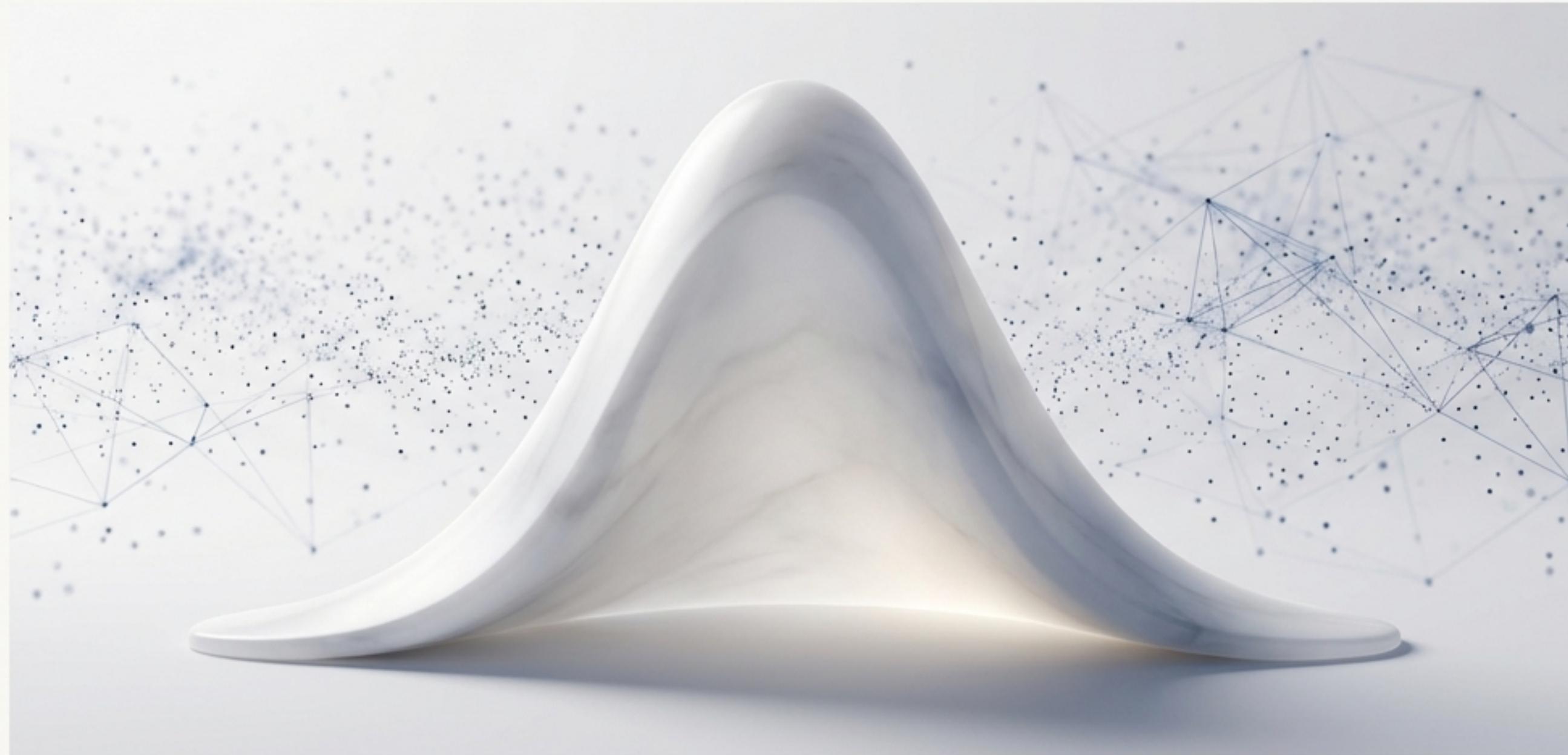


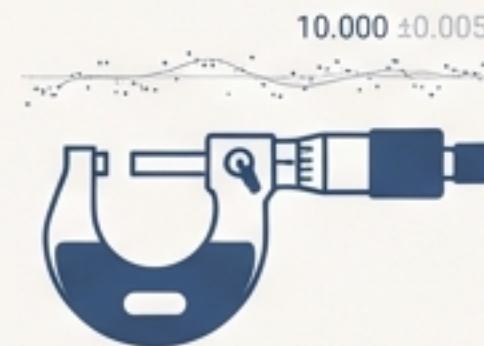
The Ubiquitous Bell

The Elegant Mathematics of Aggregate Randomness and the Cornerstone of Machine Learning



When many small, random forces combine, a predictable pattern emerges.

What do these phenomena have in common?



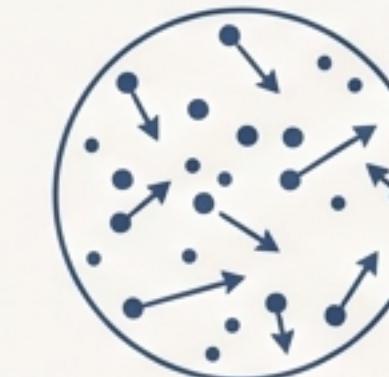
Measurement errors in a scientific experiment.



The sum of many dice rolls.



The heights of a large population.



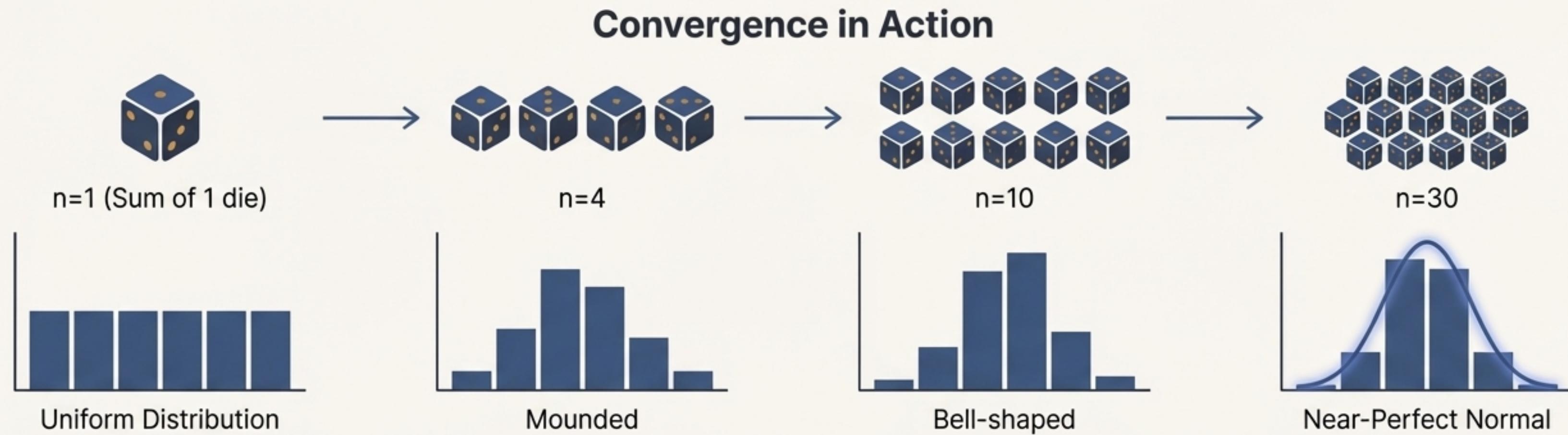
The velocity distribution of molecules in a gas.

