



190F Foundations of Data Science

Spring 2020

Lecture 15

The Normal Curve

Standard Units (review)

Standard Units

- How many SDs above average?
 - **$z = (\text{value} - \text{average})/\text{SD}$**
 - Negative z : value below average
 - Positive z : value above average
 - $z = 0$: value equal to average
 - When values are in standard units: average = 0, SD = 1
 - Chebyshev: At least 96% of the values of z are between -5 and 5
- (Demo)
-

Discussion Question

Find whole numbers
that are close to:

(a) the average age

(b) the SD of the ages

(Demo)

Age in Years	Age in Standard Units
27	-0.0392546
33	0.992496
28	0.132704
23	-0.727088
25	-0.383171
33	0.992496
23	-0.727088
25	-0.383171
30	0.476621
27	-0.0392546

... (1164 rows omitted)

The SD and the Histogram

- Usually, it's not easy to estimate the SD by looking at a histogram.
 - But if the histogram has a bell shape, then you can.
-

The SD and Bell-Shaped Curves

If a histogram is bell-shaped, then

- the average is at the center
- the SD is the distance between the average and the points of inflection on either side

(Demo)

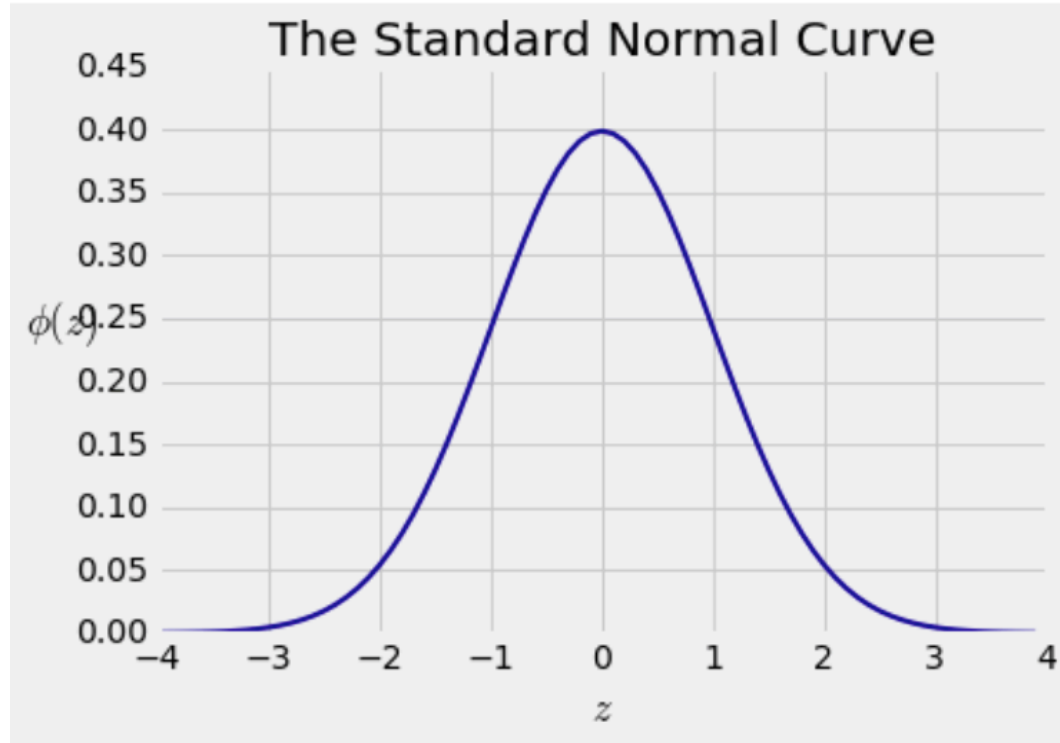
The Normal Distribution

The Standard Normal Curve

A beautiful formula that we won't use at all:

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}, \quad -\infty < z < \infty$$

Bell Curve



Normal Proportions

How Big are Most of the Values?

