

Error analysis in biology

Marek Gierliński
Division of Computational Biology

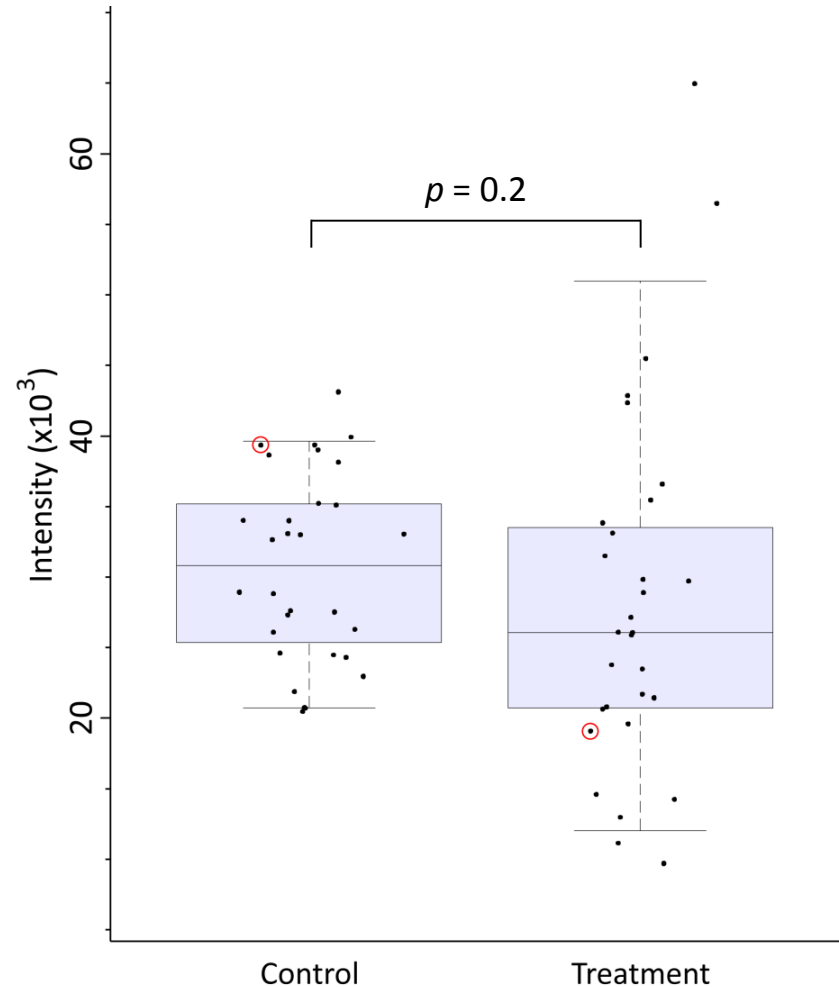
Hand-outs available at <http://is.gd/statlec>

Why do we need error analysis?

- Consider a microarray experiment
- Comparing control and treatment
- Expression level of FLG
 - control = 41,723
 - treatment = 19,786
- There is a 2-fold change in intensity
- Great! Gene is repressed in our treatment!

Why do we need error analysis?

- Consider a microarray experiment
- Comparing control and treatment
- Expression level of FLG
 - control = 41,723
 - treatment = 19,786
- There is a 2-fold change in intensity
- Great! Gene is repressed in our treatment!
- Repeat the experiment in 30 replicates
 - control = $(31.5 \pm 1.6) \times 10^3$
 - treatment = $(27.7 \pm 2.4) \times 10^3$
- Reveal **variability** of expression
- No difference between control and treatment



“A measurement without error is meaningless”

My physics teachers

Data Analysis Group



Chris Cole



Stuart MacGowan



Pietà Schofield



Marek Gierliński

Computational Biology
Barton Group
Level 2 CTIR

<http://www.compbio.dundee.ac.uk/dag.html>

$\nabla \cdot \mathbf{F} = \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 F_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta F_\theta) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \phi} (F_\phi \sin \theta)$
 $\nabla \cdot \mathbf{F} = \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 F_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta F_\theta) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \phi} (F_\phi \sin \theta)$

Magnetic Field & Vector Potentials:
 $\mathbf{A} = \frac{\mu_0}{4\pi} \int \frac{d\mathbf{j}(\mathbf{r}')}{r}$
 $\mathbf{B} = \nabla \times \mathbf{A}$
 $\nabla \times \mathbf{A} = \frac{\mu_0}{4\pi} \nabla \times \int \frac{d\mathbf{j}(\mathbf{r}')}{r}$
 $\mathbf{B} = \frac{\mu_0}{4\pi} \int \frac{d\mathbf{j}(\mathbf{r}') \times \mathbf{r} - \mathbf{r}' \times \mathbf{j}(\mathbf{r}')}{r^3}$

Quantum Mechanics:
 $\psi(\mathbf{r}, t) = \psi(\mathbf{r}) e^{-iEt/\hbar}$
 $\hat{H} \psi = E \psi$
 $\hat{H} = -\frac{\hbar^2}{2m} \nabla^2 + V(\mathbf{r})$
 $\psi(r, \theta, \phi) = R(r) Y(\theta, \phi)$
 $\psi(r, \theta, \phi) = R(r) Y(\theta, \phi) e^{-iEt/\hbar}$

Atomic Physics:
 $E_n = -13.6 \text{ eV} \frac{Z^2}{n^2}$
 $\lambda = \frac{hc}{E} = \frac{hc}{13.6 \text{ eV} \frac{Z^2}{n^2}} = \frac{1.215 \times 10^{-7} \text{ m}}{Z^2} \left(\frac{1}{n_f^2} - \frac{1}{n_i^2} \right)^{-1}$
 $\lambda = 109.7 \text{ nm} \frac{1}{Z^2} \left(\frac{1}{n_f^2} - \frac{1}{n_i^2} \right)^{-1}$

Electron Properties:
 $v = \frac{h}{m \lambda} = \frac{1240 \text{ eV nm}}{m \lambda} = \frac{1240 \text{ eV nm}}{m \cdot 109.7 \text{ nm} \frac{1}{Z^2} \left(\frac{1}{n_f^2} - \frac{1}{n_i^2} \right)^{-1}}$
 $v = \frac{1240 \text{ eV nm} \cdot Z^2 \left(\frac{1}{n_f^2} - \frac{1}{n_i^2} \right)}{m \cdot 109.7 \text{ nm}}$

Wave Functions:
 $\psi(r, \theta, \phi) = R(r) Y(\theta, \phi) e^{-iEt/\hbar}$
 $\psi(r, \theta, \phi) = R(r) Y(\theta, \phi) e^{-iEt/\hbar}$

Relativistic Energy:
 $E = \gamma mc^2 = \frac{mc^2}{\sqrt{1 - v^2/c^2}}$
 $E = \frac{mc^2}{\sqrt{1 - \frac{v^2}{c^2}}} = \frac{mc^2}{\sqrt{1 - \frac{v^2}{c^2}}}$

Atomic Structure:
 $\lambda = \frac{h}{m v} = \frac{h}{m \frac{v}{c}} = \frac{h}{m \beta c} = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$
 $\lambda = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$

Electron Spin:
 $\mu_B = \frac{e \hbar}{2m} = 9.27 \times 10^{-24} \text{ J/T}$
 $\mu_B = \frac{e \hbar}{2m} = 9.27 \times 10^{-24} \text{ J/T}$

Angular Momentum:
 $L = \hbar \sqrt{l(l+1)}$
 $L_z = \hbar m_l$

Probability Density:
 $P(r, \theta, \phi) = |\psi(r, \theta, \phi)|^2$
 $P(r, \theta, \phi) = |\psi(r, \theta, \phi)|^2$

Electron Configuration:
 $1s^2 2s^2 2p^6 3s^2 3p^4$
 $1s^2 2s^2 2p^6 3s^2 3p^4$

Bohr Model:
 $r_n = n^2 a_0 = n^2 \cdot 0.529 \text{ \AA}$
 $r_n = n^2 a_0 = n^2 \cdot 0.529 \text{ \AA}$

Wave Function Normalization:
 $\int |\psi|^2 d\tau = 1$
 $\int |\psi|^2 d\tau = 1$

Energy Levels:
 $E_n = -13.6 \text{ eV} \frac{Z^2}{n^2}$
 $E_n = -13.6 \text{ eV} \frac{Z^2}{n^2}$

Wave Function Asymptotics:
 $\psi(r) \sim e^{-\kappa r}$
 $\psi(r) \sim e^{-\kappa r}$

Atomic Size:
 $r \sim \frac{\hbar}{m v} = \frac{\hbar}{m \frac{v}{c}} = \frac{\hbar}{m \beta c} = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$
 $r \sim \frac{\hbar}{m v} = \frac{\hbar}{m \frac{v}{c}} = \frac{\hbar}{m \beta c} = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$

