



# Statistics in Data Analysis

*All you ever wanted to know about statistics but never dared to ask*

*part 3*

Pawel Brückman de Renstrom  
(pawel.bruckman@ifj.edu.pl)

March 19, 2025

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

## Question from the previous lecture

- 1 Suppose two independent measurements of the same quantity gave the following results:

$$x_1 \pm \sigma_1 \quad \text{and} \quad x_2 \pm \sigma_2$$

Take the weighted mean to be  $\bar{x} = wx_1 + (1 - w)x_2$ . Find the  $w$  which minimizes the error on the mean, hence provide expressions for the weighted mean  $\bar{x}$  and its variance  $\sigma_{\bar{x}}^2$ .

Solution to be sent to me before the next lecture

## Solution

We have to express the variance of the weighted mean

$$\bar{x} = wx_1 + (1 - w)x_2$$

using the recipe for error propagation:

$$\begin{aligned} \text{Var}(\bar{x}) &= \left( \frac{\partial \bar{x}}{\partial x_1} \right)^2 \sigma_1^2 + \left( \frac{\partial \bar{x}}{\partial x_2} \right)^2 \sigma_2^2 \\ &= w^2 \sigma_1^2 + (1 - w)^2 \sigma_2^2 \end{aligned}$$

and minimise it w.r.t. the weight  $w$ .

$$\begin{aligned} \frac{\partial \text{Var}(\bar{x})}{\partial w} &= 2w\sigma_1^2 - 2(1 - w)\sigma_2^2 = 0 \\ \implies w &= \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2} \end{aligned}$$

Hence we get:

$$\bar{x} = \frac{\sigma_2^2 x_1 + \sigma_1^2 x_2}{\sigma_1^2 + \sigma_2^2} \quad \text{and} \quad \text{Var}(\bar{x}) = \frac{\sigma_1^2 \sigma_2^2}{\sigma_1^2 + \sigma_2^2} \quad \therefore$$

# Boost transformation

NOT a unitary transformation!

$$V = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \quad (1)$$

$$A = \begin{pmatrix} \cosh \theta & \sinh \theta \\ \sinh \theta & \cosh \theta \end{pmatrix} \quad (2)$$

$$\sinh(x) = \frac{e^x - e^{-x}}{2} \quad (3)$$

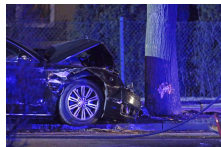
$$\cosh(x) = \frac{e^x + e^{-x}}{2} \quad (4)$$

**NOTE:** Correlation is introduced starting from uncorrelated variables!

# Accidents happen...

## Exponential distribution

- Imagine a fleet of governmental limousines circulating daily. For any of them there is a probability  $\lambda$  to be crashed in an accident in a day. We start with  $N_0$  limousines. What is the time p.d.f. of the accidents?
- For many circulating cars, accident rate is simply proportional to their number:



$$\frac{dN}{dt} = -\lambda N \quad \Rightarrow \quad \frac{dN}{N} = -\lambda dt \quad \Bigg/ \int$$

$$\ln N = -\lambda t + C \quad \Rightarrow \quad N(t) = N_0 e^{-\lambda t} \quad \Rightarrow \quad \frac{dN(t)}{dt} = -\lambda N_0 e^{-\lambda t} \quad (5)$$

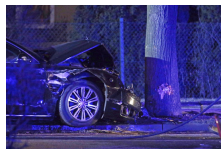
...so we observe an exponential decay of the fleet.

# Accidents happen...

## Exponential distribution

- Now consider just a single limousine of the PM. What is the time p.d.f. for its accident?

Let  $t_{1/2}$  (half-life) be the time of 50% survival probability:



$F_s(t_{1/2}) = (1 - \varepsilon)^n = 0.5$ ,  $n\delta = t_{1/2}$ ,  $k\delta = t$ ,  $\delta$  is an infinitesimal time interval.

$$n = \frac{\ln(0.5)}{\ln(1 - \varepsilon)} \simeq \frac{-\ln(0.5)(1 - \varepsilon)}{\varepsilon} \xrightarrow{\varepsilon \rightarrow 0} \frac{\ln(2)}{\varepsilon}$$

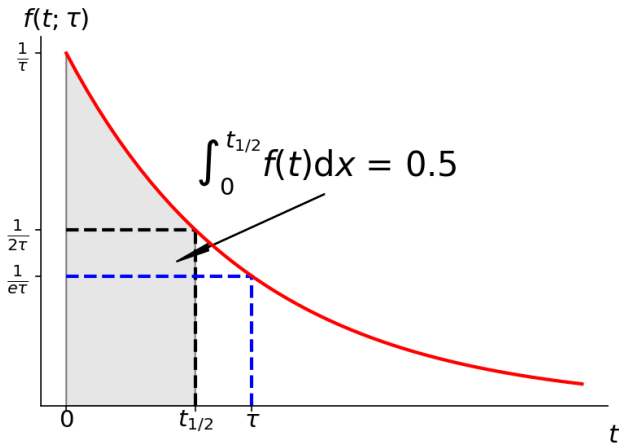
$$F_s(t) = (1 - \varepsilon)^k = (1 - \varepsilon)^{\frac{1}{\varepsilon} \frac{t}{t_{1/2}} \ln(2)} = \left| \lim_{\varepsilon \rightarrow 0} (1 - \varepsilon)^{\frac{\alpha}{\varepsilon}} = e^{-\alpha} \right| =$$
$$= e^{-\frac{t}{t_{1/2}} \ln(2)} \implies F_a(t) = 1 - e^{-\frac{t}{t_{1/2}} \ln(2)} \quad (6)$$

- $F_a$  is the cumulative accident probability. Hence the the p.d.f.:

$$f_a(t) = F'_a(t) = \frac{1}{\tau} e^{-\frac{t}{\tau}}, \quad \text{with } \tau = \frac{t_{1/2}}{\ln 2} \approx 1.44 t_{1/2} \quad (7)$$

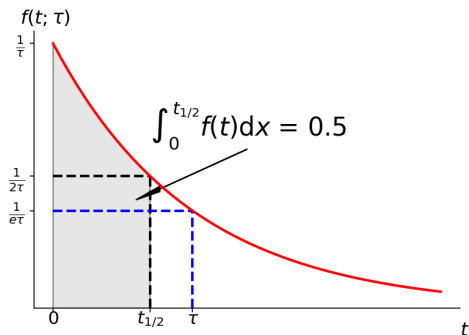
$$E[t] = \tau = \text{mean lifetime}, \quad V[t] = \tau^2. \quad \text{show these!} \quad (8)$$

# Exponential distribution



You are most likely to damage a brand new limousine!!!

