

CSCI 3022

intro to data science with probability & statistics

October 10, 2018

1. The normal distribution

Stuff & Things

- **Homework 3** due Friday.
- **Midterm** tonight.
 - Start of the course up through variance.
 - One Letter-size sheet of paper, handwritten notes.
 - Bring a calculator. Cannot bring your phone.
- **Extra office hours** today: 4-6 PM

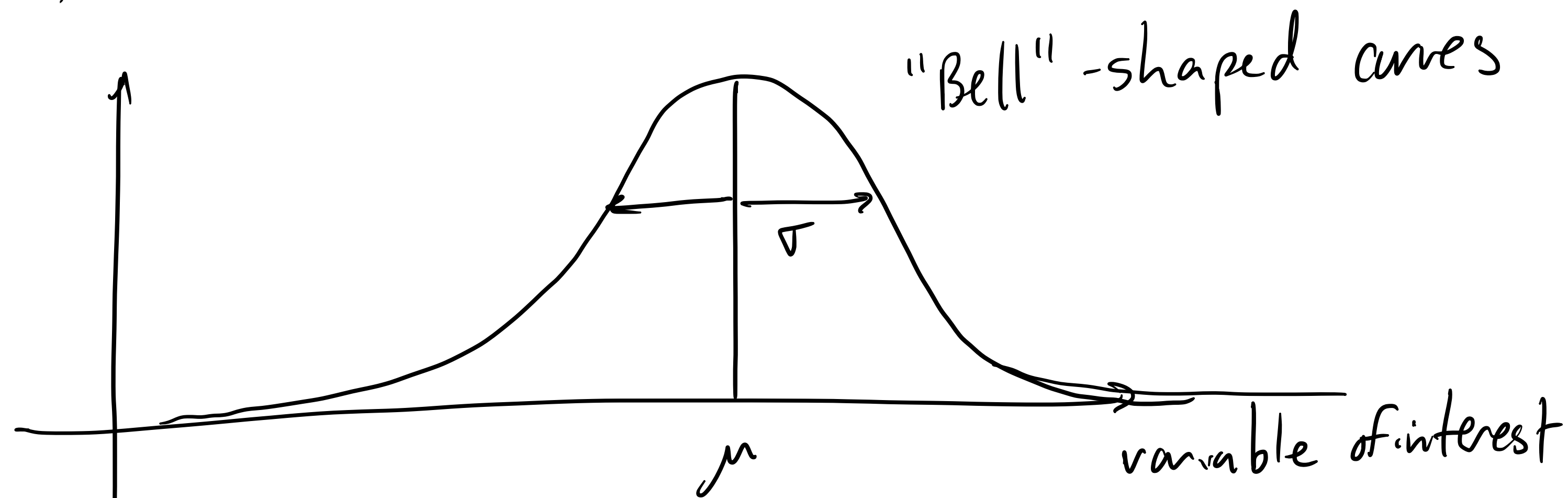
no phones as
calculators.



The normal distribution

wiki Gauss

- The **normal distribution** (aka the Gaussian distribution) is probably the most important distribution in all of probability and statistics!
- Many populations have distributions that are well-approximated by an appropriate normal distribution.
- **Examples:** height, weight, and other physical characteristics, scores on various tests, etc.



The normal distribution

- **Definition:** A continuous random variable X is said to have a **normal distribution** with parameters μ and $\sigma > 0$ (or μ and σ^2) if the pdf of X is:

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

variable parameters

- If a random variable is normally distributed, we say: $X \sim N(\mu, \sigma^2)$

<https://academo.org/demos/gaussian-distribution/>

The Standard Normal Distribution

- **Definition:** A normal distribution with parameter values $\mu = 0, \sigma^2 = 1$ is called the **standard normal distribution**.
*mean 0 $\sigma = 1$
variance, stdev = 1*
- A random variable with this distribution is called a standard normal random variable, and is often denoted by Z . Its PDF is:

$$p(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \quad (\text{we use } z \text{ for standard normal})$$

- We use a special notation to denote the CDF of the standard normal curve:

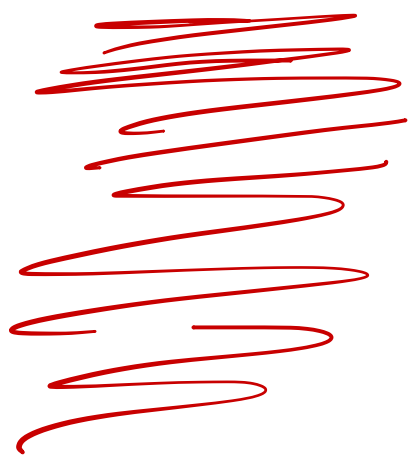

$$F(z) = \int_{-\infty}^z p(t) dt = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \Phi(z) = F(z)$$

The Standard Normal Distribution

- The standard normal distribution *rarely* occurs naturally.
- Instead, it's a **reference distribution** that allows us to learn about *other* (non-standard) normal distributions using a simple formula.
- Recall that for computing probability integrals, having the CDF is just as good as having the PDF. (Can you recall why?) In the past, when we used paper books, the values of the standard normal CDF could be found in "normal tables" in the back of any probability textbook.

