

# ENM 53 I: Data-driven Modeling and Probabilistic Scientific Computing

## *Lecture #2: Primer on Probability and Statistics*

Paris Perdikaris  
January 20, 2022



# Discrete random variables

- A ***discrete random variable*** is one which may take on only a countable number of distinct values such as 0,1,2,3,4,..... Discrete random variables are usually (but not necessarily) counts. If a random variable can take only a finite number of distinct values, then it must be discrete. Examples of discrete random variables include the number of children in a family, the Friday night attendance at a cinema, the number of patients in a doctor's surgery, the number of defective light bulbs in a box.
- The ***probability distribution*** of a discrete random variable is a list of probabilities associated with each of its possible values. It is also sometimes called the probability function or the probability mass function.



# Continuous random variables

- A *continuous random variable* is one which takes an infinite number of possible values. Continuous random variables are usually measurements. Examples include height, weight, the amount of sugar in an orange, the time required to run a mile.
- A continuous random variable is not defined at specific values. Instead, it is defined over an *interval* of values, and is represented by the *area under a curve* (in advanced mathematics, this is known as an *integral*). The probability of observing any single value is equal to 0, since the number of values which may be assumed by the random variable is infinite.
- Suppose a random variable  $X$  may take all values over an interval of real numbers. Then the probability that  $X$  is in the set of outcomes  $A$ ,  $P(A)$ , is defined to be the area above  $A$  and under a curve. The curve, which represents a function  $p(x)$ , must satisfy the following:
  - *1: The curve has no negative values ( $p(x) \geq 0$  for all  $x$ )*
  - *2: The total area under the curve is equal to 1.*
  - A curve meeting these requirements is known as a *density curve*.

# Probability density functions

$$\Pr[a \leq X \leq b] = \int_a^b f_X(x) dx.$$

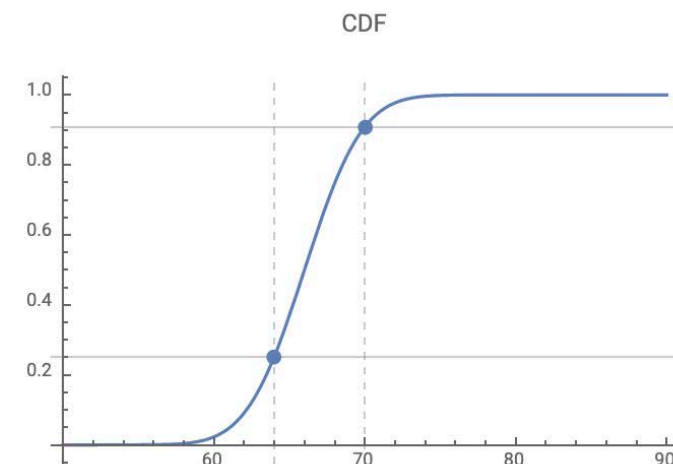
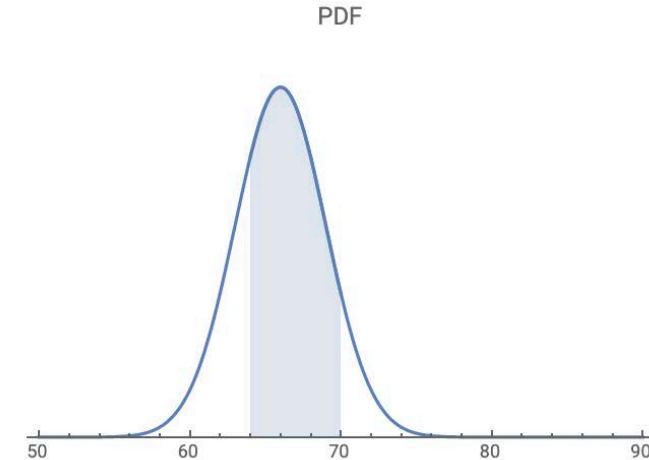
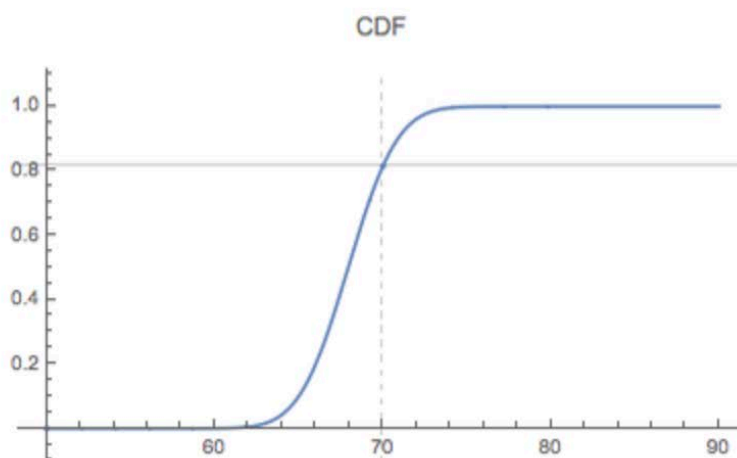
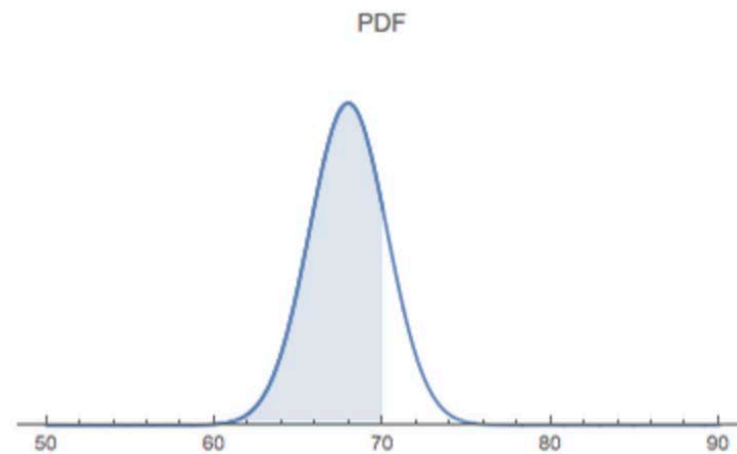
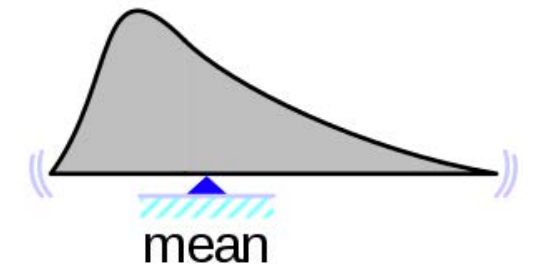
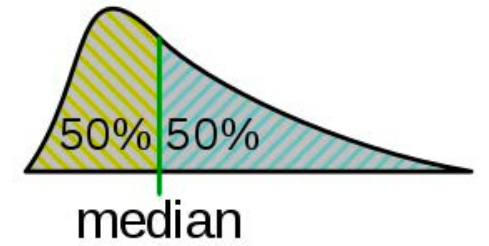
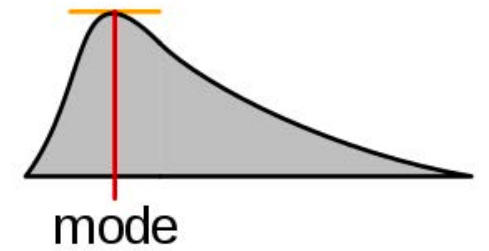
Hence, if  $F_X$  is the **cumulative distribution function** of  $X$ , then:

$$F_X(x) = \int_{-\infty}^x f_X(u) du,$$

and (if  $f_X$  is continuous at  $x$ )

$$f_X(x) = \frac{d}{dx} F_X(x).$$

Intuitively, one can think of  $f_X(x) dx$  as being the probability of  $X$  falling within the infinitesimal **interval**  $[x, x + dx]$ .