Quantum Tunneling:
 $\psi \sim e^{-\kappa x}$
 $\psi \sim e^{-\kappa x}$

Angular Momentum Addition:
 $\mathbf{L} = \mathbf{L}_1 + \mathbf{L}_2$
 $\mathbf{L} = \mathbf{L}_1 + \mathbf{L}_2$

Probability Distributions:
 $P(r) = 4\pi r^2 |\psi(r, \theta, \phi)|^2$
 $P(r) = 4\pi r^2 |\psi(r, \theta, \phi)|^2$

Bohr Model Constants:
 $a_0 = 0.529 \text{ \AA}$
 $a_0 = 0.529 \text{ \AA}$

Atomic Masses:
 $m_p = 1.67 \times 10^{-27} \text{ kg}$
 $m_e = 9.11 \times 10^{-31} \text{ kg}$
 $m_p = 1.67 \times 10^{-27} \text{ kg}$
 $m_e = 9.11 \times 10^{-31} \text{ kg}$

Atomic Structure Diagrams:
 $1s^2 2s^2 2p^6 3s^2 3p^4$
 $1s^2 2s^2 2p^6 3s^2 3p^4$

Atomic Size Comparison:
 $r \sim \frac{\hbar}{m v} = \frac{\hbar}{m \frac{v}{c}} = \frac{\hbar}{m \beta c} = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$
 $r \sim \frac{\hbar}{m v} = \frac{\hbar}{m \frac{v}{c}} = \frac{\hbar}{m \beta c} = \frac{2.19 \times 10^{-10} \text{ m}}{\beta}$

Atomic Structure Diagrams:
 $1s^2 2s^2 2p^6 3s^2 3p^4$
 $1s^2 2s^2 2p^6 3s^2 3p^4$



Table of contents

- | | |
|-------------------------------|-----|
| 1. Probability distribution | ① |
| 2. Random errors | ② |
| 3. Statistical estimators | |
| 4. Confidence intervals | ③ ④ |
| 5. Error bars | ⑤ |
| 6. Quoting numbers and errors | |
| 7. Error propagation | ⑥ |
| 8. Linear regression errors | |

Example

- Experiment: estimate bacterial concentration using a spectrophotometer
- 6 replicates
- Find the following OD600
0.37 0.34 0.41 0.40 0.30 0.33
- Experimental result is a **random variable**
- It follows a certain **probability distribution**



1. Random variables and probability distributions

“Misunderstanding of probability may be the greatest of all general impediments to scientific literacy”

Stephen Jay Gould

Random variable: random numbers



12
9
12
11
4
6
7
8
3
5

Discrete and continuous random variables

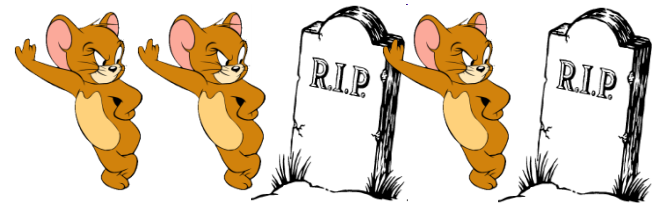
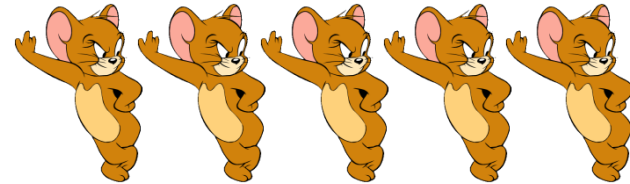
- Discrete variables:

- sum of 2 dice (2, 3, 4, ..., 12)
- categorical outcome
- number of mice (5, non random?)
- number of mice in survival experiment (random)



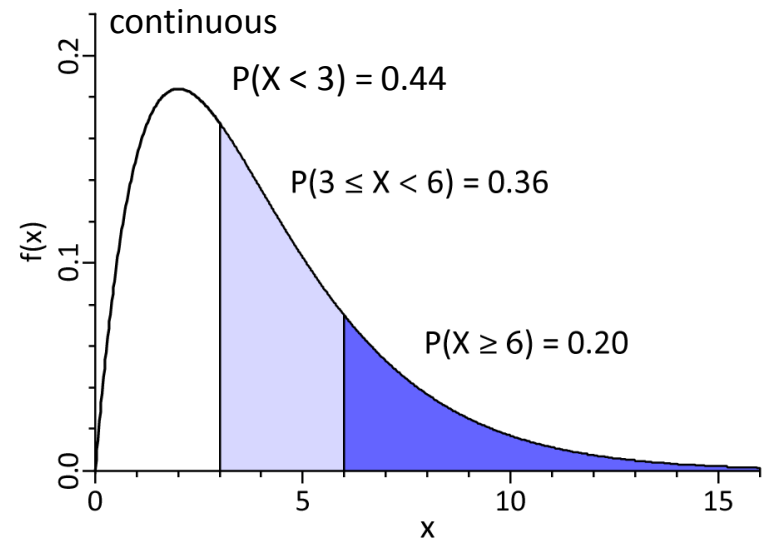
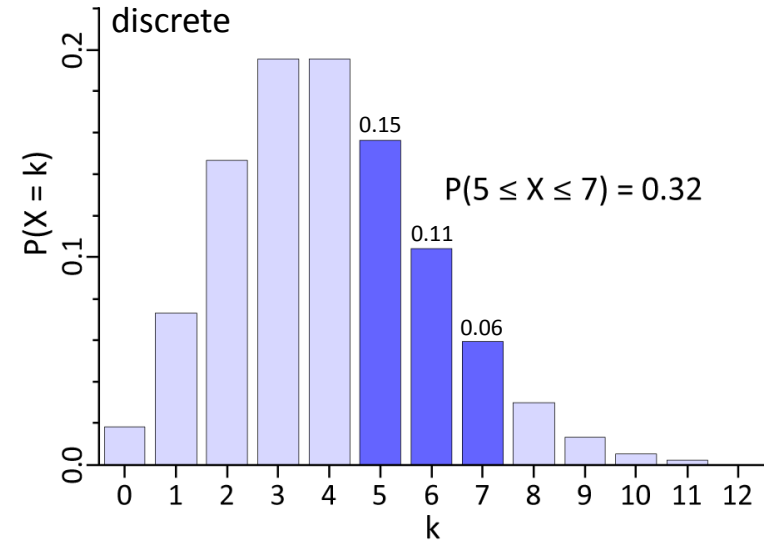
- Continuous variables:

- weight of a mouse
- height of a person
- fluorescent marker luminosity
- protein abundance



Probability distribution

- Assigns a probability to each of the possible outcomes
- X – random variable
- $P(X = 5)$ – probability of X being 5
- $P(5 \leq X \leq 7)$ – probability of X between 5 and 7 (sum of probabilities)
- $f(x)$ – probability density function
- $P(X < 3)$ – area under the curve $f(x)$
- $P(X = 5) = 0$

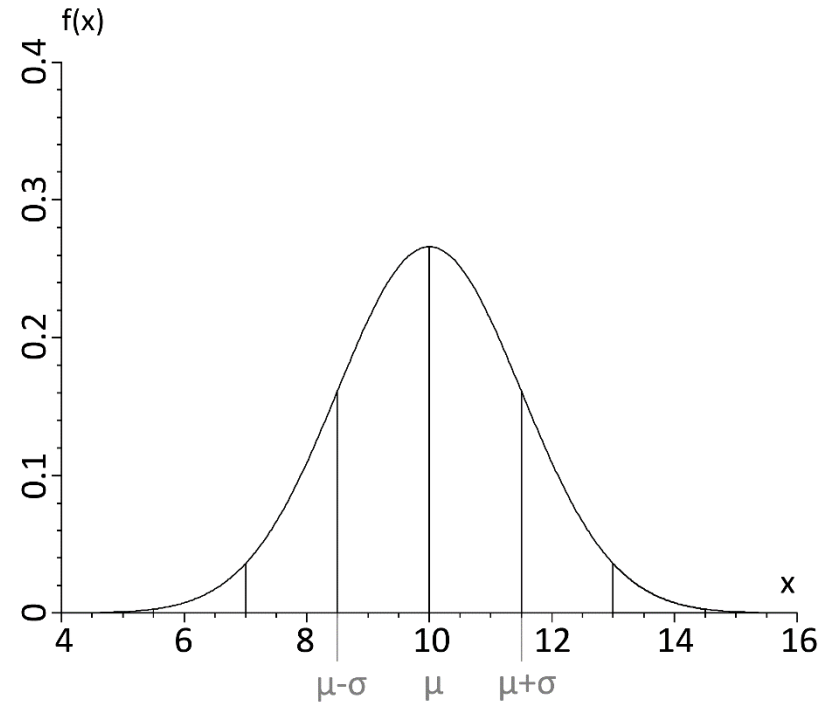


Gaussian distribution

- Gaussian (or normal) probability distribution

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

- μ - mean
 - σ - standard deviation
 - σ^2 - variance
- It is called “normal” as it often appears in nature

