

# Lecture 2:

- Central Limit Theorem
- Properties of Normal Distributions
- Trials and Tribulations!
- Regression to the Mean
- Correlations

As a consequence of the Central Limit Theorem, many **(but not all!)** physical processes often tend towards the Normal Distribution shape.

**However, few achieve this exactly !!**

Although calculated probabilities are often couched in terms of an ideal Normal Distribution to give a rough intuition of the scaling

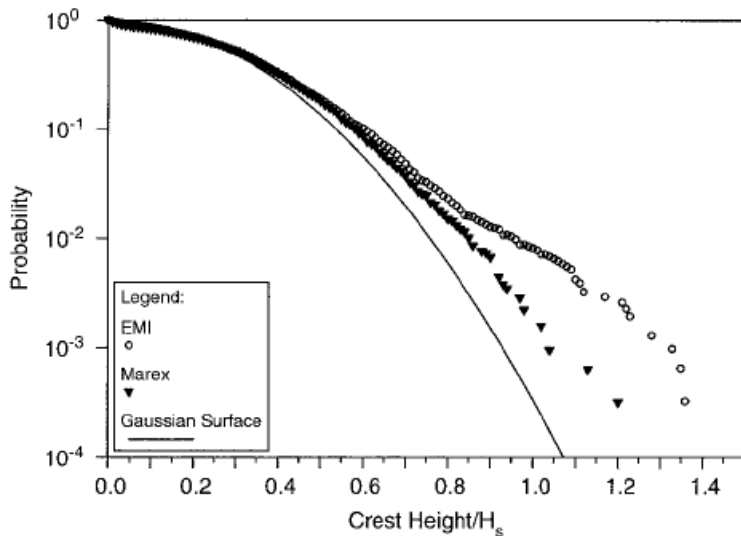
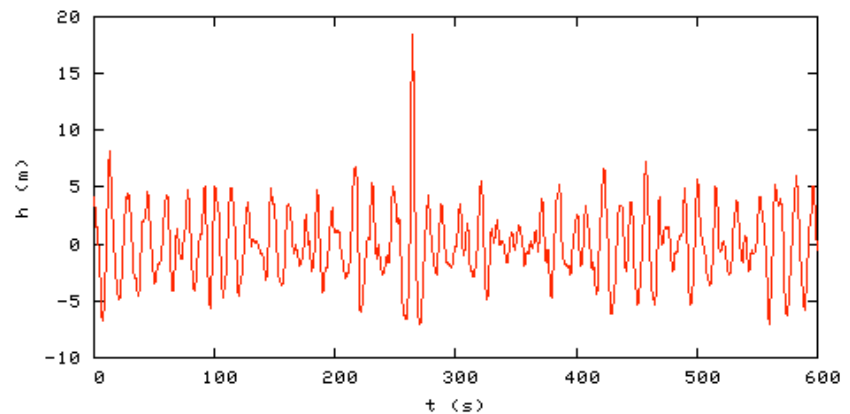
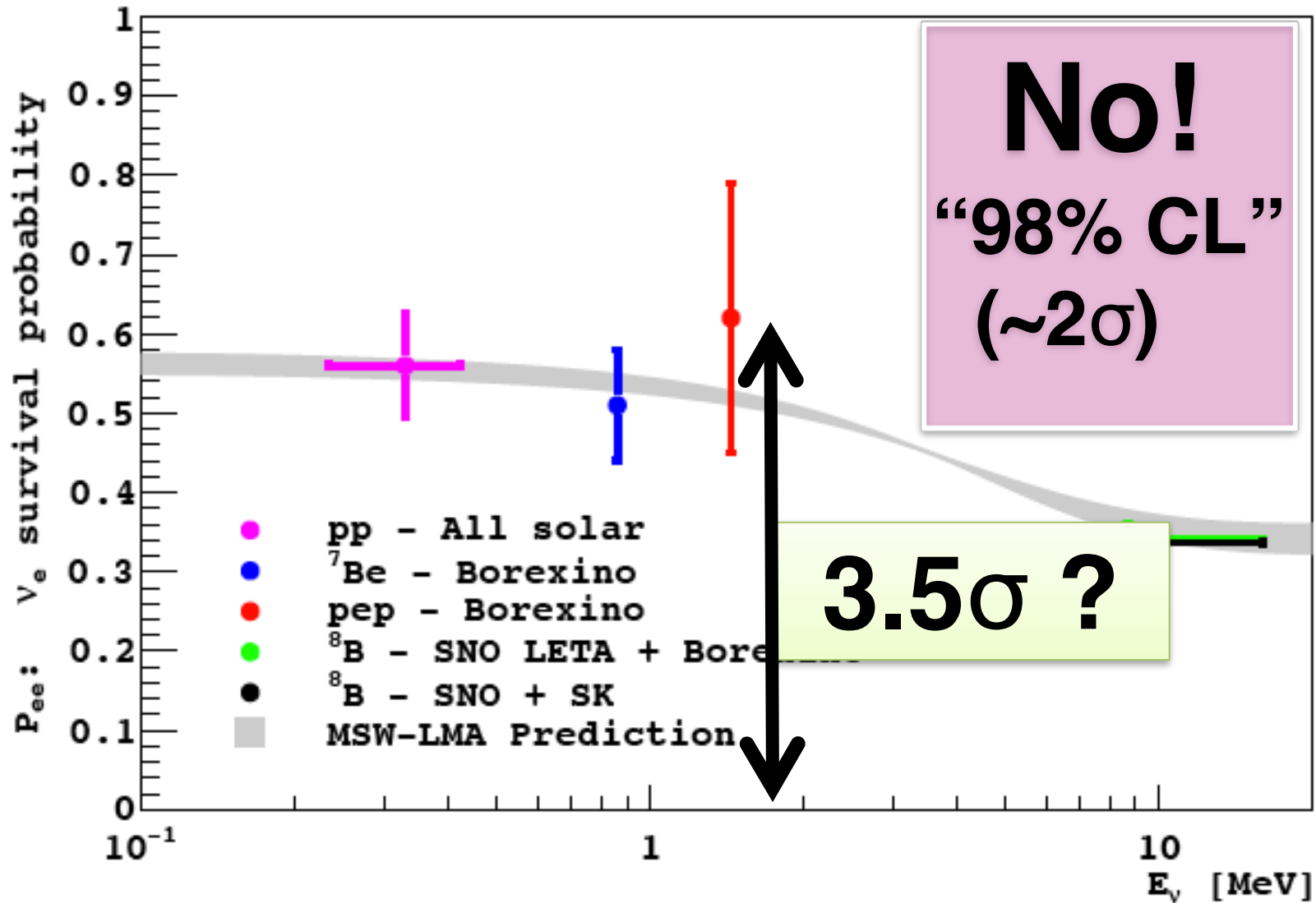