# Basic rules of probability

*Sum rule*

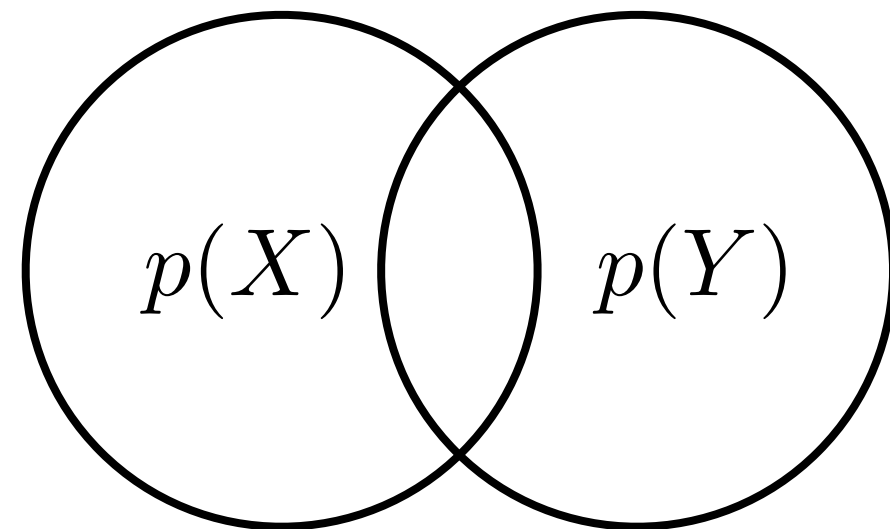
$$p(X) = \sum_Y p(X, Y)$$

*Product rule*

$$p(X, Y) = p(Y|X)p(X)$$

*Bayes rule*

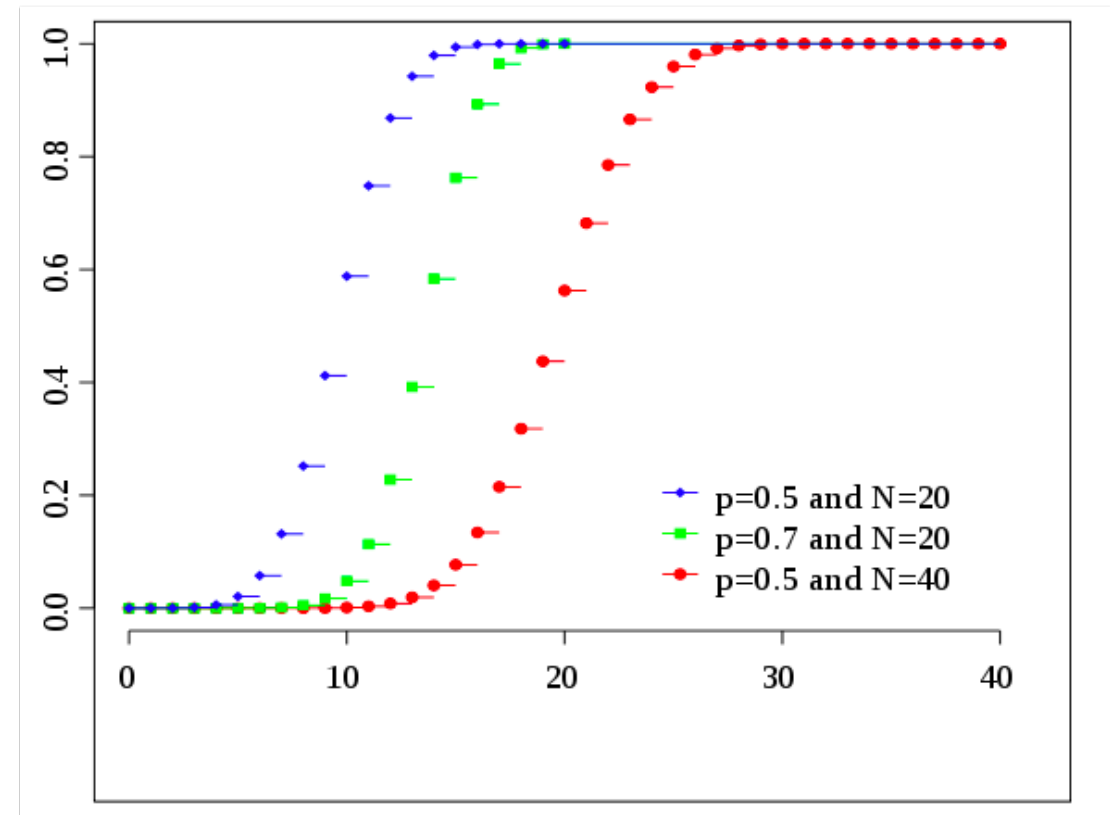
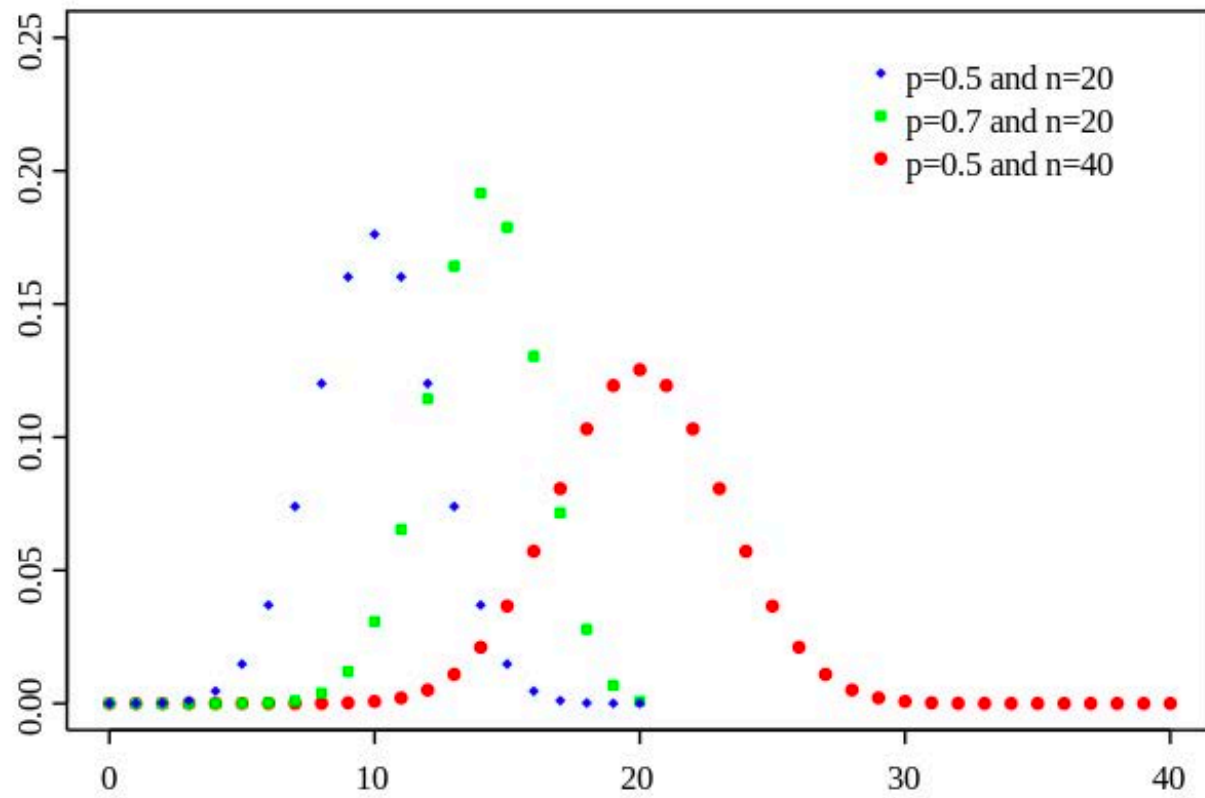
$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$



*Venn diagrams*

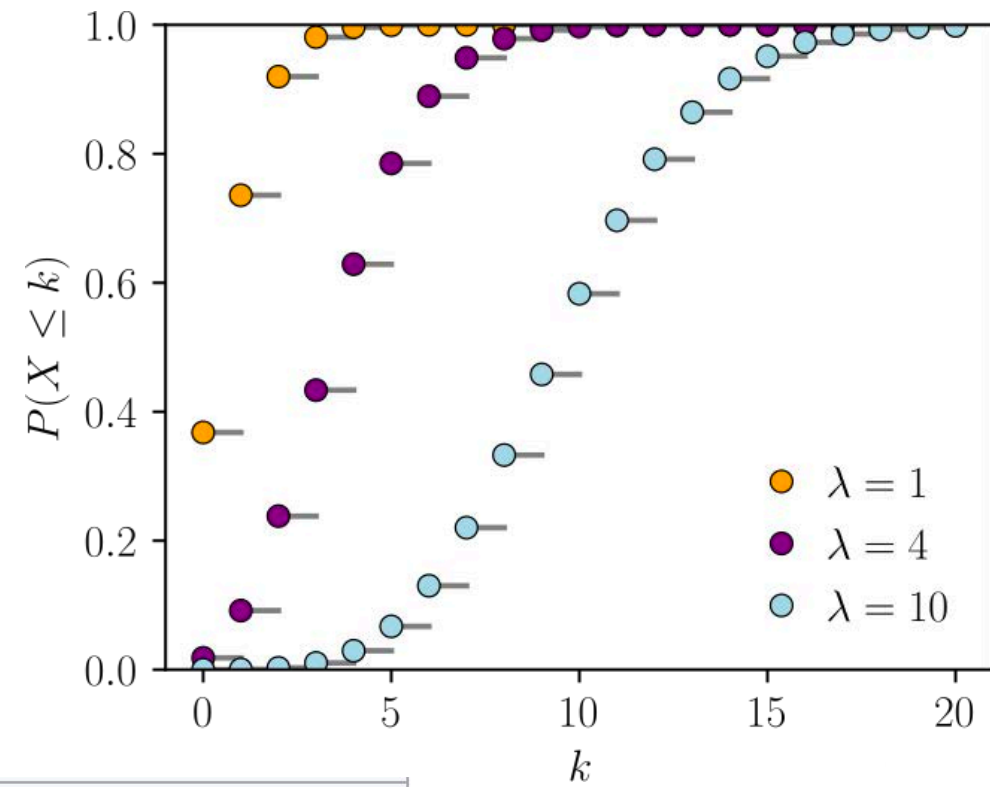
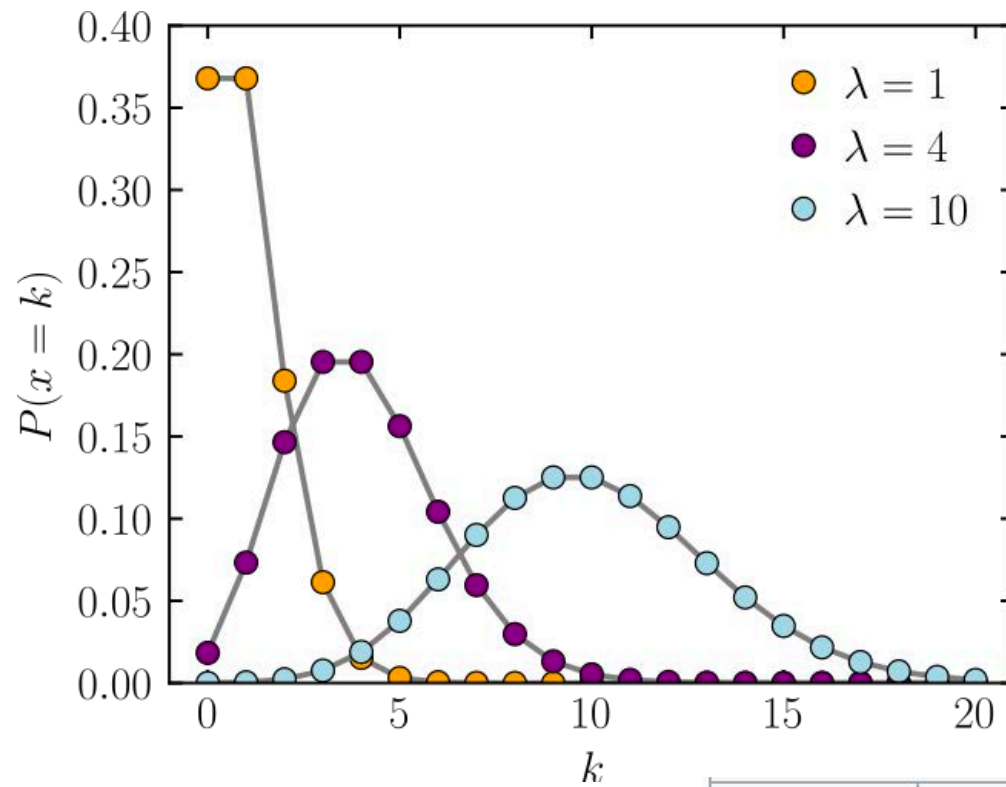