# Exponential distribution



$$\begin{aligned} f_a(t|t_0) &= f_a(t)/F_s(t_0) = \\ \frac{1}{\tau} e^{-\frac{t}{\tau}} / e^{-\frac{t_0}{\tau}} &= \\ \frac{1}{\tau} e^{-\frac{t-t_0}{\tau}} &= f_a(t-t_0). \end{aligned}$$

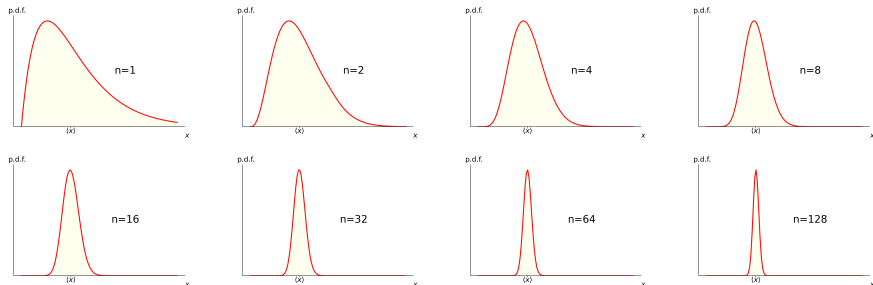
Do not be fooled! Probability of crashing a limo any day remains constant provided it has survived this far (conditional probability!).



# Mean of a random variable ensemble

## Central Limit Theorem

Imagine a measurement being a sum of many  $n$  independent ones, or an average of  $n$  random numbers drawn from an **arbitrary distribution** (sampling distribution).



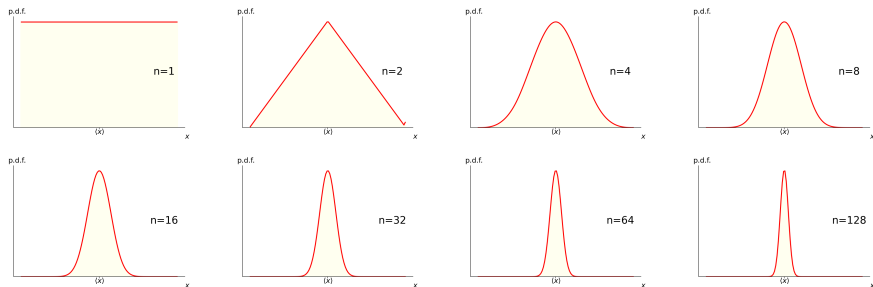
The mean  $\langle x \rangle$  converges on the initial distribution mean while the shape tends to a...

...**Gaussian** with ever decreasing width as  $n \nearrow$ .

# Mean of a random variable ensemble

## Central Limit Theorem

Ok, that was a well behaved distribution. Let's try something a bit less "gaussian" to start with:



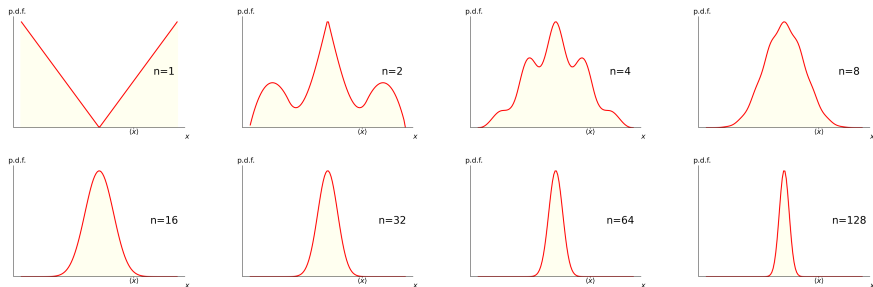
The mean  $\langle x \rangle$  converges on the initial distribution mean while the shape tends to a...

...**Gaussian** with ever decreasing width as  $n \nearrow$ .

# Mean of a random variable ensemble

## Central Limit Theorem

Ok, that was not austere enough. Let's try being bolder:



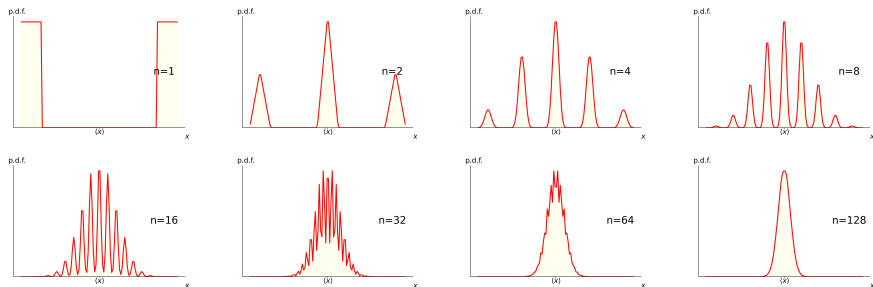
The mean  $\langle x \rangle$  converges on the initial distribution mean while the shape tends to a...

...**Gaussian** with ever decreasing width as  $n \nearrow$ .

# Mean of a random variable ensemble

## Central Limit Theorem

And again. Something manifestly non-Gaussian:



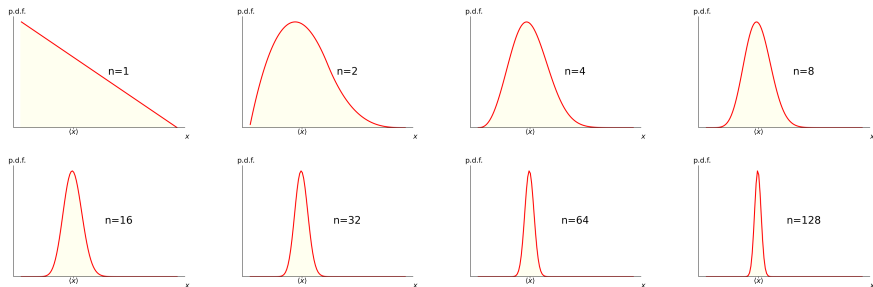
The mean  $\langle x \rangle$  converges on the initial distribution mean while the shape tends to a...

...**Gaussian** with ever decreasing width as  $n \nearrow$ .

# Mean of a random variable ensemble

## Central Limit Theorem

Finally, give up the symmetry:



The mean  $\langle x \rangle$  converges on the initial distribution mean while the shape tends to a...

...**Gaussian** with ever decreasing width as  $n \nearrow$ .

# Central Limit Theorem

Sum of  $n$  random variables drawn from a probability distribution function of a finite variance,  $\sigma^2$ , tends to be Gaussian distributed about the expectation value for the sum, with variance  $n\sigma^2$ .