No matter what the shape of the distribution,
the bulk of the data are in the range “average \pm a few SDs”

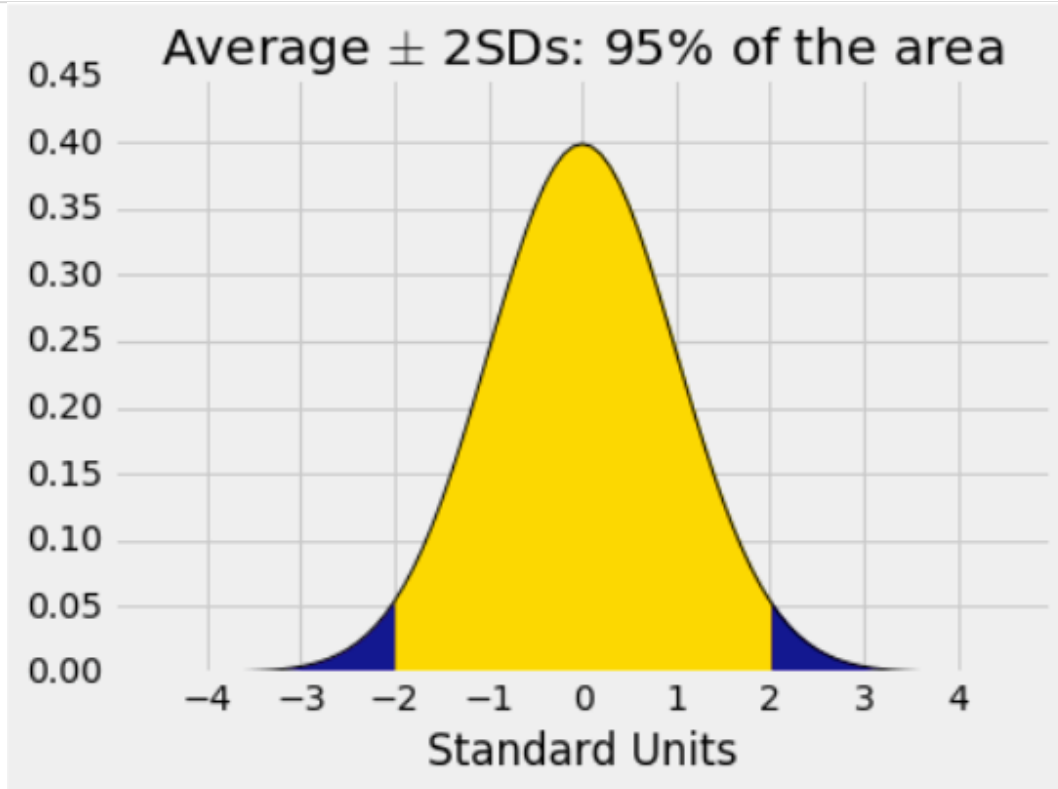
If a histogram is bell-shaped, then

- Almost all of the data are in the range
“average \pm 3 SDs”