# The binomial distribution



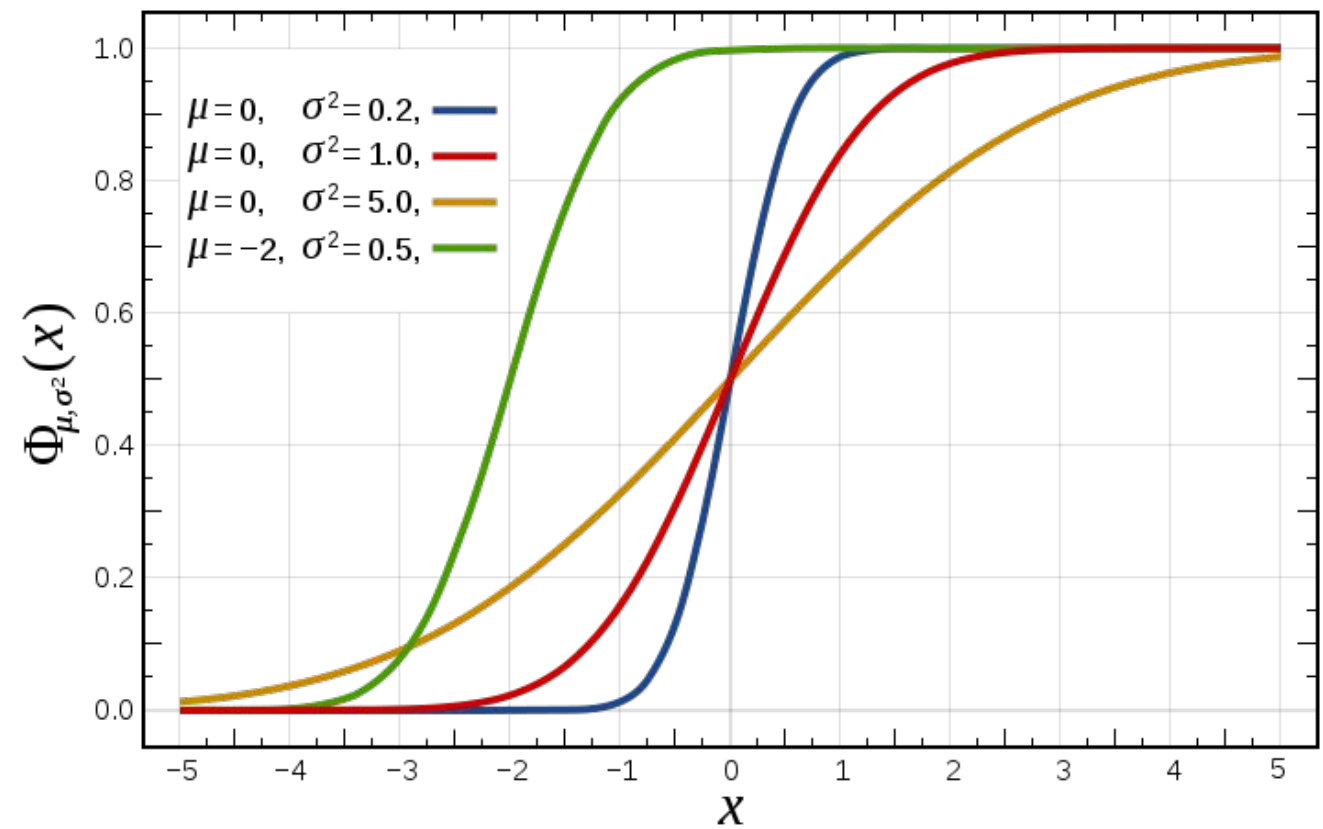
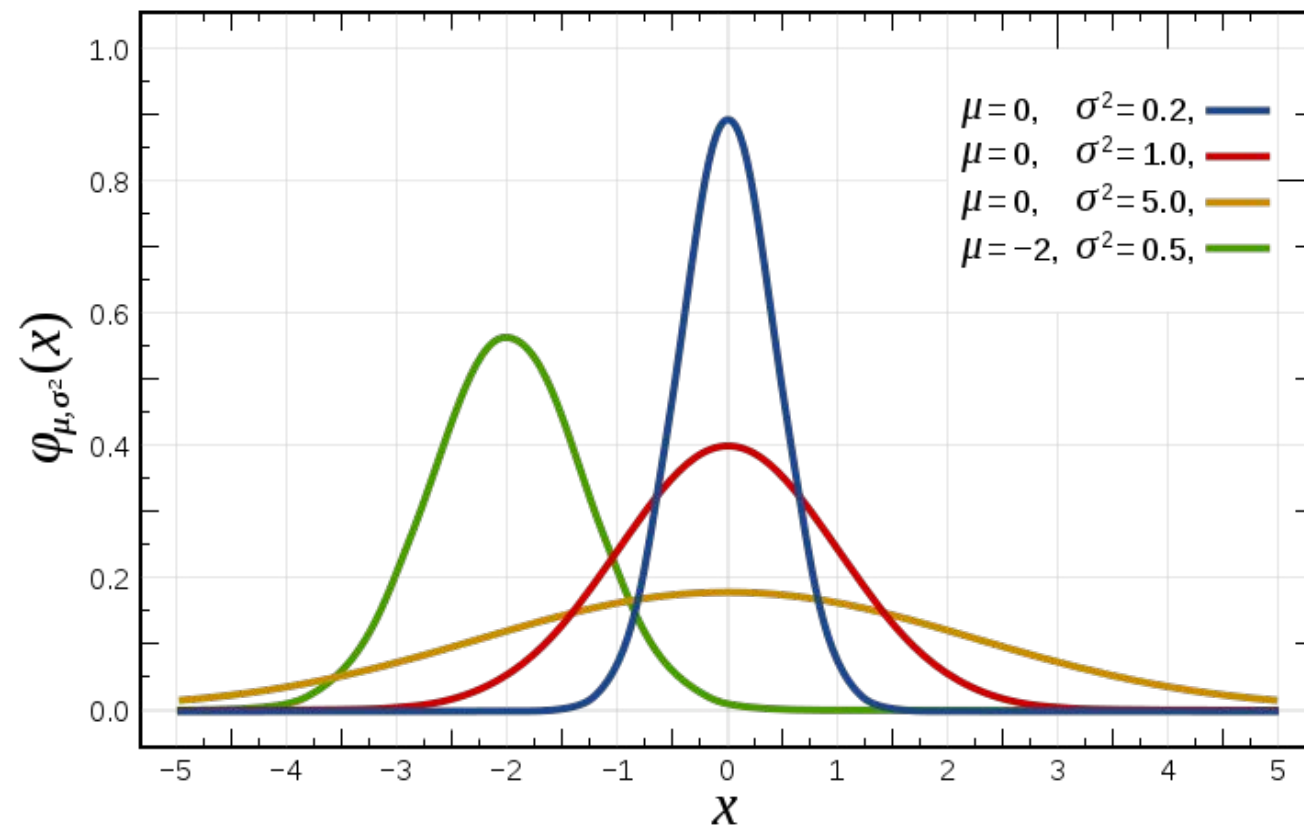
<b>Notation</b>	$B(n, p)$
<b>Parameters</b>	$n \in \{0, 1, 2, \dots\}$ – number of trials $p \in [0, 1]$ – success probability for each trial $q = 1 - p$
<b>Support</b>	$k \in \{0, 1, \dots, n\}$ – number of successes
<b>PMF</b>	$\binom{n}{k} p^k q^{n-k}$
<b>CDF</b>	$I_q(n - k, 1 + k)$
<b>Mean</b>	$np$
<b>Median</b>	$\lfloor np \rfloor$ or $\lceil np \rceil$
<b>Mode</b>	$\lfloor (n + 1)p \rfloor$ or $\lceil (n + 1)p \rceil - 1$
<b>Variance</b>	$npq$