FIG. 4. Probability distribution of normalized crest heights measured at Tern during the storm on 4 Jan 1993. The crest heights are normalized by the significant wave height during each hour of the measurements. Nine hours of measurements with an average significant wave height of about 12 m were combined to produce the observed distribution.



The Draupner wave, a single giant wave measured on New Year's Day 1995, finally confirmed the existence of freak waves, which had previously been considered near-mythical

# Borexino "1 $\sigma$ error" on solar pep flux (2011)



## Example: Search for Episodic X-Ray Emission

Over the course of a year, 36000 x-rays are observed to come from a particular astrophysical object. However, on one particular day, 130 events are observed. What is the statistical significance of this observed burst?

$$\langle x \rangle = \frac{36000}{365} = 98.6 \quad \mu \simeq \langle x \rangle \quad \sigma = \sqrt{\mu}$$

$$s \simeq \frac{(130 - 98.6)}{\sqrt{98.6}} = 3.16\sigma$$

odds of getting at least this many events by a chance fluctuation from the average rate of emission

$$P = 8 \times 10^{-4}$$

**Is this sufficient to claim the observation of a burst from this object?**

## Correct question:

What is the chance of seeing at least one burst with an excess at least as large given the number of independent tests I've done ?

# Binomial !!

**N** Bernoulli trials where the chance of each success is **P**

$$\sum_{i=1}^{\infty} \binom{N}{i} P^i (1-P)^{N-i} = 1 - \binom{N}{0} P^0 (1-P)^{N-0}$$

$$P_{\text{post-trial}} = 1 - (1-P)^N \quad (\sim NP \text{ for } NP \ll 1)$$

$$P = 8 \times 10^{-4}, N = 365 \quad \rightarrow \quad P_{\text{post-trial}} = 25\%$$

How many timescales were considered? How many objects examined?

An appreciation of trials factors (“look elsewhere effect”) is hugely important... an improper handling of this can lead to incorrect conclusions and opens the door to biased analyses!

**This is not trivial !** A full accounting for this can be tricky:

- How many hypotheses have you actually tested?
- How many different ways have you tested each hypothesis?
- How many other things would have caught your eye?
- In general, how many ways have you looked at the data?

**At the same time, the data needs to be thoroughly checked to look for possible problems and confirm how well it’s understood**

