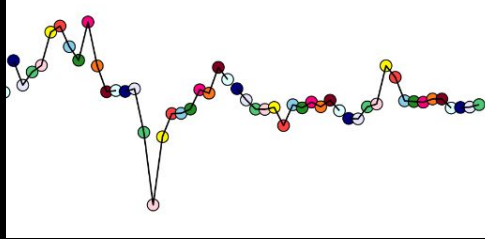


# Statistical Modelling - 3



Saket Choudhary

[saketc@iitb.ac.in](mailto:saketc@iitb.ac.in)

Introduction to Public Health Informatics

DH 302

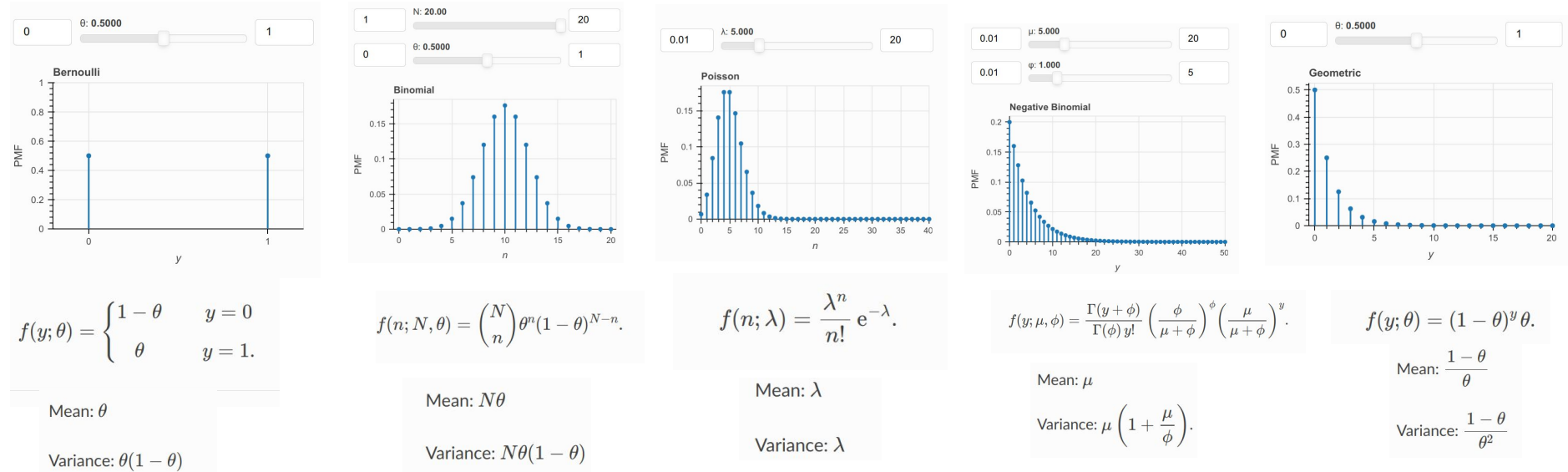
Lecture 04 || Friday, 17<sup>th</sup> January 2025

# Goals for today

- How to assess the observations for “extremeness” under a given model
- “Visualizing” likelihood, maximum likelihood and model selection
- How to decide what model best explains the observation

