

# New normality test based on the conditional second moments and the 20-60-20 rule

Damian Jelito

Based on the joint work with Marcin Pitera

`damian.jelito@im.uj.edu.pl`

Jagiellonian University in Kraków

27-th November, 2018

# Outline

Motivation - different aspects of normality

The 20-60-20 Rule

The new test statistic

Numerical experiments

Derivation of the asymptotic distribution

Possible generalisations

## Several tests for normality

- ▶ Cramer-von Mises family

$$CvM := n \int_{-\infty}^{+\infty} (F_n(x) - F(x))^2 w(x) dF(x)$$

- ▶ Special case - Anderson-Darling

$$w_{AD}(x) := \frac{1}{F(x)(1 - F(x))}$$

- ▶ Computational version of the AD test statistic

$$AD := -n - \frac{1}{n} \sum_{i=1}^n (2i - 1) (\ln(\Phi(Y_{(i)})) - \ln(1 - \Phi(Y_{(n-i+1)})))$$

with

$$Y_{(i)} := \frac{X_{(i)} - \bar{X}}{\hat{\sigma}}$$

# Several tests for normality

## ► Shapiro-Wilk test

$$SW := \frac{(\sum_{i=1}^n a_i X_{(i)})^2}{\sum_{i=1}^n (X_{(i)} - \bar{X})^2}$$

where

$$(a_1, \dots, a_n) = \frac{m^T V^{-1}}{(m^T V^{-1} V^{-1} m)^{\frac{1}{2}}}, \quad m_i = \mathbb{E}(X_{(i:n)}), \quad V_{ij} = \text{Cov}(X_{(i:n)}, X_{(j:n)}).$$

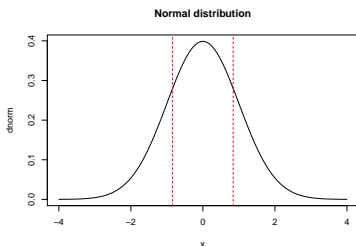
## ► Jarque-Bera test

$$JB := \frac{n}{6} \left( \hat{S}^2 + \frac{1}{4}(\hat{C} - 3)^2 \right)$$

with

$$\hat{S} := \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \hat{X})^3}{\hat{\sigma}^3}, \quad \hat{C} := \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \hat{X})^4}{\hat{\sigma}^4}$$

# The 20-60-20 Rule



$$\sigma_L^2 = \sigma_M^2 = \sigma_R^2$$

- ▶ We split the normal random variable into three disjoint conditioning sets: left (L), middle (M), and right (R):

$$L := \left( -\infty, F_X^{-1}(\tilde{q}) \right],$$

$$M := \left( F_X^{-1}(\tilde{q}), F_X^{-1}(1 - \tilde{q}) \right),$$

$$R := \left[ F_X^{-1}(1 - \tilde{q}), +\infty \right)$$

- ▶ For a unique  $\tilde{q} \approx 0.2$  the conditional variances coincide
- ▶ This property might be linked to the statistical phenomenon known as *The 20-60-20 Rule*.

# Mathematical formulation of the 20-60-20 Rule

- Recall the conditional variance  $\sigma_A^2 := \mathbb{E}((X - \mathbb{E}(X|A))^2|A)$

## Theorem

If  $X \sim \mathcal{N}(\mu, \sigma)$ , then

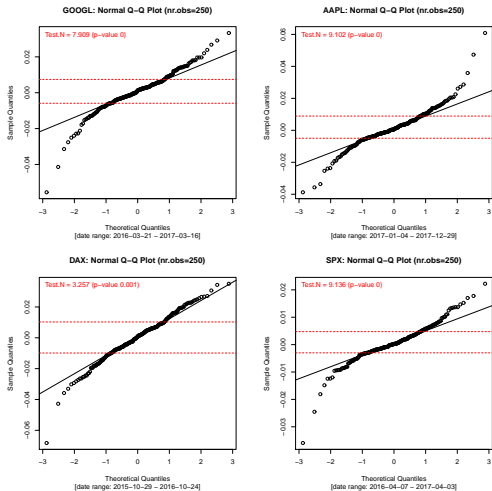
$$\sigma_L^2 = \sigma_M^2 = \sigma_R^2,$$

where  $\tilde{q} = \Phi(x) \approx 0.19809$  and  $x$  is a unique negative solution of the equation

$$-x\Phi(x) - \phi(x)(1 - 2\Phi(x)) = 0.$$

- The property holds true for an arbitrary number of conditioning sets as well as in multivariate and elliptical case (for covariance matrices).