**This is why physicists set the bar high in terms of significance level in order to claim a discovery**

**But it’s easy to get carried away...**

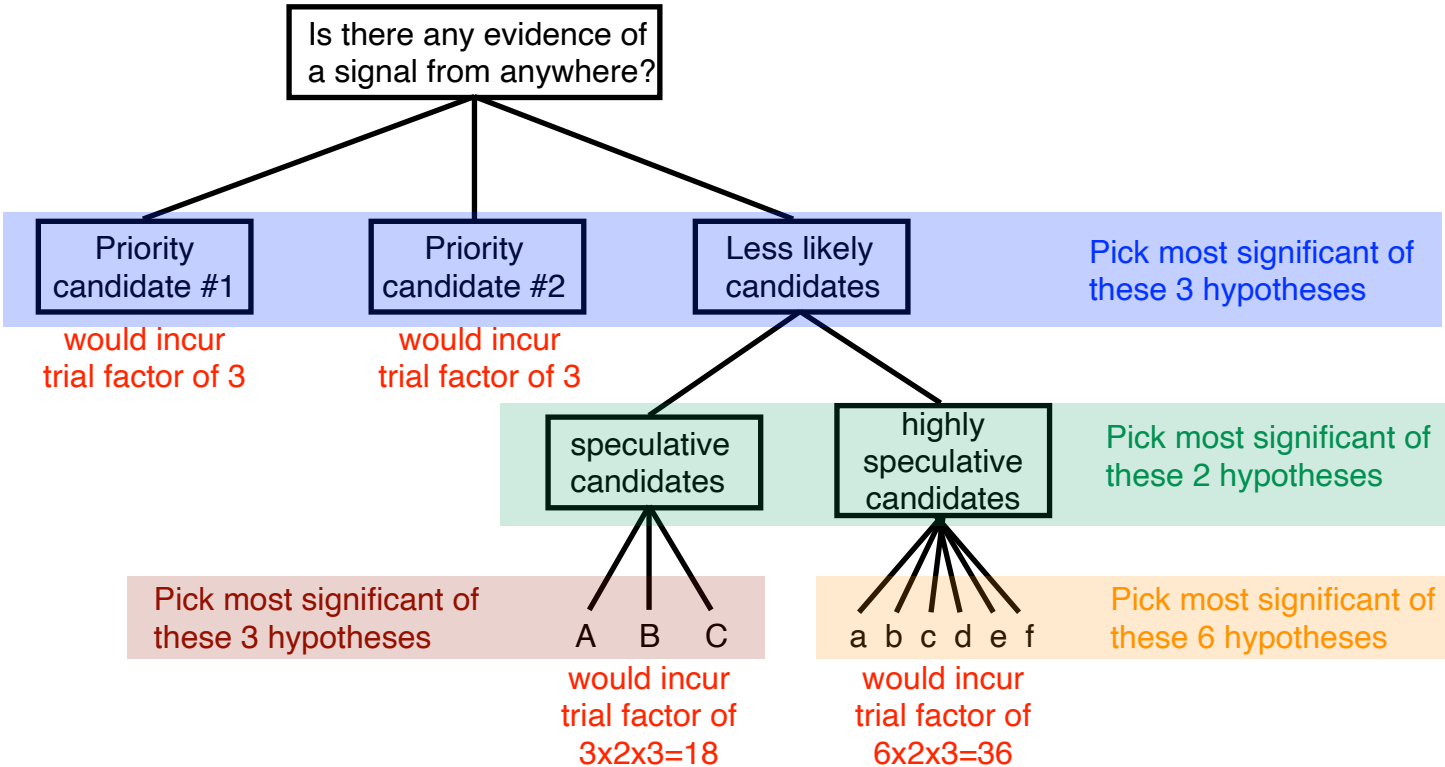


1) Trials factors apply to observations that would potentially lead to making a meaningful claim.