They are all the result of many independent (or nearly independent) random effects acting together. This additive process leads to an astonishingly consistent shape.

The Central Limit Theorem: Nature's Gravitational Pull Towards the Normal

The Law: The CLT states that the distribution of the *sum* (or average) of a large number of independent, identically distributed random variables approaches a normal distribution, *regardless of the original variable's distribution*.

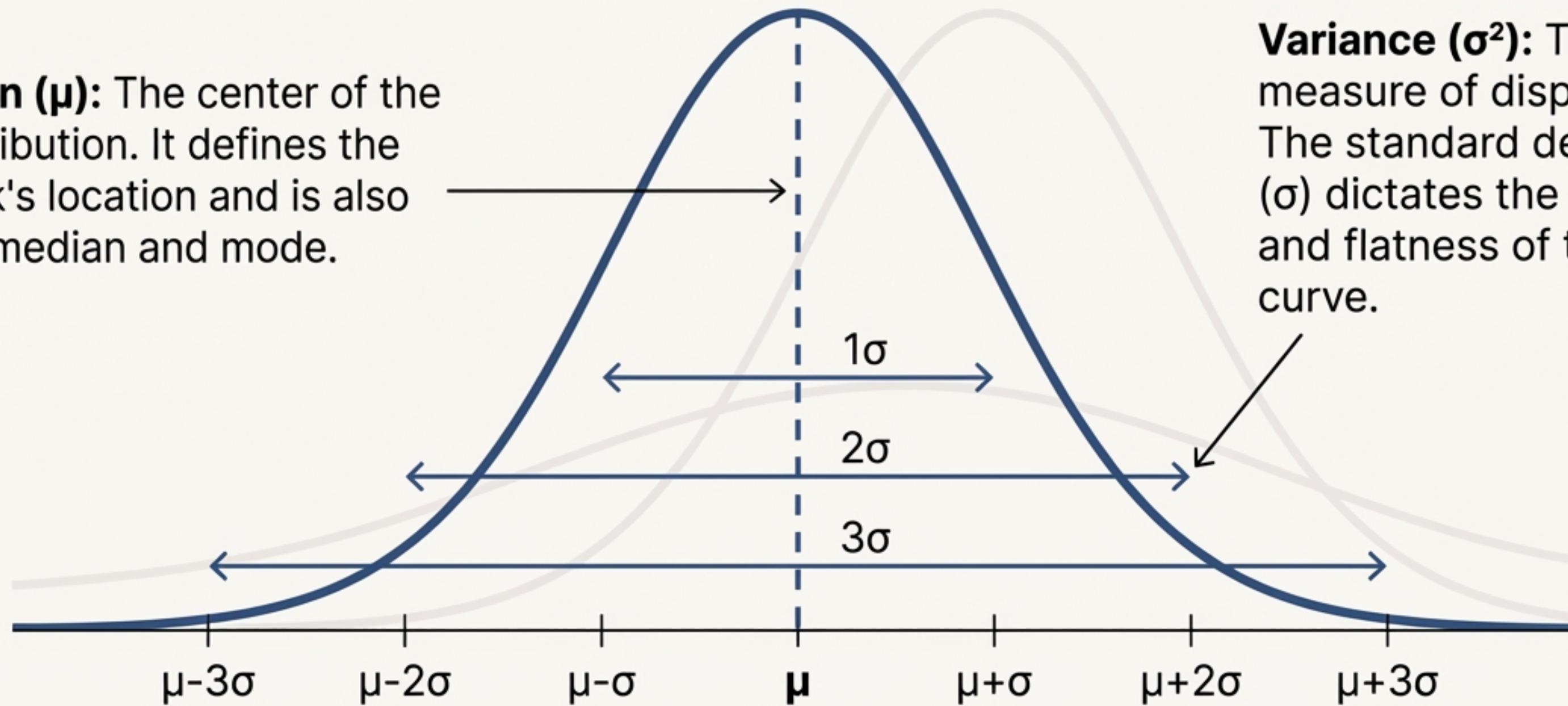
Implication: This is why physical quantities that are the sum of many independent processes, like measurement errors, so often have nearly normal distributions. It also allows us to approximate other distributions (like Binomial or Poisson) under certain conditions.



The Blueprint of the Bell: Defined by Center and Spread

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

Mean (μ): The center of the distribution. It defines the peak's location and is also the median and mode.



Variance (σ^2): The measure of dispersion. The standard deviation (σ) dictates the width and flatness of the curve.

