

# Multivariate model assessment without $\chi^2$

**Sara Algeri**

salgeri@umn.edu

School of Statistics, University of Minnesota

PHYSTAT 2-Samples,  
Virtual meeting.

June 2, 2023.

**To truly understand what are the main objects involved in our analysis, we begin with the one-dimensional setting...**

# Goodness-of-fit vs 2-Samples tests

## Goodness-of-fit (GOF)

- Inputs: A sample  $x_1, \dots, x_n$  from  $P$  and a postulated distribution  $Q$

- Test:

$$H_0 : P = Q \quad \text{vs} \quad H_1 : P \neq Q$$

- Test statistics: Functionals of the empirical process

$$v_n(x) = \sqrt{n} \left[ \hat{P}(x) - Q(x) \right]$$

with  $\hat{P}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq x\}}$ .

**Note:** we can incorporate unknown parameters through  $Q(x)$ .

## 2-Samples

- Inputs: Samples  $x_{11}, \dots, x_{1n_1}$  from  $P_1$  and  $x_{21}, \dots, x_{2n_2}$  from  $P_2$

- Test:

$$H_0 : P_1 = P_2 = P \quad \text{vs} \quad H_1 : P_1 \neq P_2$$

- Test statistics: Functionals of the empirical process

$$v_n(x) = \sqrt{\frac{n_1 n_2}{n}} \left[ \hat{P}_1(x) - \hat{P}(x) \right]$$

with  $\hat{P}(x) = \frac{n_1}{n} \hat{P}_1(x) + \frac{n_2}{n} \hat{P}_2(x)$   
 $\hat{P}_j(x) = \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbb{1}_{\{x_{ji} \leq x\}}$ ,  $j = 1, 2$ ,  
and  $n = n_1 + n_2$ .

## A very famous example

Recall that our empirical process specifies as

- $v_n(x) = \sqrt{n} \left[ \hat{P}(x) - Q(x) \right]$ , for GOF problems.
- $v_n(x) = \sqrt{\frac{n_1 n}{n_2}} \left[ \hat{P}_1(x) - \hat{P}(x) \right]$ , for 2-Samples problems.

Then the Kolmogorov-Smirnov statistic specifies as

$$KS = \sup_x |v_n(x)|. \quad (1)$$

In the one-dimensional setting and if, in the context of GOF,  $Q$  does not depend on unknown parameters,  $KS$  is **asymptotically distribution-free**.

### Distribution-freeness in GOF

We have *distribution-freeness* whenever the null distribution of the test statistic considered does depend on the distribution  $Q$  being tested.

## Another desirable property of GOF

- Let's keep in mind that distribution-freeness is not all we need.
- We also want sensitivity (power) against “all” alternatives.

### Note

The latter is essentially what differentiates GOF tests from tests of hypotheses (e.g., Neyman-Pearson) where the power is concentrated towards the specific alternative model specified under  $H_1$ .

### But what does “all” actually mean?

...there exist alternatives that cannot be detected even by Neyman-Pearson, so there is no hope we can detect those via GOF.

# Power against (converging) contiguous alternatives

## A more sensible criterion...

We require our GOF test to have some power against all (converging) **contiguous** alternatives.

## What is a (converging) contiguous alternative?

Heuristically...

- They are alternatives that get progressively closer to the null as the sample size increases.
- They are detectable via Neyman-Pearson\*.
- In the limit, we can identify the direction from which they approach the null.

\*See “Oosterhoff J., van Zwet W.R. A note on contiguity and Hellinger distance. *Springer New York*, 2012.”

## ...a little more formally

### Converging contiguous alternatives

Let  $Q$  be the null distribution postulated for the underlying continuous data generating process, and let  $q$  be its density. We say that the distribution  $\tilde{P}_n$  is a contiguous alternative to  $Q$  if its density can be specified as

$$\tilde{p}_n(x) = q(x) \left[ 1 + \frac{h_n(x)}{\sqrt{n}} \right],$$

with  $\|h_n(x)\|_Q^2 < \infty$  and  $\|h_n(x) - h\|_Q^2 \rightarrow 0$ . This last condition is what makes them “converging” and the function  $h$  which corresponds to the direction from which  $\tilde{p}_n$  approaches  $q$ .