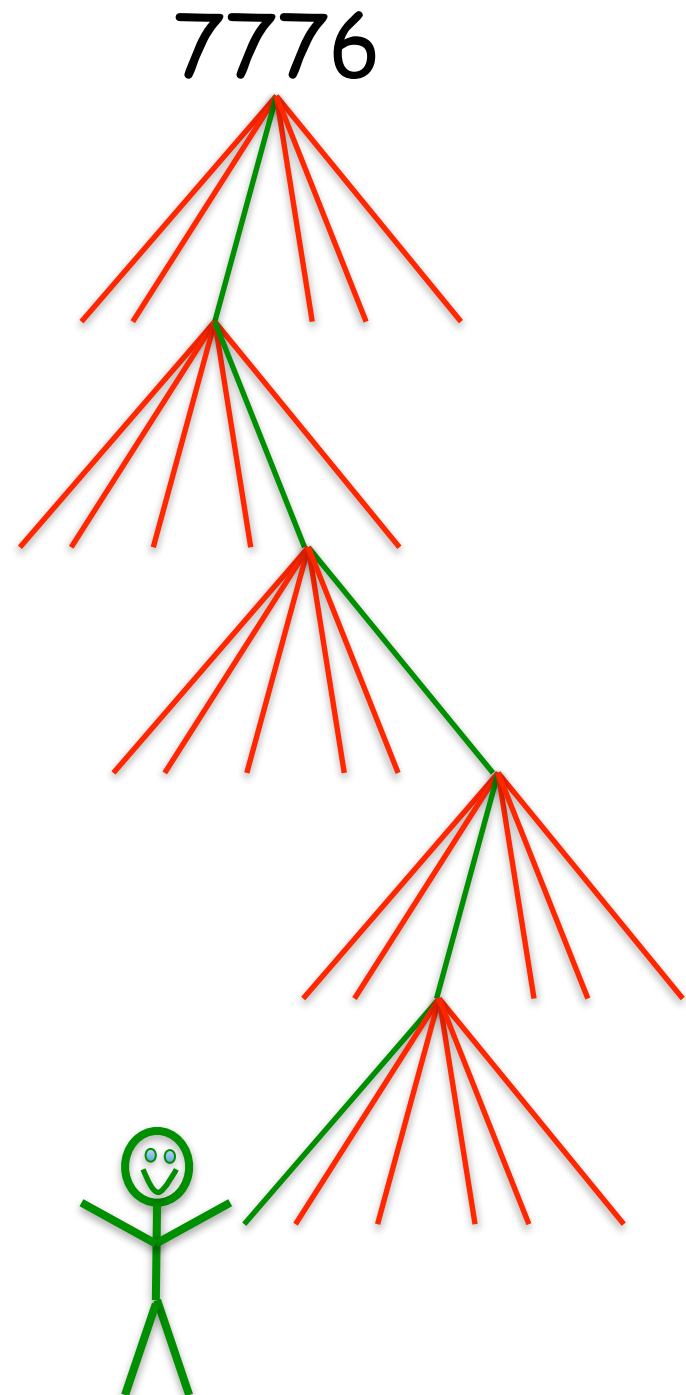
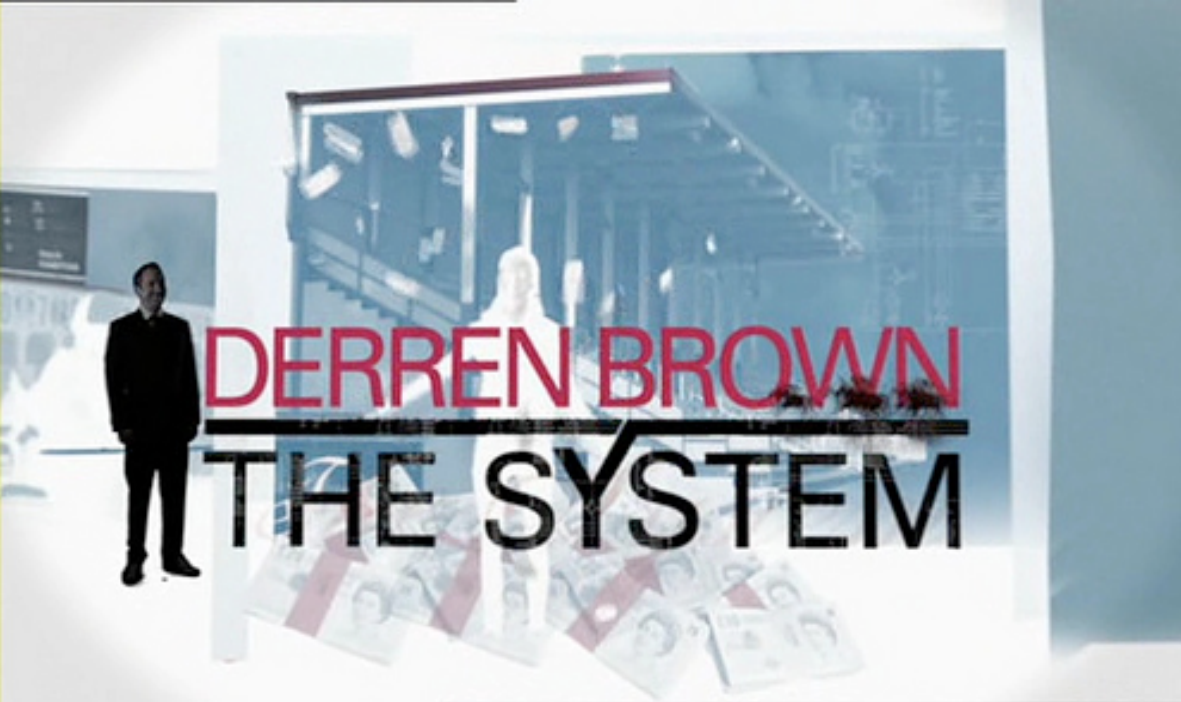
2) Verification based on applying the same analysis to an independent set of data is a good way to avoid misinterpretation of statistical fluctuations.

How do you deal with trial factors in the context of an open-ended search when an independent data set may not be available?

It's possible to structure trial factors based on an a priori ranking of hypothesis plausibility:







# “Regression to the Mean”

## Pop Quiz:

100 true/false questions on details of 17<sup>th</sup> century Swedish architecture.

100 true/false questions on details of 17<sup>th</sup> century Danish architecture.

**What an improvement! This particular group of students must know much more about Danish architecture!!**

**ONLINE EXAM:**

**17th Century Swedish Architecture**

1) Nicodemus Tessin the Elder designed the cathedral in Kalmar in 1690.  true  false

2) Adolf Fredriks kyrka replaced an old wooden chapel in central Stockholm.  true  false