The Standard Normal: Creating a Common Language for Probability

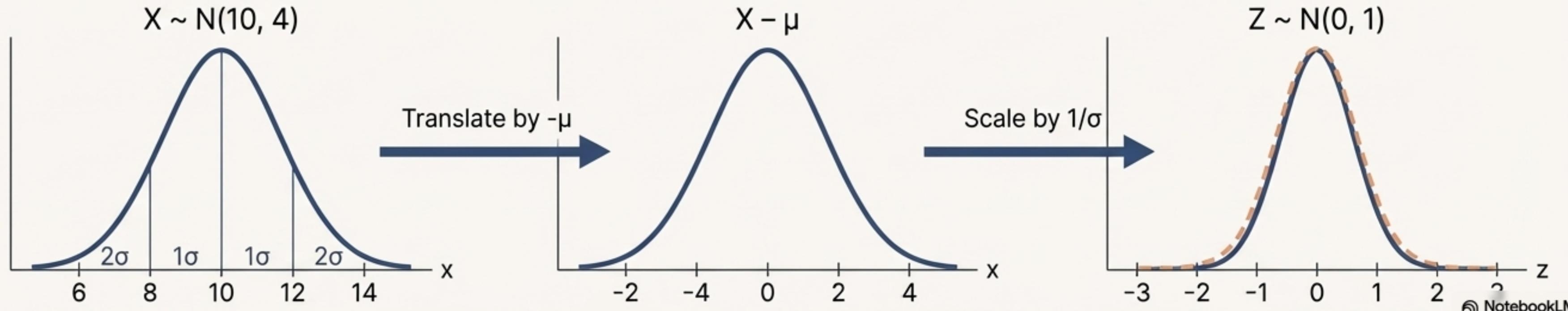
The **Standard Normal Distribution** is the specific case where $\mu=0$ and $\sigma=1$. Its PDF is often denoted by $\phi(z)$:

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

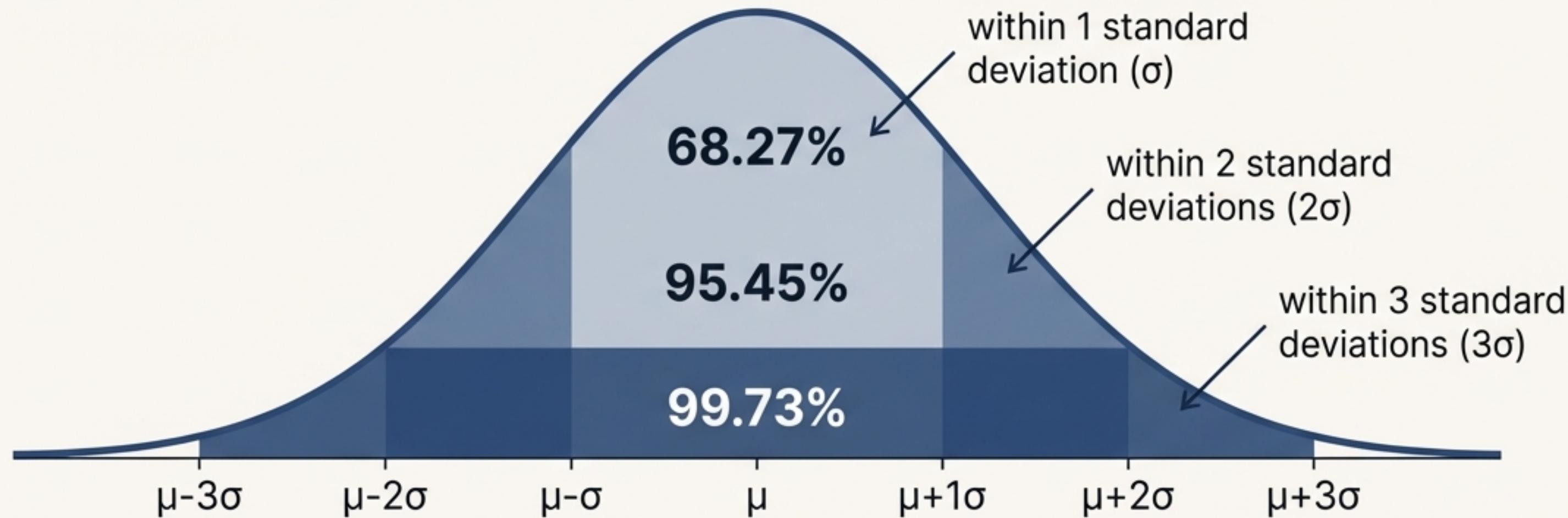
Any normally distributed variable $X \sim N(\mu, \sigma^2)$ can be converted to the standard normal $Z \sim N(0, 1)$ by calculating its **Z-score**:

$$Z = \frac{X - \mu}{\sigma}$$

The Z-score measures how many standard deviations an observation is from the mean, allowing us to compare values from different normal distributions.



The Empirical Rule: A Practical Shorthand for a Normally Distributed World



For any normal distribution, we can approximate the data coverage:

- **~68.27%** of data falls within **1 standard deviation** (σ) of the mean.
- **~95.45%** of data falls within **2 standard deviations** (2σ) of the mean.
- **~99.73%** of data falls within **3 standard deviations** (3σ) of the mean.

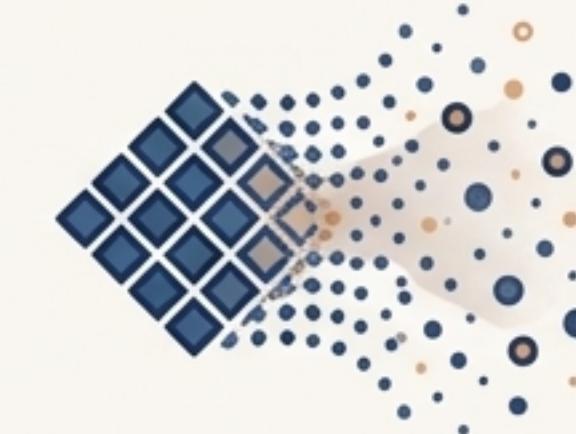
This is also known as the '3-sigma rule'.

The Properties That Make Gaussians Analytically Powerful



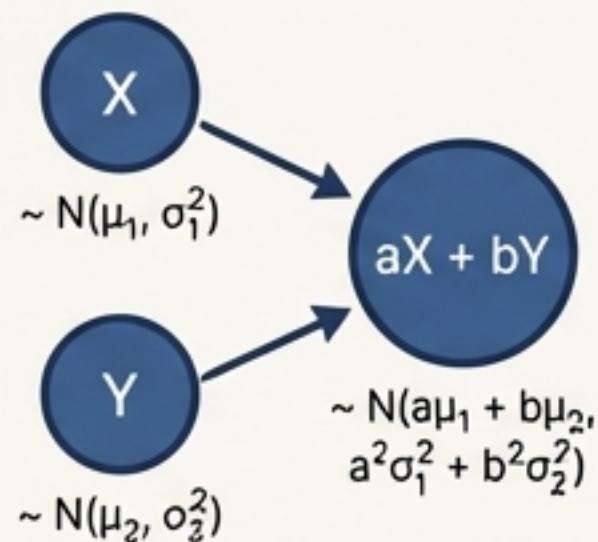
Symmetry & Moments

Perfectly symmetric around the mean. Skewness is 0, and all odd central moments are 0. Excess Kurtosis is 0.