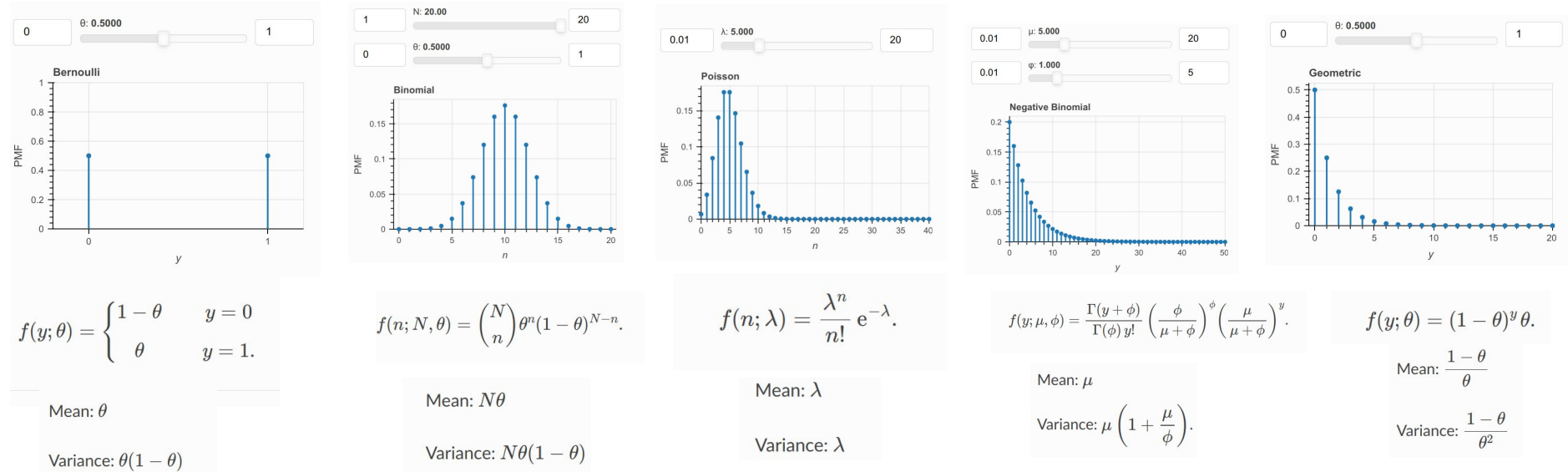
# R demo

# Previously



- What are the parameters of a distribution?
  - “Knobs” of a distribution - center and spread
- What is likelihood function?
  - Probability of seeing the observations under a model
- How are the maximum likely parameters learned?
  - Maximum likelihood estimation, we maximize log of likelihood function

# Previously



- What are the parameters of a distribution?
  - “Knobs” of a distribution - center and spread
- What is likelihood function?
  - Probability of seeing the observations under a model
- How are the maximum likely parameters learned?
  - Maximum likelihood estimation, we maximize log of likelihood function

# Compendium of Indian Data



## Compendium of Datasets and Registries in India, 2024

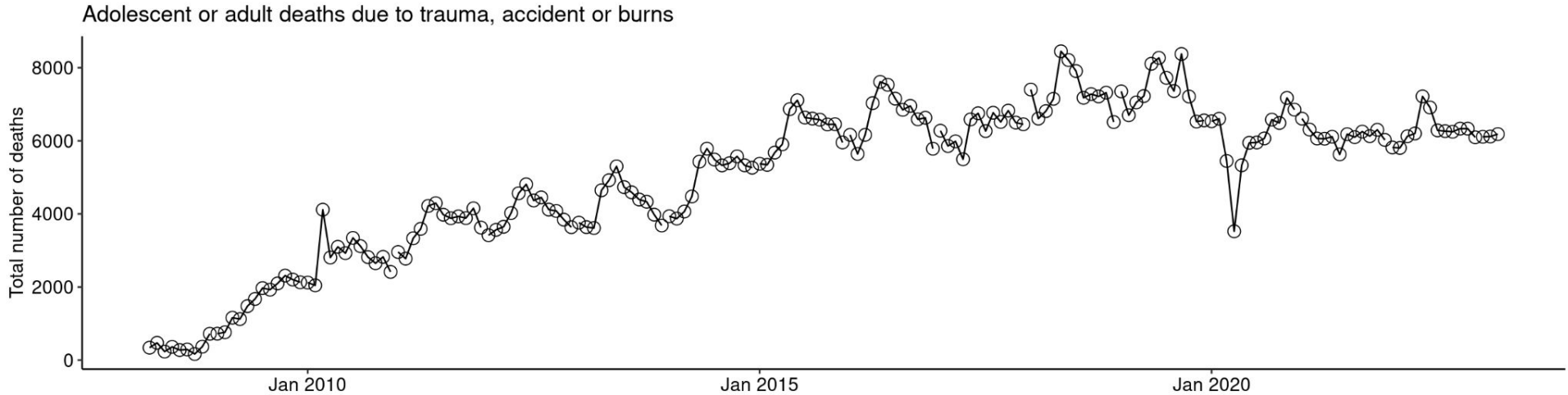
Data Informatics and Innovation Division  
Ministry of Statistics & Programme  
Implementation

