algorithm), and Factor analysis (to reduce correlation by splitting random factors of the covariance matrix). The main reference books for the course are Anderson (2003) and Johnson and Wichern (1992). It is felt hard to explain the computation using just the SAS outputs as they are usually printed as isolated "magic" numbers. The Excel tool described in the following was initially written to help explain the computations but it has later accumulated to be a standalone tool.

## 2. Examples

In this section, we give a self-contained explanation of selected statistical methods with examples worked out on the Excel tool. We will not be able to cover every aspect implemented but will try to cover the most important ones.

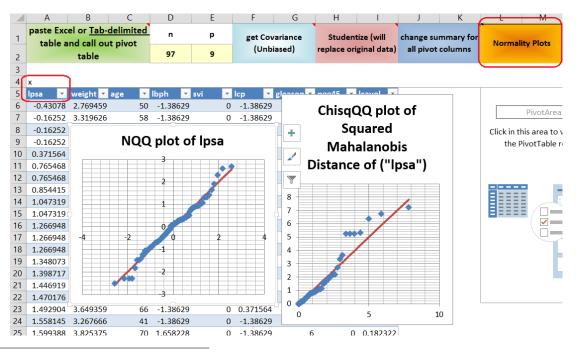
### 2.1. Regression with Excel Tool

This example demonstrates multiple regression, i.e. a single y-variate to be regressed on a number of x-variates. We use sheet "Correl" (**Correl**) to perform multiple regression and a

manual variable selection for the prostate cancer data<sup>1</sup> already stored on the "Data" ( Data ) sheet under name "Prostate", which can be found in the drop-down menu of cell Data!H2. The original source of this dataset is Stamey, Kabalin, McNeal, Johnstone, Freiha, Redwine, and Yang (1989). The column "lpsa" is the y-variate; all other columns are x-variates.

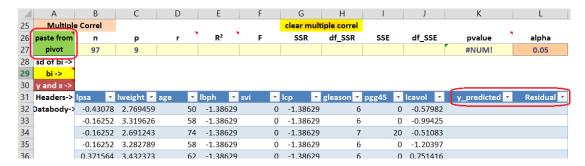
	Α	В	С	D	Е	F	G	Н	1	J	K
1	Paste	from		Register	Selected		Regis	Registered Data table			
2	Clipk	oard		Data	table			Prostate			
2209											
2210		Icavol	lweight	age	lbph	svi	lcp	gleason	pgg45	lpsa	
2211		-0.57982	2.769459	50	-1.38629	0	-1.38629	6	0	-0.43078	
2212		-0.99425	3.319626	58	-1.38629	0	-1.38629	6	0	-0.16252	
2213		-0.51083	2.691243	74	-1.38629	0	-1.38629	7	20	-0.16252	
2214		-1.20397	3.282789	58	-1.38629	0	-1.38629	6	0	-0.16252	
2215		0.751/16	2 //22272	62	-1 38620	0	-1 38620	6	n	0.37156/	

- 1. Copy the Prostate dataset from the sheet "Data" ( Data") to the system clipboard then immediately switch to sheet "Pivot" (Pivot) and double-click the top-left green cell at Pivot! A1 to create a working copy of the dataset on the "Pivot" (Pivot) sheet. All subsequent operations will be performed on this copy.
- 2. Select column "lpsy" by putting an "x" above the header (a reordering of columns will be triggered to prioritize the selected) and double-click the orange button "Normality Plots" (Normality Plots) to have a visual check of the response variable's normality condition. The plots show that the data have a little bit excess kurtosis over that of a normal distribution. Nevertheless, we will proceed for demonstration purpose.

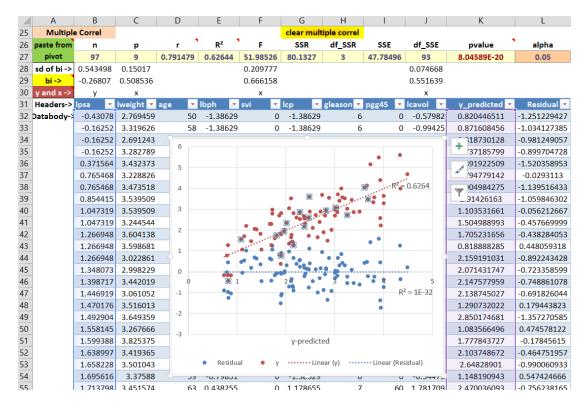


<sup>&</sup>lt;sup>1</sup>The dataset is extracted from R's lasso2 package under name "Prostate". Schematic summary about the dataset can be found at http://www.biostat.jhsph.edu/~ririzarr/Teaching/649/prostate.html and http://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/prostate.info.txt

- 3. Activate sheet "Correl" (Correl).
- 4. On sheet "Correl" (**Correl**), double-click on the green cell at **Correl!A26** to copy-paste the dataset from the "Pivot" (**Pivot**) sheet and augment it with two additional columns: the fitted response "y\_predicted" and the residuals of fitting. The sheet "Correl" (**Correl**) may now appear as



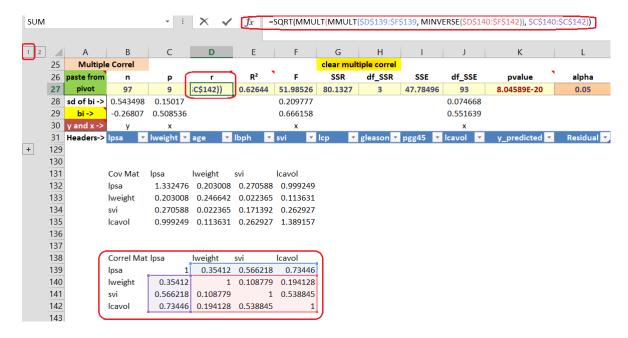
5. Enter "y" in cell Correl!B30 and "x" to any subset of cells Correl!C30:J30 while mointoring the  $R^2$  at cell Correl!E27. After a few trials, one may quickly settle to the subset of "lweight", "svi", "lcavol" giving an  $R^2 = 0.626439681190266$ . One can then quickly make some visualizations of the numbers. The "Correl" (Correl) sheet may now appear as



One may quickly verify the 3 ways of computing  $\mathbb{R}^2$  in multiple regression:

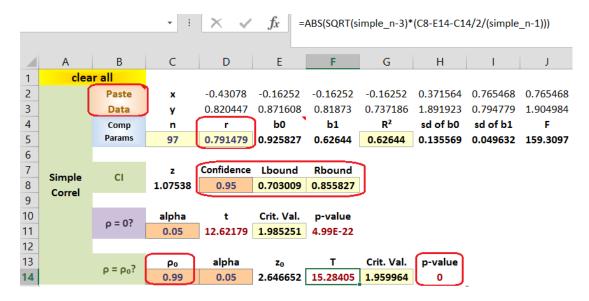
$$R^{2} = \frac{SSR}{SSR + SSE} = \mathbf{R}_{yx}\mathbf{R}_{xx}^{-1}\mathbf{R}_{xy} = \operatorname{corr}(y, \hat{y})^{2}$$

The second way is coded into the formula for r at cell Correl!D27. It interprets multiple regression as a process of maximizing the squared correlation between the response and a vector in the linear space spanned by the regressors. And the correlation-maximizing vector is the  $\hat{y}$ .