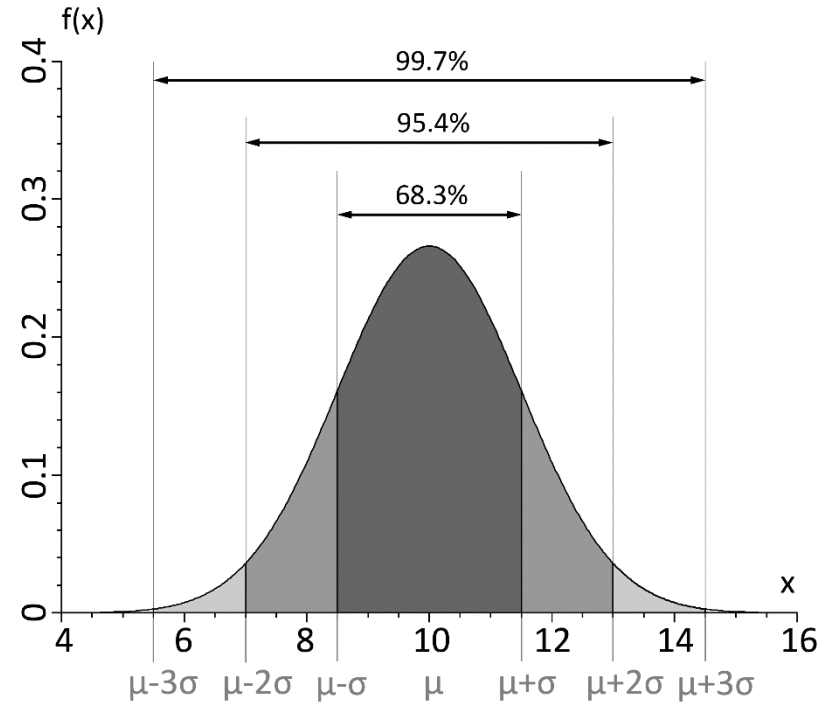


$\mathcal{N}(10, 1.5)$ - normal distribution with
 $\mu = 10$ and $\sigma = 1.5$

Gaussian distribution: a few numbers

- Area under the curve = probability
- Probability within one sigma of the mean is about $\frac{2}{3}$ (68.3%)
- 95% confidence intervals are traditionally used: correspond to about 1.96σ

	In	Out	Odds of out
$\pm 1\sigma$	68.3%	31.7%	1:3
$\pm 2\sigma$	95.4%	4.6%	1:20
$\pm 3\sigma$	99.7%	0.3%	1:400
$\pm 4\sigma$	99.994%	0.006%	1:16,000
$\pm 5\sigma$	99.99993%	0.00007%	1:1,700,000
$\pm 1.96\sigma$	95.0%	5.0%	1:20

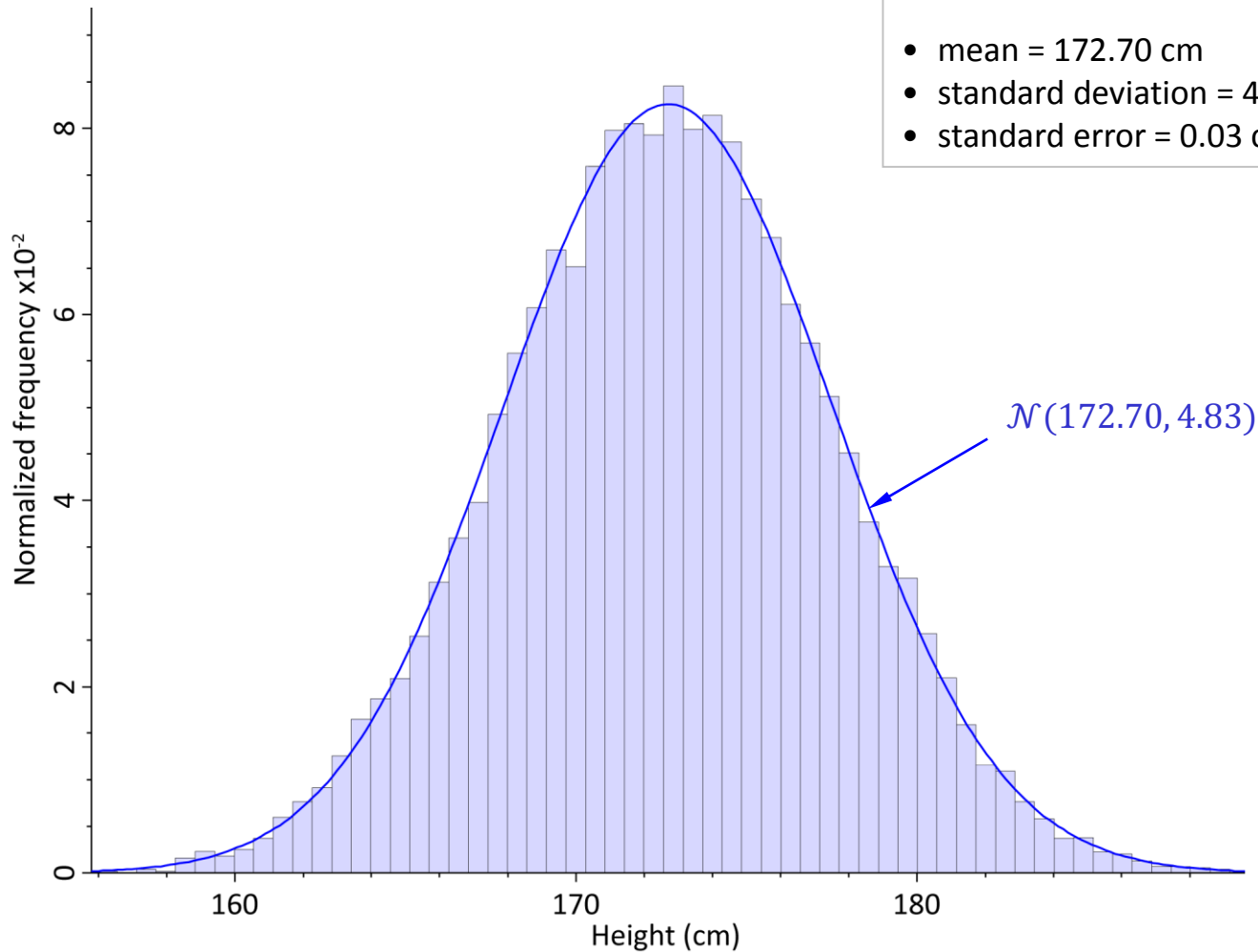


$\mathcal{N}(10, 1.5)$ - normal distribution with $\mu = 10$ and $\sigma = 1.5$

Example: Gaussian distribution

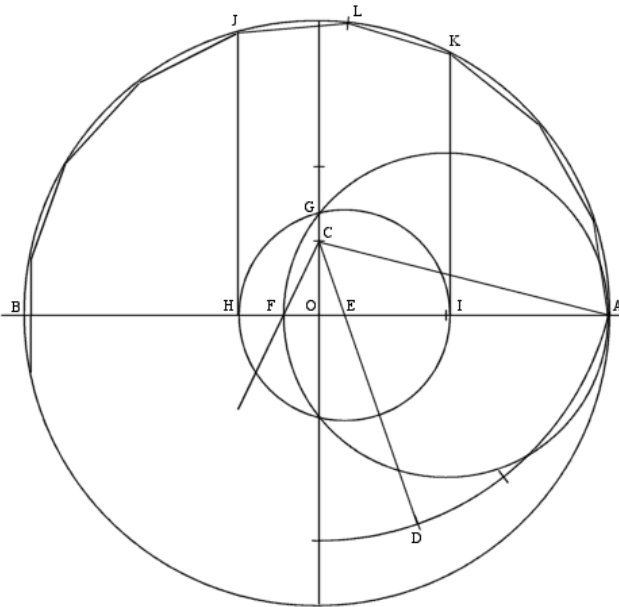
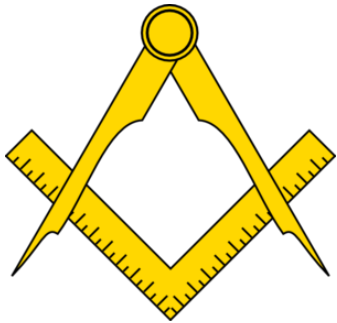
Height of 25,000 individuals from Hong Kong

