

Probability & Normal Distribution

Portland State University
USP 634 Data Analysis I
Spring 2017

Introduction to Probability

Slides developed by Mine Çetinkaya-Rundel of OpenIntro

The slides may be copied, edited, and/or shared via the [CC BY-SA license](#)

Some images may be included under fair use guidelines (educational purposes)

Probability

There are several possible interpretations of probability but they (almost) completely agree on the mathematical rules probability must follow.

- $P(A)$ = Probability of event A
- $0 \leq P(A) \leq 1$

Frequentist interpretation:

- The probability of an outcome is the proportion of times the outcome would occur if we observed the random process an infinite number of times.

Bayesian interpretation:

- A Bayesian interprets probability as a subjective degree of belief: For the same event, two separate people could have different viewpoints and so assign different probabilities.
- Largely popularized by revolutionary advance in computational technology and methods during the last twenty years.

Probability as proportion

- Probability is a proportion that can be predicted over “the long run” (many trials)
- An important consideration for insurers, gamblers, brokers...



- You win 18 of 38 (47.4%); house wins 20 of 38 (52.6%)

Practice

Which of the following events would you be most surprised by?

- (a) exactly 3 heads in 10 coin flips
- (b) exactly 3 heads in 100 coin flips
- (c) exactly 3 heads in 1000 coin flips

Practice

Which of the following events would you be most surprised by?

- (a) exactly 3 heads in 10 coin flips
- (b) exactly 3 heads in 100 coin flips
- (c) exactly 3 heads in 1000 coin flips*

Probability is the foundation of inferential statistics

- Probability theory is what enables statements like “Those who participate in the workfare program have higher household income” to be *statistically* accurate:

“The probability of observing a 30% difference in promotion rates b/w genders is less than 5% if there is no gender discrimination”

Law of large numbers

Law of large numbers states that as more observations are collected, the proportion of occurrences with a particular outcome, \hat{p}_n , converges to the probability of that outcome, p .

Disjoint and non-disjoint outcomes

Disjoint (mutually exclusive) outcomes: Cannot happen at the same time.

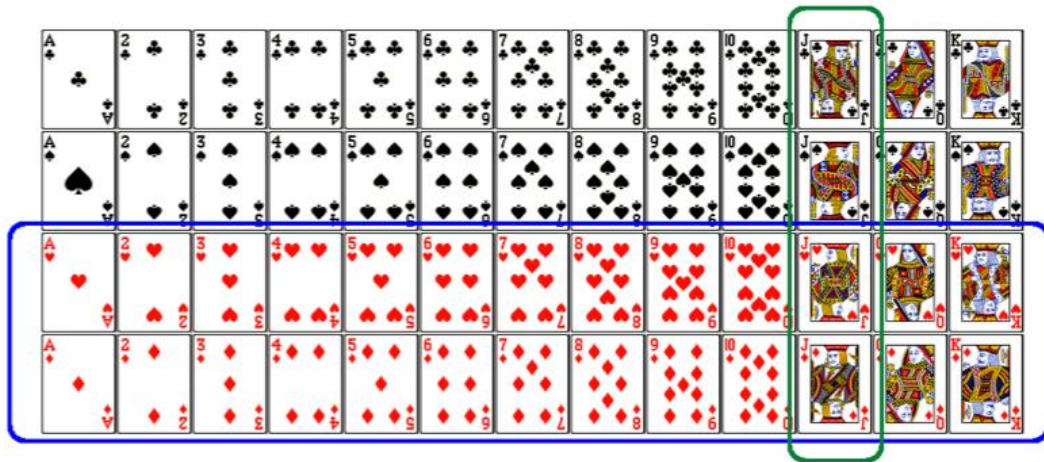
- The outcome of a single coin toss cannot be a head and a tail.
- A student both cannot fail and pass a class.
- A single card drawn from a deck cannot be an ace and a queen.

Non-disjoint outcomes: Can happen at the same time.

- A student can get an A in Stats and A in Econ in the same semester.

Union of non-disjoint events

What is the probability of drawing a jack or a red card from a well shuffled full deck?



$$P(\text{jack or red}) = P(\text{jack}) + P(\text{red}) - P(\text{jack and red})$$

$$= \frac{4}{52} + \frac{26}{52} - \frac{2}{52} = \frac{28}{52}$$

Figure from <http://www.milefoot.com/math/discrete/counting/cardfreq.htm>

Recap

General addition rule

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

Note: For disjoint events $P(A \text{ and } B) = 0$, so the above formula simplifies to $P(A \text{ or } B) = P(A) + P(B)$

Independence

Two processes are **independent** if knowing the outcome of one provides no useful information about the outcome of the other.