### 4.9.5: Dataset: Health Management Information System\*

1	Name of dataset/ indicator	Health Management Information System
2.	Source Ministry/ Department/ Organization	Ministry of Health & Family Welfare
3.	Data/ indicators compiled are based on a <i>survey data, administrative data, multiple data sources, macro-aggregates</i> or any other method	Administrative Data
4.	Themes/ Categories under which data is collated	Health
5	List of <i>key</i> variables and their units of measurement (in case of datasets)	<p>Data related to Service Delivery and Infrastructure/ Human Resource is captured in HMIS on monthly basis. The data for Infrastructure/ Human Resource is also entered on monthly basis. Various modules for which data is captured in HMIS are as follows:</p> <p><b>Service Delivery:</b></p> <ul style="list-style-type: none"><li>• Maternal Health, Child-health &amp; Immunization, Family Planning,</li><li>• Vector Borne Disease, Tuberculosis, Morbidity and Mortality</li><li>• OPD, IPD Services, Surgeries etc.</li></ul> <p><b>Infrastructure:</b></p> <ul style="list-style-type: none"><li>• Manpower, Equipment</li><li>• Cleanliness, Building</li><li>• Availability of Medical Services such as Surgery etc.</li><li>• Super Specialties services such as Cardiology etc.</li><li>• Diagnostics</li><li>• Para Medical and Clinical Services etc.</li></ul>
6.	List of <i>key</i> variables used for computation and	