# The 20-60-20 Rule and the Q-Q plot



1. We take return rates (based on adjusted daily close prices)
2. We make a simple Q-Q plot with theoretical normal distribution
3. We check if 20/60/20 division leads to accurate clustering

# The test statistic

1. We introduce the test statistic

$$N := \frac{\sqrt{n}}{\rho} \left( \frac{\hat{\sigma}_L^2 - \hat{\sigma}_M^2}{\hat{\sigma}^2} + \frac{\hat{\sigma}_R^2 - \hat{\sigma}_M^2}{\hat{\sigma}^2} \right)$$

where  $\rho \approx 1.8$  is a fixed normalising constant

2. Under the normality assumption  $N$  is a pivotal quantity
3.  $N$  can be seen as a measure of tail fatness (cf. Anderson-Darling test)
4.  $N$  is based on the conditional second moments while Jarque-Bera test uses the third and fourth moments

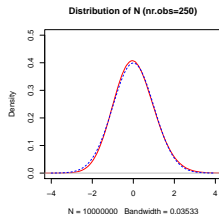
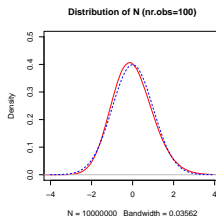
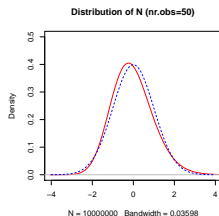
```

1 Test.N <- function(x){
  q1 <- quantile(x,0.2)
3  q2 <- quantile(x,0.8)
  n <- length(x)
5  x.low <- x[x <= q1]
  x.med <- x[x > q1 & x < q2]
7  x.high <- x[x >= q2]
  N <- var(x.low)+var(x.high)-2*var(x.med)
9  N <- N * sqrt(n)/(var(x)*1.8)
  return(N)}

```



# Asymptotic distribution of the test statistic



$\alpha$	$n$	$\Phi^{-1}(1 - \alpha)$	$F_n^{-1}(1 - \alpha)$
1.0%	50	2.33	2.64
	100		2.53
	250		2.47
2.5%	50	1.96	2.14
	100		2.08
	250		2.06
5.0%	50	1.64	1.74
	100		1.71
	250		1.71

## Theorem

Let  $X \sim N(\mu, \sigma)$ . Then,

$$N \xrightarrow{d} \mathcal{N}(0, 1), \quad n \rightarrow \infty.$$

Moreover,  $\rho$  is independent of  $\mu$ ,  $\sigma$  and  $n$ .

We will come back to the asymptotics later.

## Market data case study - overview

- ▶ We take S&P500 stocks returns from 01.2000 to 05.2018 (4610 daily adjusted close price returns for 381 stocks)
- ▶ For a given stock, the sample is split into disjoint sets of length  $n$  with  $n \in \{50, 100, 250\}$ .
- ▶  $N$  is compared with Jarque–Bera test, Anderson–Darling test, and Shapiro–Wilk test.
- ▶ Normality hypothesis is checked at confidence level  $\alpha \in \{1.0\%, 2.5\%, 5.0\%\}$ .
- ▶ Non-normality of returns is a well known fact, hence we expect the null hypothesis to be rejected.
- ▶ We compute three supplementary metrics
  - ▶ **Statistic T - total rejection ratio** of a given test at confidence level  $\alpha$  - for what proportion of all subsets the normality assumption was rejected.
  - ▶ **Statistic U - unique rejection ratio** of a given test at confidence level  $\alpha$  - for what proportion of all subsets the normality assumption was rejected only by a given test (among all four tests).
  - ▶ **Statistic A - acceptance ratio** of a given test at confidence level  $\alpha$  - for what proportion of all subsets the normality assumption was not rejected by any tests if it was not rejected by a given test.