- Knowing that the coin landed on a head on the first toss does not provide any useful information for determining what the coin will land on in the second toss.
>> Outcomes of two tosses of a coin are independent.
- Knowing that the first card drawn from a deck is an ace does provide useful information for determining the probability of drawing an ace in the second draw.
>> Outcomes of two draws from a deck of cards (without replacement) are dependent.

Gender Discrimination

At a first glance, does there appear to be a relationship between promotion and gender?

		<i>Promotion</i>		
		Promoted	Not Promoted	Total
<i>Gender</i>	Male	21	3	24
	Female	14	10	24
	Total	35	13	48

% of promoted: $35 / 48 = 0.729$

% of males promoted: $21 / 24 = 0.875$

% of females promoted: $14 / 24 = 0.583$

Checking for independence

If $P(A \text{ occurs, given that } B \text{ is true}) = P(A | B) = P(A)$,
then A and B are independent.

$$P(\text{promoted}) = 35 / 48 = 0.729$$

$P(\text{promoted, given that the gender is male})$

$$= P(\text{promoted} | \text{male}) = 21 / 24 = 0.875$$

$$P(\text{promoted} | \text{female}) = 14 / 24 = 0.583$$

$P(\text{promoted}) \neq P(\text{promoted} | \text{male}) \neq P(\text{promoted} | \text{female})$

$P(\text{promoted})$ varies by gender, therefore promotion and gender are most likely dependent.

Determining dependence based on sample data

- If conditional probabilities calculated based on sample data suggest dependence between two variables, the next step is to conduct a hypothesis test to determine if the observed difference between the probabilities is likely or unlikely to have happened by chance.
- If the observed difference between the conditional probabilities is large, then there is stronger evidence that the difference is real.
- If a sample is large, then even a small difference can provide strong evidence of a real difference.

Product rule for independent events

$$P(A \text{ and } B) = P(A) \times P(B)$$

You toss a coin twice, what is the probability of getting two tails in a row?

$$\begin{aligned} &P(\text{T on the first toss}) \times P(\text{T on the second toss}) \\ &= (1 / 2) \times (1 / 2) = 1 / 4 \end{aligned}$$

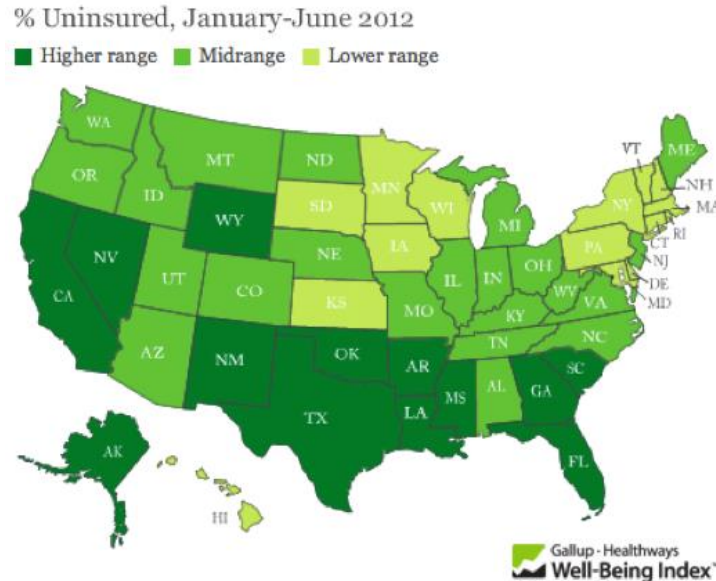
Probability of independent Events

- $\Pr(\text{Event A or Event B}) = \Pr(\text{Event A}) + \Pr(\text{Event B})$
 - $\Pr(4 \text{ or } 6) = \Pr(4) + \Pr(6) = 1/6 + 1/6 = 1/3$
- $\Pr(\text{Event A and Event B}) = \Pr(\text{Event A}) * \Pr(\text{Event B})$
 - $\Pr(4 \text{ followed by } 6) = \Pr(4) * \Pr(6) = 1/6 * 1/6 = 1/36$
 - Prob of having 7 daughters in a row = ?