3) The *Riddarhuset* was commissioned by Axel Oxentjema.  true  false

**ONLINE EXAM:**

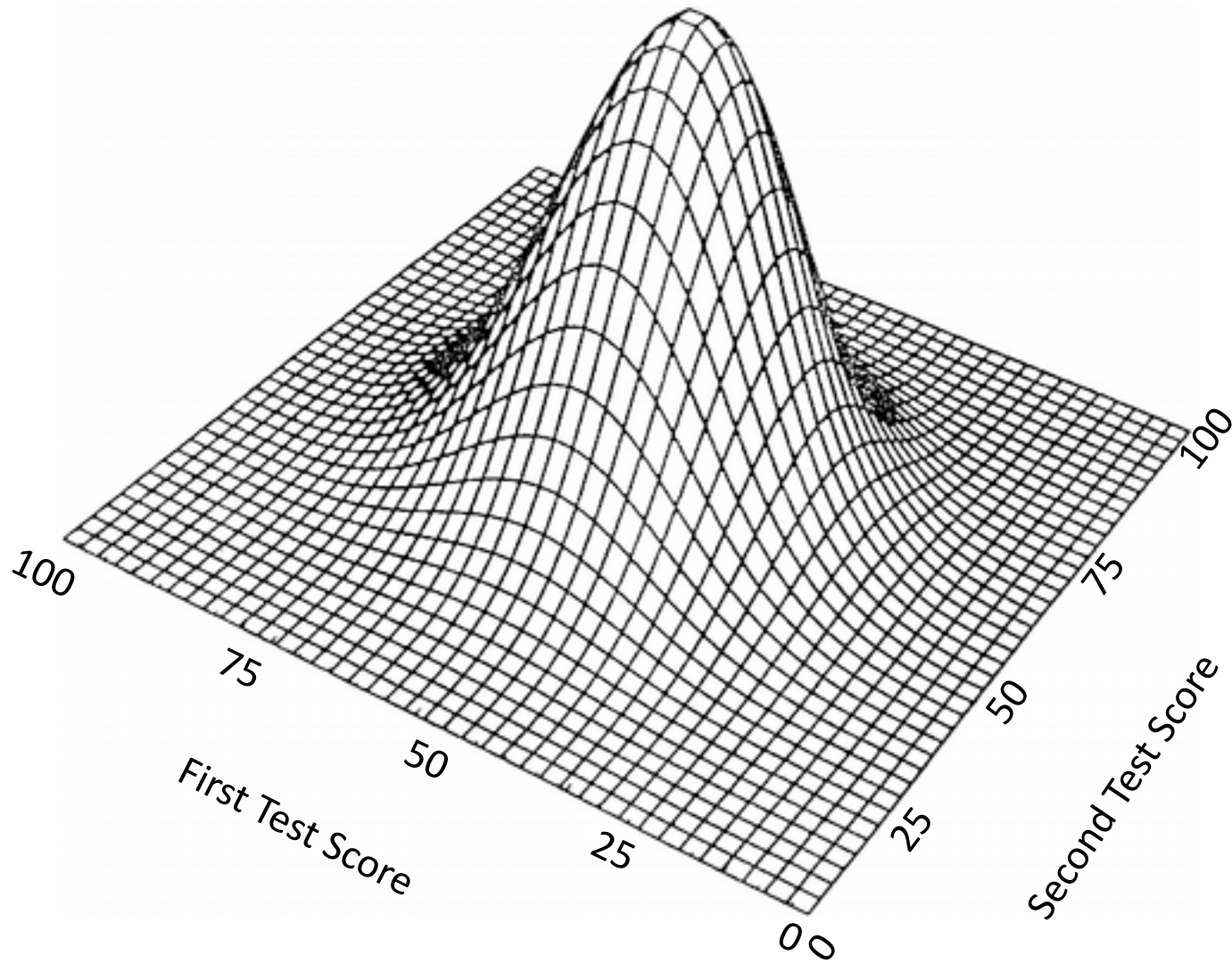
**17th Century Danish Architecture**

1) The Kunstforeningen building on Gammel Strand was built in 1690.  true  false

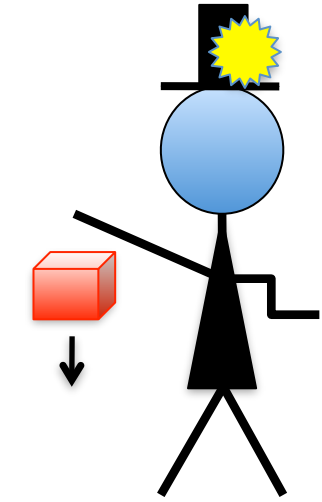
2) Knippelsbro was constructed to link Copenhagen with Christianshavn.  true  false

3) Timber-framed houses in Køge are typical of the area north of Copenhagen.  true  false

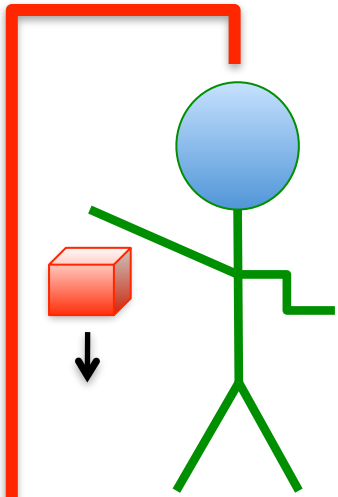
# Bi-variate Distribution with Identical Marginal Distributions (i.e. uncorrelated)



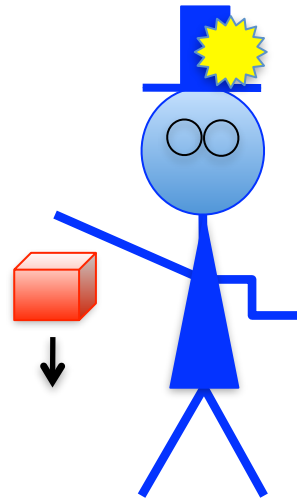
“The Effect of Hats on the Measurement of Gravity”