- mean = 172.70 cm
- standard deviation = 4.83 cm
- standard error = 0.03 cm



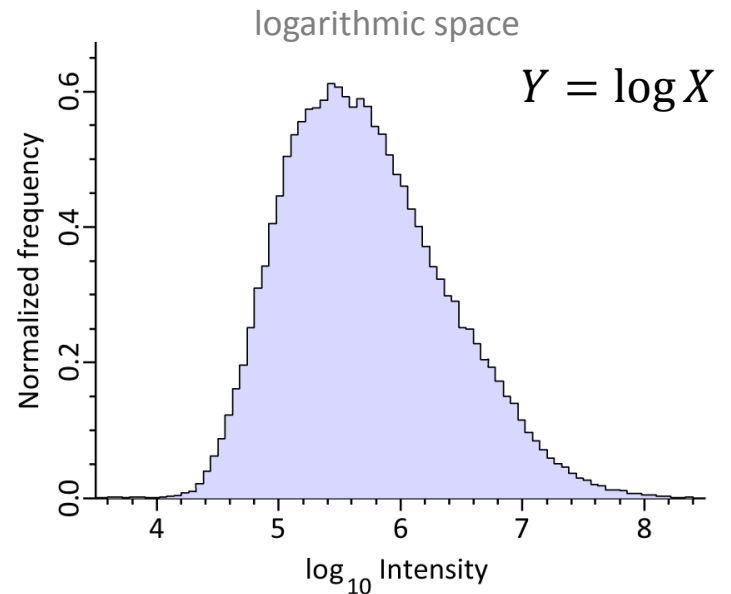
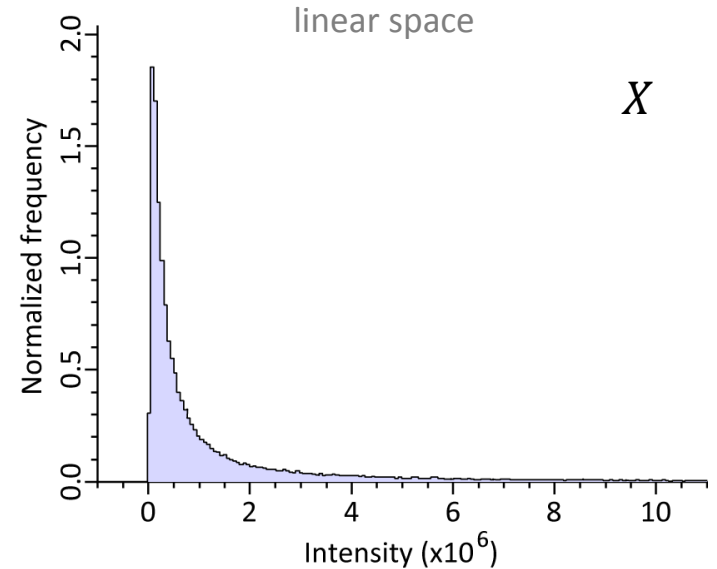
Carl Friedrich Gauss (1777-1855)

- Brilliant German mathematician
- Constructed a regular heptadecagon with a ruler and a compass
- He requested that a regular heptadecagon should be inscribed on his tombstone
- However, it was Abraham de Moivre (1667-1754) who first formulated “Gaussian” distribution



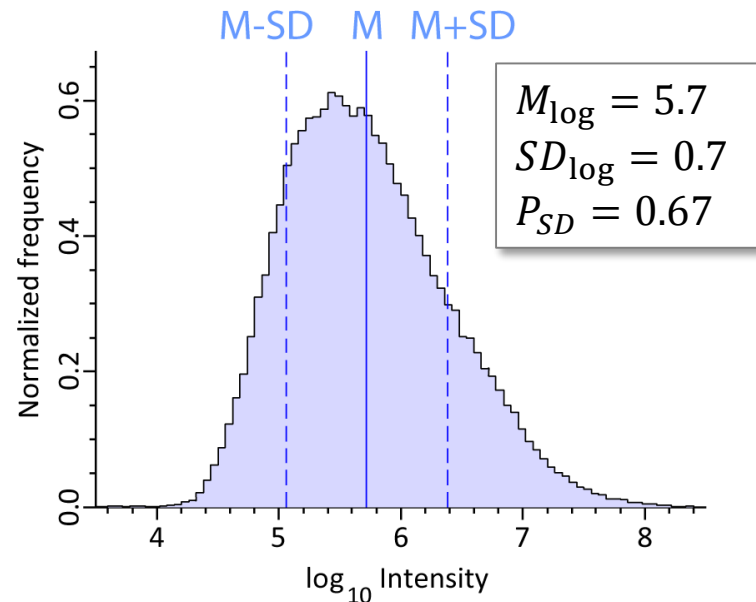
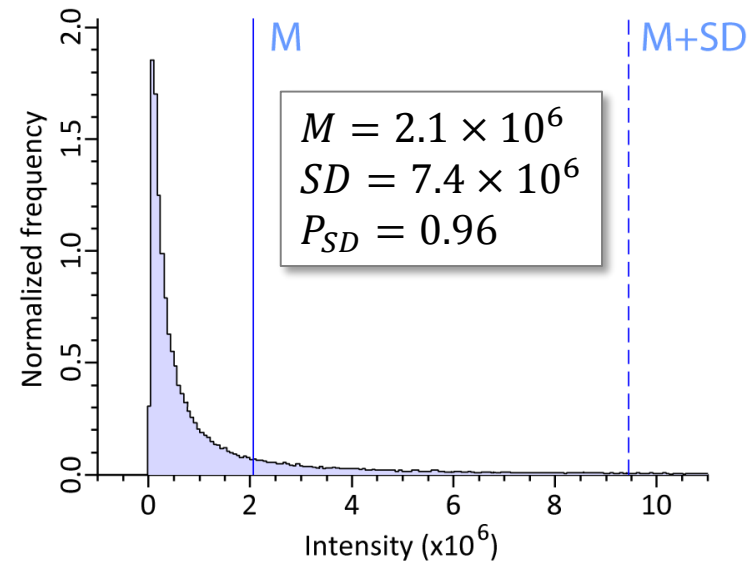
Log-normal distribution

- Probability distribution of a random variable whose logarithm is normally distributed
- Log-normal distribution can be very asymmetric!



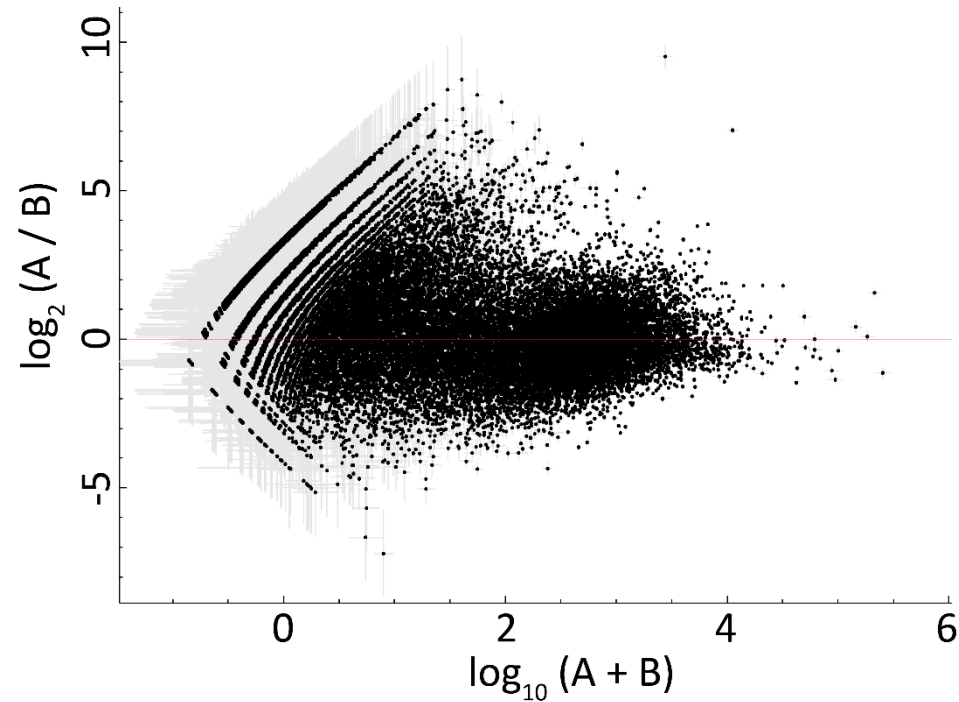
Example: log-normal distribution

- Peptide intensities from a mass spectrometry experiment
- P_{SD} - fraction of data within $M \pm SD$
- Data look better in logarithmic space
- Always plot the distribution of your data before analysis
- About two-thirds of data points are within one standard deviation from the mean **only** when their distribution is approximately Gaussian