# An important note on data binning

## Can't we just bin the data and rely Pearson?

In multidimensional settings and/or when  $Q$  depends on unknown parameters, it is common practice to bin the data and use Pearson  $\chi^2$  (or asymptotic equivalent) to perform GOF.

### Warning

It can be shown\* that Pearson (and many other similar statistics) have no power against infinitely many converging contiguous alternatives

\*See "Algeri S. and Khmaladze E.V. When Pearson  $\chi^2$  and other divisible statistics are not goodness-of-fit tests. *In preparation.*"



# First, let's extend what we know for 1D...

## 1D

- Sample:  $x_1, \dots, x_n$ . Each observation  $x_i$  is a scalar.
- Distribution function under  $H_0$  evaluated at  $x$ :

$$Q(x) = P(X \leq x | H_0)$$

- Empirical distribution function under  $H_0$  evaluated at  $x$ :

$$\hat{P}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq x\}}$$

- Empirical process:

$$v_n(x) = \sqrt{n} \left[ \hat{P}(x) - Q(x) \right]$$

## multi-D

- Sample:  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . Each observation  $\mathbf{x}_i = (x_{1i}, \dots, x_{Di})$  is a vector.
- Distribution function under  $H_0$  evaluated at  $\mathbf{x} = (x_1, \dots, x_D)$ :

$$Q(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_D \leq x_D | H_0)$$

- Empirical distribution function under  $H_0$  evaluated at  $\mathbf{x}$ :

$$\hat{P}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_{1i} \leq x_1, \dots, x_{Di} \leq x_D\}}$$

- Empirical process:

$$v_n(\mathbf{x}) = \sqrt{n} \left[ \hat{P}(\mathbf{x}) - Q(\mathbf{x}) \right]$$

## ...and then let's get parametric...

Given a set of  $n$  observations from an unknown cumulative distribution function (cdf)  $P(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_d \leq x_d)$ , we are interested in testing

$$H_0 : P(\mathbf{x}) = Q(\mathbf{x}, \theta) \quad \text{versus} \quad H_1 : P(\mathbf{x}) \neq Q(\mathbf{x}, \theta)$$

for some postulated distribution  $Q(\mathbf{x}, \theta)$ . To perform the test above, we consider the parametric empirical process  $v_Q(\mathbf{x}, \theta)$

$$v_Q(\mathbf{x}, \theta) = \sqrt{n} \left[ \hat{P}(\mathbf{x}) - Q(\mathbf{x}, \theta) \right] \quad (2)$$

## ...we now have access to an entire family of tests!

Recall that our empirical process is  $v_n(\mathbf{x}, \theta) = \sqrt{n}[P_n(\mathbf{x}) - Q(\mathbf{x}, \theta)]$  and let  $\hat{\theta}$  be the Maximum Likelihood Estimate (MLE) of  $\theta \Rightarrow$  we can construct tests statistics as functionals of  $v_n(\mathbf{x}, \hat{\theta})$ . E.g.,

• Kolmogorov-Smirnov statistic:  $KS = \sup_{\mathbf{x}} v_n(\mathbf{x}, \hat{\theta})$ .

• Cramer-von Mises statistic:  $CvM = \int |v_n(\mathbf{x}, \hat{\theta})|^2 dQ(\mathbf{x}, \hat{\theta})$ .

• Anderson-Darling statistic:  $AD = \int \left| \frac{v_n(\mathbf{x}, \hat{\theta})}{\sqrt{Q(\mathbf{x}, \hat{\theta})(1-Q(\mathbf{x}, \hat{\theta}))}} \right|^2 dQ(\mathbf{x}, \hat{\theta})$ .

(where all the integrals are multivariate).

They are not distribution-free but we can simulate their distribution via the parametric bootstrap.

**Cons:** Computational complexity may be high + simulations must be repeated on a case-by-case basis.



In the remaining of the talk we will see two approaches which will help us to overcome these two limitations.