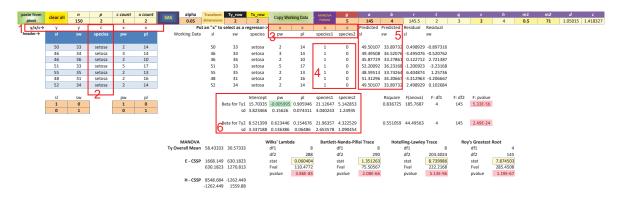
One may want to test whether the hypotheses that  $corr(y, \hat{y})$  is close enough to 1.

- 1. Select the two columns of "y" and "y\_predicted" holding the ctrl key
- 2. Copy the selection
- 3. Double click cell Correl!B2
- 4. Enter 0.99 in cell Correl!C14



#### 2.2. Multivariate Regression involving categorical variables

In a general regression setup, one frequently encounters more than a single response variable and categorical variables in the regressors. The sheet "LM" ( implemented for this task. LM stands for Linear Model. The initiation step is similar as before: after the data is pasted to sheet "Pivot" ( Pivot ) and preprocessed there, one double-clicks the green paste-from-pivot button at the top-left corner to bring data to sheet "LM" ( implemented for this task. One then specifies a y ahead each response column, an x ahead each continuous regressor column and a c ahead each categorical regressor column. A second specification, regarding Rectangle 3 of Fig[], is needed to indicate which of the x and c columns will finally be used with an x in the cells above the Working Data. Note that the categorical variables in rectangle 2 are auto-encoded into dummy variables of rectangle 4 in Fig[]. The mutlivariate regression's coefficients and standard deviation estimates are output in rectangle 6. In addition to estimation of the regression coefficients, the sheet also implements Multivariate-ANOVA tests, a SAS proc glm code generator macro, and transformation matrices on both continuous responses and continuous regressors.



## 2.3. Multivariate Hypothesis Testing with Excel Tool

Many hypotheses for analysis of differences among *correlated* variables can be tested using the Hotelling's  $T^2 \in \mathbb{R}^1_+$  statistic. The  $T^2$  statistic adopts a quadratic form, is a multivariate generalization of the Student's t statistic, and has a sampling distribution linked to the F-distribution (Hotelling 1931).

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \sim T^2(p, n - 1) = \frac{p(n - 1)}{n - p} F(p, n - p)$$

where  $\mathbf{x} \in \mathbb{R}^{p \times n}$  is the data vector,  $\bar{\mathbf{x}} \in \mathbb{R}^p$  is the sample mean vector,  $\boldsymbol{\mu} \in \mathbb{R}^p$  is the true mean vector, n is the sample size, and  $\mathbf{S}$  is the sample covariance matrix

$$\mathbf{S} = \frac{1}{n-1} \left( \mathbf{x} - \bar{\mathbf{x}} \mathbf{1}^{\mathsf{T}} \right) \left( \mathbf{x} - \bar{\mathbf{x}} \mathbf{1}^{\mathsf{T}} \right)^{\mathsf{T}}.$$

If the sample covariance  $\mathbf{S}$  is replaced by the true covariance  $\mathbf{\Sigma}$  then the resulting pivotal quantity  $n(\bar{\mathbf{x}} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu})$  is  $\chi^2$ -distributed. Since any linear transformation of a normal random vector remains normal, the test statistic remains  $T^2$ . One then focuses on the design of a linear transform T that represents the hypothesis of interest:

$$H_0: T\bar{\mathbf{x}} - \boldsymbol{\xi}_0 = \mathbf{0}.$$

Since the design of the linear transform is essentially a process of choosing the proper basis for a re-coordinatization of data, the whole process is intuitively understood as finding the most direct "angle" to view the data such that the hypothesis is settled by judging whether the distance between a pair of points is too much for them to be considered the same point. We demonstrate the methodology using the classic example of Rao (1948) with our Excel tool.

Bark deposit from 4 directions of the trunk of 28 Oak trees.

Rao (1948) exhibits a data set containing measurements of the weights of cork borings taken from 4 directions on the trunk of 28 Oak trees. The hypothesis of interest, suggested by Prof. Mahalanobis to Prof. Rao as well as other studies at the time, was that the deposit is uniform in the North-South directions and also uniform but *less* in the East-West directions.

1. Locate the dataset under name "Cork borings" on the "Data" ( Data ) sheet using drop-down menu at cell Data! A2 and copy it to clipboard

Regis	teredList			<b>-</b>	× <	fx Co	rk borings			
4	Α	В	С	D	E	F	G	Н	1	
1	Paste	from		Register	Selected		Regis	tered Data	table	
2	Clipb	oard		Data	table		(	ork boring	gs	~
374							Caterpillar			^
375		N	Е	S	W		Census Coleman Rep	ort 6th Grade	irc	
376		72	66	76	77		Cork borings			
377		60	53	66	63		Crowder and Crude Oil	Hand (1990)		
878		56	57	64	58		Diabetes			
379		41	29	36	38		Film			~
380		32	32	35	36					
881		30	25	3/1	26					

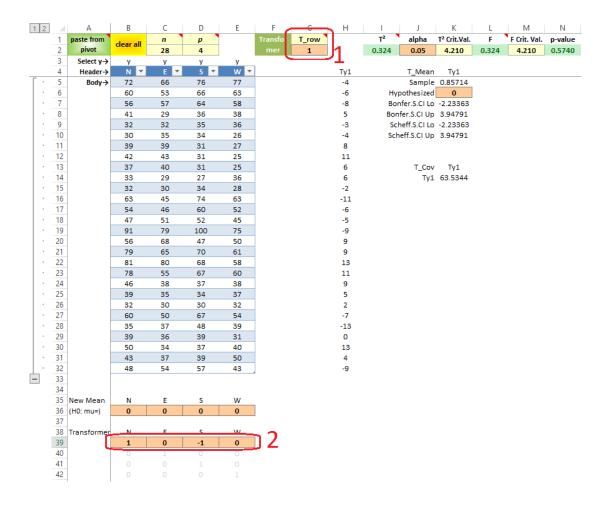
- 2. Activate the "Pivot" (Pivot) sheet and directly double-click the top-left green cell Pivot!A1 to make a copy of the dataset and equip it with the Excel table format.
- 3. Activate the "Tsquare" (Tsquare") sheet and directly double-click the top-left green cell Tsquare!A1 to copy the data over.



4. Perform the following 3 steps exactly: Select cells Tsquare!B3:E3 | press y on the keyboard | Windows user: press ctrl + Enter on the keyboard; Mac user: Press command + Enter. By doing these steps, you have entered 4 "y"s simultaneously.