A few notes on log-normal distribution

- Examples of log-normal distributions
 - gene expression (RNA-seq, microarrays)
 - mass spectrometry data
 - drug potency IC_{50}
- It doesn't matter if you use \log_2 , \log_{10} or \ln , as long as you are consistent
- \log_{10} is easier to understand in plots
 - $10^5 = 100,000$
 - $2^{10} = 1024$



John Napier (1550-1617)

- Scottish mathematician and astronomer
- Invented logarithms and published first tables of natural logarithms
- Created “Napier’s bones”, the first practical calculator
- Had an interest in theology, calculated the date of the end of the world between 1688 and 1700
- Apparently involved in alchemy and necromancy



Merchiston Castle, Edinburgh

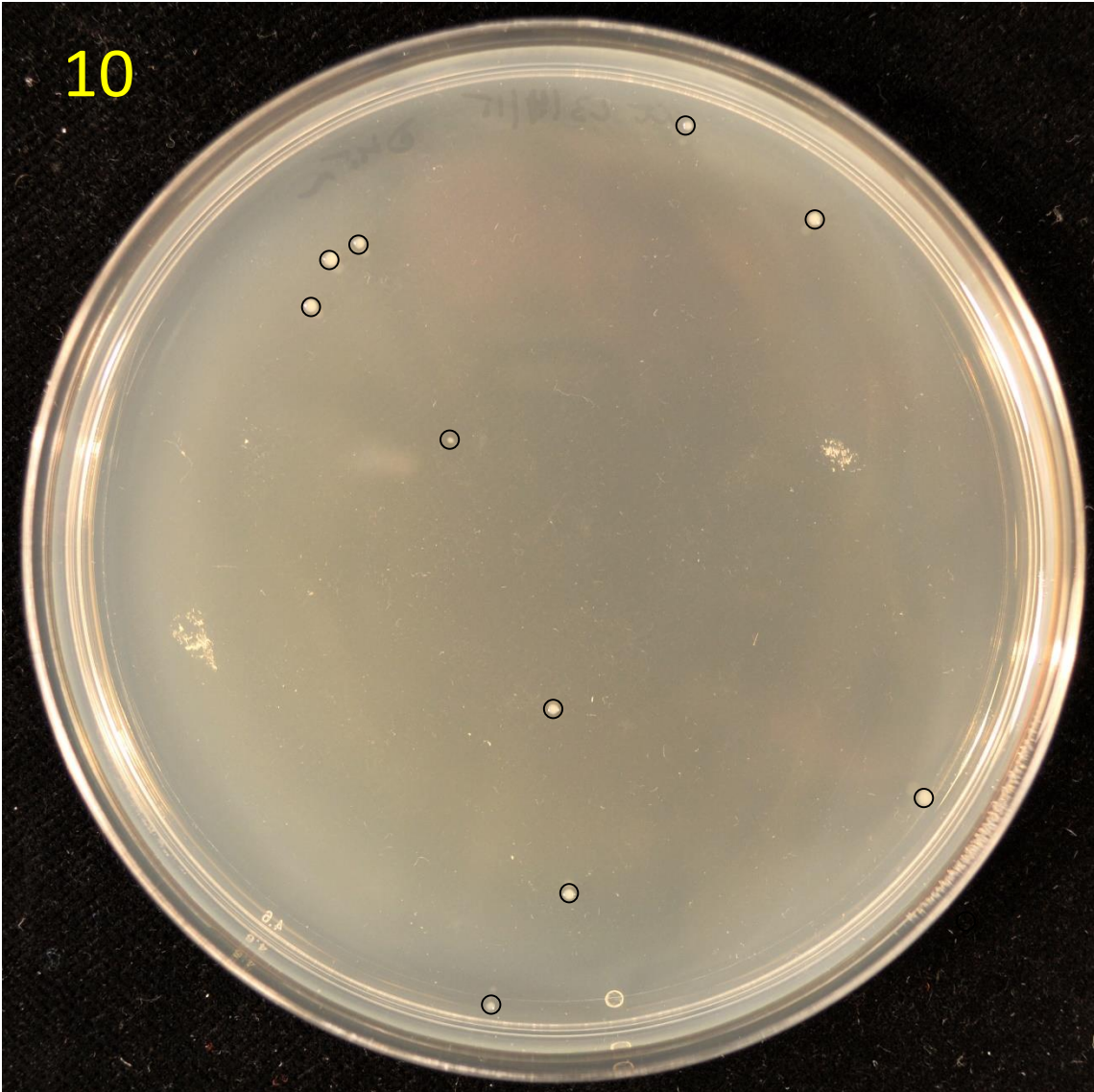




What's in the box?

Counting bacterial colonies

10

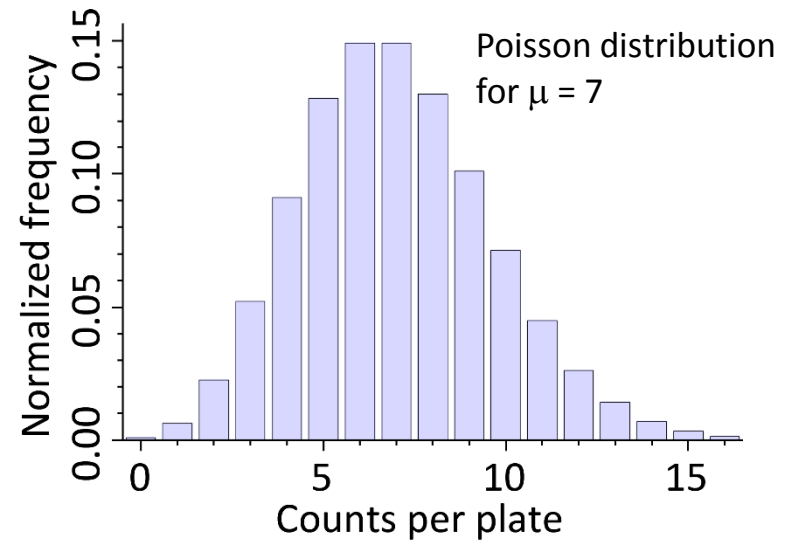
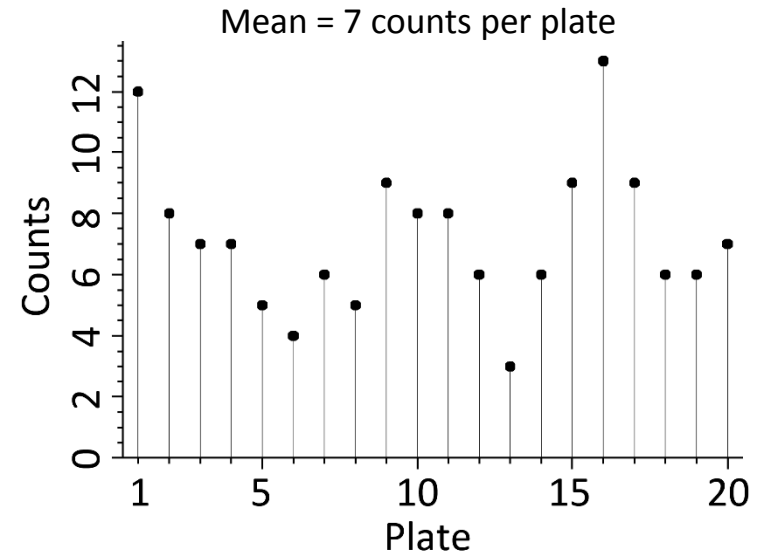