Maximum Entropy

For a given mean and variance, the normal distribution is the continuous distribution with the maximum possible entropy. It represents the “most random” choice under these constraints.



Closure Under Linearity

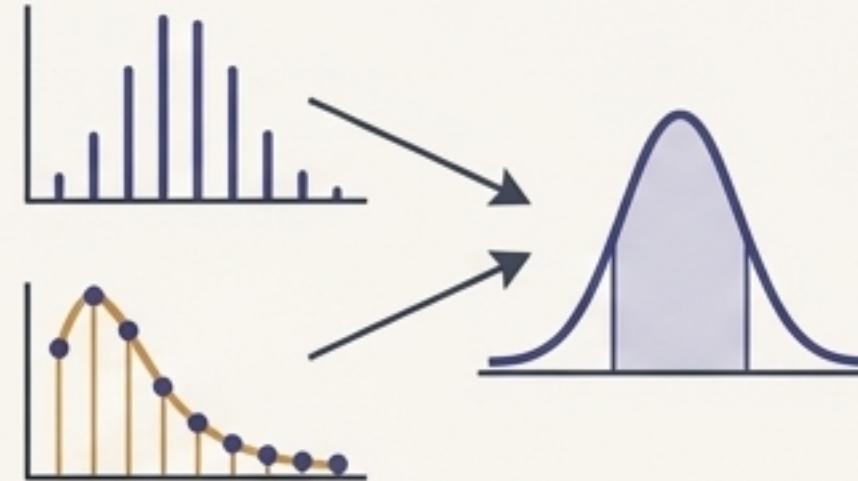
Any linear combination of independent normal variables is also a normal variable. If $X \sim N(\mu_1, \sigma_1^2)$ and $Y \sim N(\mu_2, \sigma_2^2)$ are independent, then $aX + bY \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2)$.



Independence

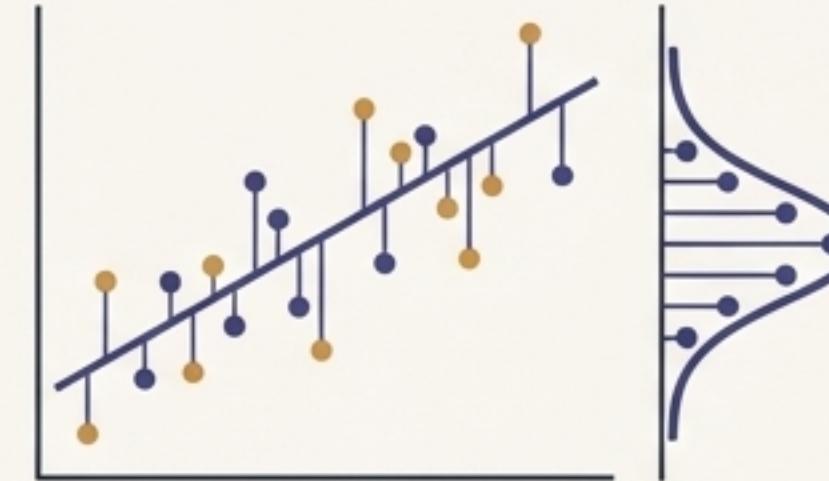
Uncorrelated *jointly* normal variables are independent. This is a unique and powerful property not shared by other distributions.

The Three Realms of Application



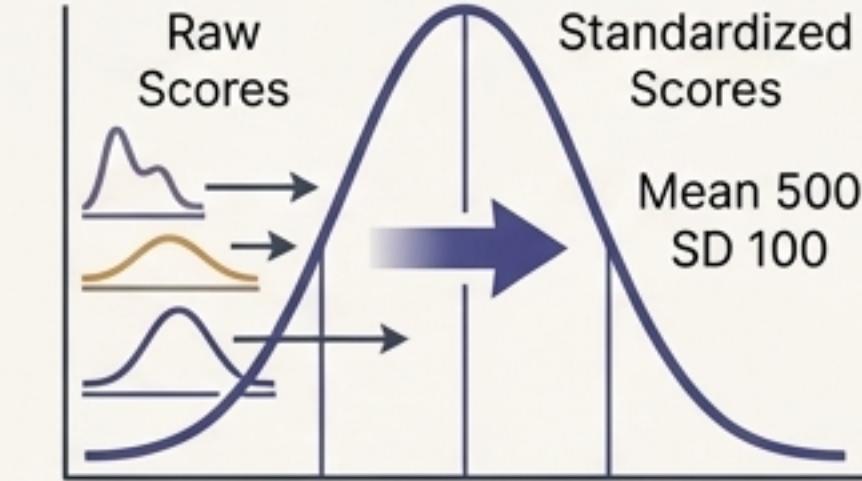
1. Approximate Normality

Used to model real-world phenomena where the Central Limit Theorem applies (e.g., population statistics, noise in signals). Also used to approximate other distributions distributions like Binomial for large n and Poisson for large λ .