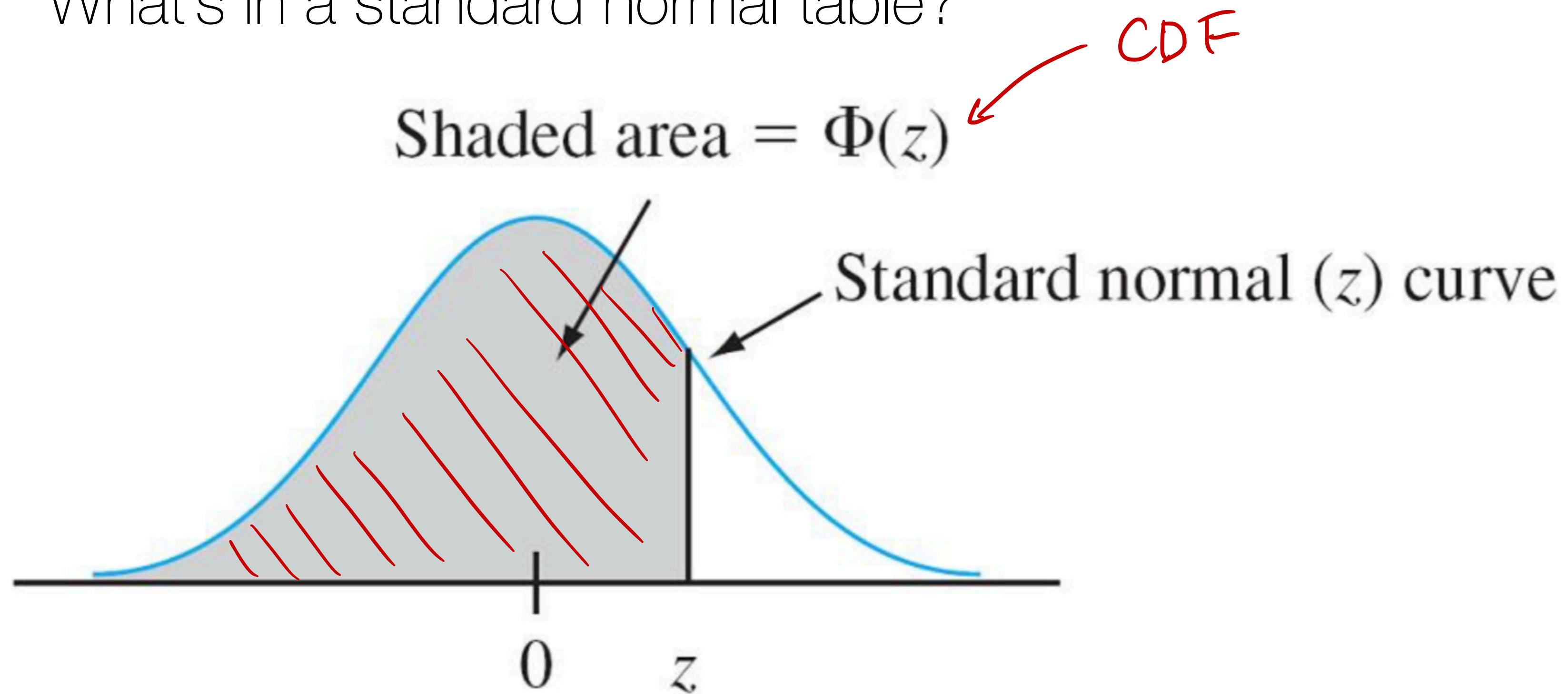
$$\int_a^b f(x) dx = F(b) - F(a) = \phi(b) - \phi(a)$$

↑
if std. normal

a	$\phi(a)$
	

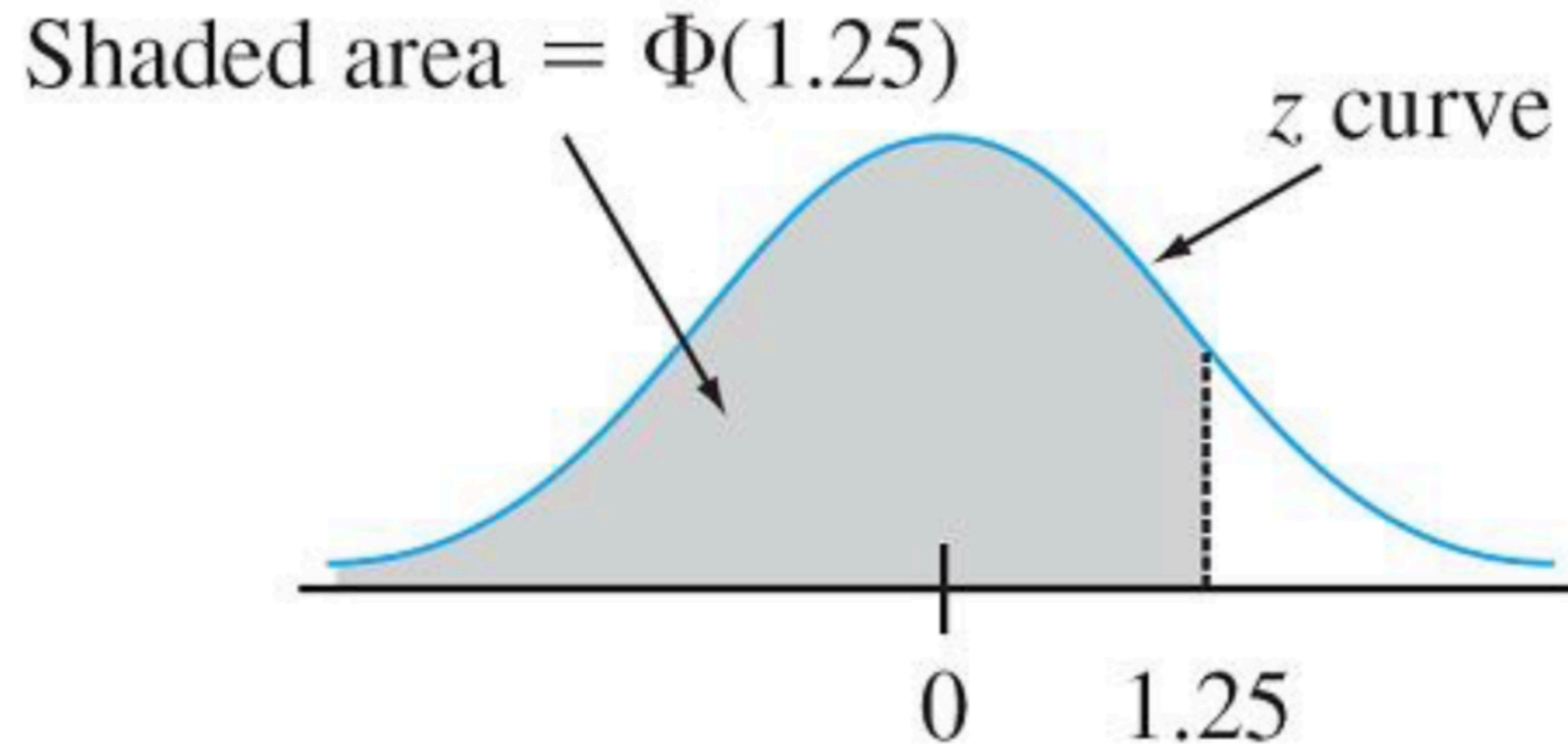
The Standard Normal Distribution

- What's in a standard normal table?



The Standard Normal Distribution

- Example: What is $P(Z \leq 1.25)$?