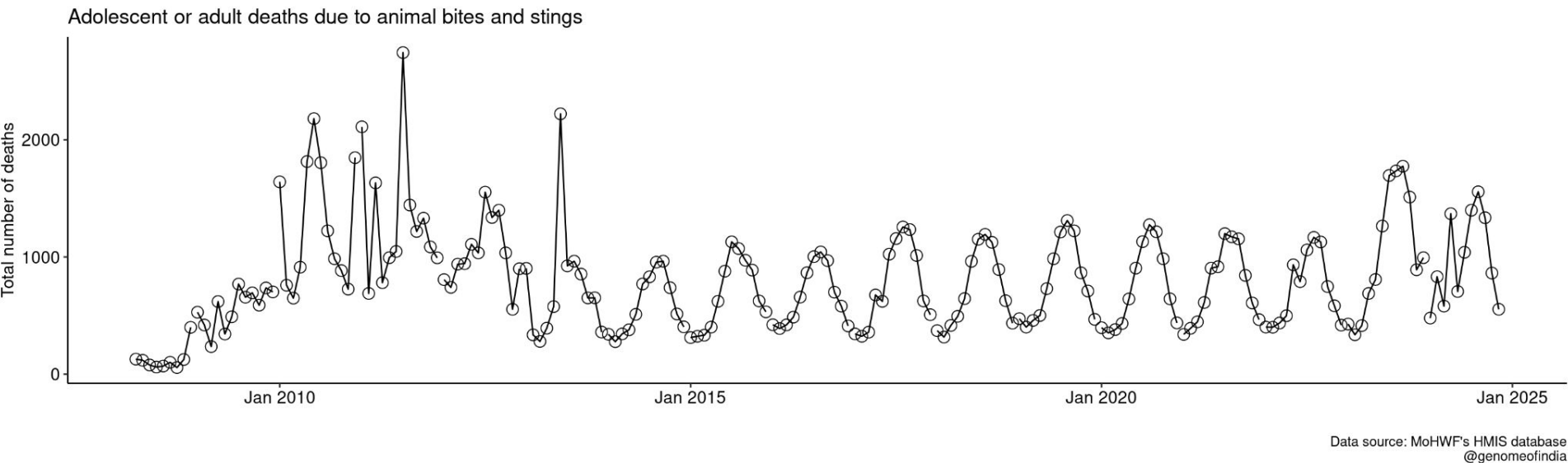
<https://hmis.mohfw.gov.in/#/>

# Question: What is going on this plot?



Data source: MoHWF's HMIS database  
@genomeofindia

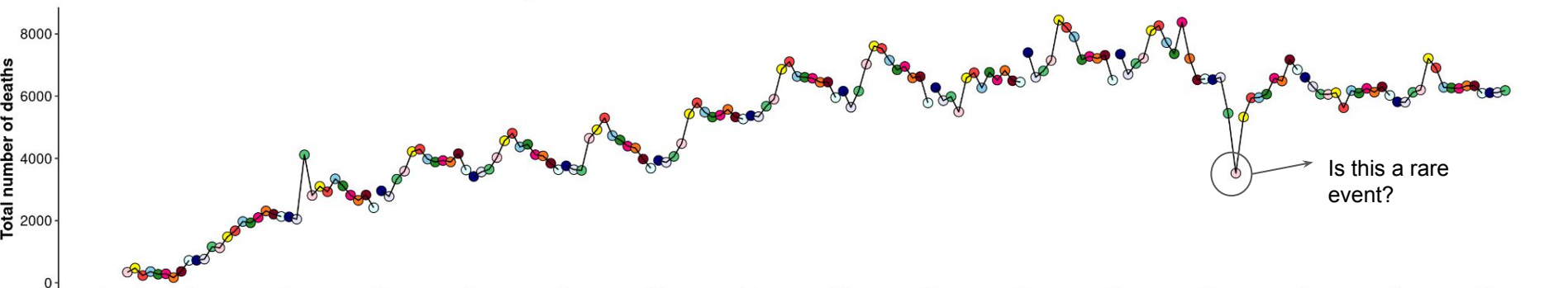
# Question: What is going on this plot?



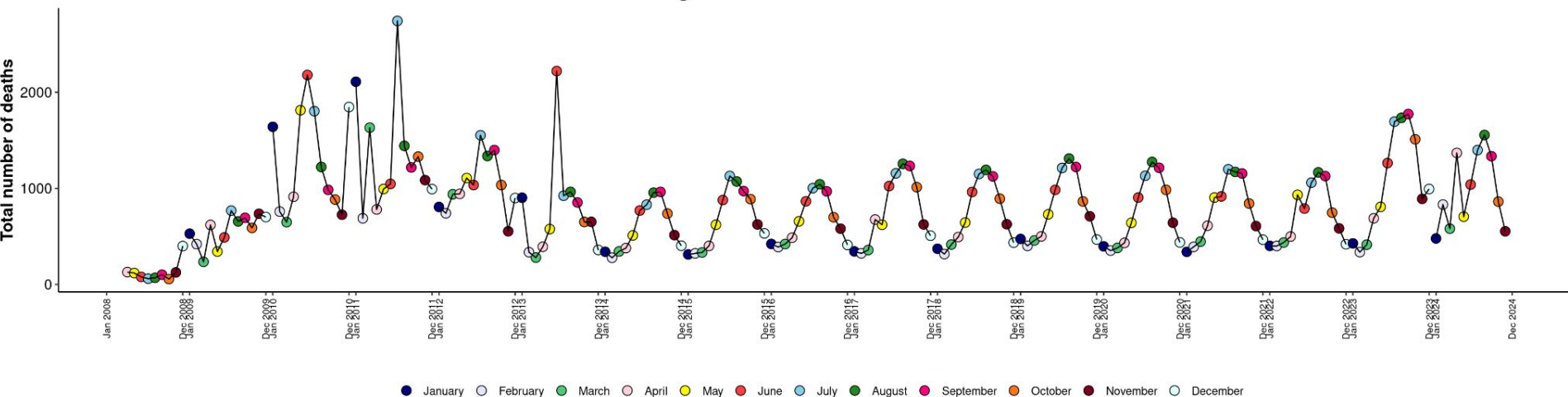


# Trauma and bite related deaths are seasonal

## Adolescent or adult deaths due to trauma, accident or burns



## Adolescent or adult deaths due to animal bites and stings



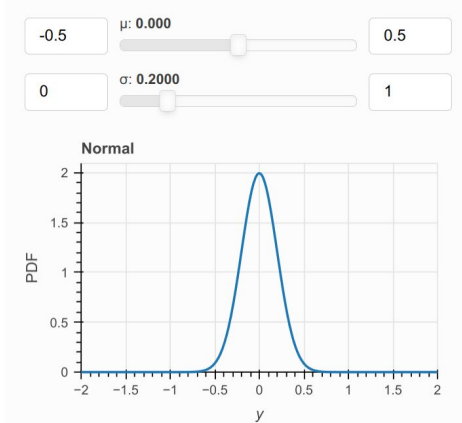
# Some potential questions we would like to answer w.r.t the plots

- Are certain points an “outlier”?
- Has the trend of deaths changed “significantly” with time?

To be able to answer this question, we first need to answer: What is the best model that presents our data?