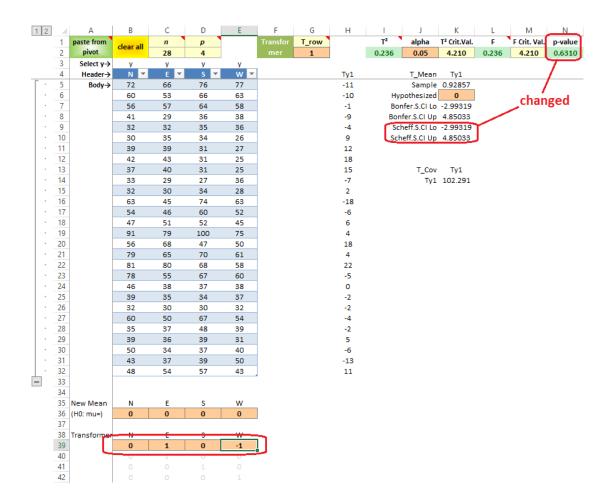
Now that the data has been loaded onto the "Tsquare" (Tsquare") sheet, we proceed to test some hypotheses. We start with the univariate null hypothesis  $H_0: N = S$ . The following step show the p-value of this test in Tsquare!N2, which will be 0.5740 and the null hypothesis therefore accepted.

5. Enter -1 in cell Tsquare!D39 such that the first row of the matrix of our linear transform becomes (1,0,-1,0). Enter 1 in cell Tsquare!G2 to indicate we use only the first row of the transformation matrix. The "Tsquare" (Tsquare) sheet should now appear as



Next we test the another univariate hypothesis  $H_0: E = W$ . The following step will result in a p-value=0.6310 and therefore null hypothesis accepted.

6. Modify the first row of the transformation matrix near Tsquare!B39:E39 into (0, 1, 0, -1). The "Tsquare" (Tsquare) sheet should now appear as



Next we test the two hypotheses jointly:  $H_0: N = S$  and E = W. The following step will result in an increased p-value=0.8194 and therefore null hypothesis accepted.

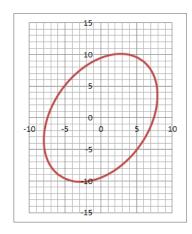
7. Change Tsquare!G2 to 2 and modify the first 2 rows of the transformation matrix near Tsquare!B39:E40 into

$$\left[\begin{array}{cccc} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \end{array}\right]$$

The "Tsquare" (Tsquare) sheet should now appear as

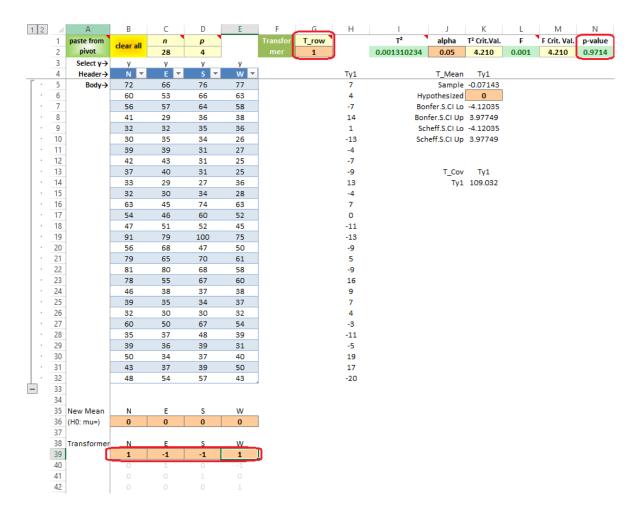
1 2 1 2 3 4 5 6 7 8 9 10 11 11 12		y 72 60 56 41 32 30 39	n 28 y 66 53 57 29 32 35	P 4 y 76 66 64 36 35	y v v 77 63 58 38	Transfor Trow mer 2	Ty1 -4 -6	T <sup>2</sup> 0.417  Ty2 -11	alpha 0.05	T <sup>2</sup> Crit.Val. 6.997  T_Mean	F 0.201 Ty1	Ty2	p-value 0.8194 ased!
3 4 7 · 6 · 7 · 8 · 9 · 10	Select y-> Header-> Body->	y N 72 60 56 41 32 30	y 66 53 57 29 32 35	y 76 66 64 36 35	77 63 58	mer 2	-4	Ty2		T_Mean	Ty1	incre	
4 - 5 - 6 - 7 - 8 - 9 - 10 - 11	Header→ Body→	72 60 56 41 32 30	E ▼ 66 53 57 29 32 35	S ▼ 76 66 64 36 35	77 63 58		-4	-11				Ty2	ased!
5 · 6 · 7 · 8 · 9 · 10 · 11	Body→	72 60 56 41 32 30	66 53 57 29 32 35	76 66 64 36 35	77 63 58		-4	-11				Ty2	asca.
· 6 · 7 · 8 · 9 · 10 · 11		60 56 41 32 30	53 57 29 32 35	66 64 36 35	63 58								
· 7 · 8 · 9 · 10 · 11		56 41 32 30	57 29 32 35	64 36 35	58		-6			Sample		0.92857	
· 8 · 9 · 10 · 11		41 32 30	29 32 35	36 35				-10		pothesized	0	0	
· 9 · 10 · 11		32 30	32 35	35	38		-8	-1		nfer.S.CI Lo		-3.60785	
· 10		30	35				5	-9		fer.S.CI Up		5.465	
. 11					36		-3	-4		heff.S.CI Lo			1
		39		34	26		-4	9	Sch	neff.S.CI Up	4.84177	5.98451	J
· 12			39	31	27		8	12					
		42	43	31	25		11	18					
. 13		37	40	31	25		6	15		T_Cov		Ty2	
· 14		33	29	27	36		6	-7		Ty1	63.5344	28.3968	
· 15		32	30	34	28		-2	2		Ty2	28.3968	102.291	
· 16		63	45	74	63		-11	-18					
. 17	-	54	46	60	52		-6	-6					
· 18		47	51	52	45		-5	6					
· 19		91	79	100	75		-9	4					
· 20		56	68	47	50		9	18					
· 21		79	65	70	61		9	4					
· 22		81	80	68	58		13	22					
· 23		78	55	67	60		11	-5					
· 24		46	38	37	38		9	0					
. 25		39	35	34	37		5	-2					
· 26		32	30	30	32		2	-2					
· 27	_	60	50	67	54		-7	-4					
· 28	_	35	37	48	39		-13	-2					
· 29		39	36	39	31		0	5					
. 30	_	50	34	37	40		13	-6					
. 31	_	43	37	39	50		4	-13					
. 32	_	48	54	57	43		-9	11					
- 33	-												
34	-												
	New Mean	N	E	S	W								
36	(H0: mu=)	0	0	0	0								
37	T		-		144								
			-			<b>S</b>							
						1							
	ļ Ļ	U	1	Ü	-1	ע							
41	-	0	0	0	1								
39 40 41	Į	1 0	0 1	-1 0	0 -1								

8. To understand why the p-value has increased in the previous step, we plot the covariance matrix of the transformed data. Now perform the following steps: copy L14:M15, the covariance of the transformed data (Ty1, Ty2) -> switch to sheet "Cov2Correl" (Cov2Correl) -> double-click on the wide green cell Cov2Correl!A3 -> double-click on the orange cell Cov2Correl!Q1 -> follow the instruction on the popup to pick the covariance range at Cov2Correl!B4:C5 and click done.The following covariance plot should appear on "Cov2Correl" (Cov2Correl) now.