Courtesy of Katharina Trunk,
Molecular Biology

100 μ l of 10^{-7} dilution of $OD_{600} = 2.0$

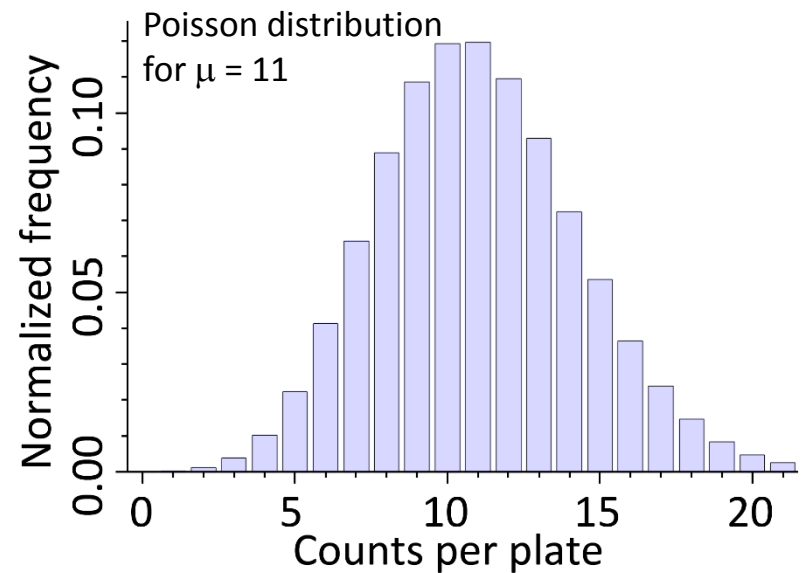
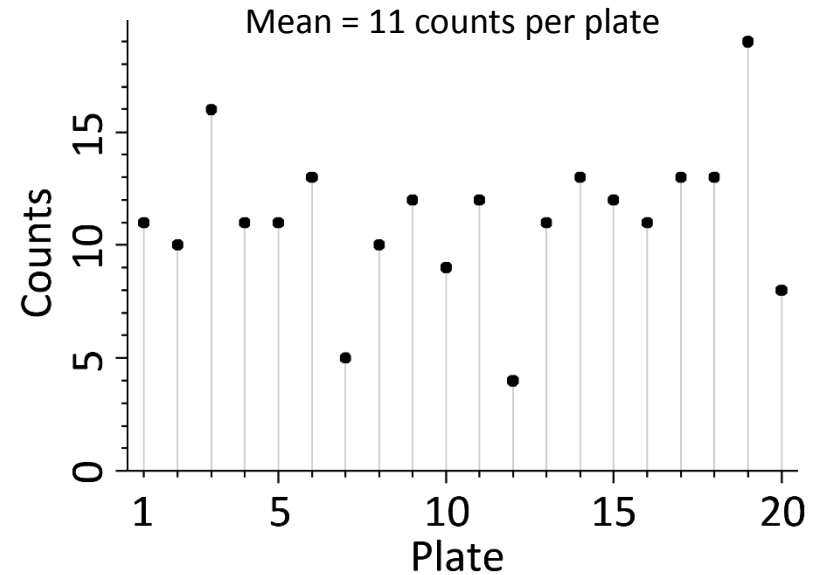
Poisson distribution

- Measure of bacterial count per unit volume
- Poisson count: always per bin
- This applies to any counts in time or space
 - radioactive decays per second
 - number of deaths in a population
 - number of cells in a counting chamber
 - number of mutations in a DNA fragment



Poisson distribution

- Measure of bacterial count per unit volume
- Poisson count: always per bin
- This applies to any counts in time or space
 - radioactive decays per second
 - number of deaths in a population
 - number of cells in a counting chamber
 - number of mutations in a DNA fragment



Poisson distribution

- *Random and independent* events
- Probability of observing exactly k events:

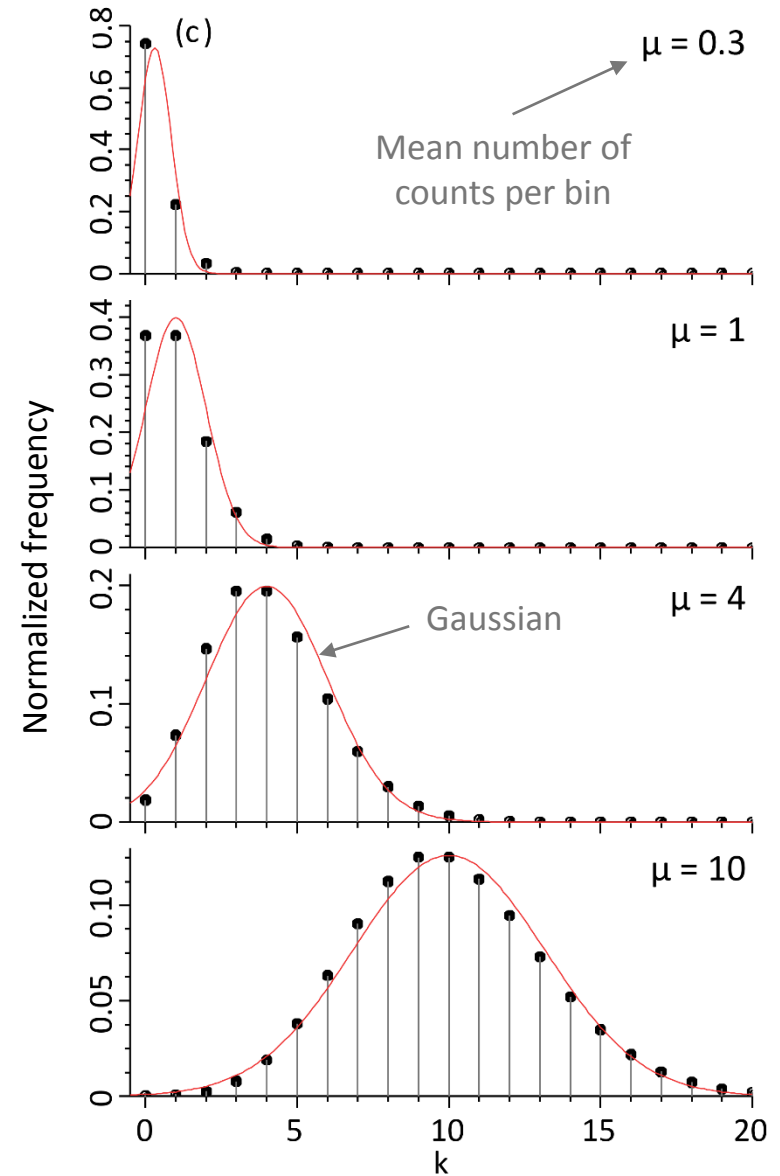
$$P(X = k) = \frac{\mu^k e^{-\mu}}{k!}$$

- One parameter: mean count rate, μ
- Standard deviation:

$$\sigma = \sqrt{\mu}$$

$$\sigma^2 = \mu$$

- For large μ Poisson distribution approximates Gaussian



Classic example: horse kicks

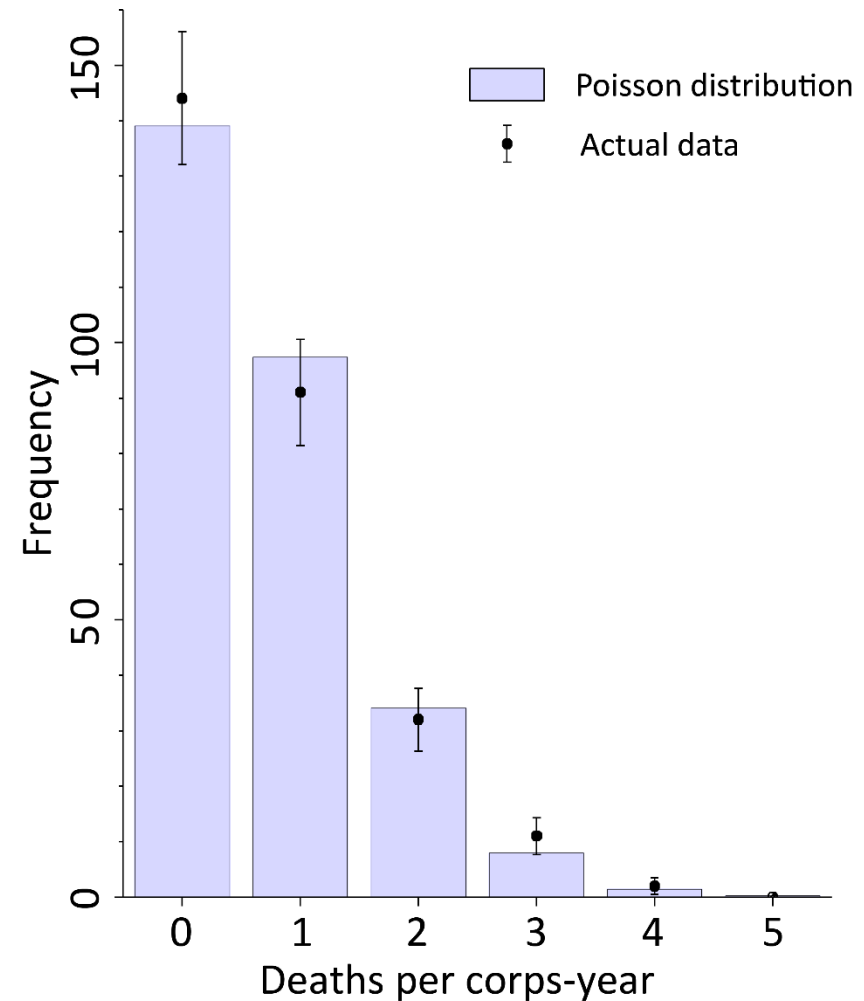
- Ladislaus von Bortkiewicz (1898) *“Das Gesetz der kleinen Zahlen”*
- Number of soldiers in the Prussian army killed by horse kicks
 - 14 army corps, 20 years of data
 - Deaths per year per army corps

In nachstehender Tabelle sind die Zahlen der durch Schlag eines Pferdes verunglückten Militärpersonen, nach Armeecorps („G.“ bedeutet Gardecorps) und Kalenderjahren nachgewiesen.¹⁾