Consequently, the mean of the same  $n$  random values will have the expectation value of the initial p.d.f. and variance  $\frac{1}{n}\sigma^2$ .

Ex: What is the probability that the mean salary of 50 randomly chosen employees of our institute exceeds 6000 pln?

**NOTE:** We don't need to know the actual distribution of salaries in the institute. All we need to know is the average and the variance (or standard dev.).

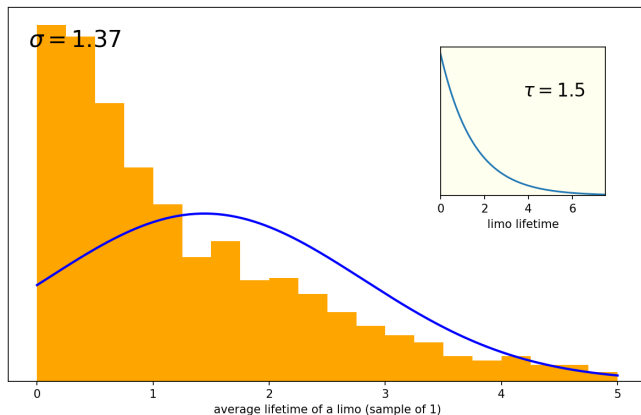
**Careful:** The *finite variance* is an important (and the only) requirement. A notable exception is the Cauchy (Breit-Wigner) distribution describing resonant states:

$$f(x) = \frac{1}{\pi} \frac{1}{1+x^2}$$

You can trivially show that the  $E[x^2]$  is divergent!

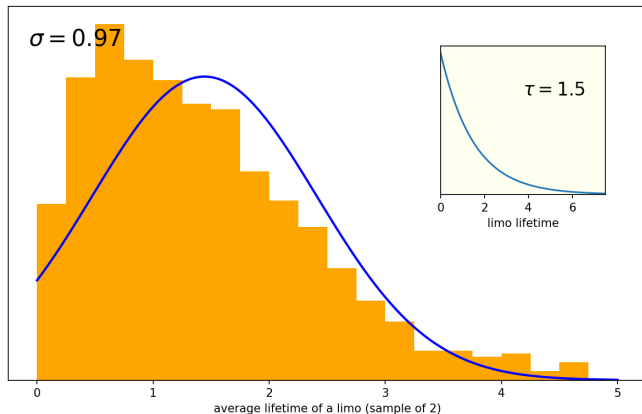
# Back to fleet of limousines...

a single limo



# Back to fleet of limousines...

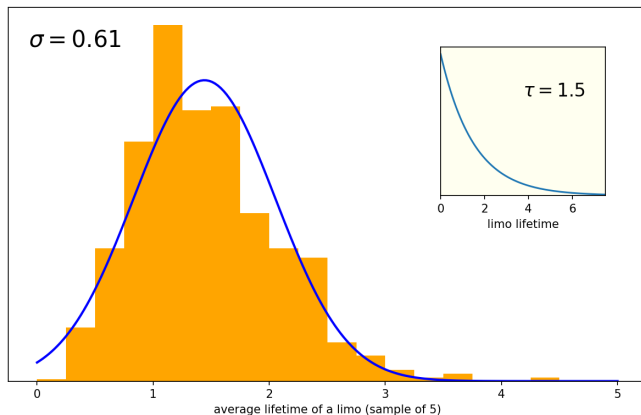
2 limo's





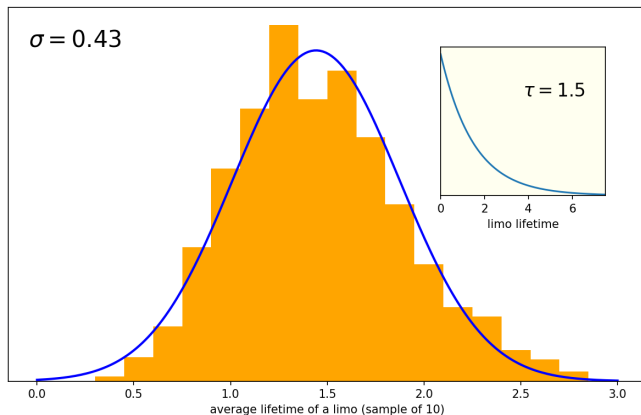
# Back to fleet of limousines...

5 limo's



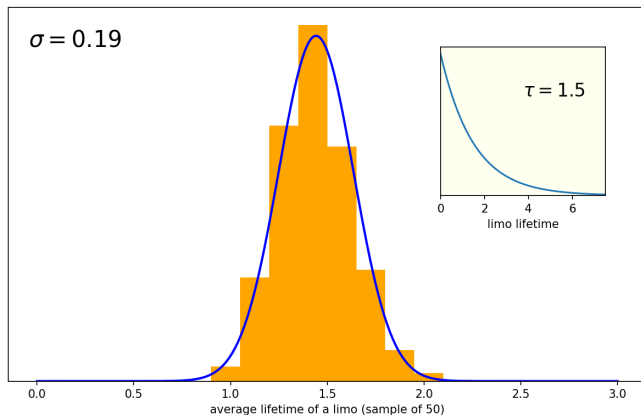
# Back to fleet of limousines...

10 limo's



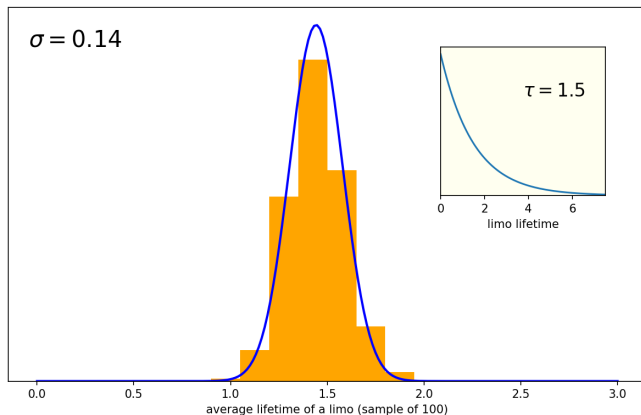
# Back to fleet of limousines...

50 limo's



# Back to fleet of limousines...

100 limo's



## Gaussian distribution

The **Gaussian** p.d.f. of the continuous random variable  $x$  with  $-\infty < x < \infty$  is defined by:

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right) \quad (9)$$

The term **normal** distribution is used when  $\mu = 0$  &  $\sigma = 1$ .