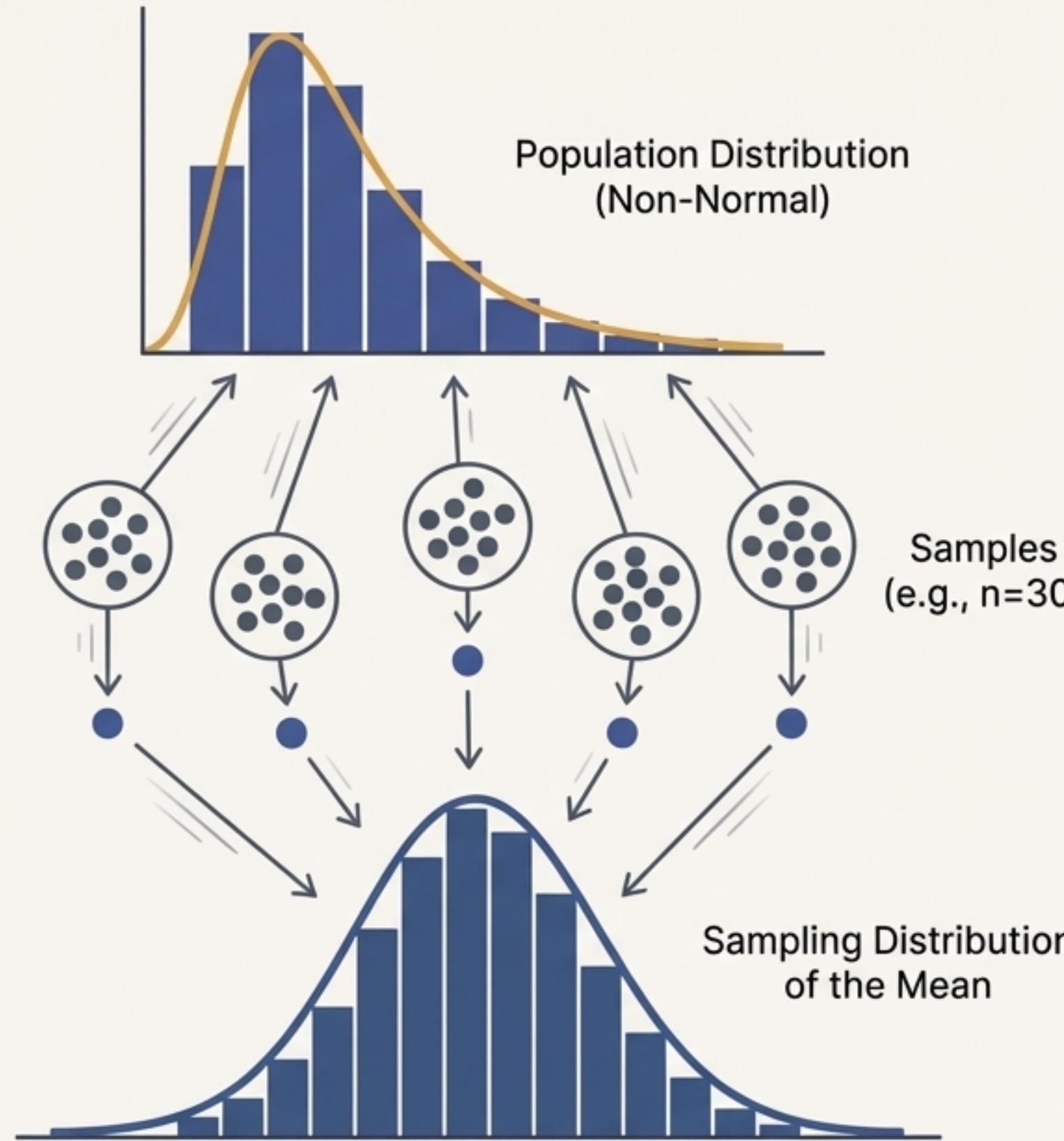
2. Assumed Normality

A powerful modeling choice for errors that are the sum of many small effects (e.g., measurement errors, residuals in linear regression). This assumption is the justification for the method of least squares.



3. Constructed Normality

In standardized testing (like IQ or SAT scores), raw scores are often transformed to fit a predefined normal distribution (e.g., mean 500, SD 100 for the SAT) for comparability and grading.



The Bedrock of Statistical Inference

Because of the Central Limit Theorem, the distribution of *sample means* from any population (with finite variance) will be approximately normal, given a large enough sample size.

This fundamental fact is what enables modern statistical inference:

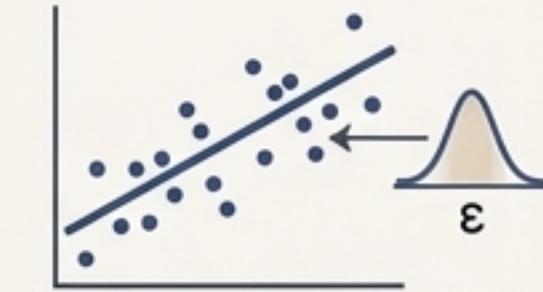
- **Confidence Intervals:** We can construct a range of plausible values for a population parameter (like the mean) because we know the distribution of our sample estimate.
- **Hypothesis Testing (z-tests, t-tests):** We can quantify the probability that an observed effect is merely due to random sampling chance, because we have a baseline (the normal curve) to compare it against.

A Fundamental Building Block for Machine Learning Models



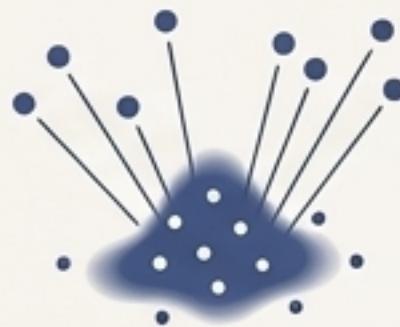
Density Estimation

Gaussian Mixture Models (GMMs) represent complex, multimodal distributions as a weighted sum of multiple simpler Gaussian distributions.



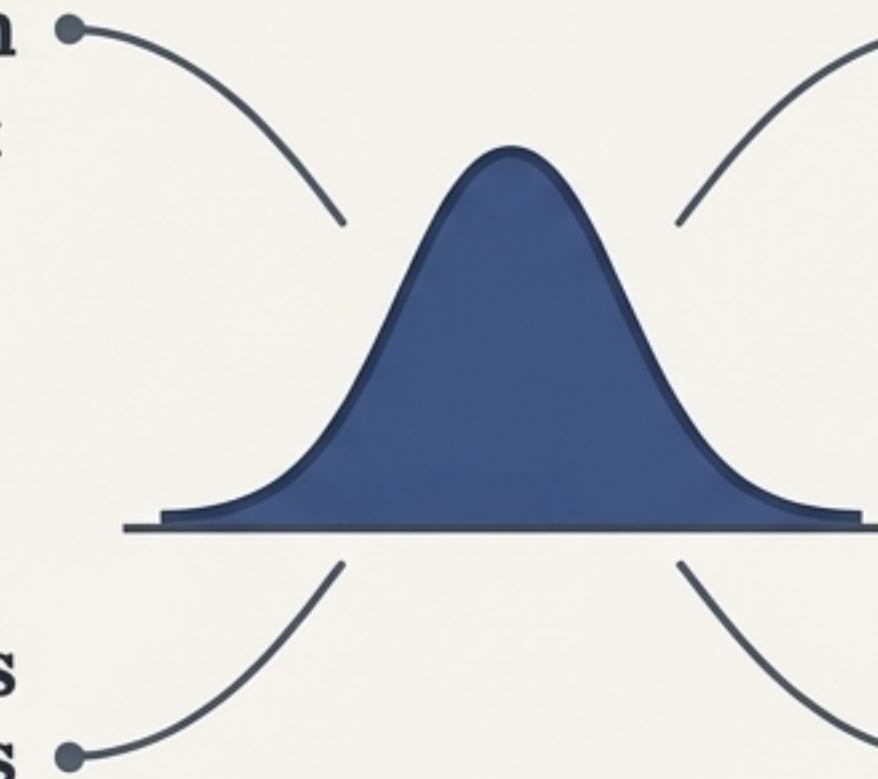
Regression

In Bayesian Linear Regression, it serves as both the prior distribution for model parameters and the likelihood function for the noise term ($y = w^T x + \epsilon$, where $\epsilon \sim N(0, \sigma^2)$)