Practice

A recent Gallup poll suggests that 25.5% of Texans do not have health insurance as of June 2012. Assuming that the uninsured rate stayed constant, what is the probability that two randomly selected Texans are both uninsured?

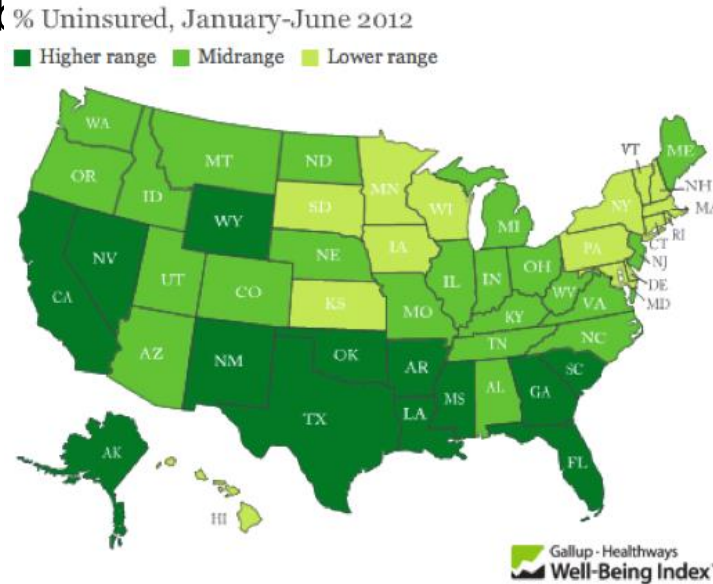
- (a) 25.5^2
- (b) 0.255^2
- (c) 0.255×2
- (d) $(1 - 0.255)^2$



Practice

A recent Gallup poll suggests that 25.5% of Texans do not have health insurance as of June 2012. Assuming that the uninsured rate stayed constant, what is the probability that both Texans are uninsured?

- (a) 25.5^2
- (b) 0.255^2
- (c) 0.255×2
- (d) $(1 - 0.255)^2$



<http://www.gallup.com/poll/156851/uninsured-rate-stable-across-states-far-2012.aspx>

Probability Distributions

Slides developed by Mine Çetinkaya-Rundel of OpenIntro

The slides may be copied, edited, and/or shared via the [CC BY-SA license](#)

Some images may be included under fair use guidelines (educational purposes)

Probability distributions

A **probability distribution** lists all possible events and the probabilities with which they occur.

- The probability distribution for the gender of one kid:

Event	Male	Female
Probability	0.5	0.5

- Rules for probability distributions:

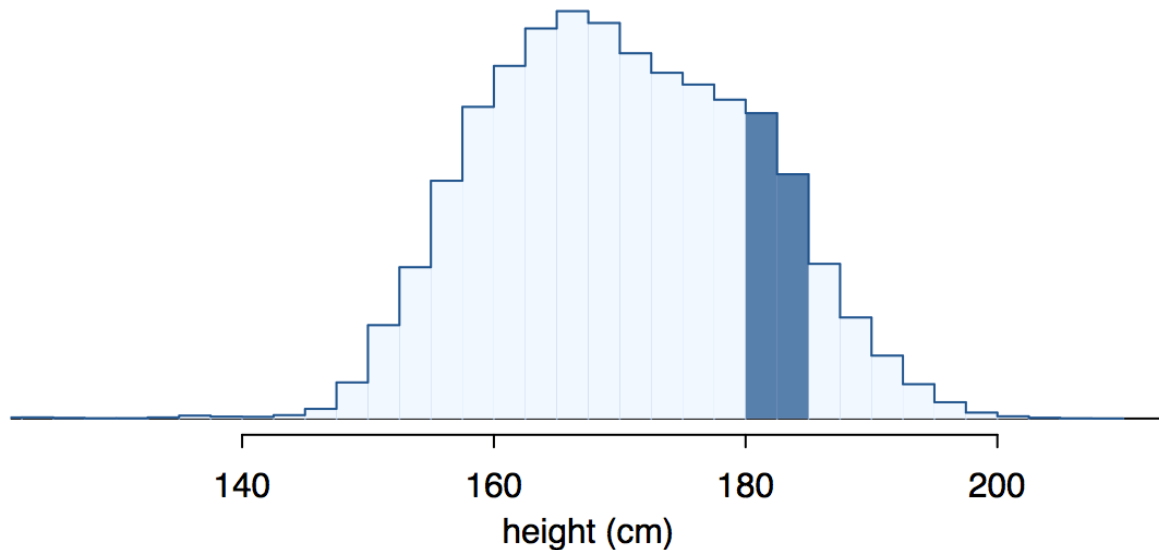
1. The events listed must be disjoint
2. Each probability must be between 0 and 1
3. The probabilities must total 1