Now we see that the covariance is elongated along the positive sloped direction, making it possible that the mean vector (0.857142857, 0.928571429), displayed at range Tsquare!L5:L6, has a smaller Mahalanobis distance from the center than both its projections on the two standard basis coordinates. This should be understood in common-sense language that the two hypotheses "cross-validate" each other. The data has indicated that it is more natural to have N=S and E=W happening together than separately, and it is rather strange to observe uniformity in only one of the directions but not in the other. The cross-validation effect would not have been captured if we were testing with only univariate procedures.

To test that the E-W direction is less uniform than the N-S direction, we use the linear combination (N-S)-(E-W) and the null hypothesis that the two directions are equally uniform so that linear combination has zero mean under the null. On the "Tsquare" (Tsquare sheet, we change back to use only the 1st row of the transform matrix by setting Tsquare!G2 to 1 and enter (1,-1,-1,1) as the 1st row at Tsquare!B39:E39. The test result accepts the null with an even bigger p-value=0.9714 (Tsquare!N2). Looking into the data, we do find a number of points where the difference between N-S is greater than that between E-W.

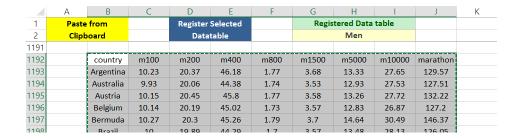


Note because (1, -1, -1, 1) = (1, 0, -1, 0) - (0, 1, 0, -1), therefore we cannot test all three hypotheses in one transformation as that would create a singular covariance matrix for the transformed data and then no  $T^2$  statistic could be constructed.

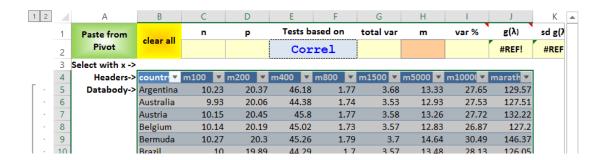
## 2.4. Principal Component Analysis with Excel Tool

Men's track data.

1. Locate the dataset under name "Men" on the "Data" (Data) sheet using drop-down menu at cell Data! A2 and copy it to clipboard

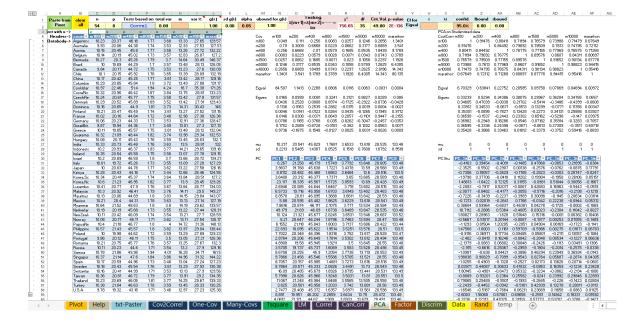


- 2. Activate the "Pivot" (Pivot) sheet and directly double-click the top-left green cell Pivot!A1 to make a copy of the dataset and equip it with the Excel table format.
- 3. Activate the "PCA" (PCA) sheet and directly double-click the top-left green cell PCA!A1 to copy the data over.



4. Perform the following 3 steps exactly: Select cells PCA!C3:J3 | press x on the keyboard | Windows user: press ctrl + Enter on the keyboard; Mac user: Press command + Enter. By doing these steps, you have entered 8 "x"s simultaneously.

The "PCA" (PCA) sheet should now appear as



The main purpose of principal component analysis is dimension reduction and, if one assumes multivariate normality of data, de-correlation as a side effect. The "PCA" (PCA) sheet implements principal component analysis with the multivariate normality assumption. Under the assumption, PCA amounts to eigen-decomposition of the covariance matrix because the eigenvectors give the principal axes of the elliptic contour of the data. The longest principal axis is the linear direction on which the data projects to maximum variance. The second longest principal axis gives the next maximum-variance linear direction, and so on. The covariance matrix is always real-symmetric and have a set of orthogonal eigenvectors with positive eigenvalues. The eigenvalues have the interpretation as the variance of the multivariate data along the corresponding eigenvector direction. There is a subjective decision to make about whether studentization of data is needed. PCA on studentized data is equivalent to eigen-decomposition of the correlation matrix. The main issue is whether one wants to retain variance ratios among the observed variables. This certainly depends on the context. Ratios among some original variables may have established interpretations and hence would be preferred to retain. In other cases, for example, one is preparing the independent variables going into the right-hand side of a regression formula, one might want to studentize the data as the regression coefficients can recover such ratio. In middle cases, a properly estimated convex combination between the two matrices might be considered. Following are some further details regarding implementation.

- 1. PCA on correlation matrix do not add back the mean vector because if we do that we must also multiply back the s.d. But we don't want to do that because the s.d. is now unity for studentized data. This reflects that PCA on correlation matrix focuses on the angular difference between coordinates and ignores the radial differences.
- 2. PCA on covariance matrix may add back the mean vector. The PCA on multivariate normal sample is essentially doing a rotation of the sample space about the mean vector, not the origin. To carry out that rotation via a rotation matrix, we need to first remove the mean before applying the rotation matrix formed by the eigenvectors, and may or may not add the mean vector back.

3. PCA is essentially a data orthogonalization routine and does not model the mean vector. It can be used to orthogonalize the covariates for regression if prediction is the main goal and interpretation of the new covariates' meaning are not a concern. A more important purpose is dimension reduction. Those eigenvectors with tiny eigenvalues indicates insufficient information in those trailing dimensions and hence can be removed to avoid over-fitting.

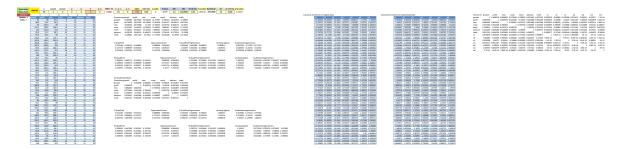
Note that the original definition of PCA is to be a method seeking the linear directions on which data projects with maximum variance. With this definition, it is not restricted to multivariate normality of data but is applicable under any sampling assumption by proper techniques. Nonetheless, if we can assume multivariate normality, the computation becomes much easier.