Bounds and Normal Approximations

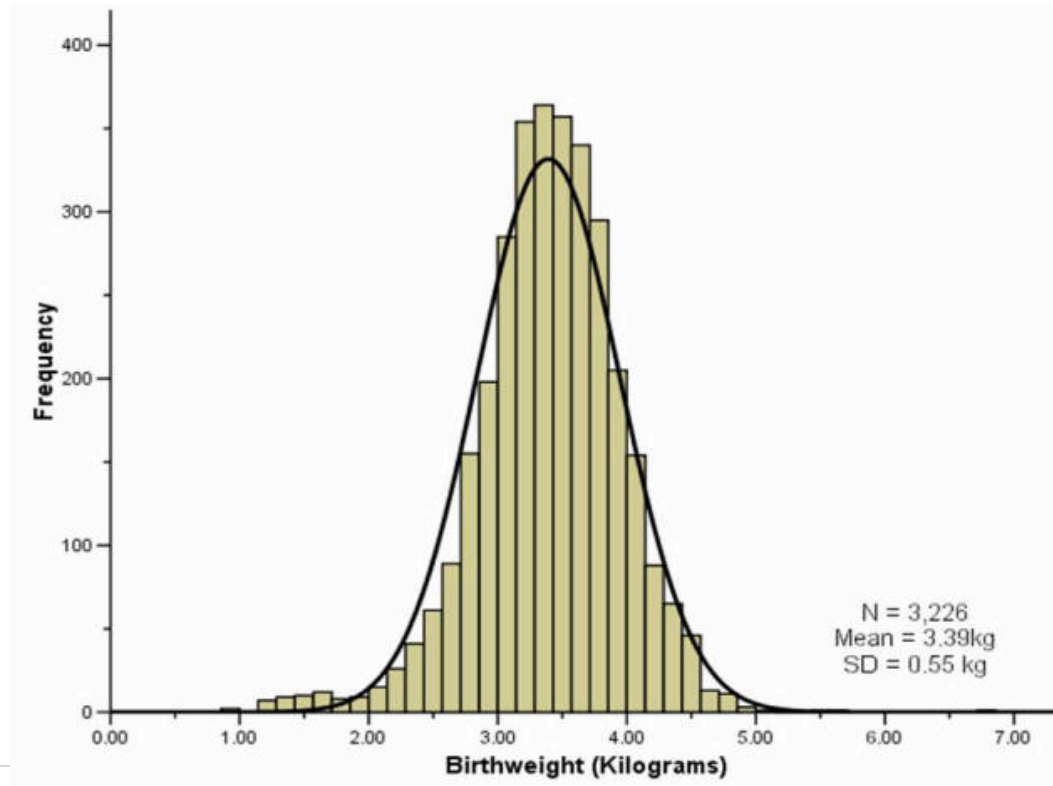
Percent in Range	All Distributions	Normal Distribution
average \pm 1 SD	at least 0%	about 68%
average \pm 2 SDs	at least 75%	about 95%
average \pm 3 SDs	at least 88.888...%	about 99.73%