# The Poisson distribution



<b>Notation</b>	$\text{Pois}(\lambda)$
<b>Parameters</b>	$\lambda \in (0, \infty)$ (rate)
<b>Support</b>	$k \in \mathbb{N}_0$ (Natural numbers starting from 0)
<b>PMF</b>	$\frac{\lambda^k e^{-\lambda}}{k!}$
<b>CDF</b>	$\frac{\Gamma(\lfloor k+1 \rfloor, \lambda)}{\lfloor k \rfloor!}, \text{ or } e^{-\lambda} \sum_{i=0}^{\lfloor k \rfloor} \frac{\lambda^i}{i!}, \text{ or}$ $Q(\lfloor k+1 \rfloor, \lambda)$ <p>(for <math>k \geq 0</math>, where <math>\Gamma(x, y)</math> is the <a href="#">upper incomplete gamma function</a>, <math>\lfloor k \rfloor</math> is the <a href="#">floor function</a>, and Q is the <a href="#">regularized gamma function</a>)</p>
<b>Mean</b>	$\lambda$
<b>Median</b>	$\approx \lfloor \lambda + 1/3 - 0.02/\lambda \rfloor$
<b>Mode</b>	$\lceil \lambda \rceil - 1, \lfloor \lambda \rfloor$
<b>Variance</b>	$\lambda$

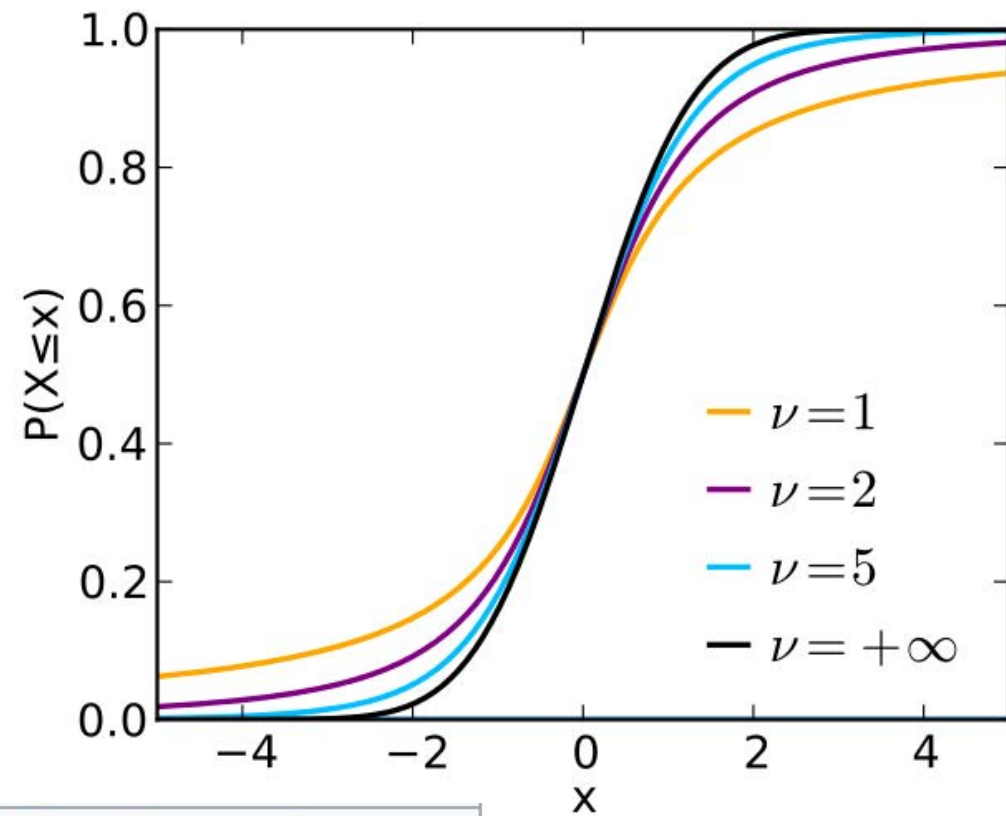
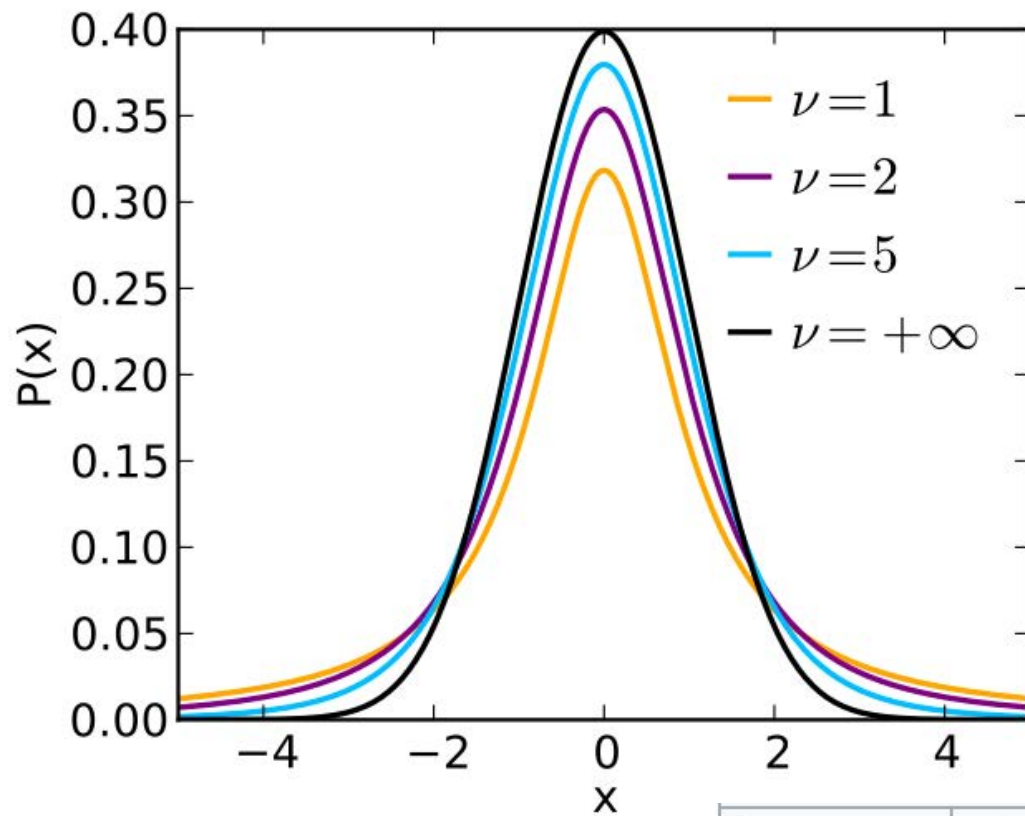
# The Gaussian distribution



<b>Notation</b>	$\mathcal{N}(\mu, \sigma^2)$
<b>Parameters</b>	$\mu \in \mathbb{R}$ = mean ( <b>location</b> ) $\sigma^2 > 0$ = variance (squared <b>scale</b> )
<b>Support</b>	$x \in \mathbb{R}$
<b>PDF</b>	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
<b>CDF</b>	$\frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right]$
<b>Quantile</b>	$\mu + \sigma\sqrt{2} \operatorname{erf}^{-1}(2F - 1)$
<b>Mean</b>	$\mu$
<b>Median</b>	$\mu$
<b>Mode</b>	$\mu$
<b>Variance</b>	$\sigma^2$