## 2.5. Correlation Analysis with Excel Tool

We use Salespeople Data to demonstrate canonical correlation analysis following Anderson (2003) The raw data can be located with name "Salespeople" on the "Data" (Data ) sheet

Paste from Clip	board	Register	Selected	F	Registered Data table				
raste from Cit	board	Data	table		Salespeople				
growth	profit	new	creat	mech	abstract	math			
93	96	97.8	9	12	9	20			
88.8	91.8	96.8	7	10	10	15			
95	100.3	99	8	12	9	26			
101.3	103.8	106.8	13	14	12	29			
102	107.8	103	10	15	12	32			
95.8	97.5	99.3	10	14	11	21			
95.5	99.5	99	9	12	9	25			

As before, the raw data should be copied first to sheet "Pivot" ( $\begin{subarray}{l} \begin{subarray}{l} \begin{subarray}{l}$ 



The overall idea of Canonical Correlation Hotelling (1936) Analysis is to construct a scalar "correlation measure" r to describe linear association between two vectors of multivariate

random variables  $\mathbf{v} \in \mathbb{R}^p$  and  $\mathbf{w} \in \mathbb{R}^q$ . The scalar is constructed by finding in each space a unit directional vector such that the usual unsigned correlation (geometrically the cosine of the angle) between the two directional vectors are maximal. The derivation of the pair of optimal directional vectors happens to become eigenvalue problems for two positive semi-definite matrices. The two matrices happen to share the a same set of non-zero eigenvalues and the largest eigenvalue is the square "correlation" measure being sought. Moreover, the eigenvectors corresponding to the largest eigenvalue for each matrix is the unit vector in the respective space.

$$r^{2}\mathbf{v} = \mathbf{S}_{vv}^{-1}\mathbf{S}_{vw}\mathbf{S}_{ww}^{-1}\mathbf{S}_{wv}\mathbf{v}$$
$$r^{2}\mathbf{w} = \mathbf{S}_{vvv}^{-1}\mathbf{S}_{wv}\mathbf{S}_{vv}^{-1}\mathbf{S}_{vw}\mathbf{w}$$

Covariance	growth	profit	new	creat	mech	abstract	math
growth	53.83664	68.79409	30.56453	16.57967	17,58522	10.58759	71.69861
profit	68.79409	2.1018	40.19508	21.65629	25.561277	10.08131	100.7442
new	30.56453	40.19508	22.205	13.03653	10.16755	6.463673	42.3351
creat	16.57967	21.65629	13.03653	15.60367	7.898367	1.241633	17.17633
mech	17.58522	25,56127	10.16755	7.898367	11,45673	2.795102	20.49306
abstract	10.58759	10.08131	6.463673	1.241633	2, 79,101	<b>1</b> 24.577959	12.7698
math	71.69861	100.7442	42.3351	17.17633	20.49306	12.7698	111.0433

V QuadProd				Eigenvalue	Cancorr		V QuadPro	d Eigenved	tors		sd along eigv	
=MMULT(MN	MULT(MIN	IVERSE(\$L\$	5:\$N\$7) <b>,</b> \$C	\$5:\$R\$7),	MMULT(MII	NVERSE(\$C	\$8:\$R\$11),	\$L\$8:\$N\$1	1))		9.780811	
0.127373	0.809859	-0.053626		0.771071	0.878107		0.20467	-0.634364	-0.189002		2.619561	
0.472958	0.145236	0.573176		0.147153	0.383606		0.765428	0.624226	-0.700056		1.825843	
W QuadProd	d				Eigenvalue	Cancorr		W QuadPro	od Eigenve	ctors		
0.393691	0.040775	0.210079	0.341437		0.988996	0.994483		0.523093	0.335881	0.617205	0.053313	
-0.109506	0.203974	-0.098968	0.646698		0.771071	0.878107		0.230529	-0.351913	-0.355207	0.950185	
0.442953 -	0.081711	0.583377	0.126459		0.147153	0.383606		0.671708	0.865515	-0.701485	0.150619	
0.116091	0.193654	0.027501	0.726178		3.64F-17	6.04F-09		0.471209	-0.119268	0.028369	-0.267618	

For the second-largest eigenvalue, it is the squared canonical correlation between the spaces  $\alpha^{\perp} \subset \mathbb{R}^{p-1}$  and  $\beta^{\perp} \subset \mathbb{R}^{q-1}$ , the orthogonal complement spaces of the two directional vectors, and recursively so doing for the other smaller non-zero eigenvalues. Note that the "Quad-Prod" matrices are real symmetric, this means that the eigenvectors are perpendicular to each other. Together with the "maximal correlation" property, the eigenvectors can be used to recoordinate the data columns, as done in range starting at cell CanCorr!AG4 and CanCorr!AP4. Fig [] shows first a few rows of re-coordinated variates.

Canonical	Transform	on Original	Data			
v1	v2	v3	w1	w2	w3	w4
15.46365	16.24594	-12.3602	3.05927	2.408177	-1.7791	2.76664
15.03552	16.29364	-13.1261	2.633711	3.263887	-2.32531	2.584792
15.7723	15.83873	-12.5111	3.366502	1.80589	-1.95766	2.184552
16.84893	17.94649	-13.488	4.233901	3.647039	-1.81534	2.821697
16.67892	16.19417	-12.1811	4.243885	2.663342	-2.66291	2.817273
15.78709	16.72753	-12.0346	3.432453	3.12063	-2.36545	3.46391
15.78675	16.1195	-12.2397	3.37342	2.066597	-1.72244	2.297155
12 42756	17 21827	-15 NAR	6 412024	3 383886	-2 02529	3 008336

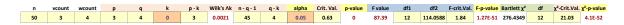
Canonica	al Trar	nsform	on Student	ized Data			
v1	v2		v3	w1	w2	w3	w4
-0.9783	8 0	.36254	0.819381	-0.97479	-0.0943	0.088519	0.06569
-1.4065	2 0.4	10239	0.053517	-1.40035	0.761407	-0.45769	-0.11616
-0.6697	4 -0	.04467	0.668475	-0.66756	-0.69659	-0.09004	-0.5164
0.40689	7 2.0	063089	-0.3084	0.19984	1.144559	0.052276	0.120748
0.23688	3 0.3	310765	0.998522	0.209824	0.160863	-0.79529	0.116324
-0.6549	5 0.8	344131	1.145015	-0.60161	0.618151	-0.49783	0.762961
-0.6552	9 0.2	236094	0.939863	-0.66064	-0.43588	0.145182	-0.40379
2 04552	גי 1	22427	-1 86845	2 383063	N 8814N6	-∩ 15767	በ ဒበ73ጸ6

The full correlation matrix including the original dataset and the constructed canonical variates is displayed in range starting at cell CanCorr!AY3.