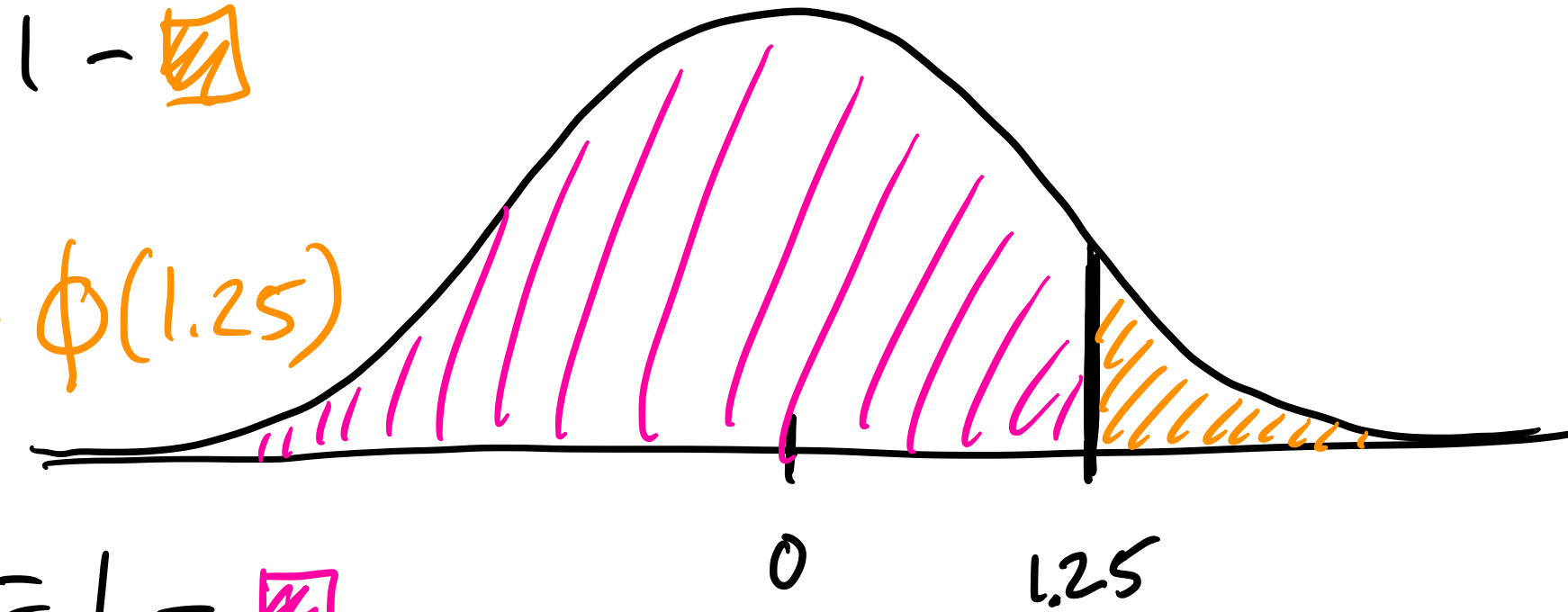
The Standard Normal Distribution

- **Example 1.** What is $P(Z \geq 1.25)$? $\neq 1 - \phi(1.25)$

table $\phi(1.25)$ \uparrow

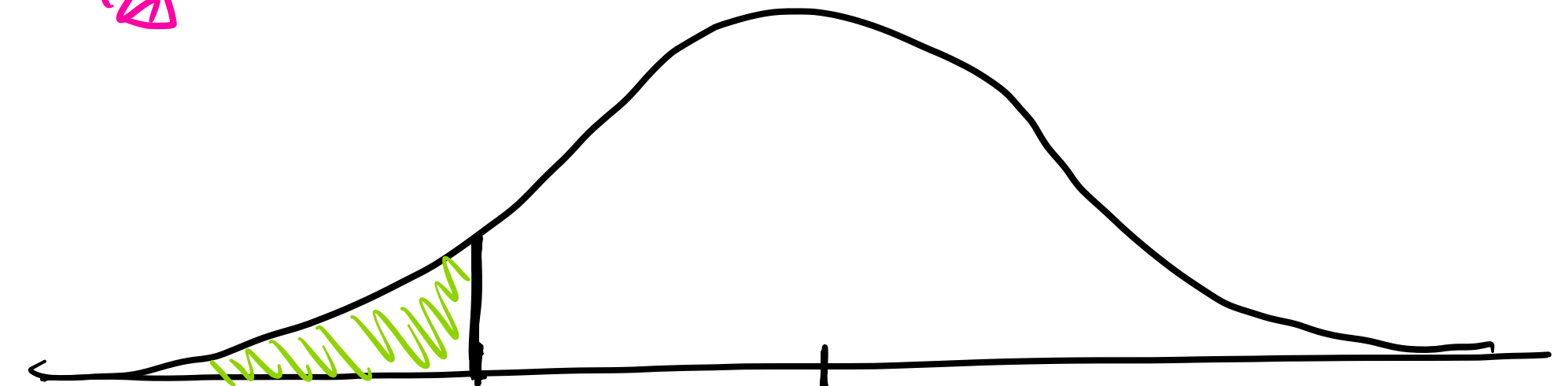
$$\text{pink box} = 1 - \text{orange box}$$