Even before we answer this question, we should learn about some (important) continuous distributions

# Some (important) continuous distributions

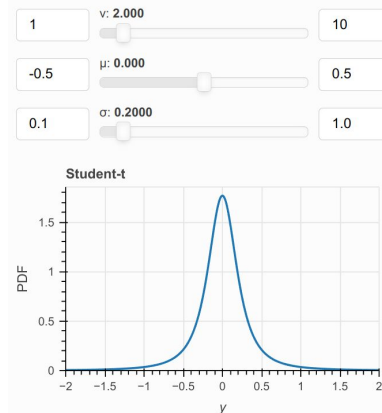


## Normal

$$f(y; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mu)^2/2\sigma^2}.$$

Mean:  $\mu$

Variance:  $\sigma^2$

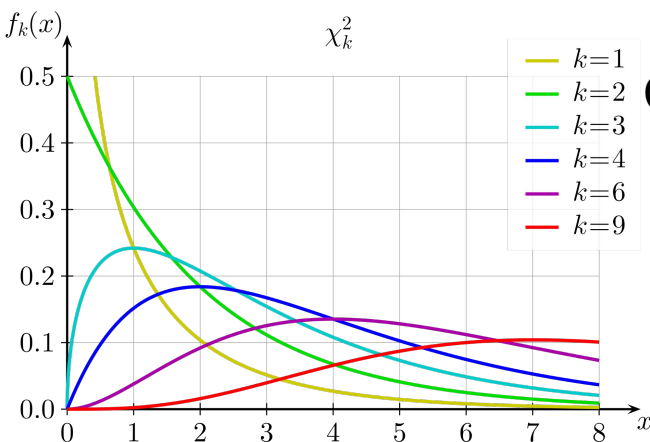


## Student's t

$$f(y; \nu, \mu, \sigma) = \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\pi\nu\sigma^2}} \left(1 + \frac{(y-\mu)^2}{\nu\sigma^2}\right)^{-\frac{\nu+1}{2}},$$

Mean:  $\mu$  for  $\nu > 1$ , otherwise undefined.

Variance:  $\frac{\nu}{\nu-2} \sigma^2$  for  $\nu > 2$ . If  $1 < \nu < 2$ , then the variance is infinite. If  $\nu \leq 1$ , the variance is undefined.

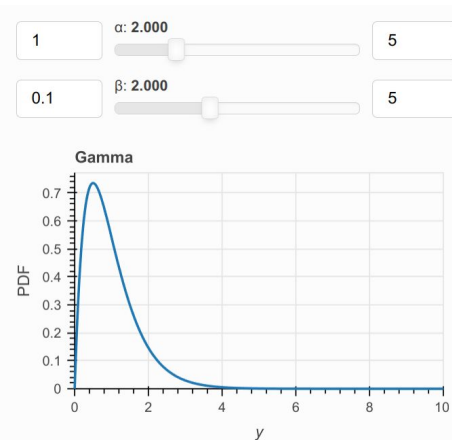


## Chi-square

$$\frac{1}{2^{k/2}\Gamma(k/2)} x^{k/2-1} e^{-x/2}$$

Mean = k

Variance = 2k



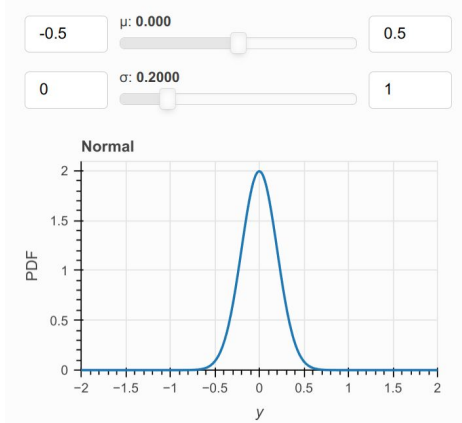
## Gamma

$$f(y; \alpha, \beta) = \frac{1}{\Gamma(\alpha)} \frac{(\beta y)^\alpha}{y} e^{-\beta y},$$

Mean:  $\frac{\alpha}{\beta}$

Variance:  $\frac{\alpha}{\beta^2}$

# Gaussian a.k.a Normal distribution



## Normal

$$f(y; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mu)^2/2\sigma^2}.$$

Mean:  $\mu$

Variance:  $\sigma^2$

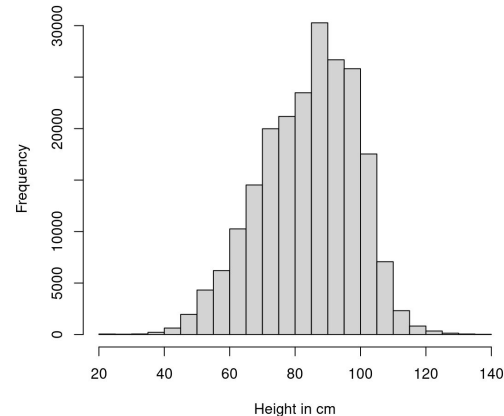
$$X_i \sim \mathcal{N}(\mu, \sigma^2)$$

Quantities that are sum of large number of subprocesses tend to be normally distributed