Full Correl	growth	profit	new	creat	mech	abstract	math	v1	v2	v3	w1	w2	w3	w4
growth	1	0.926076	0.884002	0.572036	0.708074	0.674407	0.927312	0.979878	-0.00065	0.199598	0.974471	-0.00057	-0.07657	-7.7E-15
profit	0.926076	1	0.842523	0.541508	0.74591	0.465388	0.944296	0.946409	-0.32288	-0.0075	0.941187	-0.28353	0.002879	1.45E-14
new	0.884002	0.842523	1	0.700363	0.637471	0.641089	0.852568	0.951862	0.186301	-0.24341	0.94661	0.163592	0.093375	1.89E-14
creat	0.572036	0.541508	0.700363	1	0.590736	0.146907	0.412639	0.634809	0.189406	-0.24988	0.638331	0.215698	0.65141	0.348817
mech	0.708074	0.74591	0.637471	0.590736	1	0.38595	0.574553	0.717184	-0.20861	0.025985	0.721163	-0.23756	-0.06774	0.647224
abstract	0.674407	0.465388	0.641089	0.146907	0.38595	1	0.566372	0.643678	0.440224	0.220275	0.647249	0.501333	-0.57422	-0.00096
math	0.927312	0.944296	0.852568	0.412639	0.574553	0.566372	1	0.938877	-0.17345	0.036146	0.944086	-0.19753	-0.09423	-0.24658
v1	0.979878	0.946409	0.951862	0.634809	0.717184	0.643678	0.938877	1	5.05E-13	-1.1E-13	0.994483	-5E-14	-2E-14	3.23E-14
v2	-0.00065	-0.32288	0.186301	0.189406	-0.20861	0.440224	-0.17345	5.05E-13	1	-7.9E-13	6.05E-14	0.878107	-5.4E-14	7.24E-14
v3	0.199598	-0.0075	-0.24341	-0.24988	0.025985	0.220275	0.036146	-1.1E-13	-7.9E-13	1	-1.5E-14	5.01E-15	-0.38361	-1E-13
w1	0.974471	0.941187	0.94661	0.638331	0.721163	0.647249	0.944086	0.994483	6.05E-14	-1.5E-14	1	-6.4E-14	-1.3E-15	-2.2E-15
w2	-0.00057	-0.28353	0.163592	0.215698	-0.23756	0.501333	-0.19753	-5E-14	0.878107	5.01E-15	-6.4E-14	1	-4.8E-15	-4.2E-15
w3	-0.07657	0.002879	0.093375	0.65141	-0.06774	-0.57422	-0.09423	-2E-14	-5.4E-14	-0.38361	-1.3E-15	-4.8E-15	1	-2E-15
w4	-7.7E-15	1.45E-14	1.89E-14	0.348817	0.647224	-0.00096	-0.24658	3.23E-14	7.24E-14	-1E-13	-2.2E-15	-4.2E-15	-2E-15	1

The fact that each V-canonical variate only respond to one of the W-canonical variates means that if we run a regression of all V's on all W's, then we know it is merely a bunch of 1-to-1 simple linear regression performed together.

Finally, the first two rows implements some hypothesis tests to determine how many canonical variates should be retained if dimension reduction is a concern.



#### 2.6. Factor Analysis with Excel Tool

When all observed variable are used and the residuals are still not spherical, one ponders over the existence of latent factors. The factor model is formulated as

$$y - \mu = \Lambda F + \varepsilon$$

where F is a multivariate normal vector that can be required to satisfy

$$\operatorname{var}(F) = I$$

and  $\varepsilon$  is the new residual that is hopefully more spherical than before. The coefficient matrix  $\Lambda$  is called the factor loadings. Both F and  $\Lambda$  will need to be estimated. A further simplifying assumption makes  $\Lambda$  constant so that

$$\operatorname{var}(\Lambda F) = \Lambda \Lambda^{\mathsf{T}}$$

hence

$$var(y - \mu) = \Lambda \Lambda^{\mathsf{T}} + var(\varepsilon)$$
.

This suggests expanding the real-symmetric left-hand side by eigen-decomposition

$$var(y - \mu) = \lambda_1 v_1 v_1^{\mathsf{T}} + \lambda_2 v_2 v_2^{\mathsf{T}} + \dots + \lambda_p v_p v_p^{\mathsf{T}}$$

and estimating  $\Lambda\Lambda^{\dagger}$  by first m terms of this summation. This is the method implemented in the Excel tool. After  $\Lambda$ , the factor loadings matrix, is estimated, one then proceed to estimate F by, for example, least square. In the Excel tool, we follow Anderson (2003) to implement the weighted least square method of Bartlett(1938)

$$\hat{F} = \left(\Lambda^{\mathsf{T}} \Psi^{-1} \Lambda\right)^{-1} \Lambda^{\mathsf{T}} \Psi^{-1} z$$

and the conditional expectation method of Thomson (1951)

$$\hat{F} = \Lambda^{\mathsf{T}} (\Lambda \Lambda^{\mathsf{T}} + \Psi)^{-1} z = (I + \Lambda^{\mathsf{T}} \Psi \Lambda)^{-1} \Lambda^{\mathsf{T}} \Psi^{-1} z.$$

In the following example, we analyze the olympic 88 men decathlon data using factor analysis. The dataset doesn't contain any covariants, making it suitable to demonstrate the latent factor approach.

1988 Summer Olympics Men's Decathlon data.

1. Locate the dataset under name "Olympic88" on the "Data" (Data") sheet using drop-down menu at cell Data!RegisteredList and copy it to clipboard.

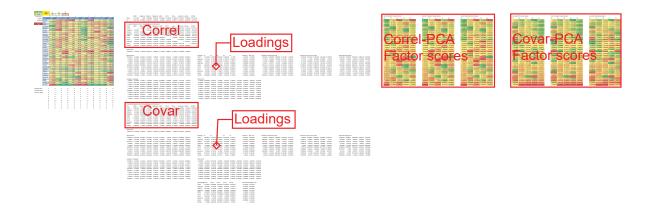
Paste fron	a Clinha	ovel	R	legister Se	lected		F	Registered	Data tab	le
Paste Iron	ii Ciipbo	oaru		Datatak	ole			Olymp	ics88	
name	m100	longjump	shotput	highjump	m400	m110	discus	polevault	javelin	m1500
Schenk	11.25	7.43	15.48	2.27	48.9	15.13	49.28	4.7	61.32	268.95
Voss	10.87	7.45	14.97	1.97	47.71	14.46	44.36	5.1	61.76	273.02
Steen	11.18	7.44	14.2	1.97	48.29	14.81	43.66	5.2	64.16	263.2
Thompson	10.62	7.38	15.02	2.03	49.06	14.72	44.8	4.9	64.04	285.11
Blondel	11.02	7.43	12.92	1.97	47.44	14.4	41.2	5.2	57.46	256.64
Plaziat	10.83	7.72	13.58	2.12	48.34	14.18	43.06	4.9	52.18	274.07

- 2. Activate the "Pivot" (**Pivot**) sheet and directly double-click the top-left green cell **Pivot!A1** to make a copy of the dataset and equip it with the Excel table format.
- 3. Activate the "Factor" (Factor) sheet and directly double-click the top-left green cell Factor! A1 to copy the data over.

Paste from	clear	n	р	m	_						
Pivot	Clear	34	11	5							
Select with x ->											
Headers->	name	m100	longjump	shotput	highjump	m400	m110	discus	polevault	javelin	m1500
Databody->	Schenk	11.25	7.43	15.48	2.27	48.9	15.13	49.28	4.7	61.32	268.95
	Voss	10.87	7.45	14.97	1.97	47.71	14.46	44.36	5.1	61.76	273.02
Run	Steen	11.18	7.44	14.2	1.97	48.29	14.81	43.66	5.2	64.16	263.2
	Thompson	10.62	7.38	15.02	2.03	49.06	14.72	44.8	4.9	64.04	285.11
	Blondel	11.02	7.43	12.92	1.97	47.44	14.4	41.2	5.2	57.46	256.64
	Plaziat	10.83	7.72	13.58	2.12	48.34	14.18	43.06	4.9	52.18	274.07

4. Perform the following 3 steps exactly: Select cells Factor!C3:L3 | press x on the keyboard | Windows user: press ctrl + Enter on the keyboard; Mac user: Press command + Enter. By doing these steps, you have entered 10 "x"s simultaneously.

The "Factor" (Factor) sheet should now appear as Fig[]. The heat-mapped regions on the right are 6 factor scores from different scoring methods and whether the input data has been studentized.