$$\text{orange box} = 1 - \text{pink box}$$

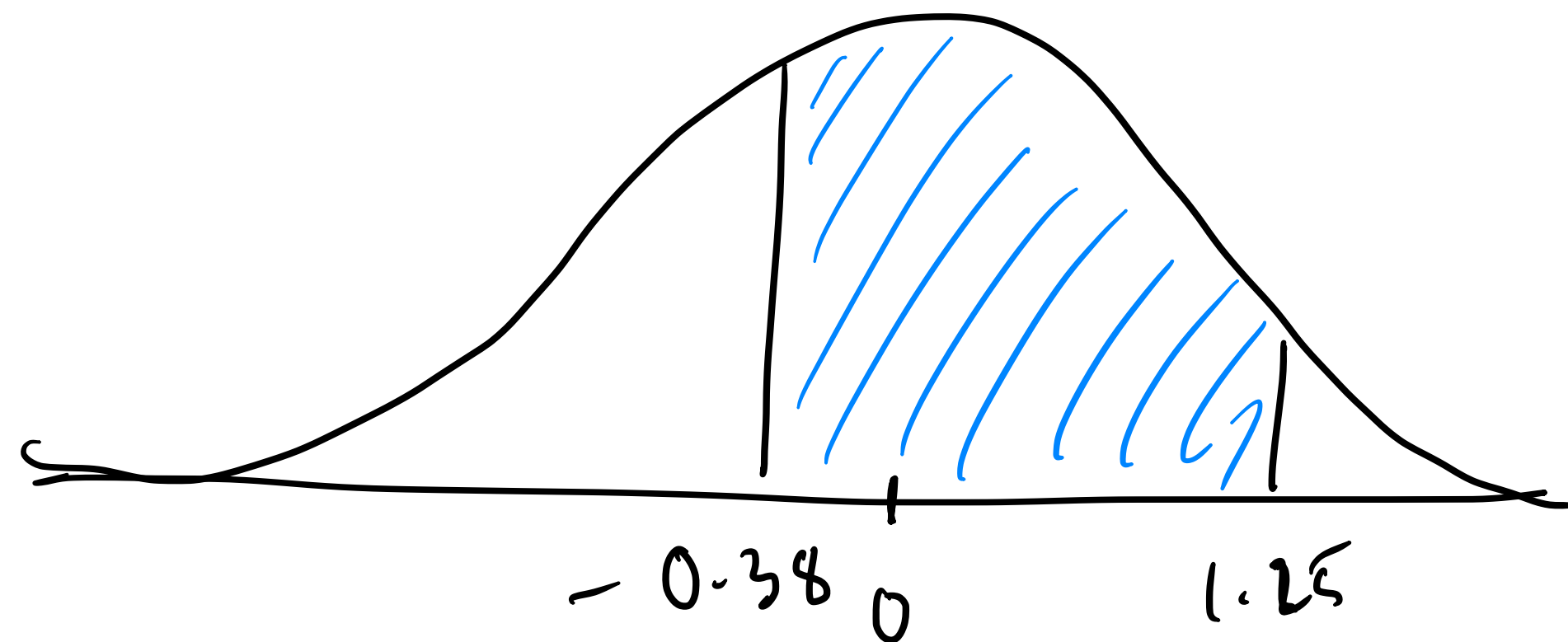


- **Example 2.** What is $P(Z \leq -1.25)$?

$$\phi(-1.25) \text{ or } 1 - \phi(1.25)$$



- **Example 3.** How can we compute $P(-0.38 \leq Z \leq 1.25)$?



$$= \phi(1.25) - \phi(-0.38)$$

\nwarrow
CDF of Z

Flip it & Reverse it

- **Recall:** what is the 99th percentile of the standard normal distribution?

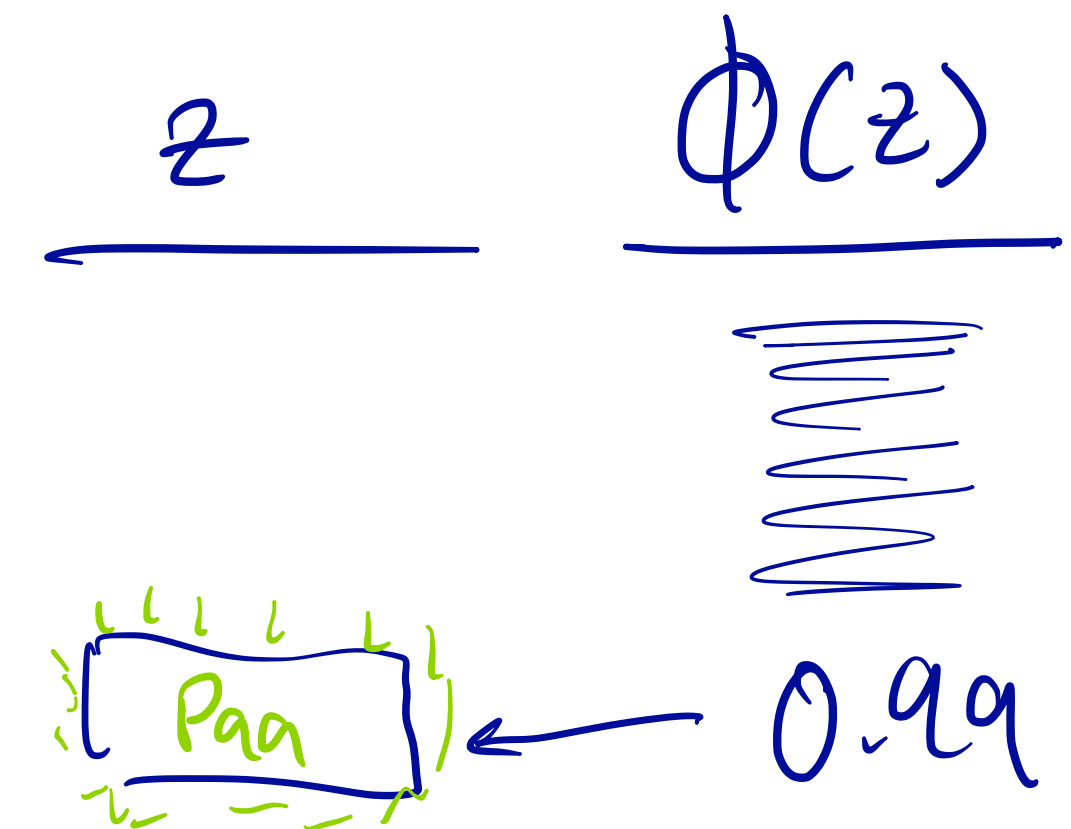
$$\int_{-\infty}^{P_{99}} f(x) dx = 0.99 \quad \longleftrightarrow \quad \Phi(P_{99}) = 0.99$$

- **Hmmm...** Tables give you a lookup from z to the area under the curve to the left of z . But we have the area and we want z .

- This is the “inverse” problem to $P(Z \leq z) = 0.99$

- How would you use a table?

- How could you sort this out in Python?