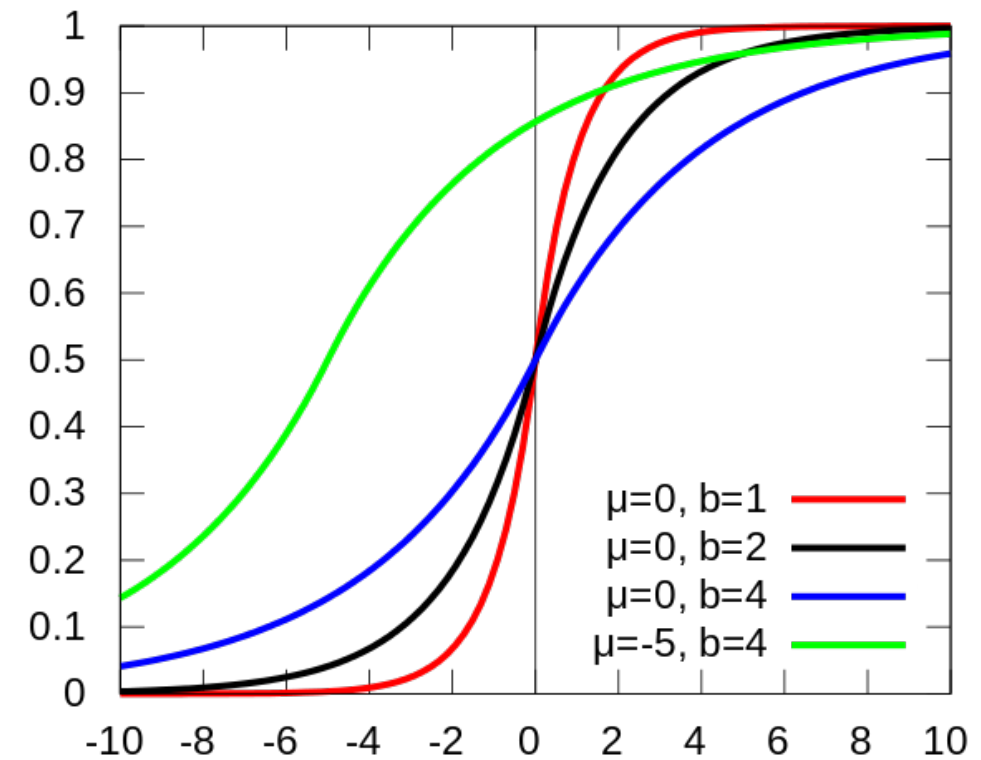
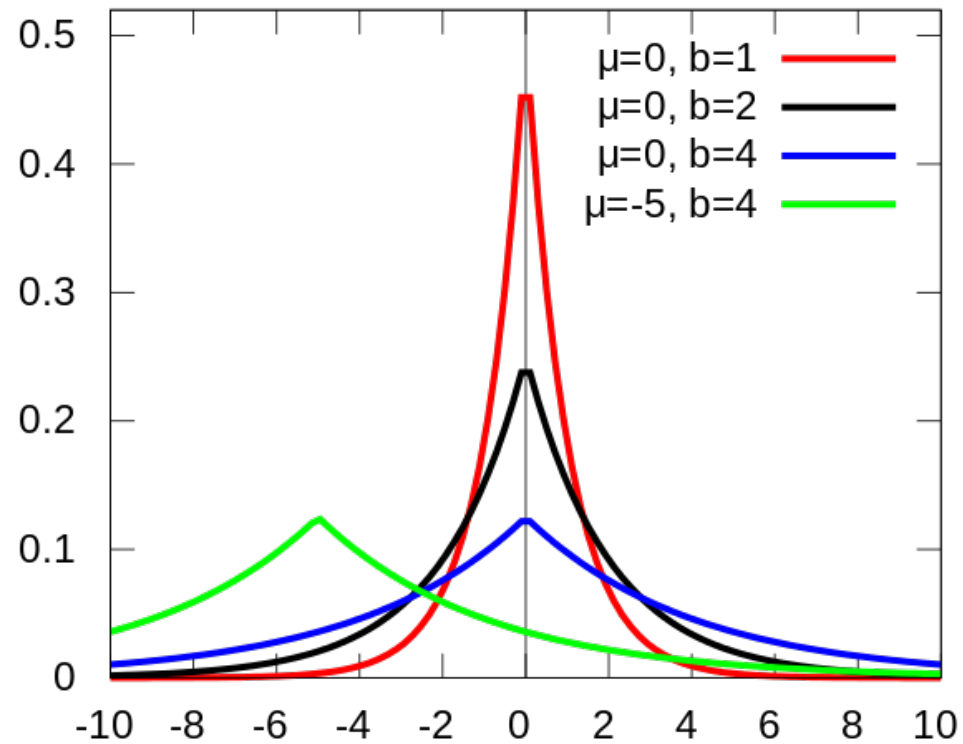


# The Student-t distribution



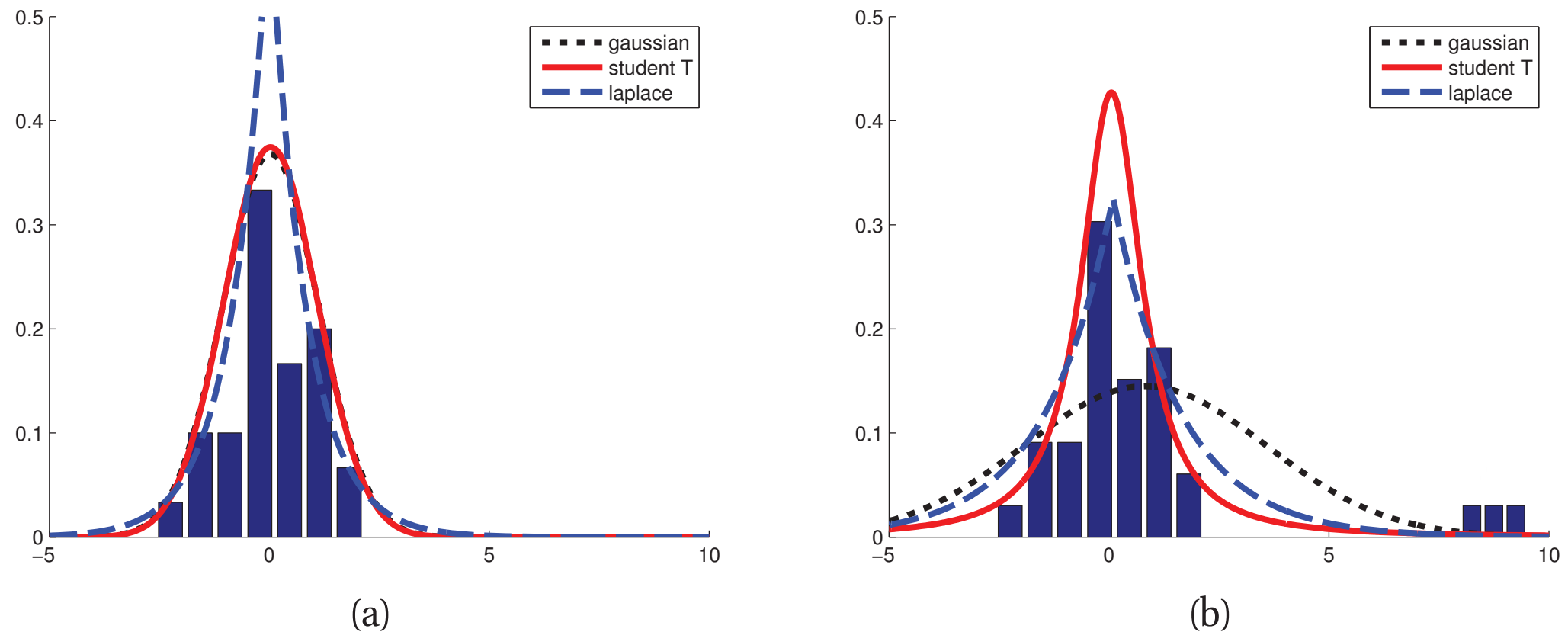
<b>Parameters</b>	$\nu > 0$ <a href="#">degrees of freedom</a> (real)
<b>Support</b>	$x \in (-\infty; +\infty)$
<b>PDF</b>	$\frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi} \Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$
<b>CDF</b>	$\frac{1}{2} + x \Gamma\left(\frac{\nu+1}{2}\right) \times$ $\frac{{}_2F_1\left(\frac{1}{2}, \frac{\nu+1}{2}; \frac{3}{2}; -\frac{x^2}{\nu}\right)}{\sqrt{\pi\nu} \Gamma\left(\frac{\nu}{2}\right)}$ <p>where <math>{}_2F_1</math> is the <a href="#">hypergeometric function</a></p>
<b>Mean</b>	0 for $\nu > 1$ , otherwise <a href="#">undefined</a>
<b>Median</b>	0
<b>Mode</b>	0
<b>Variance</b>	$\frac{\nu}{\nu-2}$ for $\nu > 2$ , $\infty$ for $1 < \nu \leq 2$ , otherwise <a href="#">undefined</a>