# Estimated and projected empirical process

We can approximate our estimated empirical process via

$$\underbrace{v_Q(x, \hat{\theta})}_{\text{Empirical process at } \hat{\theta}} = \underbrace{v_Q(x, \theta)}_{\text{Empirical process at } \theta} - \sum_{j=1}^p \int_{[-\infty, \mathbf{x}]} \underbrace{b_j(t, \theta)}_{j\text{-th normalized score function}} dt \frac{1}{\sqrt{n}} \sum_{i=1}^n b_j(x_i, \theta) + o_p(1)$$

where  $[-\infty, \mathbf{x}] = [-\infty, x_1] \times \dots \times [-\infty, x_D]$

- We denote the right hand side with  $\tilde{v}_Q(x, \theta)$ . It is a projection of  $v_Q(x, \theta)$  orthogonal to the normalized score functions  $b_j(x, \theta)$ , i.e., the components of

$$b(x, \theta) = \underbrace{\Gamma_{\theta}^{-1/2}}_{\text{Inverse sqrt of the Fisher info}} \underbrace{\frac{\partial}{\partial \theta} \log q(x, \theta)}_{\text{Score vector}}. \quad (3)$$

- The *projected empirical process*\* does not depend on  $\hat{\theta}$ !

\* See “Khmaladze, E.V. (1980). The use of  $\omega^2$  tests for testing parametric hypotheses.

*Theory of Probability & Its Applications.*”

# A toy example to assess the computational gain

We draw a sample of  $n = 100$  observations from

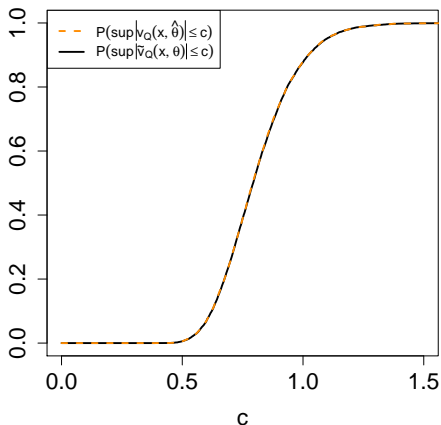
$$q(\mathbf{x}, \boldsymbol{\theta}) \propto e^{-\frac{1}{2\theta_3} [(x_1 - \theta_1)^2 + (x_2 - \theta_2)^2]} \quad \mathbf{x} \in [1, 20] \times [1, 25], \quad (4)$$

$\boldsymbol{\theta} = (-2, 5, 25)$  and its MLE is  $\hat{\boldsymbol{\theta}}_{obs} = (-0.77, 6.32, 22.02)$ .

We proceed by simulating the distribution of the KS statistic, i.e.,

1. We simulate  $\sup_{\mathbf{x}} |v_Q(\mathbf{x}, \hat{\boldsymbol{\theta}})|$  by sampling from  $Q(\mathbf{x}, \hat{\boldsymbol{\theta}}_{obs})$  via the parametric bootstrap.
2. We simulate  $\sup_{\mathbf{x}} |\tilde{v}_Q(\mathbf{x}, \boldsymbol{\theta})|$  by sampling from  $Q(\mathbf{x}, \hat{\boldsymbol{\theta}}_{obs})$  via the parametric bootstrap.

# Simulated distributions of the KS statistic



	$\sup_x  \tilde{v}_Q(x, \theta) $	$\sup_x  v_Q(x, \hat{\theta}) $
CPU time	9.429 mins	12.198 hrs

B=10,000 n=100 R=2000

**But what if we want to test another model,  $F(x, \beta)$   
for which all of this is not at all feasible?  
(Can we somehow retrieve distribution-freeness?)**

## A useful (re-)formulation

We rewrite our empirical process as,

$$\tilde{v}_Q(\mathbf{x}, \theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \left[ \mathbb{1}_{\{x_{1i} \leq x_1, \dots, x_{Di} \leq x_D\}} - Q(\mathbf{x}, \theta) \right] - b^T(x_i, \theta) \int_{[-\infty, \mathbf{x}]} b(\mathbf{t}, \theta) d\mathbf{t} \right\}$$