- The probability distribution for the genders of two kids:

Event	MM	FF	MF	FM
Probability	0.25	0.25	0.25	0.25

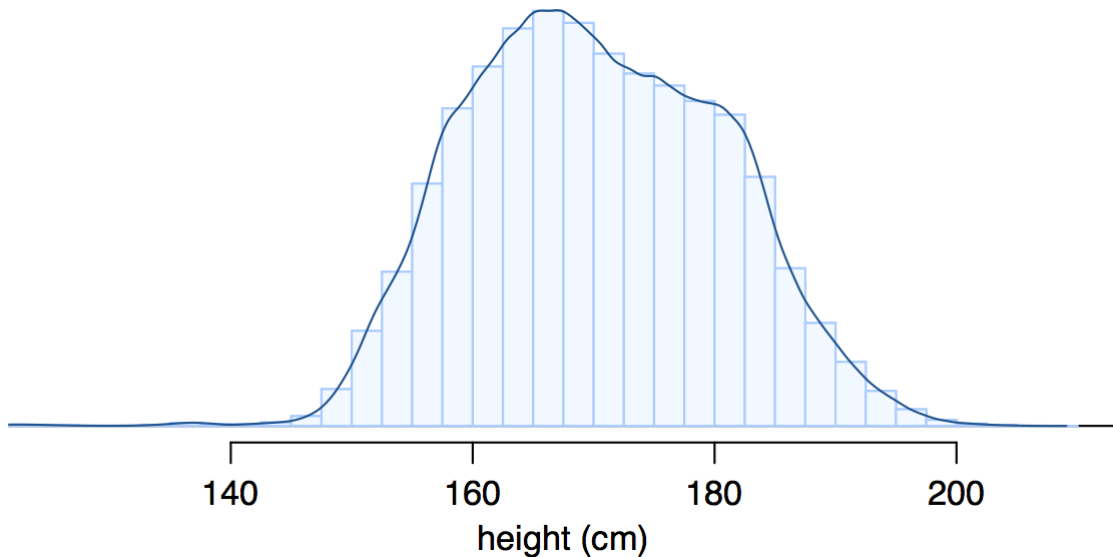
Continuous distributions

- Below is a histogram of the distribution of heights of US adults.
- The proportion of data that falls in the shaded bins gives the probability that a randomly sampled US adult is between 180 cm and 185 cm (about 5'11" to 6'1").



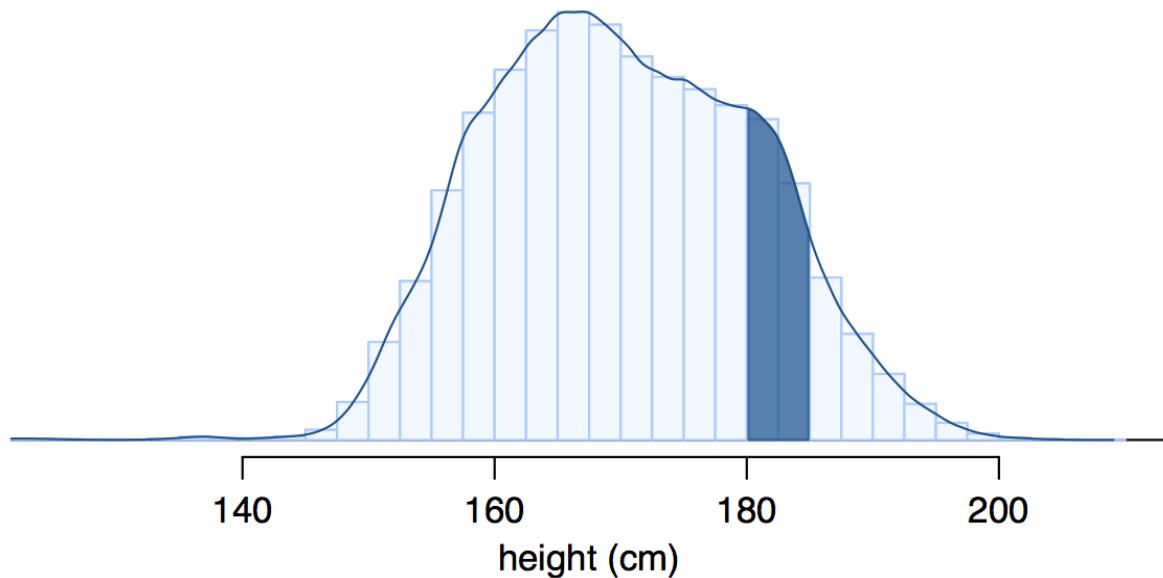
From histograms to continuous distributions