# The Laplace distribution



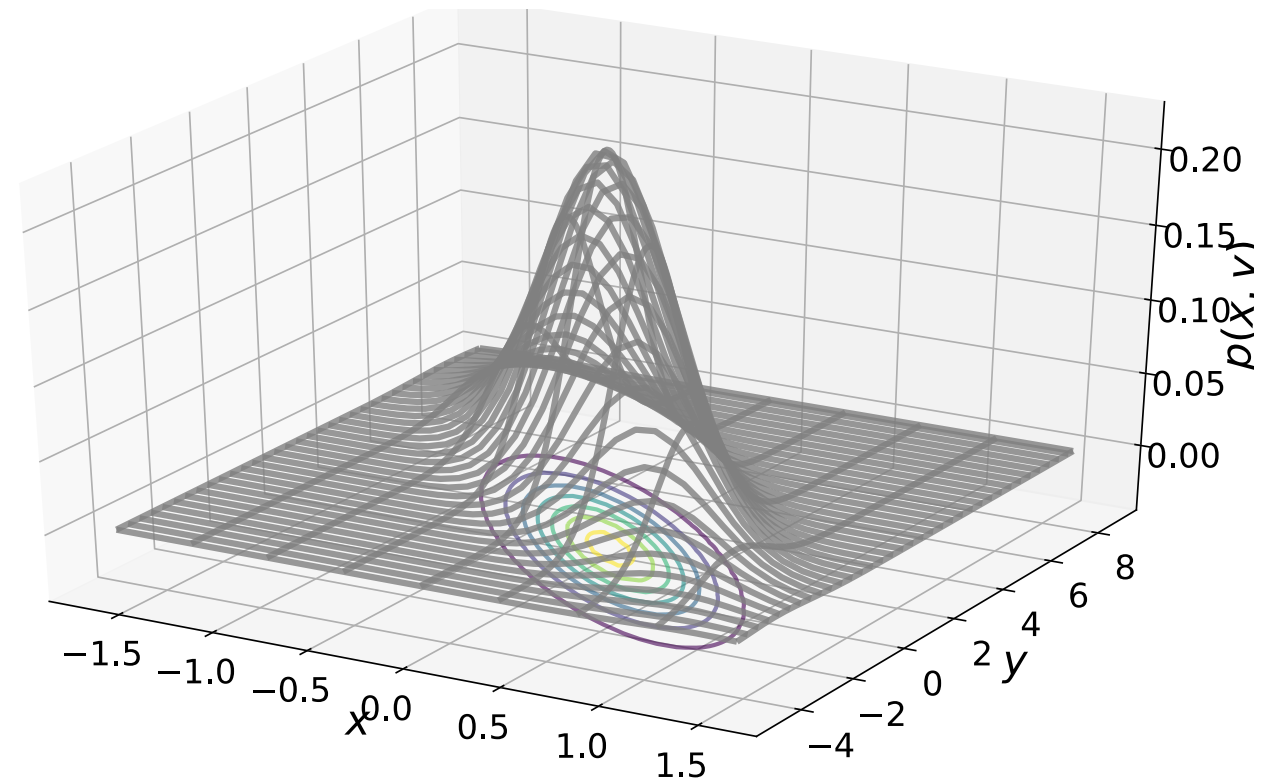
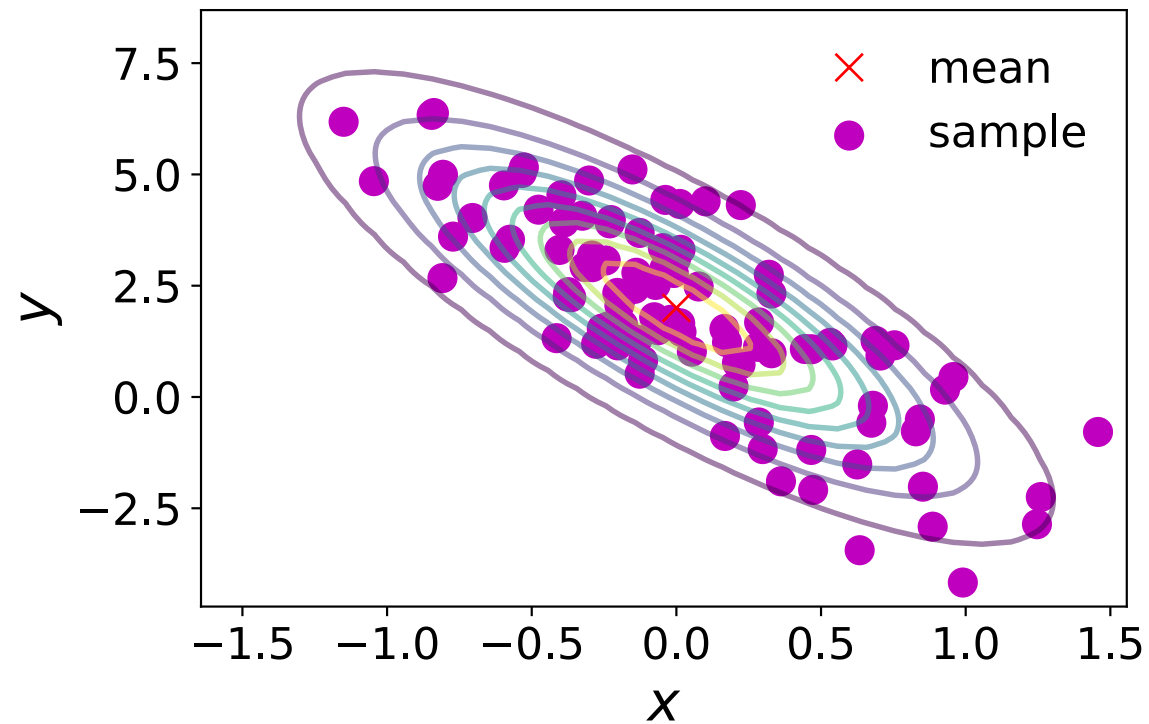
<b>Parameters</b>	$\mu$ location (real) $b > 0$ scale (real)
<b>Support</b>	$x \in (-\infty; +\infty)$
<b>PDF</b>	$\frac{1}{2b} \exp\left(-\frac{ x - \mu }{b}\right)$
<b>CDF</b>	$\begin{cases} \frac{1}{2} \exp\left(\frac{x - \mu}{b}\right) & \text{if } x < \mu \\ 1 - \frac{1}{2} \exp\left(-\frac{x - \mu}{b}\right) & \text{if } x \geq \mu \end{cases}$
<b>Mean</b>	$\mu$
<b>Median</b>	$\mu$
<b>Mode</b>	$\mu$
<b>Variance</b>	$2b^2$

# Gaussian vs Student-t vs Laplace



**Figure 2.8** Illustration of the effect of outliers on fitting Gaussian, Student and Laplace distributions. (a) No outliers (the Gaussian and Student curves are on top of each other). (b) With outliers. We see that the Gaussian is more affected by outliers than the Student and Laplace distributions. Based on Figure 2.16 of (Bishop 2006a). Figure generated by `robustDemo`.

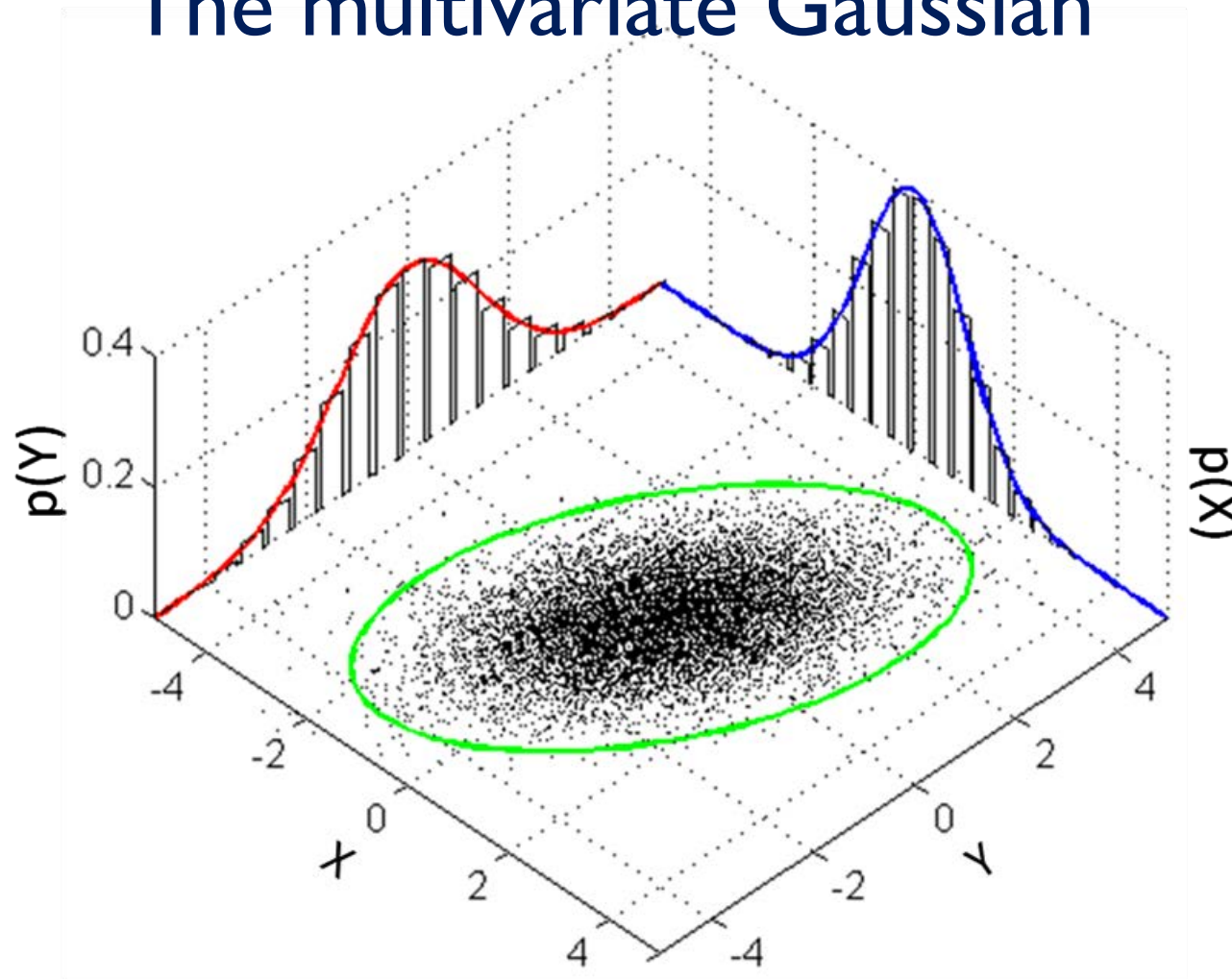
# The multivariate Gaussian



$$p(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$



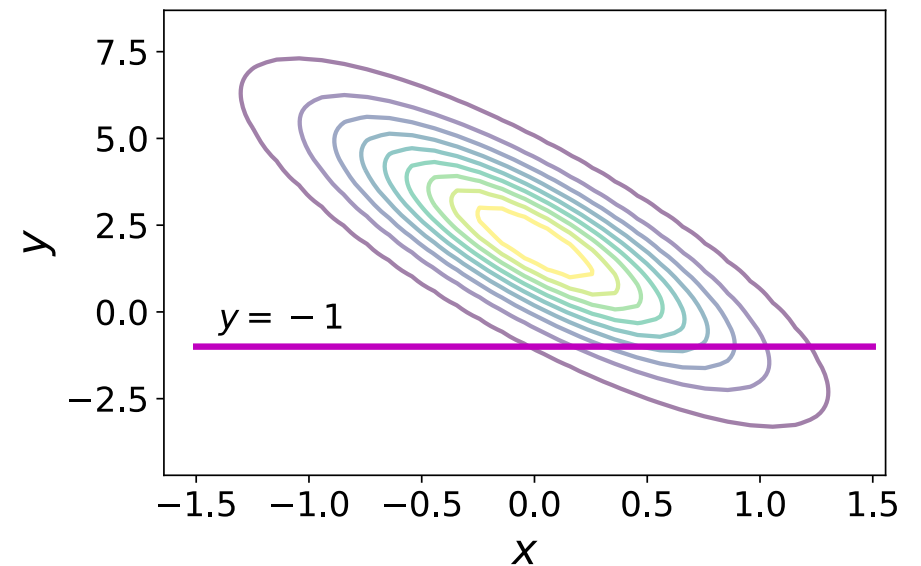
# The multivariate Gaussian



<b>Notation</b>	$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
<b>Parameters</b>	$\boldsymbol{\mu} \in \mathbf{R}^k$ — location $\boldsymbol{\Sigma} \in \mathbf{R}^{k \times k}$ — covariance (positive semi-definite matrix)
<b>Support</b>	$\mathbf{x} \in \boldsymbol{\mu} + \text{span}(\boldsymbol{\Sigma}) \subseteq \mathbf{R}^k$
<b>PDF</b>	$\det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$ , exists only when $\boldsymbol{\Sigma}$ is positive-definite
<b>Mean</b>	$\boldsymbol{\mu}$
<b>Mode</b>	$\boldsymbol{\mu}$
<b>Variance</b>	$\boldsymbol{\Sigma}$