Setting everything in the curly brackets equal to  $\psi_{\mathbf{x}}(\mathbf{x}_i, \theta)$ , we have

$$\tilde{v}_Q(\mathbf{x}, \theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{\mathbf{x}}(\mathbf{x}_i, \theta). \quad (5)$$

- One can show that the limit of  $\tilde{v}_Q(\mathbf{x}, \theta)$  is Gaussian.
  - its mean and covariance are  $E_Q[\psi_{\mathbf{x}}] = 0$  and  $E_Q[\psi_{\mathbf{x}}\psi_{\mathbf{x}'}]$
- $\Rightarrow$  the  $\psi_{\mathbf{x}}$  fully characterize the limiting distribution of  $\tilde{v}_Q(\mathbf{x}, \theta)$ .



## Towards (asymptotic) distribution-freeness

Can we construct another process whose limit, under  $F(\mathbf{x}, \beta)$ , will be the same as that of  $\tilde{v}_Q(\mathbf{x}, \theta)$  under  $Q$ ?

The key here is to “play” with our  $\psi_{\mathbf{x}}(\mathbf{x}_i, \theta)$  functions so that, by taking a suitable transformation of them, we can construct a new process that, under  $F$ , will have the same limiting distribution as  $\tilde{v}_Q(\mathbf{x}, \theta)$ , under  $Q$ .

**This can be done by means of the Khmaladze-2 (K-2) transform\*.**

\*See “Khmaladze, E.V. (2016). Unitary transformations, empirical processes and distribution free testing. *Bernoulli*, 2016.”

# The K-2 transform in a nutshell

The K-2 transform applied to the functions  $\psi_{\mathbf{x}}(\mathbf{x}_i, \boldsymbol{\theta})$  is

$$\phi_{\mathbf{x}}(\mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\beta}) = \underbrace{\mathbf{U} \left[ \mathbf{K} \left[ l_{\boldsymbol{\theta}, \boldsymbol{\beta}}(\mathbf{x}_i) \psi_{\mathbf{x}}(\mathbf{x}_i, \boldsymbol{\theta}) \right] \right]}_{\text{K-2 transform}}$$

- The isometry  $l_{\boldsymbol{\theta}, \boldsymbol{\beta}}(\mathbf{x}) = \sqrt{\frac{q(\mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{x}, \boldsymbol{\beta})}}$  ensures  $E_F[(l_{\boldsymbol{\theta}, \boldsymbol{\beta}}\psi_{\mathbf{x}})(l_{\boldsymbol{\theta}, \boldsymbol{\beta}}\psi_{\mathbf{x}'})] = E_Q[\psi_{\mathbf{x}}\psi_{\mathbf{x}'}]$ .
- The unitary operator  $\mathbf{K}$  ensures that  $E_F[\mathbf{K}l_{\boldsymbol{\theta}, \boldsymbol{\beta}}\psi_{\mathbf{x}}] = E_Q[\psi_{\mathbf{x}}] = 0$ .
- The unitary operator  $\mathbf{U}$  ensures orthogonality w.r.t. to the normalized score functions under  $F(\mathbf{x}, \boldsymbol{\beta})$ .

## A new family of test statistics

Recall that

$$\tilde{v}_F(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_{\mathbf{x}}(\mathbf{x}_i, \boldsymbol{\theta}, \boldsymbol{\beta}) \quad \text{and} \quad \tilde{v}_Q(\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{\mathbf{x}}(\mathbf{x}_i, \boldsymbol{\theta})$$