### Gaussian p.d.f.: normalisation, mean & variance

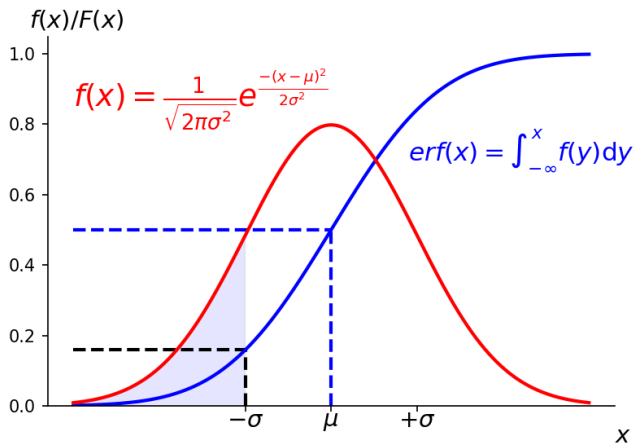
$$\int_{-\infty}^{\infty} f(x; \mu, \sigma^2) dx = 1 \quad (10)$$

$$E[x] = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right) dx = \mu, \quad (11)$$

$$V[x] = \int_{-\infty}^{\infty} (x - \mu)^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right) dx = \sigma^2. \quad (12)$$

Can you prove the above?

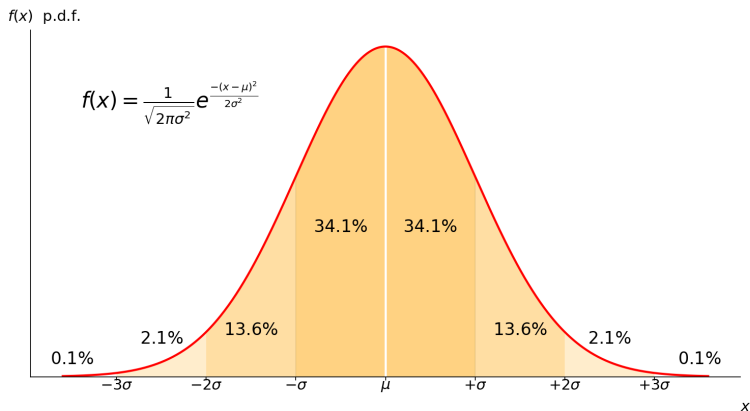
# Gaussian distribution



The cumulative distribution of the Gaussian p.d.f. is not analytically calculable. Nonetheless, quantiles of the normal distribution are of paramount importance for statistics!

# Gaussian distribution

## Quantiles



Standard deviation ( $\sigma$ ) of a Gaussian distribution has central importance for error analysis:

$$\mu \pm 1\sigma : 68.27\%, \quad \mu \pm 2\sigma : 95.45\%, \quad \mu \pm 3\sigma : 99.73\%.$$

# Characteristic function

## Fourier Transform of a p.d.f.: the **characteristic function**

$$\phi(k) = E[e^{ikx}] = \int_{-\infty}^{\infty} dx f(x)e^{ikx} \Rightarrow f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} dk \phi(k)e^{-ikx} \quad (13)$$

- $m$ 'th algebraic moment of  $f(x)$  is obtained by simple differentiation of  $\phi(k)$ :

$$\begin{aligned} (-i)^m \frac{d^m}{dk^m} \phi(k) \Big|_{k=0} &= (-i)^m \frac{d^m}{dk^m} \int_{-\infty}^{\infty} dx f(x)e^{ikx} \Big|_{k=0} = \\ &= (-i^2)^m \int_{-\infty}^{\infty} dx x^m f(x) = E[x^m] \end{aligned} \quad (14)$$

- Let  $z = \sum_i x_i$ , where  $x_1, \dots, x_n$  are  $n$  independent random variables:

$$\phi_z(k) = \int \dots \int e^{ik \sum_i x_i} f_1(x_1) \dots f_n(x_n) dx_1 \dots dx_n = \quad (15)$$

$$= \int e^{ikx_1} f_1(x_1) dx_1 \dots \int e^{ikx_n} f_n(x_n) dx_n = \phi_1(k) \dots \phi_n(k). \quad (16)$$



# Central Limit Theorem

Derivation of...

Let  $z = \frac{1}{\sqrt{n}}(x_1 + \dots + x_n) = \sum_{j=1}^n \frac{x_j}{\sqrt{n}}$ . For a single variable  $u \equiv x/\sqrt{n}$ , the characteristic function is given by:

$$\begin{aligned}\phi_u(k) &= \int_{-\infty}^{\infty} du f(u)e^{iku} = 1 + iE[u]k - \frac{1}{2}E[u^2]k^2 + O(k^3) = \\ &= 1 + iE[x]\frac{k}{\sqrt{n}} - \frac{1}{2}E[x^2]\frac{k^2}{n} + O\left(\frac{k^3}{\sqrt{n}}\right)\end{aligned}\quad (17)$$

Without any loss of generality, we can assume that  $E[x] = 0$  and  $E[x^2] = \sigma^2$  (otherwise use  $\bar{x} \equiv x - E[x]$ ):

$$\begin{aligned}\lim_{n \rightarrow \infty} \phi_z(k) &= \lim_{n \rightarrow \infty} \prod_{j=1}^n \phi_{u_j}(k) = \lim_{n \rightarrow \infty} \prod_{j=1}^n \left(1 - E[x^2]\frac{k^2}{2n} + O\left(\frac{k^3}{n^{3/2}}\right)\right) \simeq \\ &\simeq \lim_{n \rightarrow \infty} \left(1 - \frac{\sigma^2 k^2}{2n}\right)^n = e^{-\sigma^2 k^2/2}\end{aligned}\quad (18)$$

# Central Limit Theorem

... and the Gaussian distribution

So far we have found the characteristic function of the  $z$ . The p.d.f. is given by its inverse Fourier transform:

$$\begin{aligned} f_z(z) &= \frac{1}{2\pi} \int_{-\infty}^{\infty} dk \phi_z(k) e^{-ikz} = \frac{1}{2\pi} \int_{-\infty}^{\infty} dk e^{-\sigma^2 k^2/2} e^{-ikz} = \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} dk e^{-(\sigma k/\sqrt{2} + iz/(\sigma\sqrt{2}))^2 - z^2/(2\sigma^2)} = \frac{1}{\sqrt{2\pi}\sigma} e^{-z^2/(2\sigma^2)} \end{aligned} \quad (19)$$

## We have derived the **Central Limit Theorem**

The sum of independent random variables, sampled from the same distribution, will tend towards a **Gaussian** distribution, independently of the initial distribution.

**Note:** In the proof we used the strong assumption that all moments were finite. In fact, it is sufficient that the second moment ( $\sigma^2$ ) is finite, but we shall leave it without a proof. This holds for most well-behaved p.d.f.'s, but not all!