## Market data case study - results

Desc	nr runs	$\alpha$	n	rejects	JB	AD	SW	N
T	35052	1.0%	50	31.5%	<b>25.9%</b>	17.3%	23.2%	<b>25.9%</b>
U					1.9%	0.6%	0.3%	<b>3.1%</b>
A					<b>94.4%</b>	85.8%	91.7%	<b>94.4%</b>
T	35052	2.5%	50	39.6%	32.4%	22.9%	28.3%	<b>32.5%</b>
U					2.4%	0.9%	0.3%	<b>3.9%</b>
A					<b>92.8%</b>	83.3%	88.7%	<b>92.8%</b>
T	35052	5.0%	50	47.9%	38.6%	29.1%	33.7%	<b>39.6%</b>
U					2.4%	1.3%	0.4%	<b>5.1%</b>
A					90.6%	81.1%	85.8%	<b>91.6%</b>
T	17526	1.0%	100	52.8%	45.2%	31.8%	41.3%	<b>46.1%</b>
U					2.2%	0.6%	0.3%	<b>4.4%</b>
A					92.5%	79.1%	88.6%	<b>93.4%</b>
T	17526	2.5%	100	61.3%	52.9%	38.8%	47.6%	<b>54.3%</b>
U					2.2%	0.7%	0.2%	<b>5.1%</b>
A					91.6%	77.5%	86.3%	<b>93.1%</b>
T	17526	5.0%	100	68.4%	59.7%	45.7%	53.4%	<b>61.3%</b>
U					2.2%	0.8%	0.2%	<b>5.3%</b>
A					91.3%	77.3%	84.9%	<b>92.9%</b>
T	6858	1.0%	250	88.5%	82.1%	71.2%	79.3%	<b>85.4%</b>
U					1.0%	0.4%	0.1%	<b>3.8%</b>
A					93.6%	82.7%	90.8%	<b>96.9%</b>
T	6858	2.5%	250	91.8%	86.8%	77.7%	83.7%	<b>89.4%</b>
U					0.7%	0.2%	0.1%	<b>3.0%</b>
A					95.1%	85.9%	91.9%	<b>97.6%</b>
T	6858	5.0%	250	93.9%	89.7%	82.4%	86.9%	<b>92.0%</b>
U					0.5%	0.2%	0.0%	<b>2.5%</b>
A					95.7%	88.5%	92.9%	<b>98.1%</b>

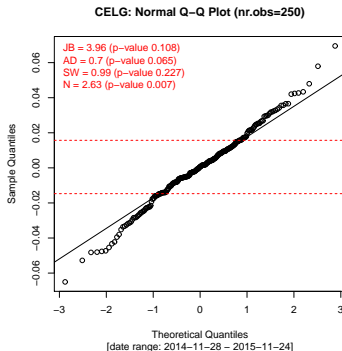
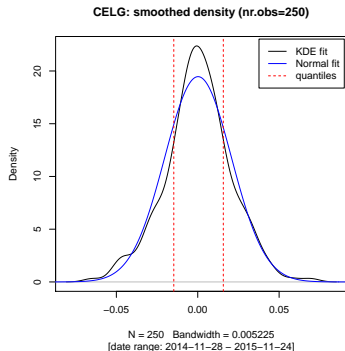
Metrics for the  $i$ -th test

$$T_i = \frac{\#\{\text{i-th test rejected}\}}{\#\text{samples}},$$

$$U_i = \frac{\#\{\text{only i-th test rejected}\}}{\#\text{samples}},$$

$$A_i = \frac{\#\{\text{no test rejected}\}}{\#\{\text{i-th test didn't reject}\}}$$