## 3. Storing Data

The data import/export of the tool is delegated to Excel's own data i/o utilities. The user can add a blank sheet to import the dataset from various original sources. Next, all datasets need to be transformed into the *data frame* format, i.e., a matrix of data with column header texts. A row in the data frame is a multivariate sample vector jointly observed for all the variates and a column is a univariate sample observed repeatedly for a single variate. The number of rows in the data frame is the sample size. This format should be familiar to both SAS (sas7bdat) and R (data frame) users.

After importing and transforming into the data frame format, the dataset should be registered on the sheet "Data" ( Data ) and archived there for future usage. The sheet "Data" ( Data ) can be navigated via a drop-down menu near cell Data!A2, which lists all registered datasets. A working copy of a dataset should be put on the sheet "Pivot" ( Pivot ) .

Following is an example of generating multivariate normal random sample using the "Rand" (Rand) sheet and then registering and storing it on the "Data" (Data) sheet.

- 1. Activate sheet "Rand" (Rand).
- 2. Put 3000 to Rand!C2 to specify the sample size.
- 3. Enter the mean vector as a column vector right below cell Rand!D1:  $[0,0,0]^{\mathsf{T}}$
- 4. Enter the covariance matrix as a symmetric positive-definite matrix right below cell Rand!E1 and extend to the right:

$$\left[\begin{array}{ccc}
9 & 5 & -5 \\
5 & 8 & 0 \\
-5 & 0 & 7
\end{array}\right]$$

Sheet "Rand" (Rand) should now appear as

	Α	В	С	D	Е	F	G	Н
1	Multiv	ariate	Size	Mean	Cov			
2	Non	mal	3000	0	9	5	-5	
3				0	5	8	0	
4	Clear s	tart ↓		0	-5	0	7	
5	\$D	\$2						

Now if you double-click on the cell Rand!A7 (Generate) some equation will be entered to the sheet by a VBA macro (shtRand.generate) triggered on the double-click event you just performed to the cell Rand!A7. The  $3000 \times 3$  range Rand!E7:G3006 is also automatically selected so that you can directly press the scatter plot button to have a visual check as I am doing. Sheet "Rand" (Rand) should now appear as

E7			- i	× <	<i>f</i> x {=N	IMULT(\$N\$	7:\$P\$3006	, TRANSPO	DSE(\$N\$2:\$P	\$4)) + TF	RANSPOSE(\$D	\$2:\$D\$4	)}		
	АВ	С	D	Е	F	G	Н	1	J	К	L	М	N	0	Р
1	Multivariate	Size	Mean	Cov		U	- 11	'	,	K	L	IVI	IN	0	
2	Normal	3000	0		5	-5					15.37495		0.74205	0.070511	0.66662
3			0	5	8	0					7.507461		0.503088	-0.71579	-0.484
4	Clear start ↓		0	-5	0	7					1.117591		-0.44302		
5	\$D\$2														
6															
7				0.419668	-1.48264	-1.59495							0.272105	2.198934	0.0940
8	Generate				1.960644	-1.6788							3.771528		-0.0668
9					2.260400								0.507853		
10				-3.1915	-2.5	Sca	atter plo	ot of (X2	, X3) vs X1				-5.18415	-0.79712	1.06747
11				1.227685	-1.5				, , , , , , , , , , , ,				1.8729	3.96985	-0.6630
12				2.540642	10.0		15						5.669984	-9.0092	-1.5473
13				1.058948	1.00								0.859103	-1.32353	0.77220
14				2.909658	1.41		10	•	•				4.527777	1.783666	-0.8639
15				-2.51951	-2.1	9.0							-4.28941	-0.68511	1.06769
16				-7.11651	-2.4	-	THE PARTY NAMED IN	4.0	2,50				-8.9458	-2.52061	-0.4508
17				-0.37733	-1.1		3		Sec.				0.057659	2.286698	-0.8720
18				-0.56706	-0.9				<b>:</b>				-1.08375	0.331083	0.32069
19				4.122614			J 0		500		• X2		5.160819	2.576374	0.16706
20				-1.26798		-10	0 5 5 1	10	15	• X3		-1.15527	-1.07782	-0.502	
21				0.603142			5					1.006844	2.111382	-0.4393	
22				-0.29309	-2.0		200	4.4.4	2300				-1.04567	1.807599	0.53311
23				-0.47049	-2	200		•				-1.18845	2.496551	0.35306	
24				2.151751	-0.7		-10		•				2.681309	2.978093	-0.07186
25				-1.50274	-2.3								-1.31541	3.124838	-1.1205
26				-5.36459	-4.€		-15	1/4					-7.4995	1.058334	0.188696
27				0.182531	-3.8			X1					-0.77395	3.621897	0.75222
28				-1.3394	-1.34800	-U.U <del>444</del> /							-1.65269	0.901805	-0.26493
29				3.123384	1.204387	-2.96681							4.237966	1.419314	-0.18222
30				-3.69119	-5.56519	-0.24385							-5.4308	3.89267	0.09638
31				-2.34551	-1.03397	0.32607							-2.40512	0.348187	-0.8780
32				-2.44091	2.257429	1.496102							-1.33839	-2.82737	-1.8727
33				1.494838	3.227683	0.032101							2.718832	-2.22725	-0.5484

The quick visual check confirms that: (i) the mean location is near the origin, (ii) both data exhibits the elliptical contour consistent with the positive definite quadratic form embedded in the MVN density, and (iii) the positive correlation between X1 and X2 gives the  $+45^{\circ}$  rotation of the blue sample while the negative correlation between X1 and X2 gives the  $-45^{\circ}$  rotation of the red sample.

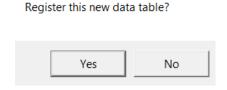
A very important feature here is that the output is a function of the input and is connected to input via a formula chain. As a result, if the user changes the covariance input, then immediately the plot will update. This reveals the many upsides of using Excel to do mathematical modeling on small-to-medium sized data: it is a functional environment; it has a robust event system; it has a lot of productive utilities to operate the data; and it lets you monitor all variables at the same time. These are all conducive to (self-)teaching core multivariate statistics.

Next we will register the generated random sample to the Data sheet. Note that the following step of storing and registering dataset on the tool is the same for any data as long as it is presented in the data frame format. One can leverage Excel's own utilities to prepare the raw data into the data frame format.

1. Add names to the 3 columns by typing into cells Rand!E6:E8 "X1", "X2", "X3". During

the process the sample may be regenerated.

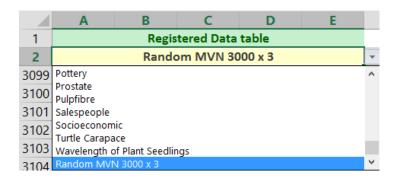
- 2. Press ctrl+a on Windows or command+a on Mac to select the entire  $3001 \times 3$  data range Rand!E6:G3006 (now with a header row)
- 3. Press ctrl+c to copy to clipboard.
- 4. Launch a simple text editor and paste the data there to make it plain text.
- 5. Copy everything in the simple text editor to clipboard
- 6. Activate sheet "Data" (Data)
- 7. Double click on Data!I1 (Paste from Clipboard) to initiate pasting and registration of a new dataset
- 8. Click Yes to confirm registration of this dataset to sheet "Data" (Data)



9. Enter "Random MVN 3000 x 3" to name the dataset being registered



10. Roll out the drop-down menu at cell Data! A2 and look for the newly registered dataset and select it. After select, the screen will auto-navigate to the dataset and select it.