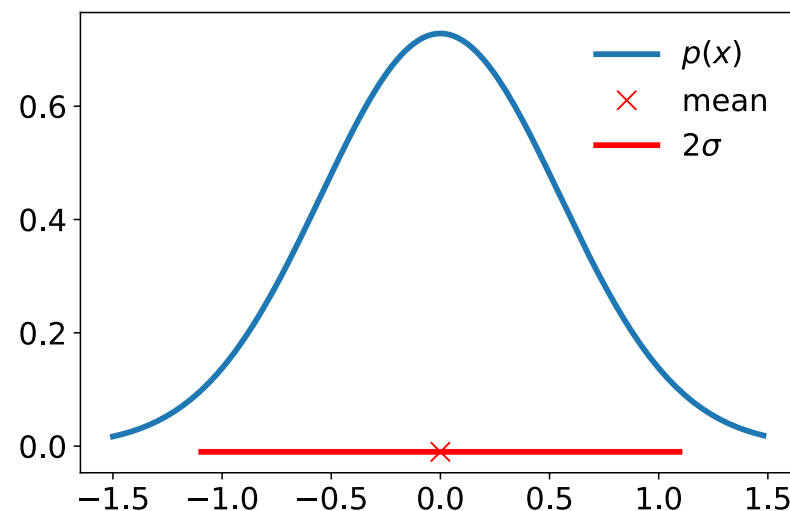
# Marginals and conditionals of a Gaussian

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left( \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$



*Marginal distribution*

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{x} \mid \mu_x, \Sigma_{xx})$$

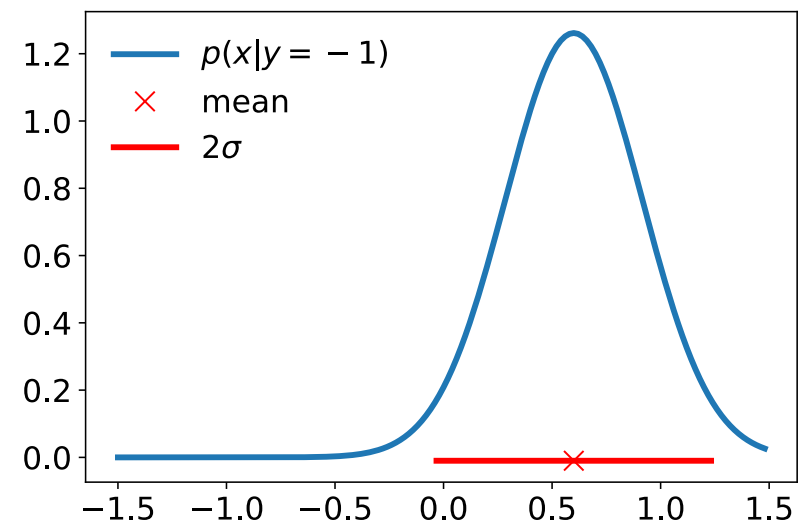


*Conditional distribution*

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(\mu_{x \mid y}, \Sigma_{x \mid y})$$

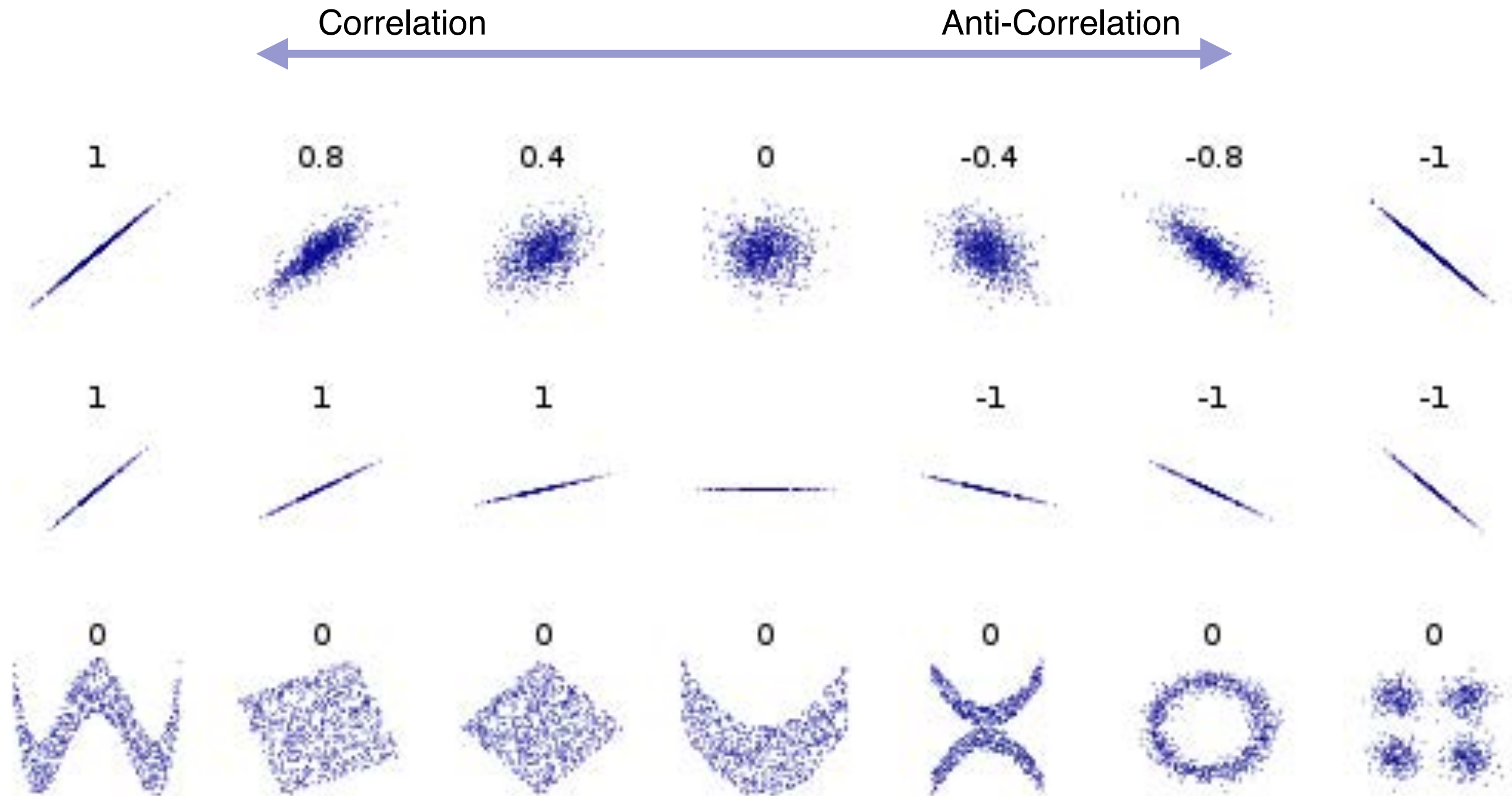
$$\mu_{x \mid y} = \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\mathbf{y} - \mu_y)$$