	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94
G	—	2	2	1	—	—	1	1	—	3	—	2	1	—	—	1	—	1	—	1
I	—	—	—	2	—	3	—	2	—	—	—	1	1	1	—	2	—	3	1	—
II	—	—	—	2	—	2	—	—	1	1	—	—	2	1	1	—	—	2	—	—
III	—	—	—	1	1	1	2	—	2	—	—	—	1	—	1	2	1	—	—	—
IV	—	1	—	1	1	1	1	—	—	—	—	1	—	—	—	—	1	1	—	—
V	—	—	—	—	2	1	—	—	1	—	—	1	—	1	1	1	1	1	1	—
VI	—	—	1	—	2	—	—	1	2	—	1	1	3	1	1	1	—	3	—	—
VII	1	—	1	—	—	—	1	—	1	1	—	—	2	—	—	2	1	—	2	—
VIII	1	—	—	—	1	—	—	1	—	—	—	—	1	—	—	—	1	1	—	1
IX	—	—	—	—	—	2	1	1	1	—	2	1	1	—	1	2	—	1	—	—
X	—	—	1	1	—	1	—	2	—	2	—	—	—	—	2	1	3	—	1	1
XI	—	—	—	—	2	4	—	1	3	—	1	1	1	1	2	1	3	1	3	1
XIV	1	1	2	1	1	3	—	4	—	1	—	3	2	1	—	2	1	1	—	—
XV	—	1	—	—	—	—	—	1	—	1	1	—	—	—	2	2	—	—	—	—

Example: Poisson distribution

- Death distribution follows Poisson law
- mean = 0.70 deaths / corps / year
- 4 deaths in a corps-year are expected to happen from time to time!
- $P(X = 4) = 0.078$ in 14 corps
- On average it should happen once in 13 years



Exercise: Poisson distribution

- Poisson law:

$$P(X = k) = \frac{\mu^k e^{-\mu}}{k!}$$

- You transfect a marker into a population of $n = 3 \times 10^5$ cells
- It functionally integrates with the genome at a rate of $r = 10^{-5}$
- What is the probability of having at least one cell with the marker?

- First calculate the mean (expected) number of marked cells:

$$\mu = nr = 3$$

- Now we can use the Poisson law to find $P(X = 0)$

$$P(X = 0) = \frac{3^0 e^{-3}}{0!} = \frac{1 \times 0.05}{1} = 0.05$$

- Hence, the solution

$$P(X > 0) = 1 - P(X = 0) = 0.95$$

Binomial distribution

- A series of n “trials”
- In each trial, the probability of of:
 - “success” = p
 - “failure” = $1 - p$
- What is the probability of having exactly k successes in n trials?

- Mean and standard deviation

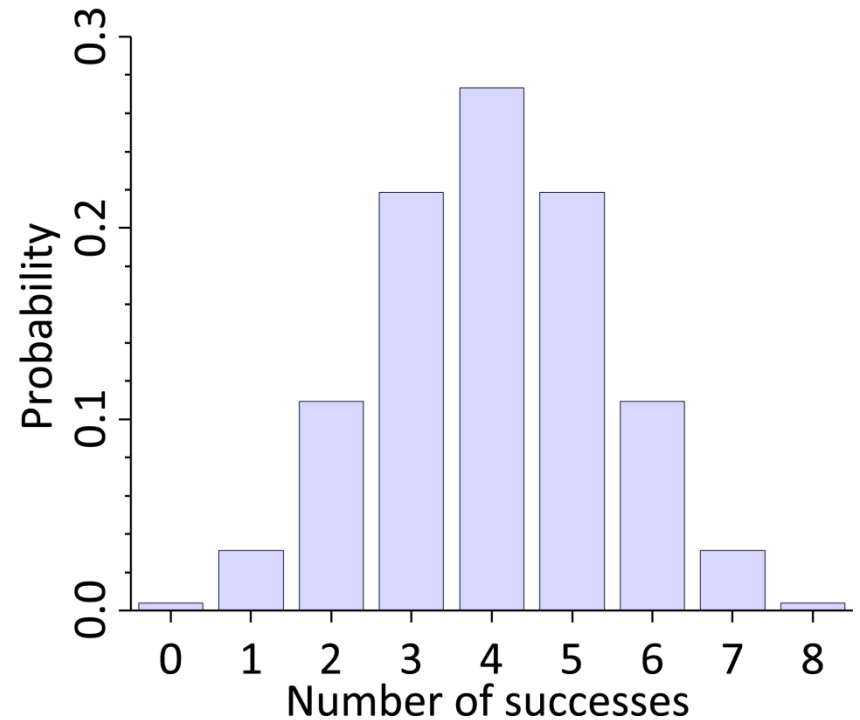
$$\mu = np$$

$$\sigma = \sqrt{np(1 - p)}$$

- For large n approximates Gaussian

- Applications:

- random errors
- error of the proportion
- error of the median

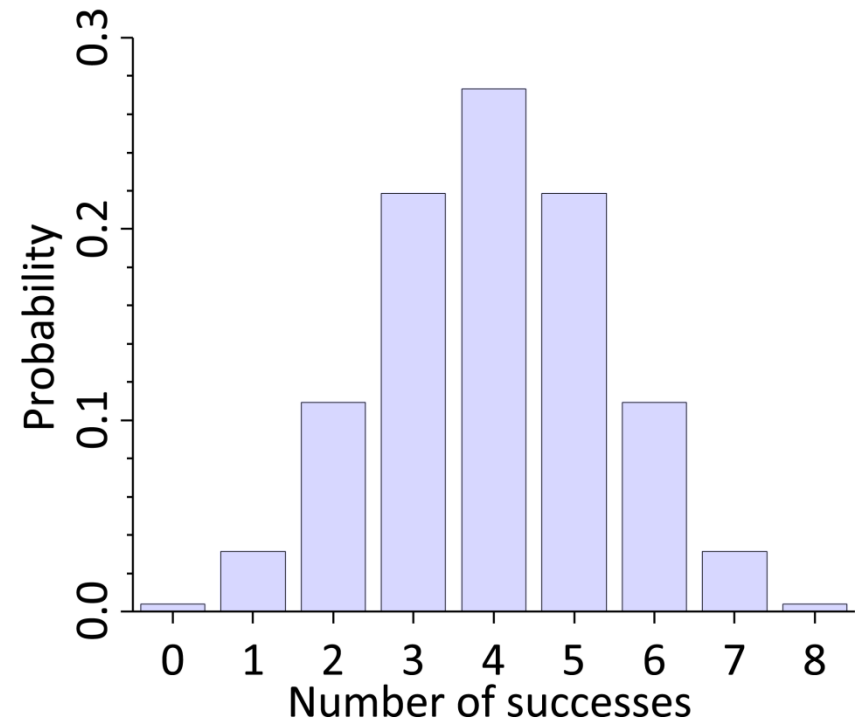


Example: toss a coin
heads = success ($p = 0.5$)
tails = failure ($1 - p = 0.5$)

What is the probability of obtaining heads k times from 8 coins?

Example: tossing a coin

- Toss 8 coins
- Question: why is the probability having heads 4 times much larger than the probability of heads 8 times?



Example: toss a coin
heads = success ($p = 0.5$)
tails = failure ($1 - p = 0.5$)

What is the probability of obtaining heads k times from 8 coins?

Example: tossing a coin

- There is only one way of having heads 8 times

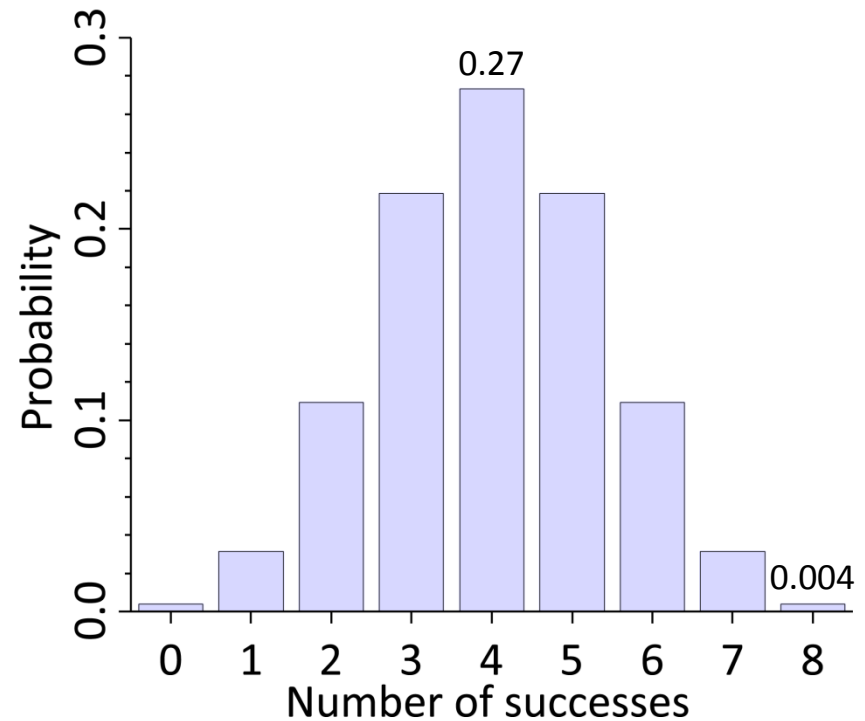


- There are many ways of getting 4 heads and 4 tails



...