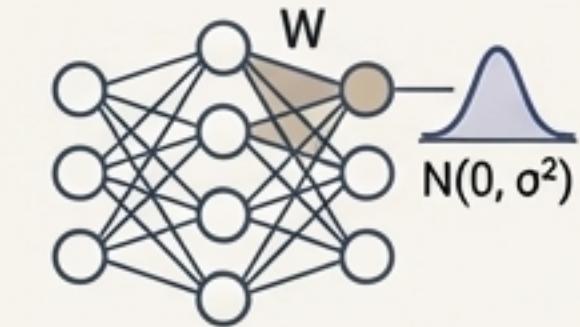
Generative Models & Latent Variables

Probabilistic PCA and Variational Autoencoders (VAEs) often use a Gaussian distribution to model the latent (unobserved) variables that generate the data.



Initialization

In deep learning, neural network weights are frequently initialized by drawing from a normal distribution (e.g., $N(0, \sigma^2)$) to break symmetry and aid optimization.

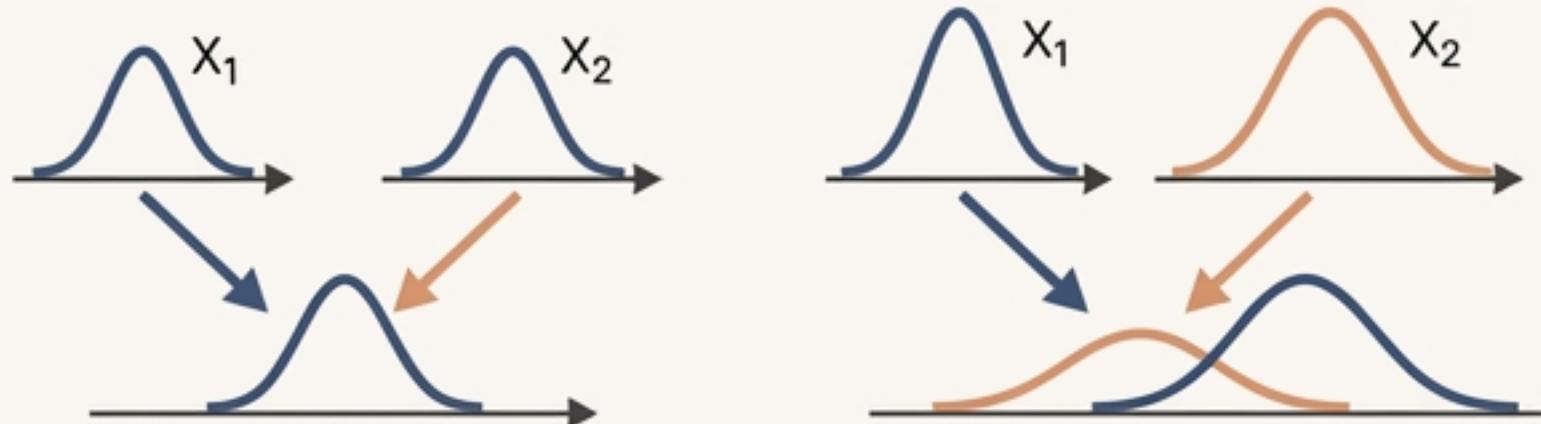


A Predictable and Stable Toolkit

Sum of Normals

If $X_1 \sim N(\mu_1, \sigma_1^2)$ and $X_2 \sim N(\mu_2, \sigma_2^2)$ are independent, then:

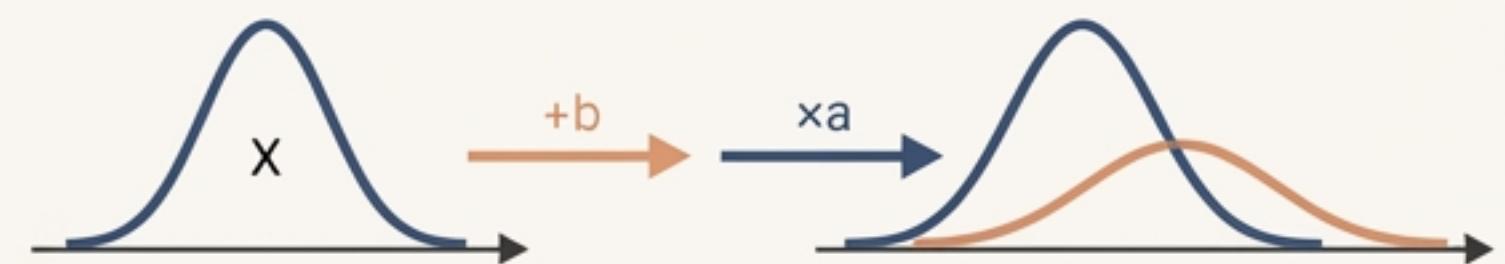
$$X_1 + X_2 \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$



Linear Transformation (Scaling & Shifting)

If $X \sim N(\mu, \sigma^2)$, then for constants a and b :