$$\Sigma_{x \mid y} = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}$$



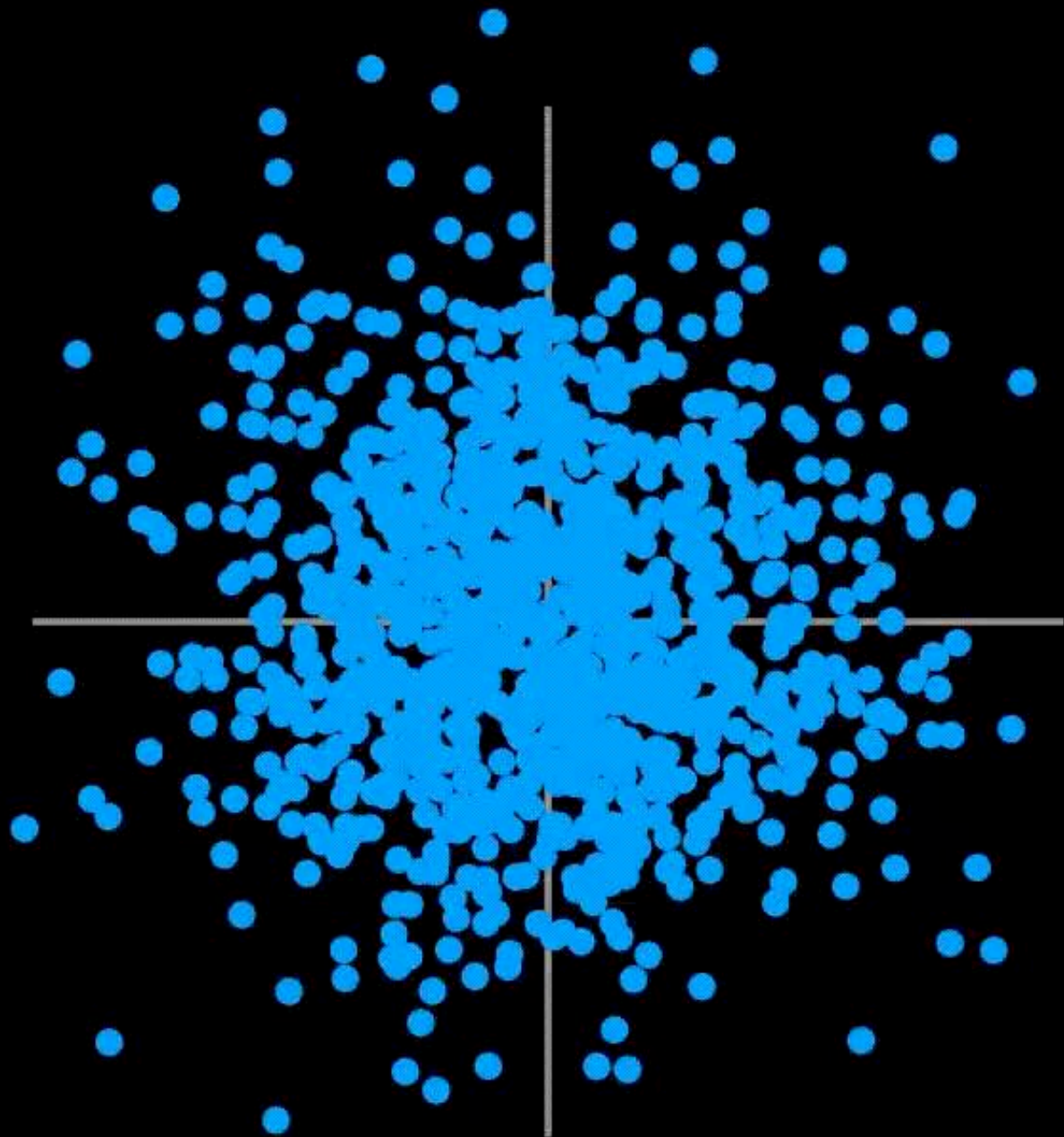
These are unique properties that make the Gaussian distribution very simple and attractive to compute with! It is essentially our main building block for computing under uncertainty.

# Correlation and linear dependence



# Covariance vs Mutual Information

$$\text{cov}(X, Y) \quad I(X; Y)$$

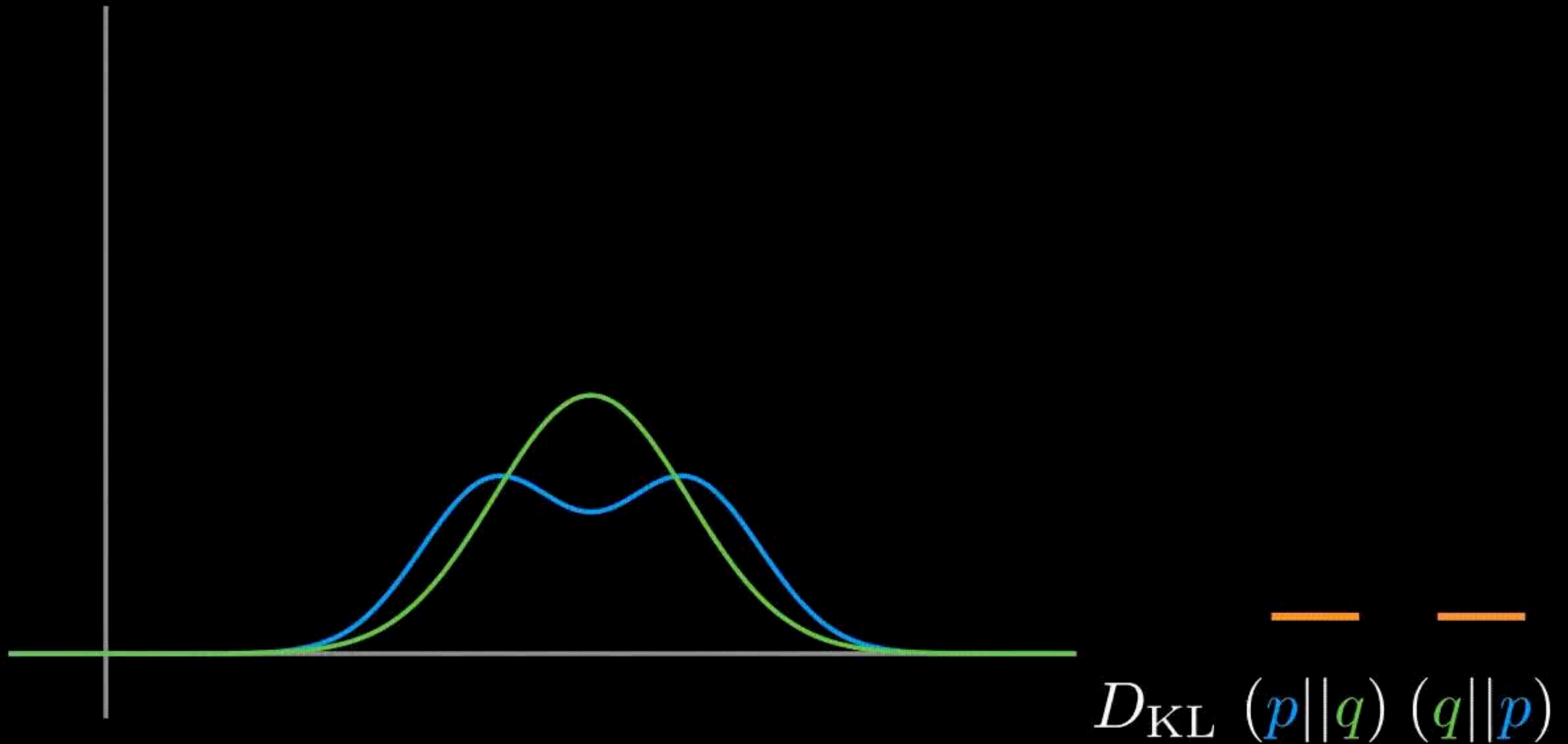


@ari\_seff

\*credit: Ari Seff (Princeton)

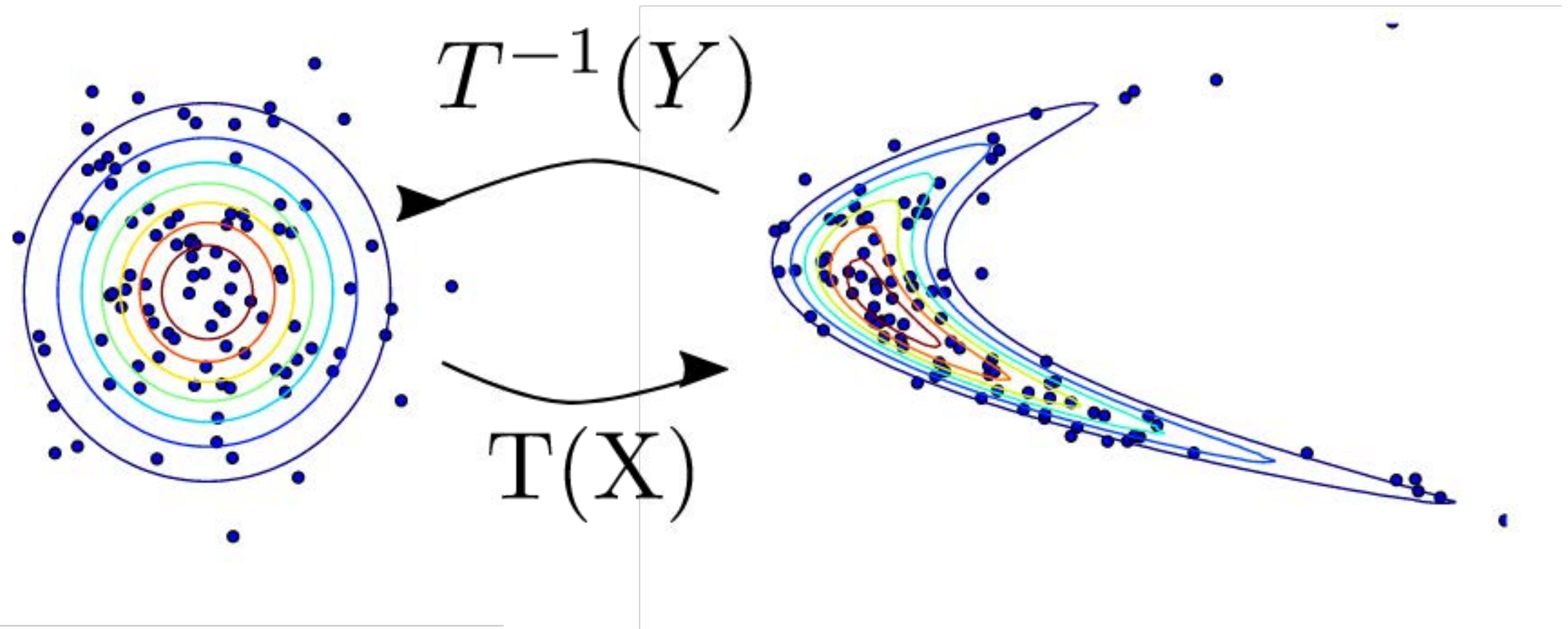


# Kullbak-Leibler divergence



\*credit: Ari Seff (Princeton)

# Transformations



# Maximum likelihood estimation

$$\theta_{\text{MLE}} = \arg \max_{\theta \in \Theta} p(\mathcal{D}|\theta)$$