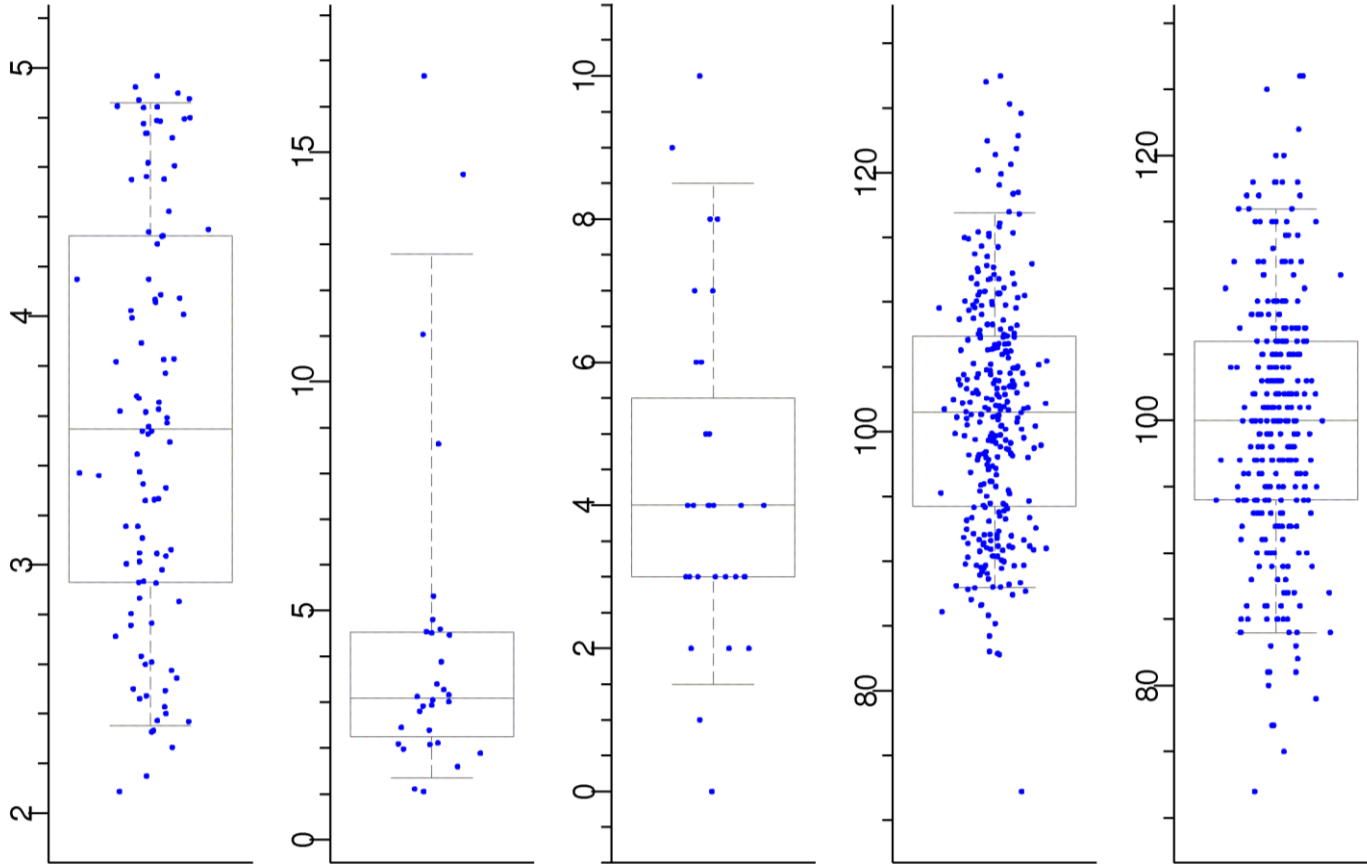
$$\binom{8}{4} = 70$$



Example: toss a coin
heads = success ($p = 0.5$)
tails = failure ($1 - p = 0.5$)

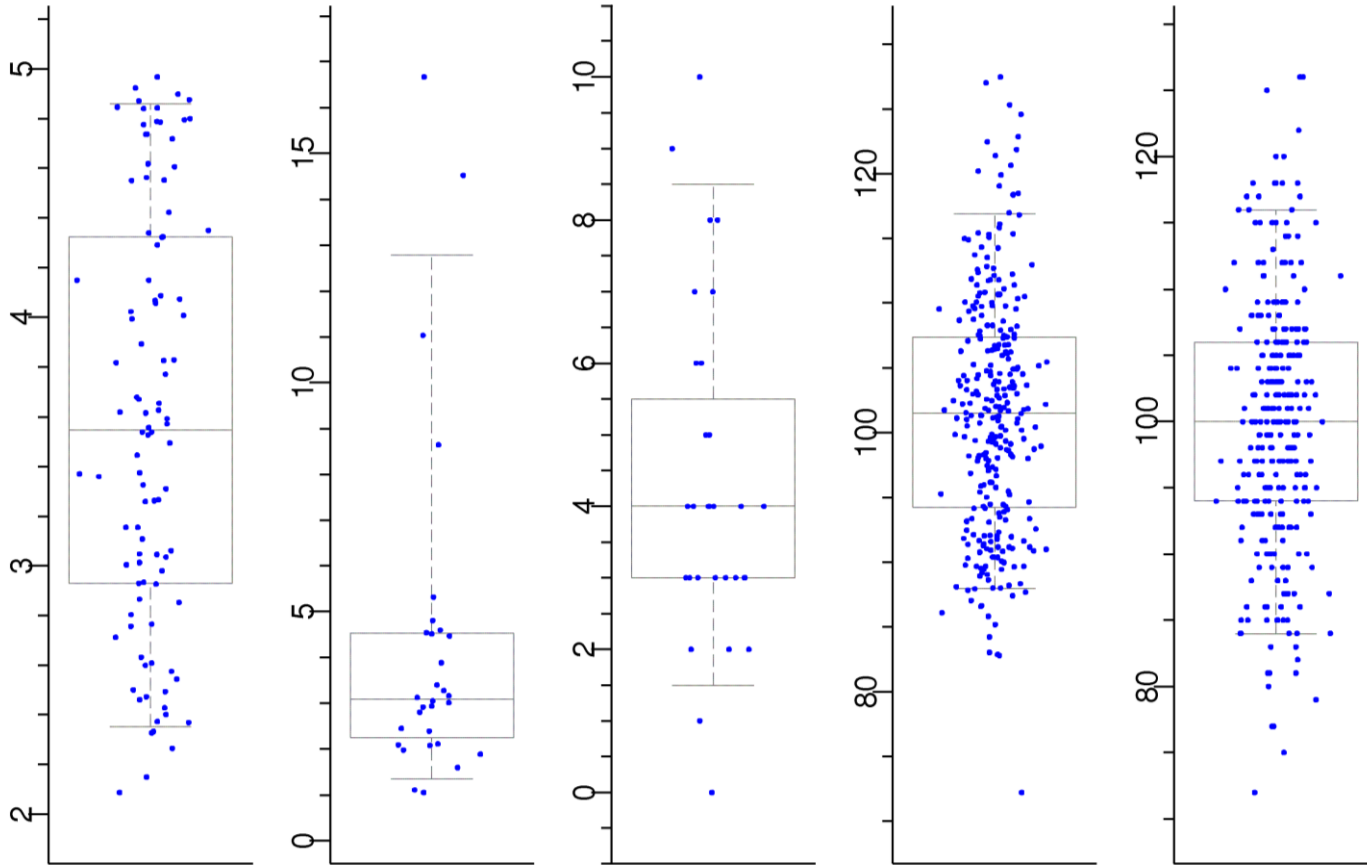
What is the probability of obtaining heads k times from 8 coins?

Exercise: recognize these distributions



Distribution					
Mean					
<i>SD</i>					

Exercise: recognize these distributions



Distribution	Uniform	Log-normal	Poisson	Gaussian	Poisson
Mean	3.5	3.5	4	100	100
<i>SD</i>	0.87	0.90	2	10	10



Hand-outs available at <http://is.gd/statlec>

Please leave your feedback forms on the table by the door