# Central Limit Theorem

## consequences

For the above derivation we used particularly normalised sum ( $z = \sum_{j=1}^n \frac{x_j}{\sqrt{n}}$ ) which led to the variance of the  $z$  being equal to the variance of  $x_i$ .

It is easy to see that:

- 1 For the algebraic sum  $z = \sum_{j=1}^n x_j$  we obtain  $\sigma_z = \sqrt{n}\sigma$ , or more generally  $\sigma_z^2 = \sum_{j=1}^n \sigma_j^2$ , ( $\langle z \rangle = \sum_{j=1}^n \langle x_j \rangle$ ).
- 2 For the algebraic mean  $z = \frac{1}{n} \sum_{j=1}^n x_j$  we obtain  $\sigma_z = \frac{1}{\sqrt{n}}\sigma$ , or more generally  $\sigma_z^2 = \frac{1}{n} \sum_{j=1}^n \sigma_j^2$ , ( $\langle z \rangle = \frac{1}{n} \sum_{j=1}^n \langle x_j \rangle$ ).

## What does it mean?

- If we estimate the mean from a sample, we will always tend towards the true mean,
- The uncertainty in our estimate of the mean will decrease as the sample gets bigger.

# Gaussian distribution

... generalisation

Let  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  be a  $n$ -dimensional sample space.

## $n$ -dimensional Gaussian distribution

$$f(\mathbf{x}; \mu, V) = \frac{1}{(2\pi)^{n/2} |V|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T V^{-1}(\mathbf{x} - \mu)\right) \quad (20)$$

$V$  is the covariance matrix of  $\mathbf{x}$  and  $V^{-1}$  is its inverse, called the *weight* matrix.  
 $|V|$  is the determinant of  $V$ .

What does the above give for independent random variables?

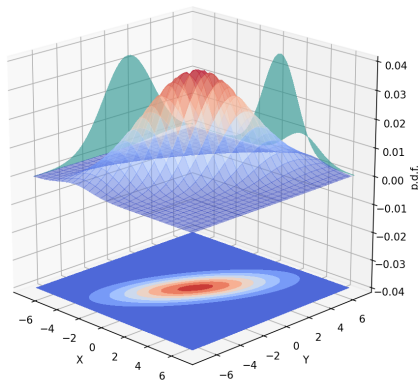
# Gaussian distribution

... 2D case

- $\sigma_1 = 2$
- $\sigma_2 = 3$
- $\rho = 0.7$

$$V = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

$$V^{-1} = \frac{1}{(1 - \rho^2)} \begin{pmatrix} \frac{1}{\sigma_1^2} & \frac{-\rho}{\sigma_1\sigma_2} \\ \frac{-\rho}{\sigma_1\sigma_2} & \frac{1}{\sigma_2^2} \end{pmatrix}$$



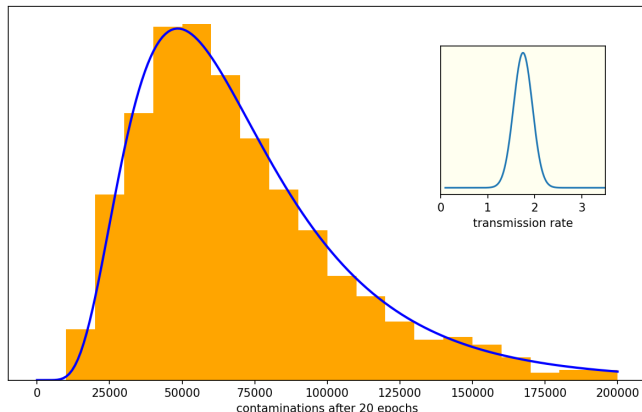
$$f(x_1, x_2; \mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[ \left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1}\right) \left(\frac{x_2 - \mu_2}{\sigma_2}\right) \right]\right) \quad (21)$$

# Spread of a pandemic

multiplicative Gaussian

Average transmission rate: 1.75 with standard deviation of 0.2.

Number of infected after 20 epochs:

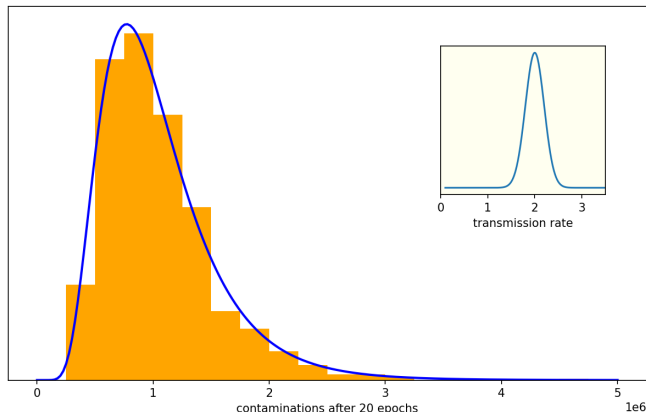


# Spread of a pandemic

multiplicative Gaussian

Average transmission rate: 2.0 with standard deviation of 0.2.

Number of infected after 20 epochs:

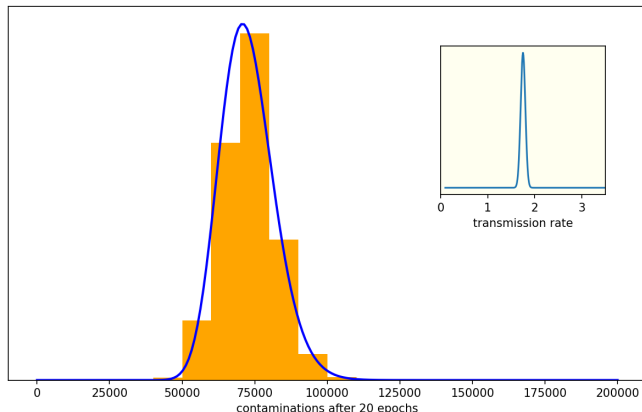


# Spread of a pandemic

multiplicative Gaussian

Average transmission rate: 1.75 with standard deviation of 0.05.