`stats.norm.ppf(0.99)`

Flip it & Reverse it

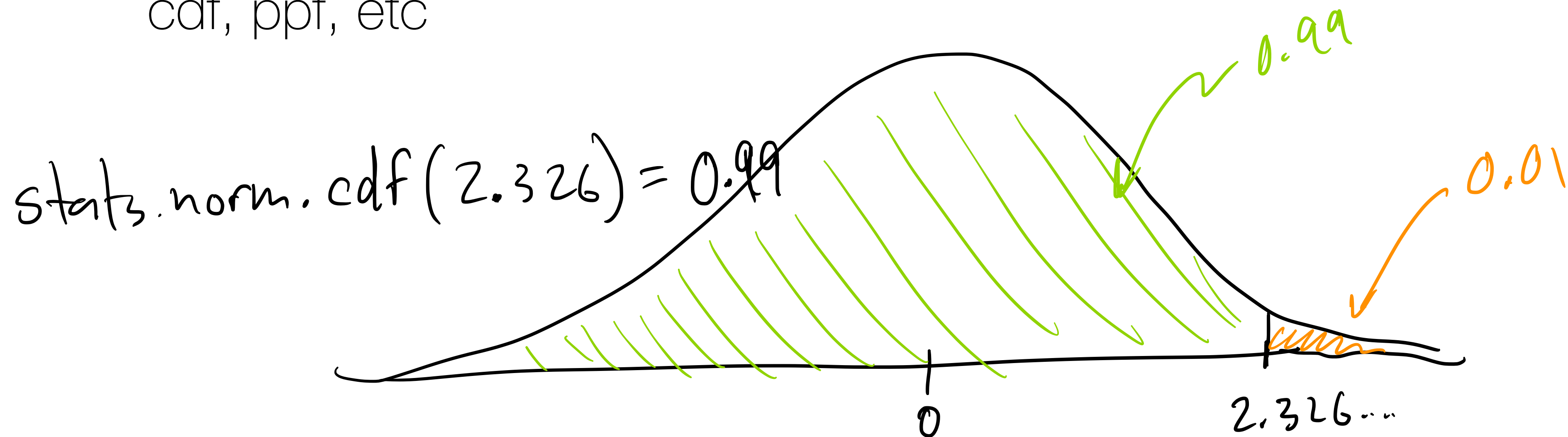
- Of course everything that we could look up in a textbook table is built into Python:

what does ppf stand for?

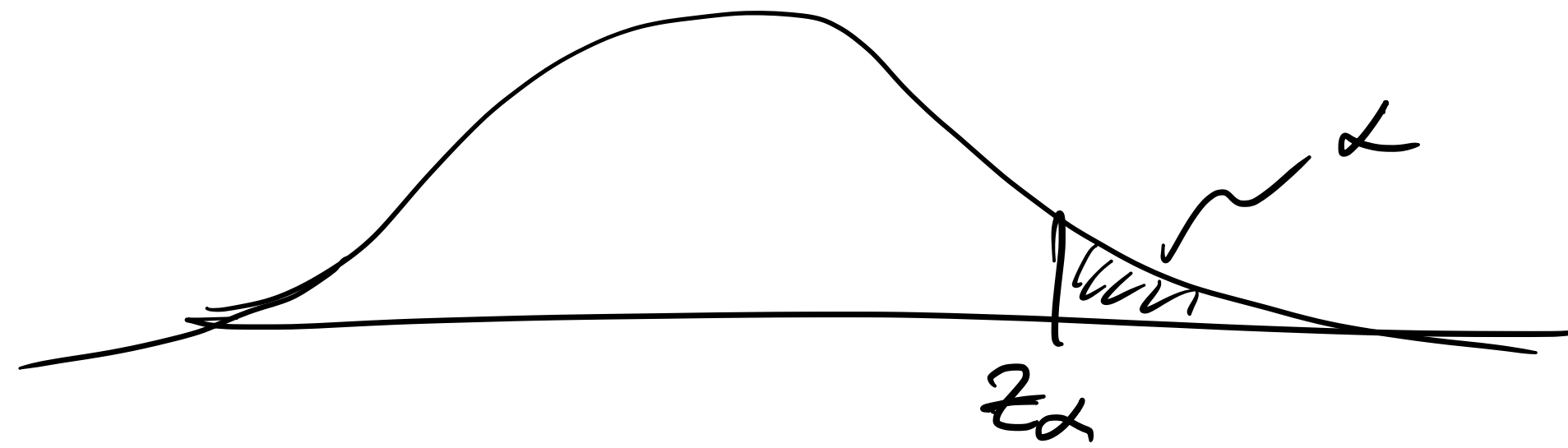
```
In [37]: 1 from scipy import stats  
        2 stats.norm.ppf(0.99)
```

```
Out[37]: 2.3263478740408408
```

- stats.norm has lots of good functions related to normal distributions: pdf, cdf, ppf, etc



... and apply it



- **Notation:** z_α is the value of z under the standard normal distribution that gives a certain “tail” area. In particular, it is the z value such that exactly α area lies to the right of z_α .
- **Hmmm...** so what is the relationship between z_α and the CDF?

$$\Phi(z_\alpha) = 1 - \alpha$$

- **Hmmm...** so what is the relationship between z_α and percentiles?

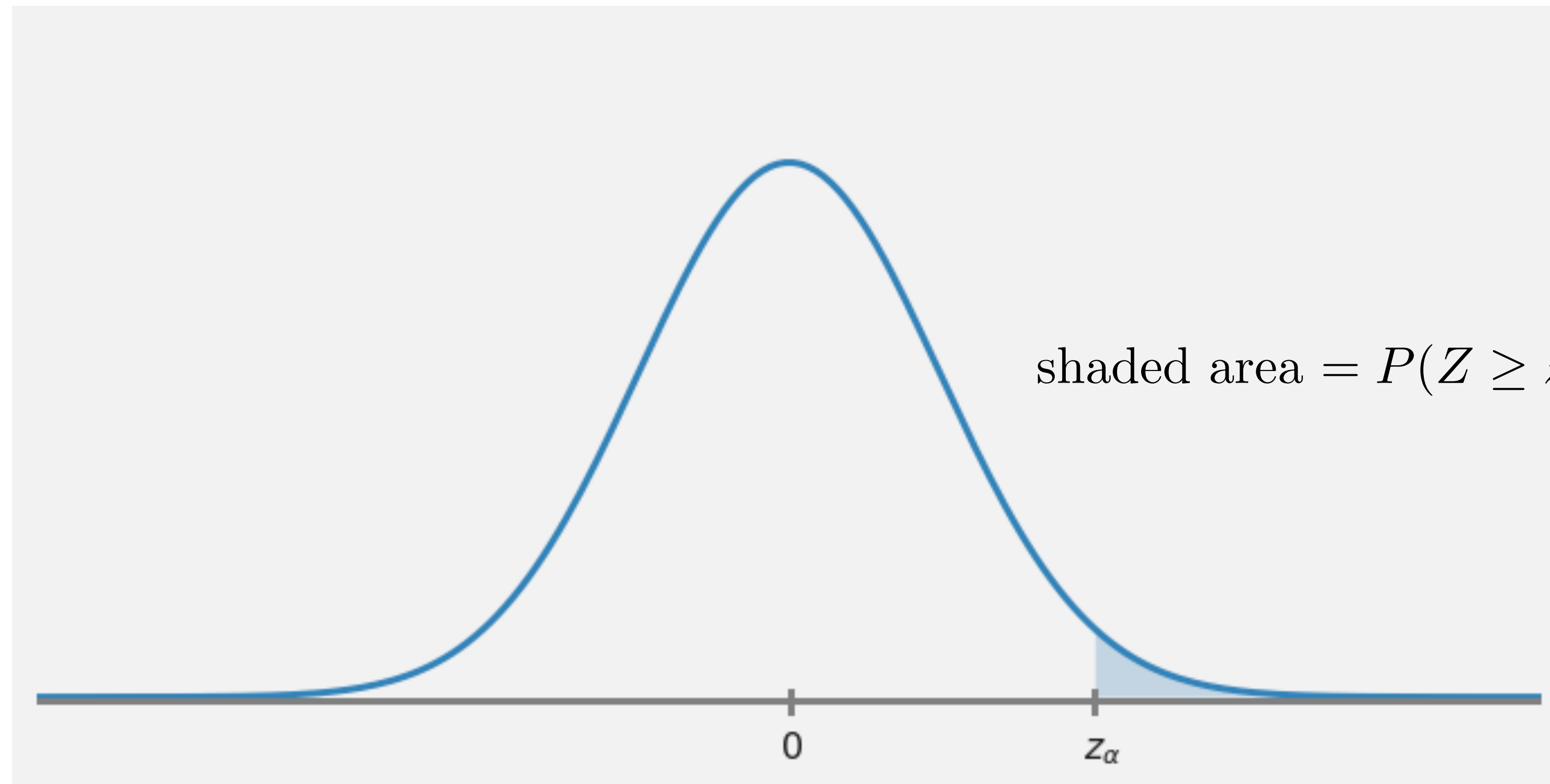
$$99^{\text{th}} \text{ percentile} = z_{0.01}$$