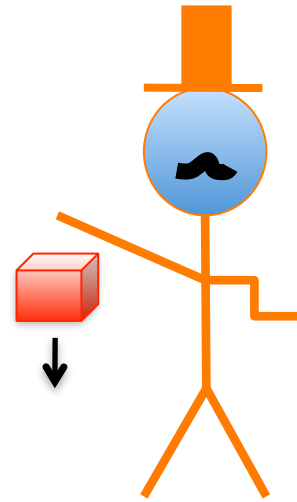
$g = 9.6$



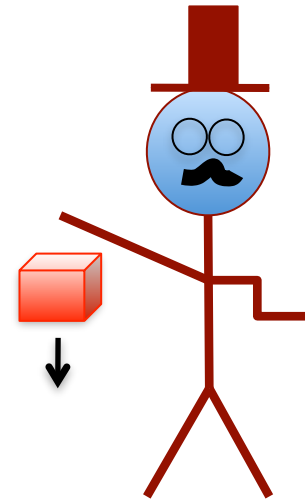
**9.3**



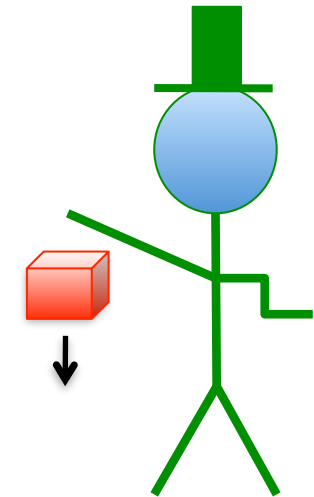
9.9



9.7



9.8

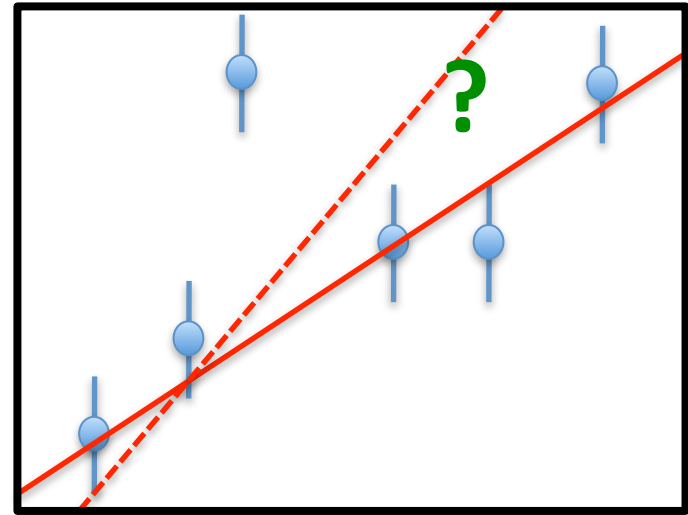


$g = 9.7$

**Much Better!!**

# So How Do You Handle Outliers?

No clear rules!



## Rules of Thumb:

- Look for possible systematic biases in the data;
- However, only reject outliers based on clear statistical/scientific criteria;
- Explicitly point out the issue and discuss the details;
- Be aware of any potential bias that could result and review the robustness of your final conclusions.



**The total number of known species is ~1.5 million**

**The number of known species that can fly is ~500,000**

$$\mathbf{P(\text{flying}) > } 5 \times 10^5 / 1.5 \times 10^6 = \mathbf{0.33}$$

**The number of plant species ~400,000**

$$\mathbf{P(\text{plant}) = } 4 \times 10^5 / 1.5 \times 10^6 = \mathbf{0.27}$$

**Thus, probability of finding a flying plant is**

$$\mathbf{P(\text{plant}) \times P(\text{flying}) = 0.089}$$

**And the expected number of flying plant species is**

$$\mathbf{(0.089)(1.5 \times 10^6) = 133,500}$$

# Correlations

And 12 points  
from Norway go to...  
**SWEDEN !!**



7		FINLAND	3		UNITED KINGDOM
6		ICELAND	2		POLAND
5		SPAIN	1		DENMARK
4		ROMANIA			



NORWAY

23 OF 37 COUNTRIES VOTING

Flying	500k	10k	400		
	<b>0.35</b>	<b>0.007</b>	<b><math>2.8 \times 10^{-4}</math></b>	<b>0</b>	<b>0</b>
Non-Flying	500k	54	6k	400k	10k
	<b>0.35</b>	<b><math>3.8 \times 10^{-5}</math></b>	<b>0.004</b>	<b>0.28</b>	<b>0.007</b>
	Insects	Birds	Mammals	Plants	Reptiles

**Joint  
PDF**

500k	54	6k	400k	10k
<b>0.54</b>	<b><math>5.9 \times 10^{-5}</math></b>	<b>0.006</b>	<b>0.44</b>	<b>0.011</b>
Insects	Birds	Mammals	Plants	Reptiles

**PDF for  
Non-Flying  
Species**

500k	54	6k	400k	10k
<b>0.70</b>	<b>0.00704</b>	<b>0.00428</b>	<b>0.28</b>	<b>0.007</b>
Insects	Birds	Mammals	Plants	Reptiles