Number of infected after 20 epochs:





## Log-normal distribution

Let  $y$  be a Gaussian-distributed random variable with mean and variance  $\mu, \sigma^2$ .  
What is the p.d.f. of  $x = e^y$ ?

$$g(x) = f(y(x); \mu, \sigma^2) \left| \frac{dy}{dx} \right| = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(\ln x - \mu)^2}{2\sigma^2}\right) \frac{d(\ln x)}{dx} \quad (22)$$

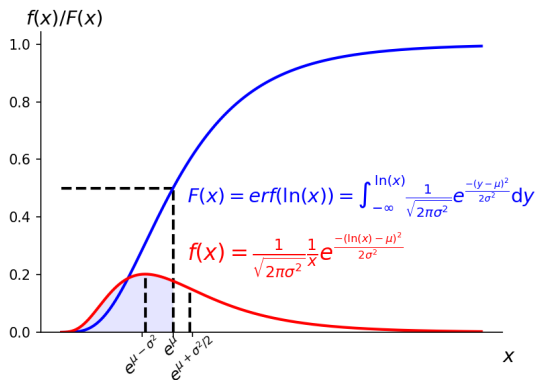
log-normal p.d.f.

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{x} \exp\left(\frac{-(\ln x - \mu)^2}{2\sigma^2}\right) \quad (23)$$

$$E[x] = e^{\mu + \frac{1}{2}\sigma^2} \quad (24)$$

$$V[x] = e^{2\mu + \sigma^2} \left[ e^{\sigma^2} - 1 \right] \quad (25)$$

# Log-normal distribution



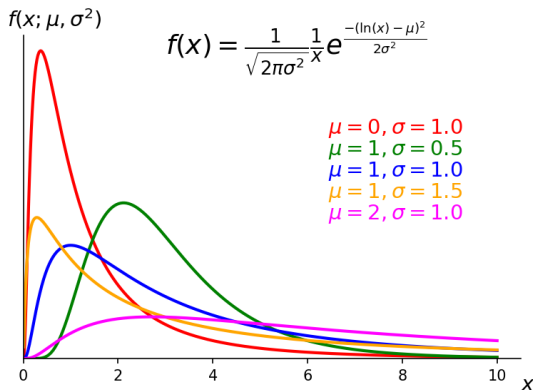
$$\int_0^X \frac{1}{x} e^{-\frac{(\ln(x)-\mu)^2}{2\sigma^2}} dx = \left| \ln(x) = y, \frac{1}{x} dx = dy \right| = \int_{-\infty}^{\ln(X)} e^{-\frac{(y-\mu)^2}{2\sigma^2}} dy = \sqrt{2\pi\sigma^2} \text{erf}(\ln(X))$$

$$\int_0^\infty x \frac{1}{x} e^{-\frac{(\ln(x)-\mu)^2}{2\sigma^2}} dx = \int_{-\infty}^\infty e^{-\frac{(y-\mu)^2}{2\sigma^2}} e^y dy = \int_{-\infty}^\infty e^{-\frac{(y-(\mu+\sigma^2/2))^2}{2\sigma^2}} e^{\mu+\frac{1}{2}\sigma^2} dy = \sqrt{2\pi\sigma^2} e^{\mu+\frac{1}{2}\sigma^2}$$

mode:  $e^{\mu-\sigma^2}$ ,    median:  $e^\mu$ ,    mean:  $e^{\mu+\frac{1}{2}\sigma^2}$ ,     $F(X) = \text{erf}(\ln(X))$

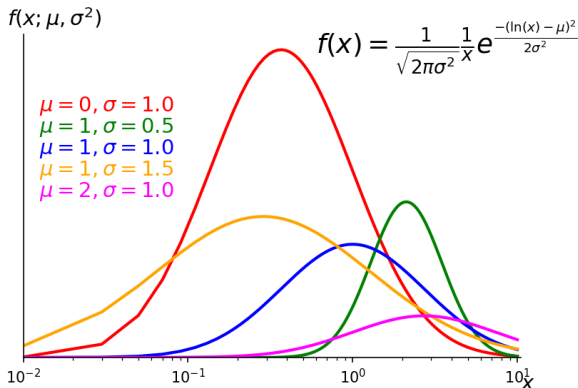
# Log-normal distribution

multiplicative factors



It becomes apparent that if  $z = \prod_{j=1}^n x_j = e^{\sum_{j=1}^n y_j}$ , the product of many random variables tends to a log-normal distribution with  $\mu = \sum_{j=1}^n \mu_j$  and  $\sigma^2 = \sum_{j=1}^n \sigma_j^2$ . Here,  $\mu_j = E[\ln x]$  and  $\sigma_j^2 = E[\ln^2 x] - E[\ln x]^2$ . Certainly,  $\forall_j x_j > 0$ .

# Log-normal distribution

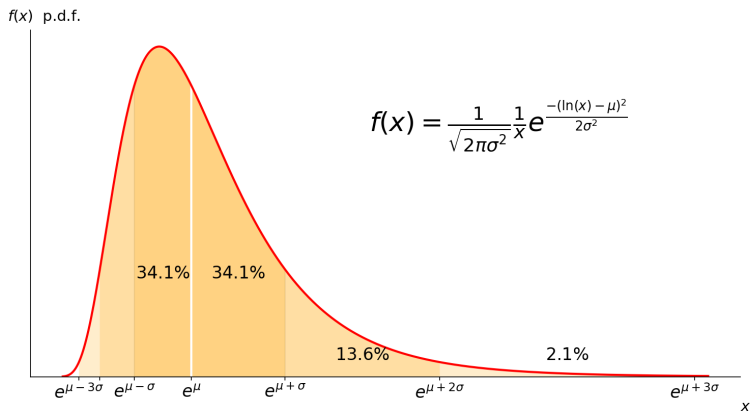


In logarithmic scale, log-norm distributions appears as Gaussian (normal).

$$y = \ln(x): \frac{1}{x} e^{-\frac{(\ln(x) - \mu)^2}{2\sigma^2}} = e^{-\frac{(y^2 - 2\mu y + \mu^2) - 2\sigma^2 y}{2\sigma^2}} = e^{-\mu + 2\sigma^2} e^{-\frac{(y - (\mu - \sigma^2))^2}{2\sigma^2}}$$

# Log-normal distribution

## Quantiles



$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{x} e^{-\frac{(\ln(x) - \mu)^2}{2\sigma^2}}$$

$$e^{\mu} \times /e^{\sigma} : 68.27\%, \quad e^{\mu} \times /e^{2\sigma} : 95.45\%, \quad e^{\mu} \times /e^{3\sigma} : 99.73\%.$$

## $\chi^2$ test statistic

Let  $x$  be a Gaussian-distributed random variable with known  $\mu$  and  $\sigma$ . We can make a simple linear transformation of this variable such, that the distribution becomes so-called *standard normal* ( $\mu = 0, \sigma = 1$ ):

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right), \quad x \rightarrow z = \frac{x - \mu}{\sigma}, \quad f(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right) \quad (26)$$

What is the distribution of  $u \equiv z^2$  ( $E[u] = E[z^2] = V[z] = 1$ )?