$$50^{\text{th}} \text{ percentile} = z_{0.5} = 0$$

(median)

Critical values

- **Notation:** z_α is the value of z under the standard normal distribution that gives a certain “tail” area. In particular, it is the z value such that exactly α area lies to the right of z_α .



Nonstandard normals

- Normal distributions that are not standard can be turned into standard normals so, so, so easily!
- **Proposition:** if X is a normal distribution with mean μ and standard deviation σ , then Z is a standard normal distribution if:

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

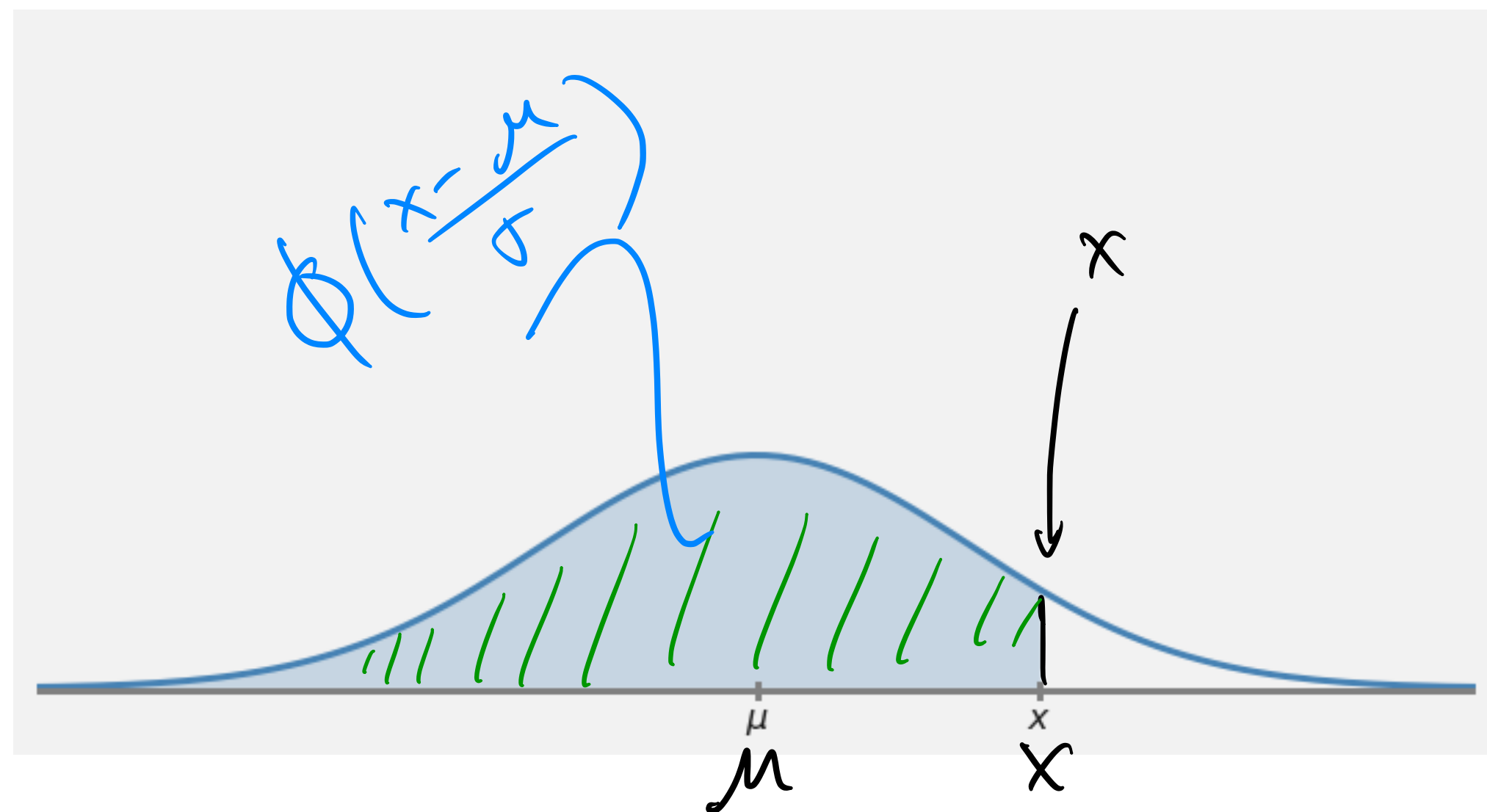
Box-Muller Transform
"Box-Muller Transform"