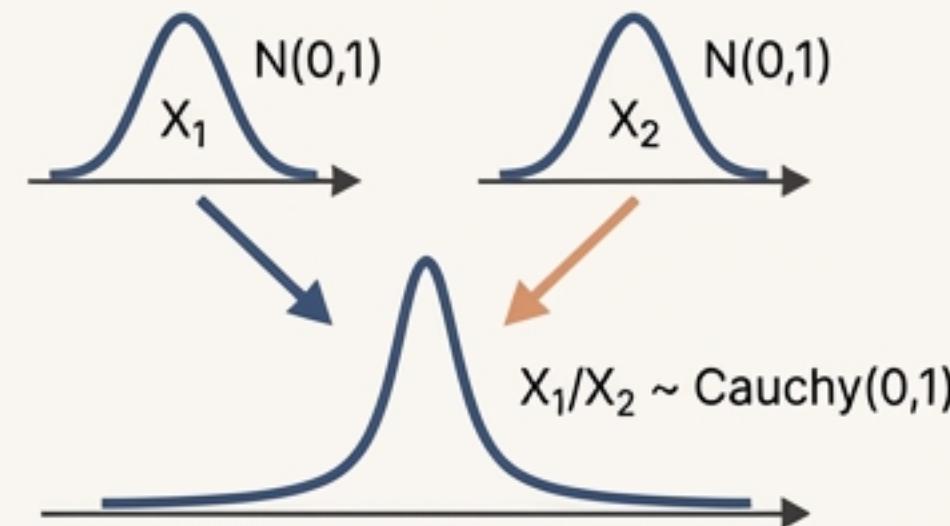
$$aX + b \sim N(a\mu + b, a^2\sigma^2)$$



Product and Ratio of Standard Normals

If X_1 and X_2 are independent standard normals, their product $Z = X_1 X_2$ follows a product distribution, while their ratio follows the standard Cauchy distribution.

$$\frac{X_1}{X_2} \sim \text{Cauchy}(0, 1)$$



Bringing the Normal Distribution to Life in Code

Generating Random Samples

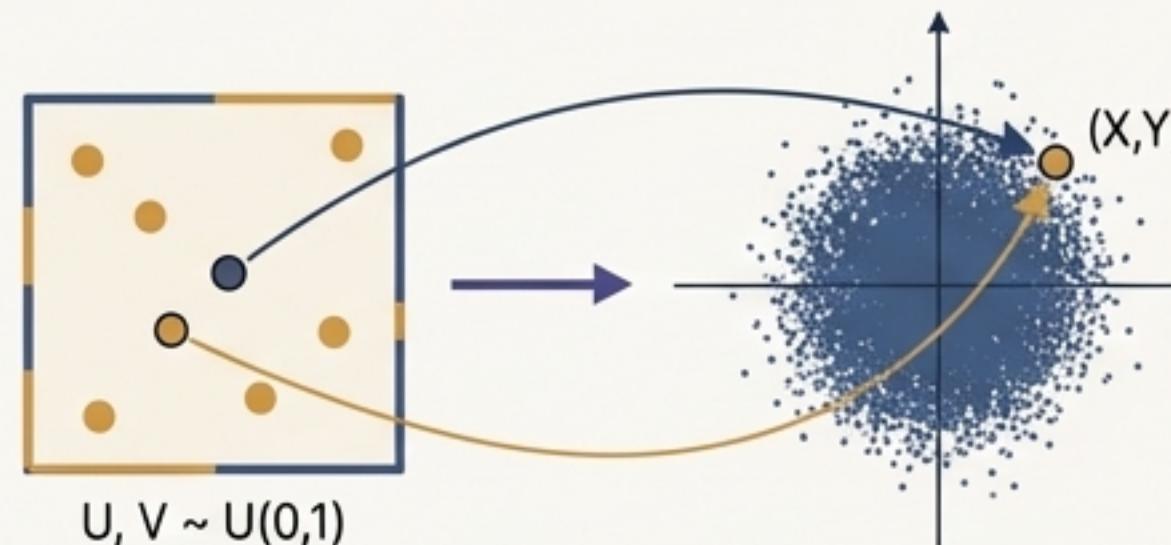
How do we generate samples?

We transform uniform random numbers.

Box-Muller Transform

A classic method using two uniform random numbers $U, V \sim U(0,1)$ to generate two independent standard normal samples.

$$X = \sqrt{-2\ln U} \cos(2\pi V)$$



Ziggurat Algorithm: A more modern, faster rejection sampling method used in many standard libraries.

Calculating the CDF

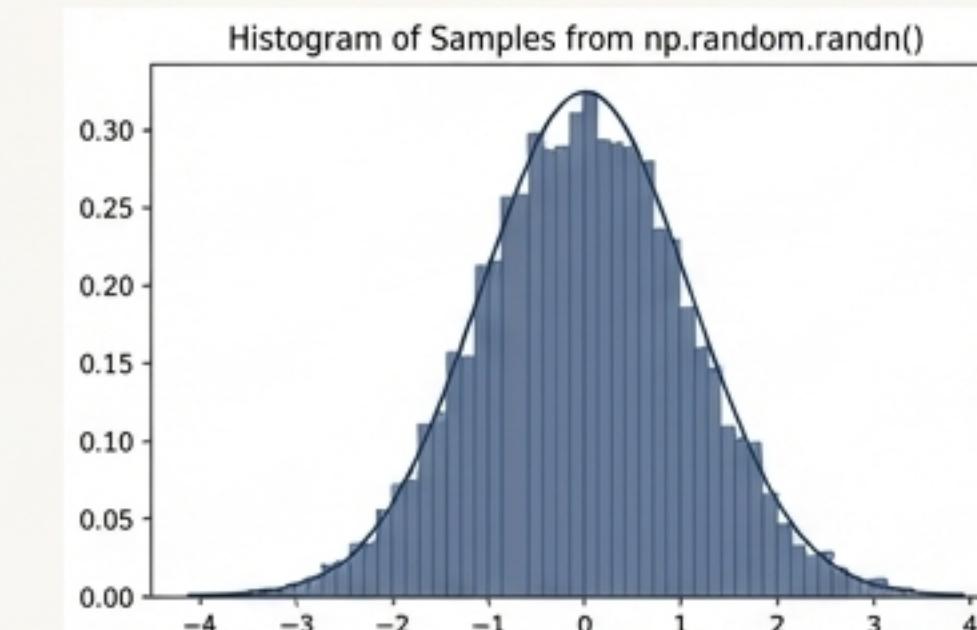
How do we calculate probabilities?