**“Marginalised”  
PDF for All  
Species**



## Just to be clear:

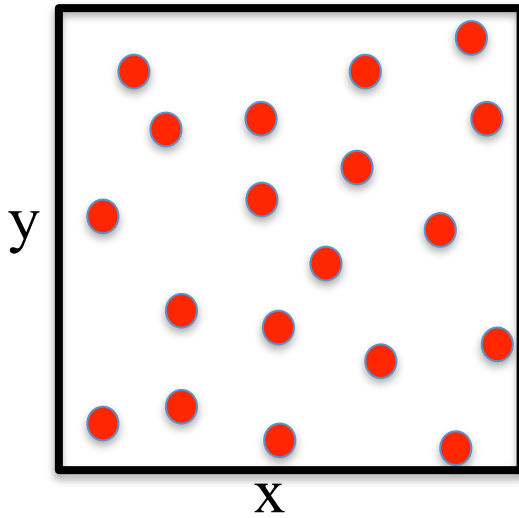
For example, if we have 2 dependent variables,  $x$  &  $y$ :

$$\int P(x, y) dx dy = 1$$

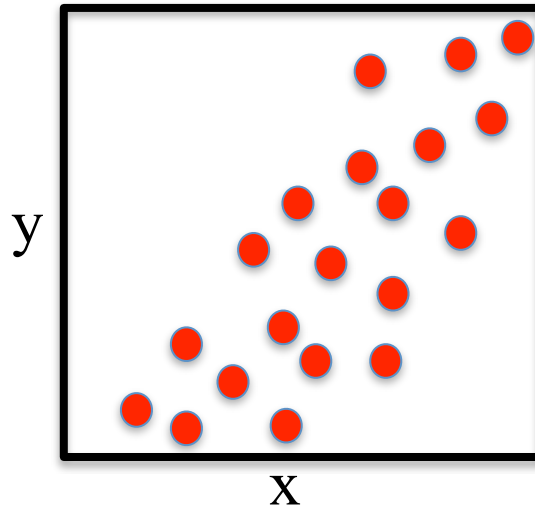
and

$$\langle f(x, y) \rangle = \int f(x, y) P(x, y) dx dy$$

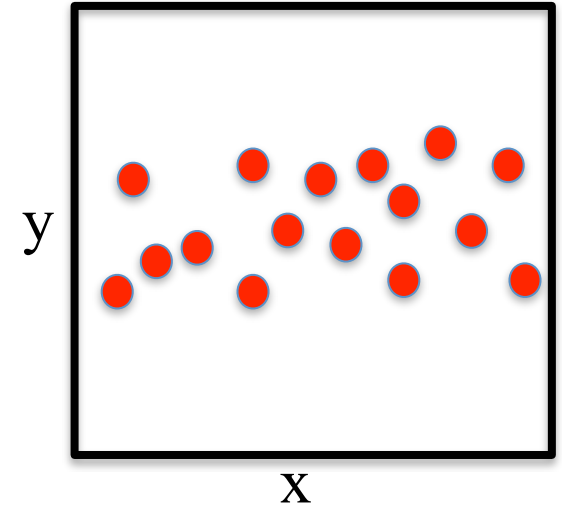
# Correlated or Uncorrelated?