$$X = \sigma Z + \mu$$

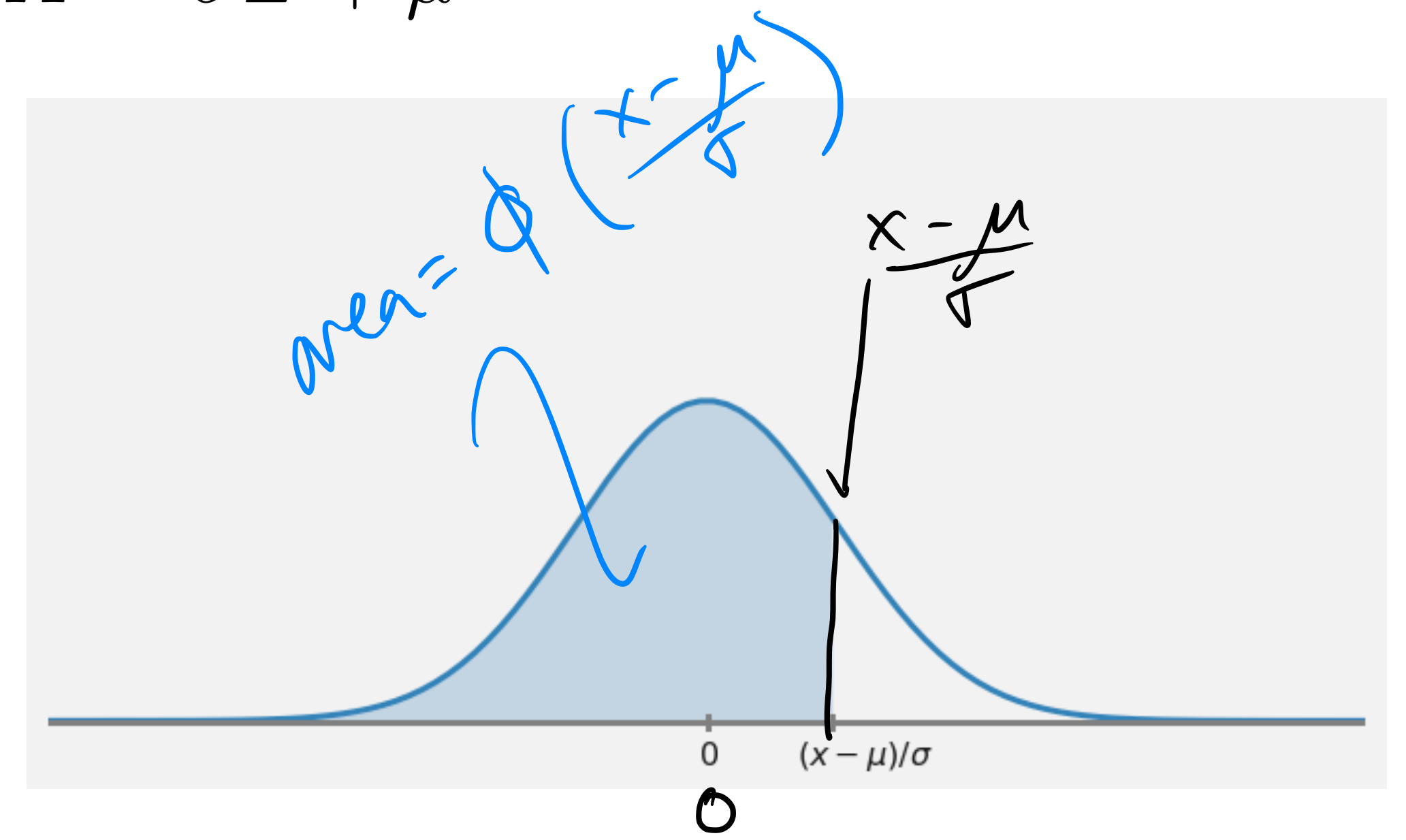
Nonstandard normals

- Normal distributions that are not standard can be turned into standard normals so, so, so easily!
- **Proposition:** if X is a normal distribution with mean μ and standard deviation σ , then Z is a standard normal distribution if:

$$Z = \frac{X - \mu}{\sigma} \quad \text{or} \quad X = \sigma Z + \mu$$



\Leftrightarrow



Area Intuition

<https://www.intmath.com/counting-probability/normal-distribution-graph-interactive.php>

Standard normals in action

- **Example:** The time that it takes a driver to react to the brake lights on a decelerating vehicle is critical in helping to avoid rear-end collisions.
- The article Fast-Rise Brake Lamp as a Collision-Prevention Device* suggests that reaction time for an in-traffic response to a brake signal from standard brake lights can be modeled with a normal distribution having **mean value 1.25 sec** and **standard deviation of .46 sec**.
- What is the probability that reaction time is between 1.00 sec and 1.75 sec?

$$\mu = 1.25$$
$$\sigma = 0.46$$

$$P(1.00 \leq X \leq 1.75) = P\left(\frac{1.00 - 1.25}{0.46} \leq Z \leq \frac{1.75 - 1.25}{0.46}\right)$$

Box-Muller

$$= P(-0.543 \leq Z \leq 1.09)$$

$$= \Phi(1.09) - \Phi(-0.543) = \boxed{0.568}$$

* (Ergonomics, 1993: 391–395)