## Market data case study - close-up



## Market data case study - p-value performance

Desc	nr runs	cond nr runs	n	JB	AD	SW	N
O	35052	16807	50	20.3%	7.8%	5.1%	28.7%
S				35.7%	17.2%	25.8%	—
O	17526	11994	100	15.0%	4.6%	3.2%	26.2%
S				26.1%	10.8%	18.7%	—
O	6858	6443	250	4.6%	1.4%	0.6%	15.4%
S				8.3%	3.4%	6.2%	—

$$O_i = \frac{\#\{\forall_j \text{p.value}_i \leq \text{p.value}_j\}}{\#\text{samples}}$$

$$S_i = \frac{\#\{\text{p.value}_i \leq \text{p.value}_N\}}{\#\text{samples}}$$

- ▶ We compare p-values between different tests.
- ▶ For brevity, we consider only samples that were rejected by at least one test at level 5%.
- ▶ We present two performance measures:
  - ▶ **Statistic O - ratio of best p-values compared to other tests** - for what number of observations the p-value for a given statistic is smaller compared to all other p-values.
  - ▶ **Statistic S - ratio of best p-values compared to single test N** - for what number of observations the p-value for a given statistic is smaller compared to  $N$  test p-value.

# Asymptotic distribution - notation

- ▶  $X \sim \mathcal{N}(\mu, \sigma)$  with mean parameter  $\mu$  and standard deviation parameter  $\sigma$
- ▶  $X_{(i)}$  -  $i$ th order statistic of the sample  $(X_1, \dots, X_n)$
- ▶ For  $0 \leq \alpha < \beta \leq 1$  we define the conditioning set

$$A[\alpha, \beta] := \{x \in \mathbb{R} : F_X^{-1}(\alpha) < x \leq F_X^{-1}(\beta)\}.$$

- ▶ Recall  $L := A[0, \tilde{q}]$ ,  $R := A[1 - \tilde{q}, 1]$ ,  $M := A[\tilde{q}, 1 - \tilde{q}]$ ,
- ▶ The conditional sample mean on set  $A$

$$\bar{X}_A := \frac{1}{[n\beta] - [n\alpha]} \sum_{i=[n\alpha]+1}^{[n\beta]} X_{(i)}$$

- ▶ The conditional sample variance on set  $A$

$$\hat{\sigma}_A^2 := \frac{1}{[n\beta] - [n\alpha]} \sum_{i=[n\alpha]+1}^{[n\beta]} (X_{(i)} - \bar{X}_A)^2$$

# Main result

- Recall that the test statistic  $N$  is given by

$$N = \frac{1}{\rho} \left( \frac{\hat{\sigma}_L^2 - \hat{\sigma}_M^2}{\hat{\sigma}^2} + \frac{\hat{\sigma}_R^2 - \hat{\sigma}_M^2}{\hat{\sigma}^2} \right) \sqrt{n}.$$

- Restate our main result:

## Theorem

Let  $X \sim N(\mu, \sigma)$ . Then,

$$N \xrightarrow{d} \mathcal{N}(0, 1), \quad n \rightarrow \infty,$$

where  $\rho$  is a fixed normalising constant independent of  $\mu$ ,  $\sigma$ , and  $n$ .

- In the proof we make use of the series of lemmas.

## Additional notation

For a fixed set  $A = A[\alpha, \beta]$ , we define

$$\mu_A := \mathbb{E}[X|X \in A],$$

$$\sigma_A^2 := \mathbb{E}[(X - \mu_A)^2|X \in A],$$

$$\kappa_A := \frac{1}{(\sigma_A^2)^2} \mathbb{E}[(X - \mu_A)^4|X \in A],$$

$$a := F_X^{-1}(\alpha) = \mu + \sigma\Phi^{-1}(\alpha),$$

$$b := F_X^{-1}(\beta) = \mu + \sigma\Phi^{-1}(\beta).$$



# Asymptotic normality of a conditional sample variance

## Lemma (Asymptotic normality of a conditional sample variance)

For any  $A = A[\alpha, \beta]$  it follows that

$$\sqrt{n} \left( \hat{\sigma}_A^2 - \sigma_A^2 \right) \xrightarrow{d} \mathcal{N}(0, \tau_A),$$

where

$$\begin{aligned} \tau_A^2 := & \frac{1}{(\beta - \alpha)^2} \left( (\beta - \alpha)(\sigma_A^2)^2(\kappa_A - 1) + \alpha(1 - \alpha) \left( (a - \mu_A)^2 - \sigma_A^2 \right)^2 \right. \\ & - \alpha(1 - \beta) \left( (a - \mu_A)^2 - \sigma_A^2 \right) \left( (b - \mu_A)^2 - \sigma_A^2 \right) \\ & \left. + \beta(1 - \beta) \left( (b - \mu_A)^2 - \sigma_A^2 \right)^2 \right).^1 \end{aligned}$$

---

<sup>1</sup>Note that for degenerate cases  $\alpha = 0$  and  $\beta = 1$ , we get  $a = -\infty$  and  $b = \infty$ , respectively. In those cases, the convention  $0 \cdot \infty = 0$  should be used.