The integral for the Cumulative Distribution Function (CDF) $\Phi(x)$ has no closed-form solution. We rely on highly accurate numerical approximations, often related to the error function (erf).

```
import numpy as np
import matplotlib.pyplot as plt

# Generate 100,000 standard normal samples
samples = np.random.randn(100000)

# Plot a histogram to visualize the distribution
plt.hist(samples, bins=100, density=True, color='#3D5A80', alpha=0.7)
plt.title('Histogram of Samples from np.random.randn()')
plt.show()
```



The Normal is Powerful, But Not Universal

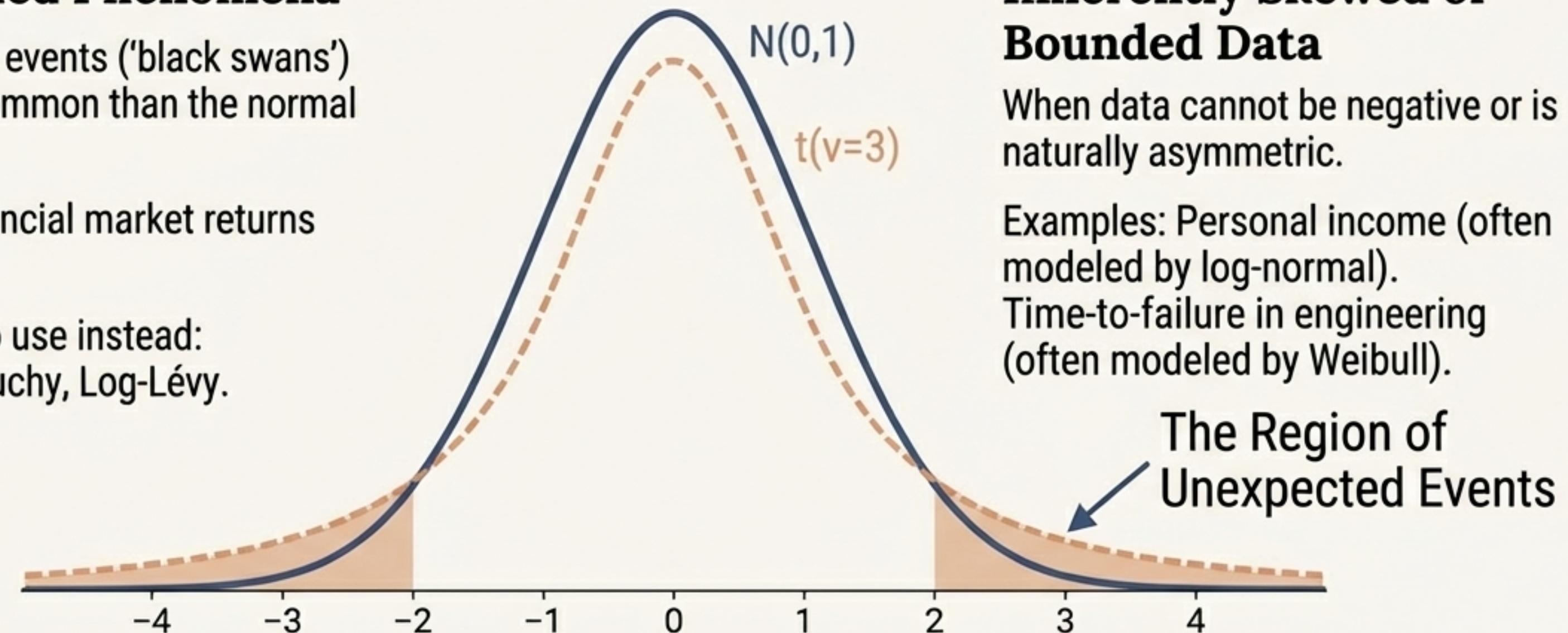
“I can only recognize the occurrence of the normal curve... as a very abnormal phenomenon.” – Karl Pearson, 1901

Heavy-Tailed Phenomena

Where extreme events ('black swans') are far more common than the normal predicts.

Examples: Financial market returns (crashes).

Distributions to use instead:
Student's t, Cauchy, Log-Lévy.



Inherently Skewed or Bounded Data

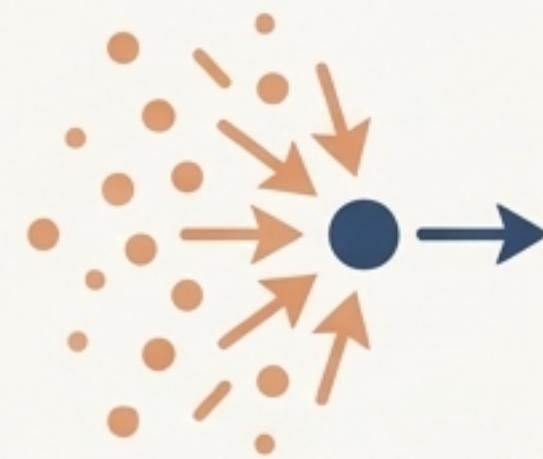
When data cannot be negative or is naturally asymmetric.

Examples: Personal income (often modeled by log-normal).
Time-to-failure in engineering (often modeled by Weibull).

The Region of Unexpected Events

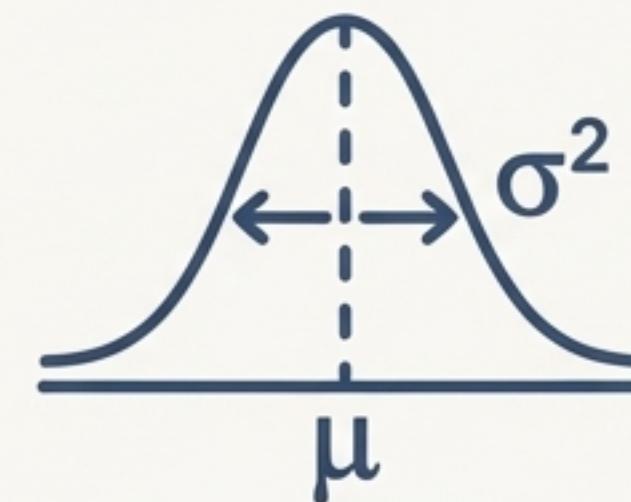
The Three Pillars of the Normal Distribution

THE LAW: A Universal Endpoint



The Central Limit Theorem establishes it as the natural result of summing many independent random effects, explaining its ubiquity.

THE BLUEPRINT: An Elegant Definition



It is perfectly described by just two parameters—mean (μ) and variance (σ^2)—and has practical, predictable properties like the 68-95-99.7 rule.

THE TOOLKIT: A Practical Foundation



It is the bedrock of classical statistical inference and serves as a fundamental building block for advanced machine learning models like GMMs and Bayesian Regression.

Further Reading

The mathematical tools needed to understand machine learning, including linear algebra, vector calculus, and probability, are essential for a deeper understanding. The concepts in this deck are explored in greater detail in the book:

Mathematics for Machine Learning

Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong

The book is freely available for personal use at:
<https://mml-book.com>