We can now construct our *K-2 transformed* test statistics as

$$\begin{aligned} \text{KS}_{F|Q} &= \sup_{\mathbf{x}} | \tilde{v}_F(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}) |, & \text{CvM}_{F|Q} &= \int \tilde{v}_F^2(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}) dQ(\mathbf{x}, \boldsymbol{\theta}), \\ \text{and } \text{AD}_{F|Q} &= \int \frac{\tilde{v}_F^2(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta})}{Q(\mathbf{x}, \boldsymbol{\theta})[1 - Q(\mathbf{x}, \boldsymbol{\theta})]} dQ(\mathbf{x}, \boldsymbol{\theta}), \end{aligned} \quad (6)$$

which have the same limiting distribution as

$$\begin{aligned} \text{KS}_Q &= \sup_{\mathbf{x}} | \tilde{v}_Q(\mathbf{x}, \boldsymbol{\theta}) |, & \text{CvM}_Q &= \int \tilde{v}_Q^2(\mathbf{x}, \boldsymbol{\theta}) dQ(\mathbf{x}, \boldsymbol{\theta}), \\ \text{and } \text{AD}_Q &= \int \frac{\tilde{v}_Q^2(\mathbf{x}, \boldsymbol{\theta})}{Q(\mathbf{x}, \boldsymbol{\theta})[1 - Q(\mathbf{x}, \boldsymbol{\theta})]} dQ(\mathbf{x}, \boldsymbol{\theta}), \end{aligned} \quad (7)$$

# Requirements on $F$ and $Q$

Can we use any  $F(\mathbf{x}, \beta)$  and any  $Q(\mathbf{x}, \theta)$ ?

- Let  $f(\mathbf{x}, \beta)$  and  $q(\mathbf{x}, \theta)$  be the densities of  $F(\mathbf{x}, \beta)$  and  $Q(\mathbf{x}, \theta)$ . We require that:
  - $f(\mathbf{x}, \beta) = 0$  iff  $q(\mathbf{x}, \theta) = 0$  (they have the same support).
  - $\theta, \beta$  are both of size  $p$  (they have the same size).
- **These are rather general criteria!**  $\Rightarrow Q(\mathbf{x}, \theta)$  can be chosen to be arbitrarily simple to ease the computations.
- We call  $Q(\mathbf{x}, \theta)$  “reference distribution” because, for any  $F_1, \dots, F_M$  satisfying these criteria, we can construct a process  $\tilde{v}_{F_m}$ ,  $m = 1, \dots, M$  with the same distribution as  $\tilde{v}_Q$ .

## An illustrative example

- **Data:** a sample of  $n = 100$  observations generated from

$$p(\mathbf{x}) \propto (2\pi)^{-1} |\Sigma|^{-1/2} [1 + (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)]^{-3/2}, \quad (8)$$

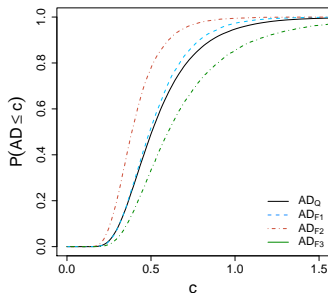
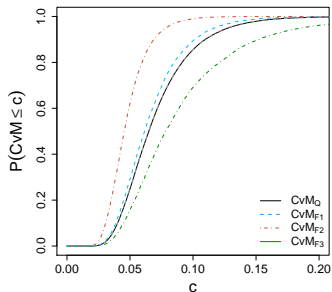
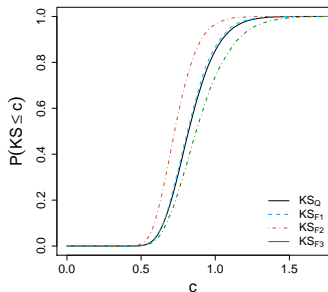
where  $\mu = (0, 3)^T$ ,  $\Sigma = \begin{bmatrix} 20 & 10 \\ 10 & 20 \end{bmatrix}$ ,  $\mathbf{x} \in [1, 20] \times [1, 25]$ .