## Additional lemmas

### Lemma (Consistency of a conditional sample mean)

For any  $A = A[\alpha, \beta]$  it follows that  $\overline{X}_A \xrightarrow{\mathbb{P}} \mu_A, \quad n \rightarrow \infty.$

### Lemma (Asymptotic normality of a conditional sample mean)

For any  $A = A[\alpha, \beta]$  it follows that

$$\sqrt{n} (\overline{X}_A - \mu_A) \xrightarrow{d} \mathcal{N}(0, \eta_A), \quad n \rightarrow \infty,$$

where  $0 < \eta_A < \infty.$

## Some remarks on quantile estimators

### Remark

We can replace  $[n\beta] - [n\alpha]$  by  $[n\beta] - [n\alpha] - 1$  in the definition of the conditional sample variance and our results remain valid.

### Remark

Consider sequences  $(\alpha_n)$  and  $(\beta_n)$  such that  $n\alpha - \alpha_n$  and  $\beta_n - n\beta$  are bounded. The corresponding conditional sample mean and variance is given by

$$\bar{X}_A^* := \frac{1}{\beta_n - \alpha_n} \sum_{i=\alpha_n+1}^{\beta_n} X_{(i)},$$

$$\hat{\sigma}_A^{2,*} := \frac{1}{\beta_n - \alpha_n} \sum_{i=\alpha_n+1}^{\beta_n} (X_{(i)} - \bar{X}_A^*)^2.$$

Then, we can replace  $\bar{X}_A$  and  $\hat{\sigma}_A^2$  by  $\bar{X}_A^*$  and  $\hat{\sigma}_A^{2,*}$  in our theorem and lemmas and their statements remain valid.

# Some possible generalisations

- ▶ Different test statistic, e.g.  $S_1 := \left( \frac{\hat{\sigma}_L^2 - \hat{\sigma}_R^2}{\hat{\sigma}_M^2} \right) \sqrt{n}$  or  $S_2 := \left( \frac{\hat{\sigma}_L^2}{\hat{\sigma}_M^2} - \lambda \right) \sqrt{n}$
- ▶ More conditioning sets, e.g.

$$N_k := \sqrt{n} \left( \frac{\hat{\sigma}_{A_1}^2 - \hat{\sigma}_{A_k}^2}{\hat{\sigma}^2} + \dots + \frac{\hat{\sigma}_{A_{k-1}}^2 - \hat{\sigma}_{A_k}^2}{\hat{\sigma}^2} \right)$$

for partitioning sets  $A_1, \dots, A_k$  such that  $\sigma_{A_1}^2 = \dots = \sigma_{A_k}^2$  (exist for any  $k \in \mathbb{N}$ )

- ▶ Multivariate case, i.e.  $\|\Sigma_A - \Sigma_B\|$ , where  $\|\cdot\|$  is a matrix norm and  $\Sigma_A$  is a conditional covariance matrix

The End

Thank you for your attention!

## References

- [1] P. Jaworski, M. Pitera, *The 20-60-20 Rule*, Discrete Cont Dyn-B, Vol. 21 (2016), No 4, pp. 1149-1166.
- [2] P. Jaworski, M. Pitera, *A note on conditional covariance matrices for elliptical distributions*, Stat Probabil Lett, Vol. 129C (2017), pp. 230-235.
- [3] D. Jelito, M. Pitera, *New fat-tail normality test based on conditional second moments with applications to finance*, preprint, arXiv: 1811.05464
- [4] H. Thode, *Testing for normality*, Marcel Dekker, 2002