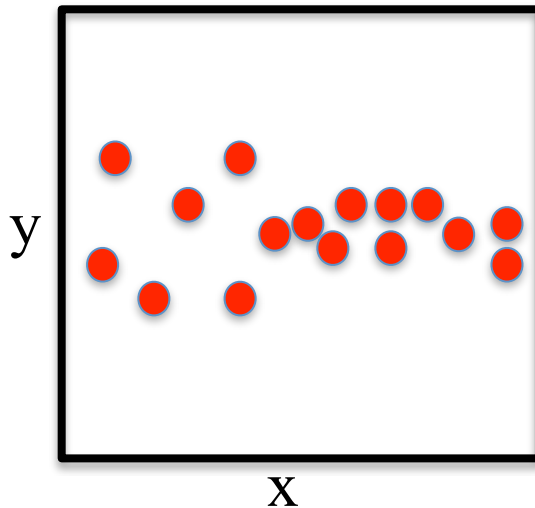
relatively uncorrelated



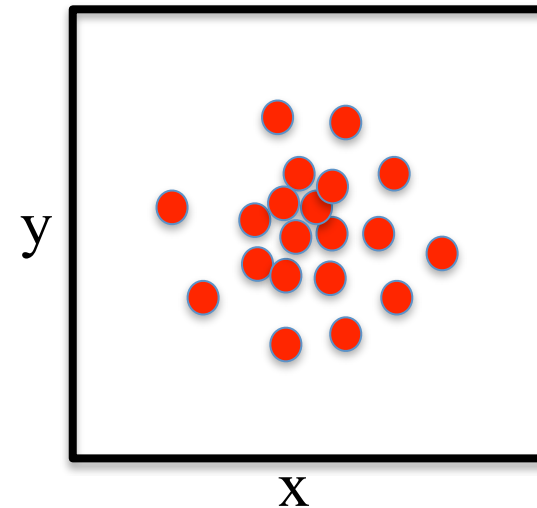
correlated



relatively uncorrelated



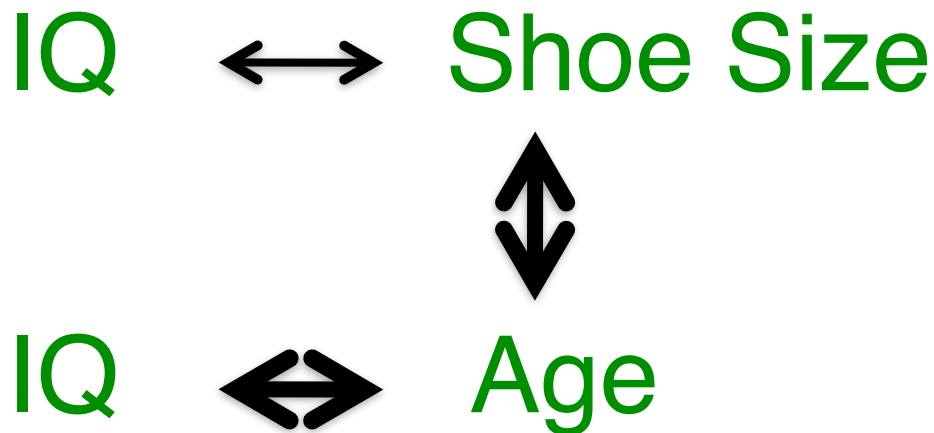
relatively uncorrelated mean,  
but correlated variance



relatively uncorrelated  
(symmetric with similar marginal distributions)

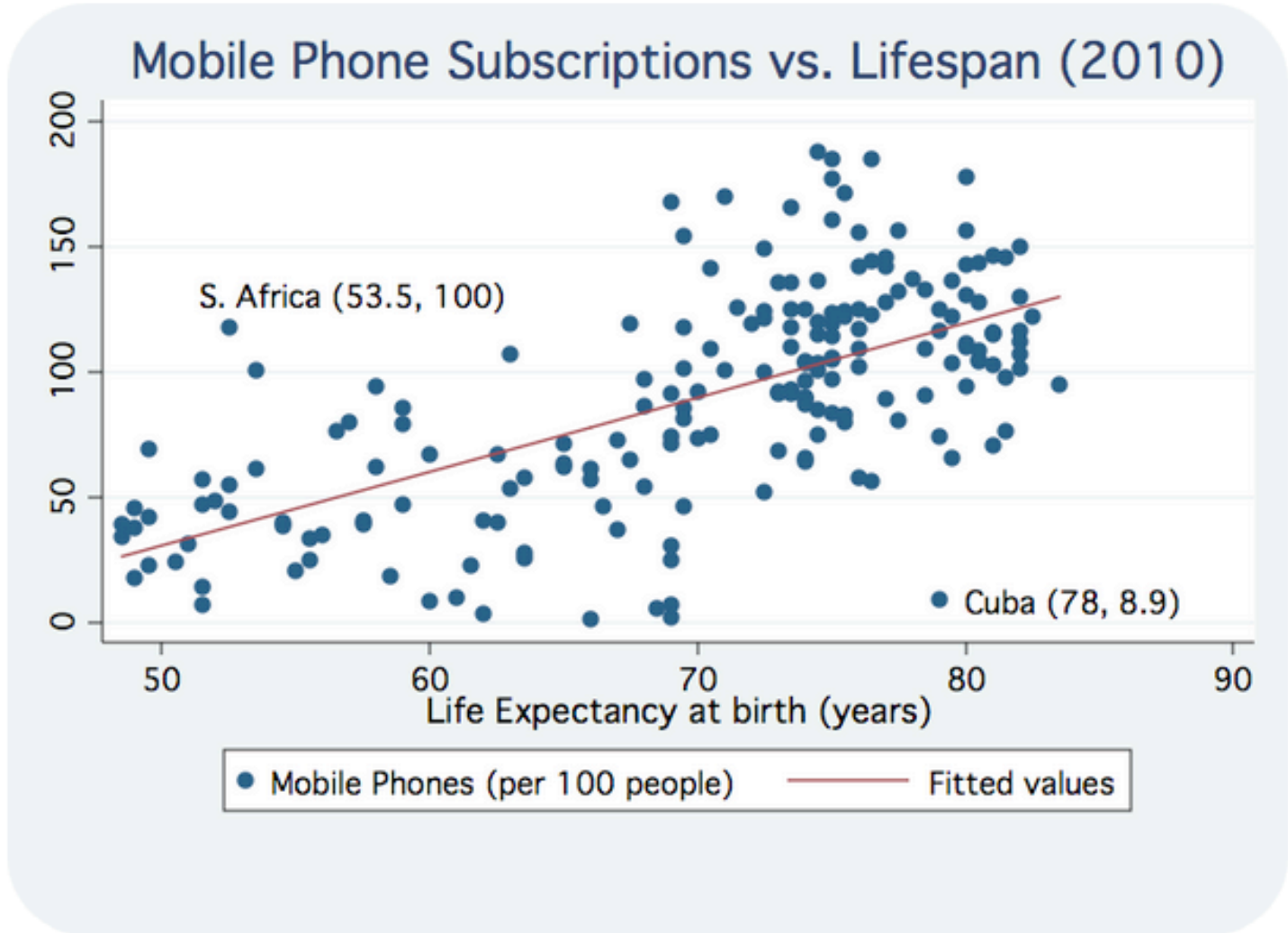


Beware of “hidden” correlations between ANY parameters that distinguish elements of your data set





# Beware of jumping to conclusion about cause and effect





# Beware of spurious correlations