- **Null models** we aim to test:

$$\begin{aligned} f_1(x; \beta) &\propto x_1^{(\beta_1-1)} x_2^{(\beta_2-1)} \exp\{-\beta_3(x_1 + x_2)\}, \\ f_2(x; \beta) &\propto \frac{\beta_3}{2\pi} [(x_1 - \beta_1)^2 + (x_2 - \beta_2)^2 + \beta_3]^{-3/2}, \\ f_3(x; \beta) &\propto e^{-\frac{1}{200} \left[ \left(\frac{x_1}{\beta_1} - 1\right)^2 + \left(\frac{x_2}{\beta_2} - 1\right)^2 - \beta_3 \left(\frac{x_1}{\beta_1} - 1\right) \left(\frac{x_2}{\beta_2} - 1\right) \right]}, \end{aligned} \quad (9)$$

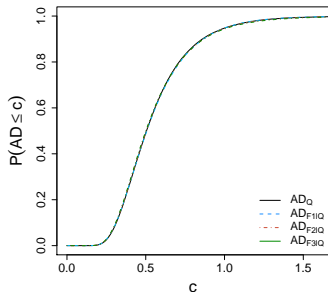
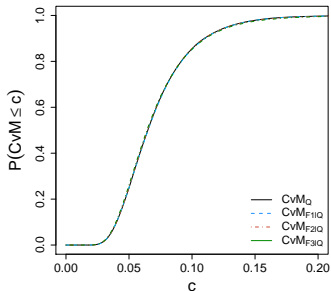
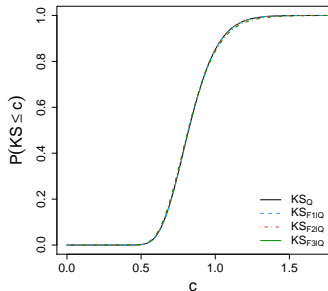
- **Reference distribution:**  $q(x, \theta) \propto e^{-\frac{1}{2\theta_3} [(x_1 - \theta_1)^2 + (x_2 - \theta_2)^2]}$ .

# Classical KS, CvM and AD: null distribution



Each simulation involves 100,000 bootstrap replicates, 100 observations, and the process is evaluated at 2000 grid points.

# Rotated KS, CvM and AD: null distribution



Each simulation involves 100,000 bootstrap replicates, 100 observations, and the process is evaluated at 2000 grid points.

## What about the power?

$H_0$	$\alpha = 0.05$					
	KS	CvM	AD	KS (K-2 rotated)	CvM (K-2 rotated)	AD
$Q$	.9331	.9817	.9382	-	-	-
$F_1$	.8623	.9529	.9092	.6971	1	1
$F_2$	.1078	.1019	.1237	.1336	.2422	.2541
$F_3$	.9528	.9820	.6356	.9153	.9746	.9470

Each simulation involves 100,000 bootstrap replicates, 100 observations, and the process is evaluated at 2000 grid points.

**Note:** We should NOT expect the K-2 rotated statistics to always dominate their classical counterparts or vice-versa!



## Key take-home ideas

When testing one continuous distribution which is multidimensional and/or depends on unknown parameters

- and we can simulate from it/evaluate it reasonably fast
  - $\Rightarrow$  we can reduce substantially the computational effort by simulating for the *projected empirical process*.
- but we cannot simulate from it/evaluate it reasonably fast
  - $\Rightarrow$  we can construct asymptotically distribution-free tests by means of the *K-2 transform*.

When testing  $M > 1$  continuous distributions,  $F_1, \dots, F_M$ , which are multidimensional and/or depend on unknown parameters

- $\Rightarrow$  we can avoid  $M$  different simulations and run just one by performing goodness-of-fit via *K-2 transform*.

# References

- **Main reference:** Algeri S. (2022). K-2 rotated goodness-of-fit for multivariate data. *Physical Review D*.
- Oosterhoff J., van Zwet W.R. A note on contiguity and Hellinger distance. *Springer New York*, 2012.
- Algeri S. and Khmaladze E.V. When Pearson  $\chi^2$  and other divisible statistics are not goodness-of-fit tests. *In preparation*.
- Khmaladze, E.V. (1980). The use of  $\omega^2$  tests for testing parametric hypotheses. *Theory of Probability & Its Applications*.
- Khmaladze, E.V. (2016). Unitary transformations, empirical processes and distribution free testing. *Bernoulli*, 2016.

**Thank you all for your time.**