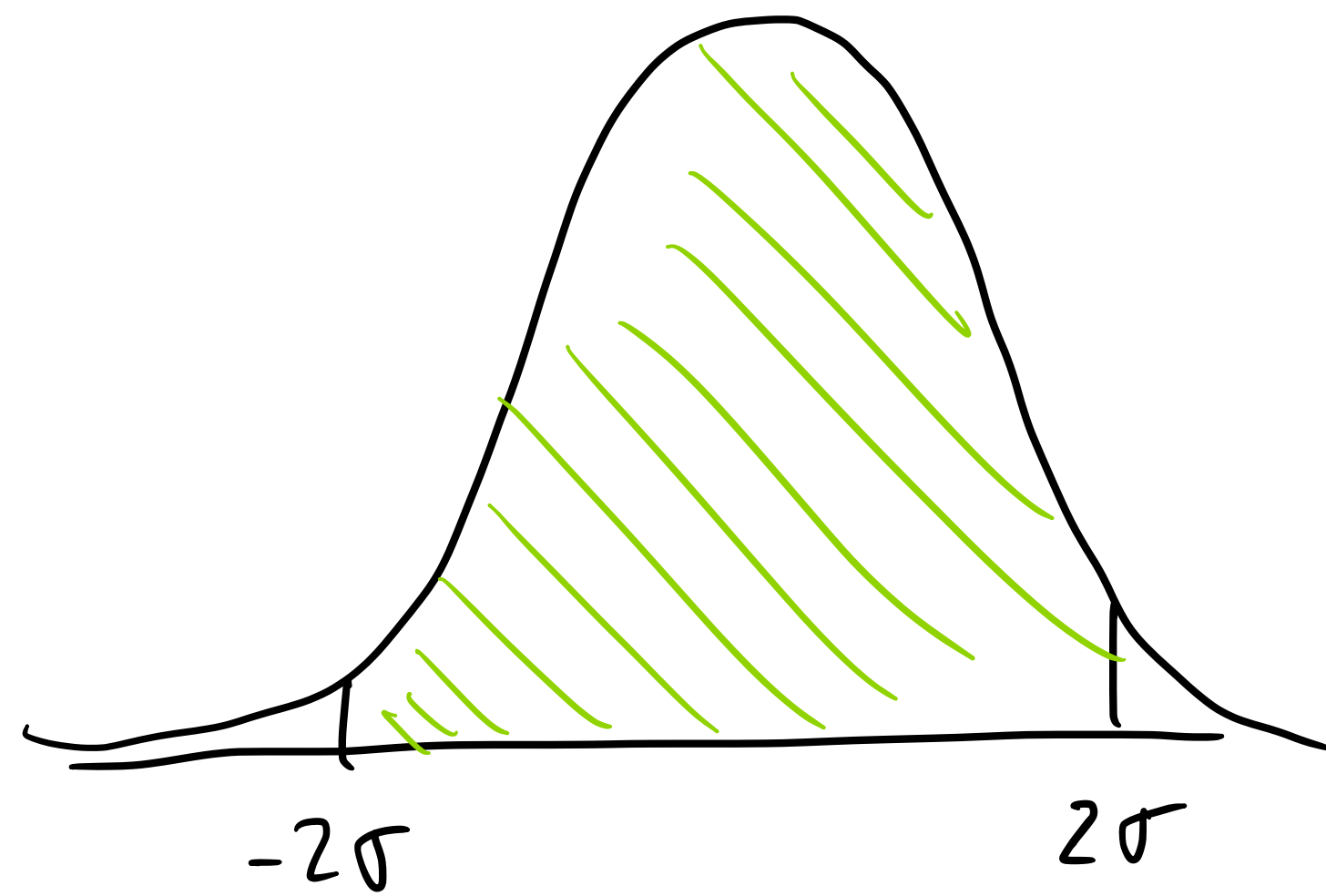
Standard Normals in action

- So what are some common things that come up with normals?



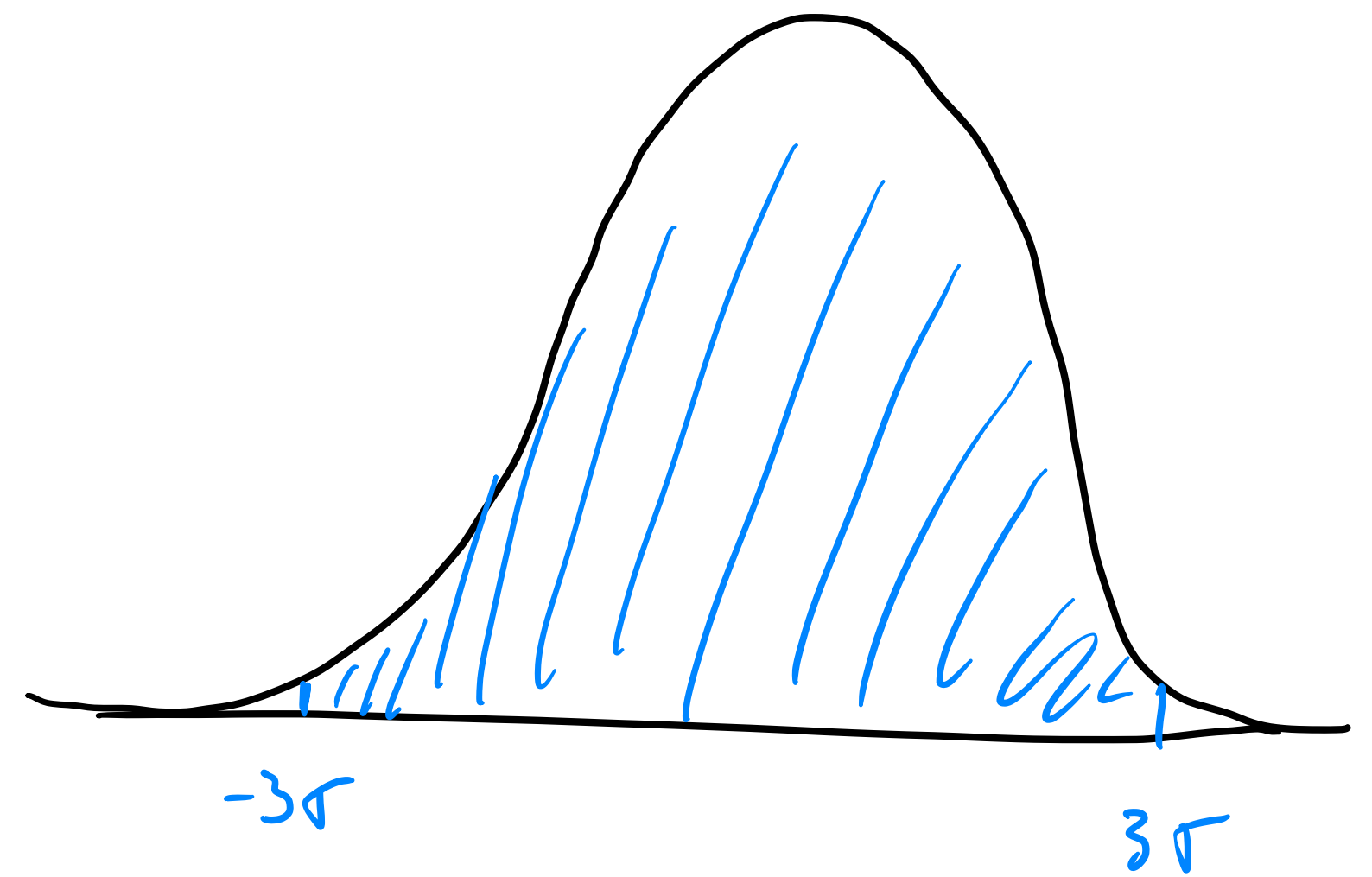
0.683

1 in 3.2



0.954

1 in 22



0.997

1 in 370.4