A “Central” Area



Normal Curve in Practice

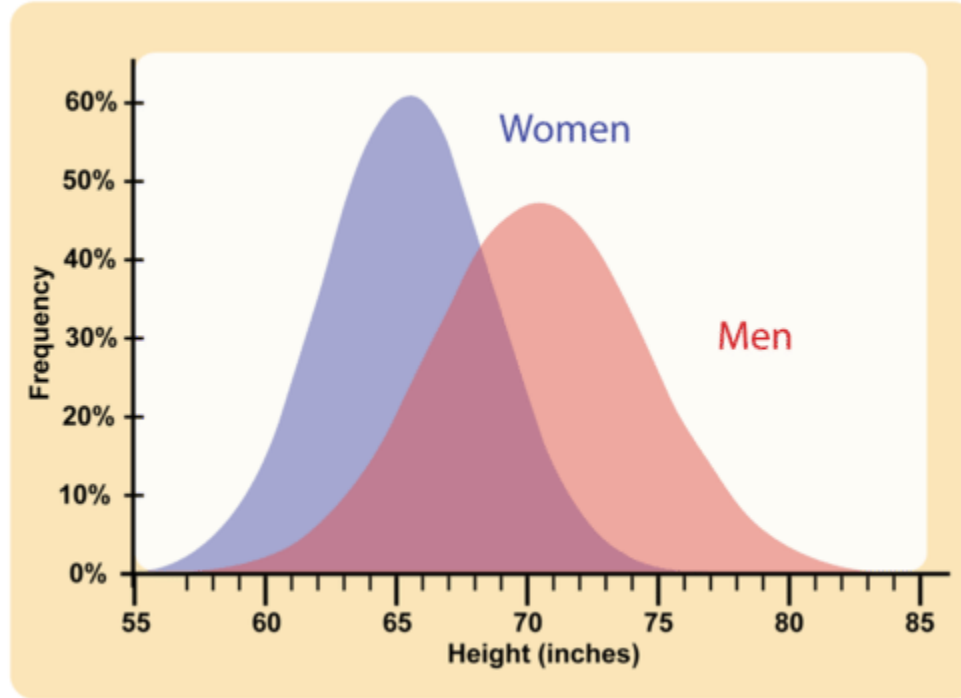
Birth weights of babies



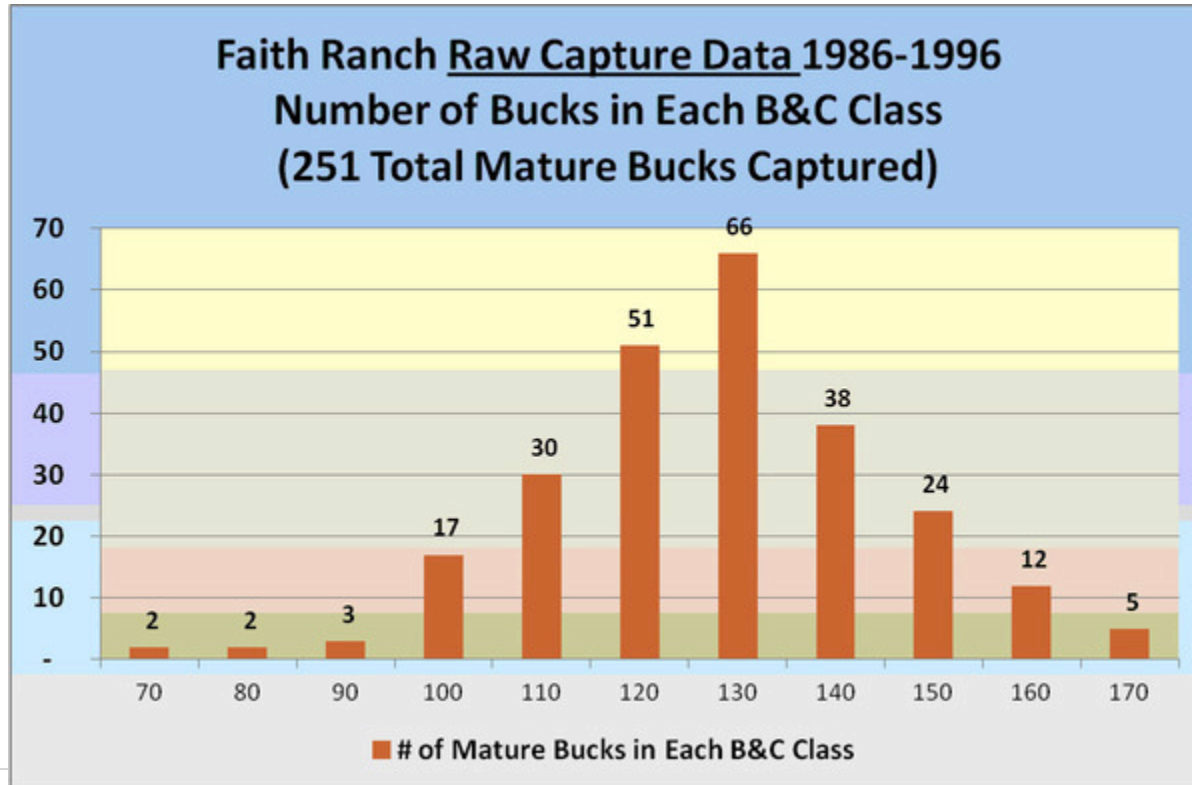
Sales of shoes for women



Heights of men and women



Lengths of Deer Antlers



The Mystery

Why is the normal distribution so common?

Central Limit Theorem

Second Reason for Using the SD

If the sample is

- large, and
- drawn at random with replacement,

Then, *regardless of the distribution of the population,*

**the probability distribution of the sample sum
(or of the sample average) is roughly normal**

(Demo)