$$\chi_1^2(u) = 2f(z(u)) \left| \frac{dz}{du} \right| = \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{u}} \exp\left(-\frac{u}{2}\right) \quad (27)$$

**Recall:**  $z \in (-\infty, \infty) \rightarrow u = z^2 \in (0, \infty)$ .

### $\chi_1^2$ : mean & variance

$$E[u] = \int_0^\infty u \chi_1^2(u) du = 1 \quad (28)$$

$$V[u] = \int_0^\infty u^2 \chi_1^2(u) du = 2 \quad (29)$$

## $\chi^2$ test statistic

$\chi_1^2$  can be extended to distribution of two **independent** normal-distributed random variables  $u = z_1^2 + z_2^2$  by means of Fourier convolution. The operation executed recurrently provides the expression for any value of  $n$  ( $u = \sum_{i=1}^n z_i^2$ ):

$$\chi_n^2(u) = \frac{1}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} u^{\frac{n}{2}-1} \exp\left(-\frac{u}{2}\right) \quad (30)$$

**Recall:**  $\Gamma(n) = (n-1)!$ ,  $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$

$\chi_n^2$ : mean & variance

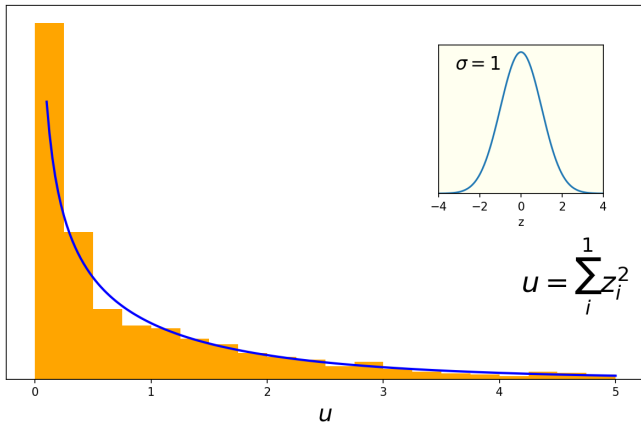
$$E[u] = \int_0^\infty u \chi_n^2(u) du = n \quad (31)$$