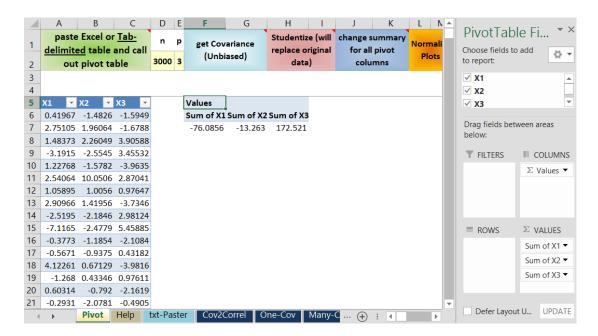
11. You can now press the keyboard shortcut to copy the selection to the system clipboard.

	Α	В	С	D	Е							
1	Registered Data table											
2	Random MVN 3000 x 3											
3099												
3100		X1	X2	X3								
3101		0.419668	-1.48264	-1.59495								
3102		2.751046	1.960644	-1.6788								
3103		1.483725	2.260492	3.905876								
3104		-3.1915	-2.55449	3.455321								
3105		1.227685	-1.57823	-3.96347								

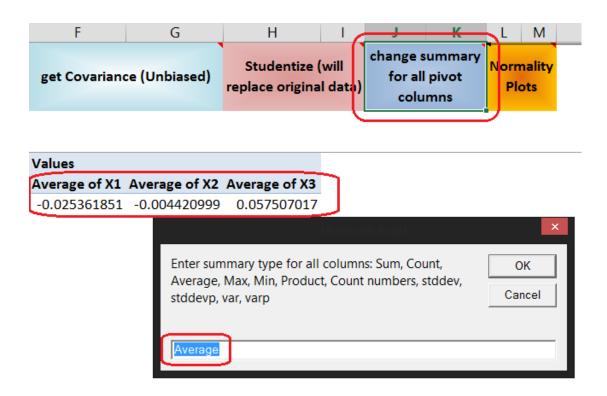
12. Once the data is on the clipboard, switch to the "Pivot" (**Pivot**) sheet and immediately double-click on the top-left green button



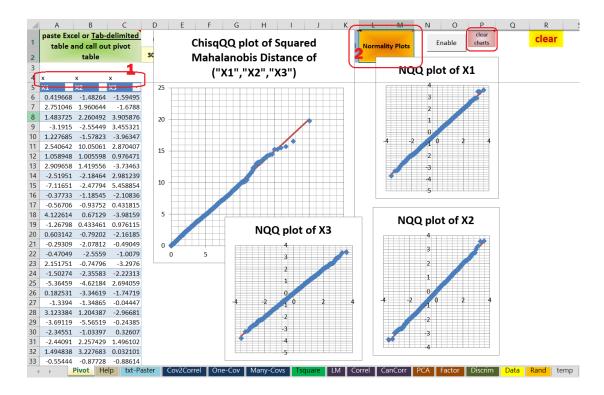
13. After double-click, the dataset will be pasted to the "Pivot" (Pivot) sheet as an Excel table and hence it enjoys all Excel table associated utilities such as filtering with complex conditions (SQL select in disguise), multi-column sort, pivot table (aggregation, a bit like R apply), and graphics. Depending on how you configure the pivot table, "Pivot" (Pivot) may now appear as



14. The cell Pivot!J1 can help you change the aggregate function to one of Sum, Count, Average, Max, Min, Product, Count numbers, stddev, stddevp, var, and varp. This is a quick way to get the mean vector and the sd vector.



15. The orange button at cell Pivot!L1 makes normal and  $\chi^2$  QQ plots to help visually check marginal and joint normality. To do this, put an "x" in cells Pivot!A4:C4 above the column headers and then double click on the orange button. The sheet "Pivot" (Pivot) may now appear as



The normal quantile-quantile plots are made by the VBA macro NormalQQplot. The chi-squared quantile-quantile plot for inspecting violation of joint normality is made by the macro MahalanobisChisqQQplot. The Mahalanobis distance is defined as

$$D = \sqrt{(\mathbf{x} - \bar{\mathbf{x}})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{x} - \bar{\mathbf{x}})}$$

Its square has an asymptotic  $\chi(p)$  distribution where p is the number of variates (p=3 here).

## 4. Discussion

The Excel tool is sheet-oriented. There four types of sheets: Data storage sheet ("Data" (Data")), Data simulation sheet ("Rand" (Rand")), Data pre-process sheet ("Pivot" (Pivot")), and Method sheet ("LM" (Data")), etc). The "Pivot" (Pivot") contains the data in analysis-ready state. The method sheets all implement a paste-from-pivot button at the top-left corner to create its own copy of the analysis-ready data and then builds formula chains to arrive at results. All sheets can use built-in Excel functionalities as well as custom addin functions written in VBA, XLL(COM), or .NET(VSTO). The transparency of the computation together with Excel's own tools around formula building, tracing, and checking allows complex models to be understood quickly. It is also a good self-documentation of an implementation elsewhere such as R. We recommend all R implementation has an equivalent Excel Tool sheet. This will solve an important problem of getting one's implementation details understood, extended with confidence, and understood again.

## References

Anderson TW (2003). An Introduction to Multivariate Statistical Analysis. 3 edition. Wiley.

Hotelling H (1931). "The Generalization of Student's Ratio." The Annals of Mathematical Statistics, 2(3), 360-378. doi:10.1214/aoms/1177732979. URL http://dx.doi.org/10.1214/aoms/1177732979.

Hotelling H (1936). "Relations between two sets of variates." *Biometrika*, **28**(3/4), 321–377. 1936.

Johnson RA, Wichern DW (1992). Applied multivariate statistical analysis. 4th edition. Prentice hall Englewood Cliffs, NJ.

Rao CR (1948). "Tests of significance in multivariate analysis." *Biometrika*, **35**(1/2), 58–79. 1948.

Stamey TA, Kabalin JN, McNeal JE, Johnstone IM, Freiha F, Redwine EA, Yang N (1989). "Prostate specific antigen in the diagnosis and treatment of adenocarcinoma of the prostate. II. Radical prostatectomy treated patients." *The Journal of wrology*, **141**(5), 1076–1083.

## Affiliation:

Fanghu Dong
Department of Statistics and Actuarial Science
Faculty of Science
University of Hong Kong
Hong Kong
E-mail: jdong@connect.hku.hk

E-mail: jdong@connect.hku.hk
URL: http://web.hku.hk/~jdong/

Guosheng Yin Department of Statistics and Actuarial Science Faculty of Science University of Hong Kong Hong Kong

E-mail: gyin@hku.hk

URL: http://web.hku.hk/~gyin/

http://www.jstatsoft.org/

http://www.amstat.org/ Submitted: yyyy-mm-dd

Accepted: yyyy-mm-dd