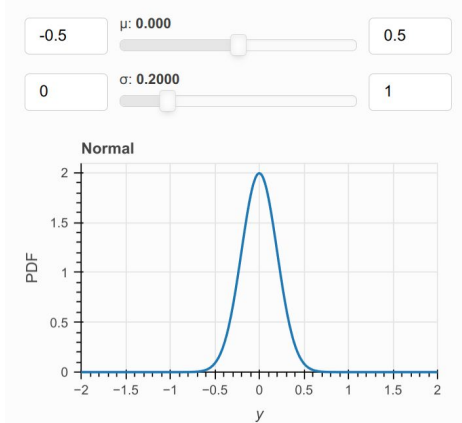
Example: Height/blood pressure distribution of a sample



Histogram of height distribution in a subset of NFHS 5



# What is normal about normal?



$$f(y; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mu)^2/2\sigma^2}.$$

Mean:  $\mu$

Variance:  $\sigma^2$

$$X_i \sim \mathcal{N}(\mu, \sigma^2)$$

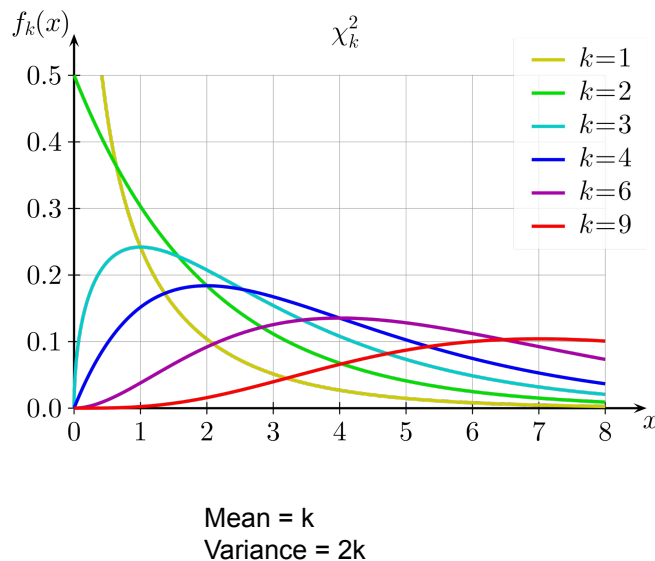
*"The literature gives conflicting evidence about the origin of the term Normal distribution". Karl Pearson (1920) claimed to have introduced it "many years ago", in order to avoid an old dispute over priority between Gauss and Legendre; but he gave no reference. Hilary Seal (1967) attributes it instead to Galton; but again fails to give a reference, so it would require a new historical study to decide this. However, the term had long been associated with the general topic: given a linear model  $y = X\beta + e$  where the vector  $y$  and the matrix  $X$  are known, the vector of parameters and the noise vector  $e$  unknown, Gauss (1823) called the system of equations which give the least squares parameter estimates, "the normal equations  $X'X\hat{\beta} = X'y$ , ellipsoid of constant probability density was called the "normal surface." It appears that somehow the name got transferred from the equations to the sampling distribution that leads to those equations"*

**Standard normal** has mean 0,  
variance 1

[Source](#)

# Chi-square distribution

Given  $Z_1, Z_2, \dots, Z_k$  are independent standard normal distribution, i.e.  $Z_i \sim \mathcal{N}(0, 1)$ , then the sum of their squares follows a  $\chi^2$  distribution with  $k$  degrees of freedom



$$X = \sum_{i=1}^k Z_i^2 \quad f(x; k) = \begin{cases} \frac{x^{k/2-1} e^{-x/2}}{2^{k/2} \Gamma\left(\frac{k}{2}\right)}, & x > 0; \\ 0, & \text{otherwise.} \end{cases}$$

$$X \sim \chi_k^2 \quad \Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt, \quad \Re(z) > 0.$$

**Example:** Estimate the parameters by curve fitting and check how “good” does it explain the observations

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

$O_i$  = an observed count for bin  $i$

$E_i$  = an expected count for bin  $i$

# Questions?