$$V[u] = \int_0^\infty u^2 \chi_n^2(u) du = 2n \quad (32)$$

**Note:**  $\chi^2$  distribution has only one parameter,  $n$ , called *number of degrees of freedom* (nDoF).

# $\chi^2$ test statistic

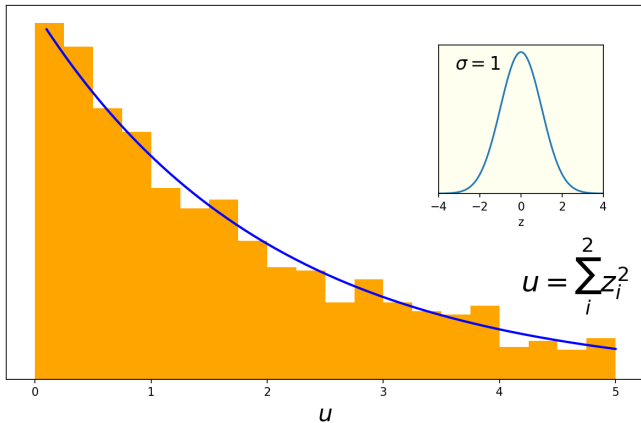
nDoF=1





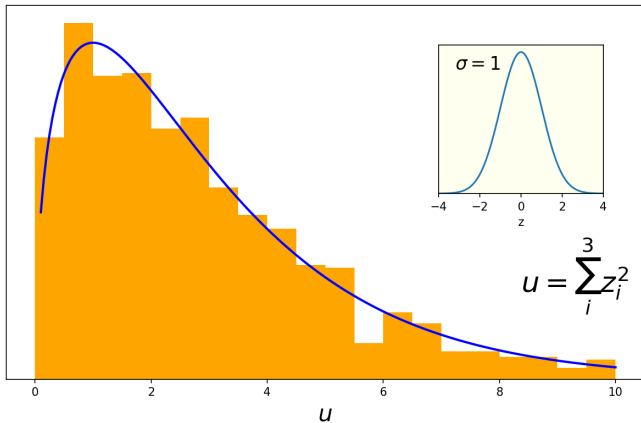
# $\chi^2$ test statistic

nDoF=2



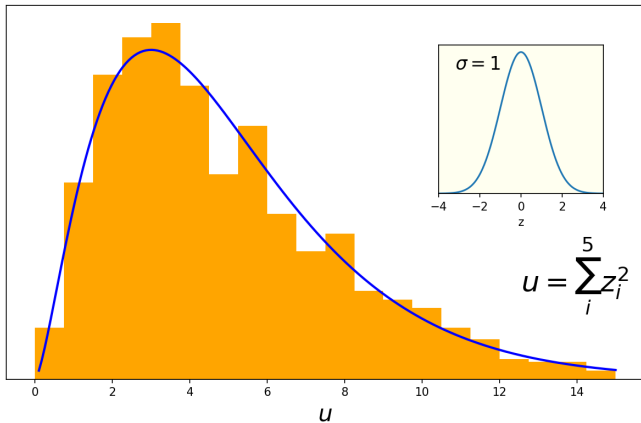
# $\chi^2$ test statistic

nDoF=3



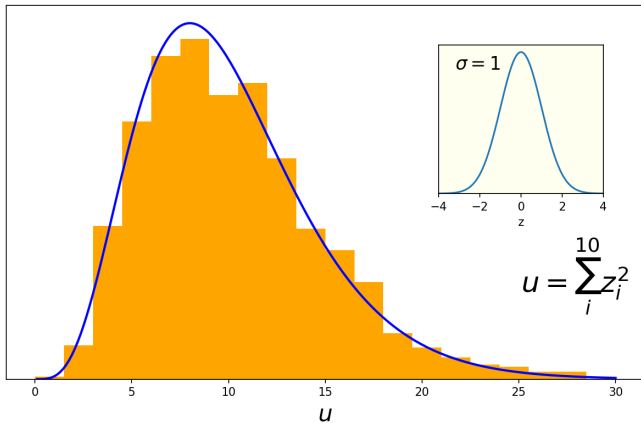
# $\chi^2$ test statistic

nDoF=5



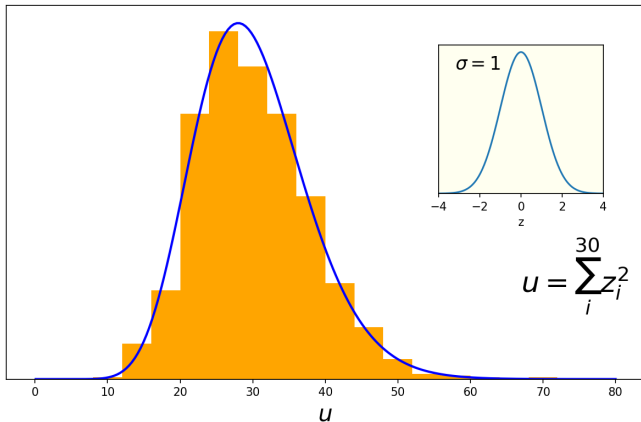
# $\chi^2$ test statistic

nDoF=10



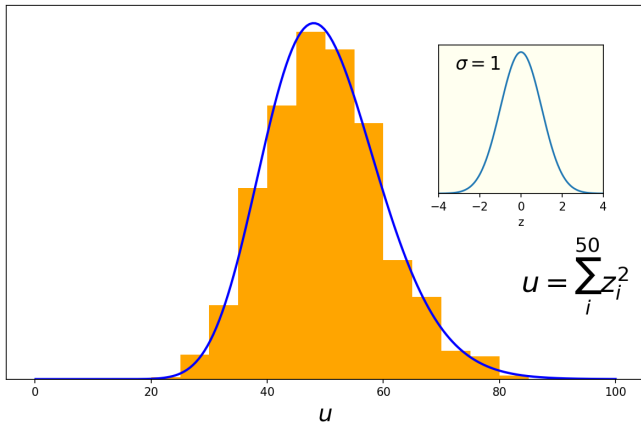
# $\chi^2$ test statistic

nDoF=30



# $\chi^2$ test statistic

nDoF=50



# $\chi^2$ test statistic

general  $n$ -dimensional case

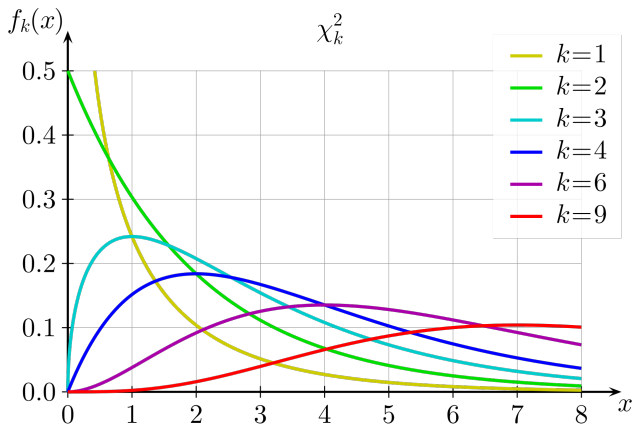
So far independence of the normal-distributed variables was as assumed. This can be generalised to  $n$ -dimensional Gaussian distribution with an arbitrary covariance matrix  $\mathbf{V}$ .

## $\chi^2$ -distributed $n$ -dimensional Gaussian

$$z = (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (33)$$

is a  $\chi_n^2$  random variable with  $n$  DoF's.

# $\chi^2$ distribution



The  $\chi_k^2$  distribution approaches a Gaussian (recall CLT!) for  $k \rightarrow \infty$ . For practical applications, it can be considered Gaussian for  $n > O(50)$  ( $\mu = k$ ,  $\sigma = \sqrt{2k}$ ).

mode:  $k - 2$ , median:  $\approx k \left(1 - \frac{2}{9k}\right)^3$ , mean:  $k$ ,  $F(X, k) = \frac{1}{\Gamma\left(\frac{k}{2}\right)} \gamma\left(\frac{k}{2}, \frac{X}{2}\right)$   
 $\gamma(s, x) = \int_0^x t^{s-1} e^{-t} dt$



# Questions

Consider the exponential p.d.f.,

$$f(x; \tau) = \frac{1}{\tau} e^{-x/\tau}, \quad x \geq 0.$$

- 1 Show that the corresponding cumulative distribution is given by

$$F(x; \tau) = 1 - e^{-x/\tau}$$

- 2 Show that the conditional probability to find a value  $x < x_0 + x'$  given that  $x > x_0$  is equal to the (unconditional) probability to find  $x$  less than  $x'$ , i.e.

$$P(x < x_0 + x' | x \geq x_0) = P(x \leq x').$$

Solutions to be sent to me before the next lecture

# Thank you

# Back-up

# Fourier convolution - revisited

$z = x + y$ , find  $f_z(z)$  given  $f_{x,y}(x, y)$

$$\begin{aligned} P(z \leq z_1) &= \int_{-\infty}^{z_1} d\kappa f_z(\kappa) = \\ &= \int_{-\infty}^{\infty} dy \int_{-\infty}^{z_1-y} dx \underbrace{f_{x,y}(x, y)}_{\text{joint p.d.f.}} = \int_{-\infty}^{\infty} dx \int_{-\infty}^{z_1-x} dy f_{x,y}(x, y) \end{aligned} \quad (34)$$

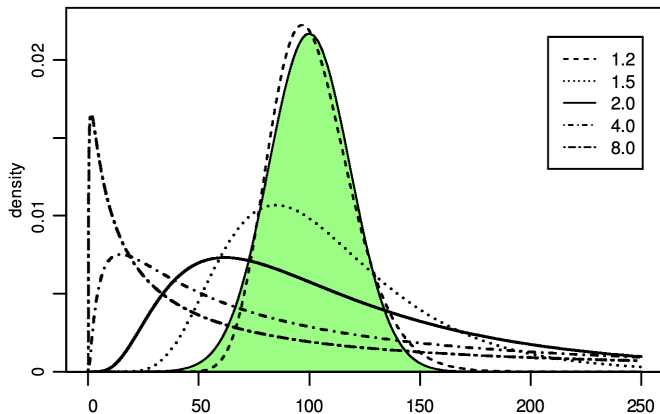
$$f_z(z) = \frac{dP}{dz} = \int_{-\infty}^{\infty} dx f_{x,y}(x, z-x) = \int_{-\infty}^{\infty} dy f_{x,y}(z-y, y) \quad (35)$$

Hence for independent variables (  $f_{x,y}(x, y) = f_x(x) * f_y(y)$  ) we obtain:

$z = x + y$  : Fourier convolution

$$f(z) = \int_{-\infty}^{+\infty} g(x)h(z-x)dx = \int_{-\infty}^{+\infty} g(z-y)h(y)dy. \quad (36)$$

# Log-normal distribution



Gaussian  $\mu, \sigma^2$  are additive, log-normal are **multiplicative**.

The log-normal distribution approaches a Gaussian for  $\sigma \rightarrow 0$ .