Since height is a continuous numerical variable, its **probability density function** is a smooth curve.



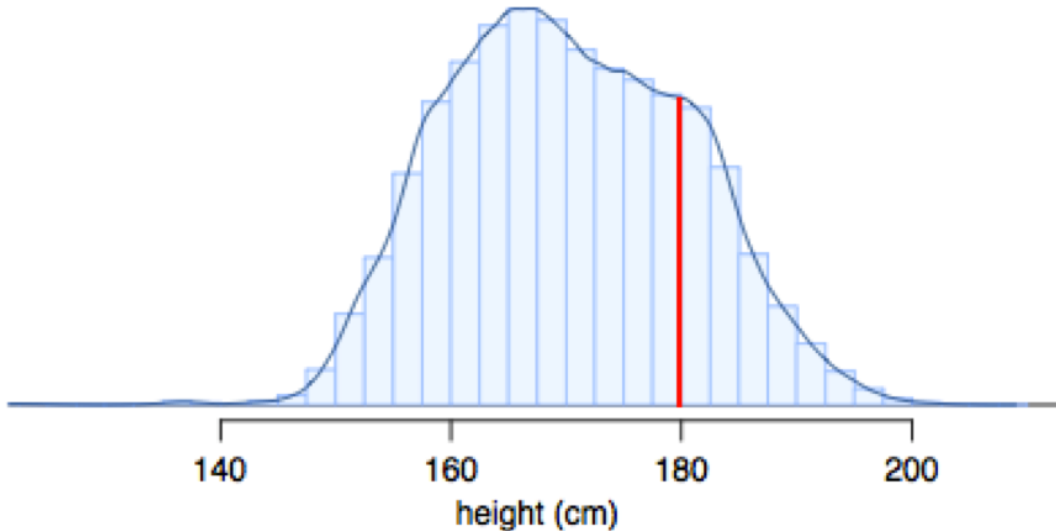
Probabilities from continuous distributions

Therefore, the probability that a randomly sampled US adult is between 180 cm and 185 cm can also be estimated as the shaded area under the curve.



By definition...

Since continuous probabilities are estimated as “the area under the curve”, the probability of a person being exactly 180 cm (or any exact value) is defined as 0.



Normal distribution

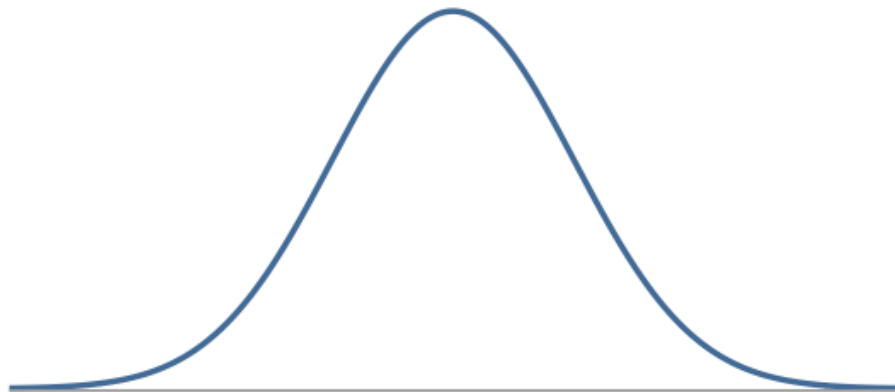
Slides developed by Mine Çetinkaya-Rundel of OpenIntro

The slides may be copied, edited, and/or shared via the [CC BY-SA license](#)

Some images may be included under fair use guidelines (educational purposes)

Normal Distribution

- Unimodal and symmetric, bell shaped curve
- Many variables are nearly normal, but none are exactly normal
- Denoted as $N(\mu, \sigma)$ → Normal with mean μ and standard deviation σ

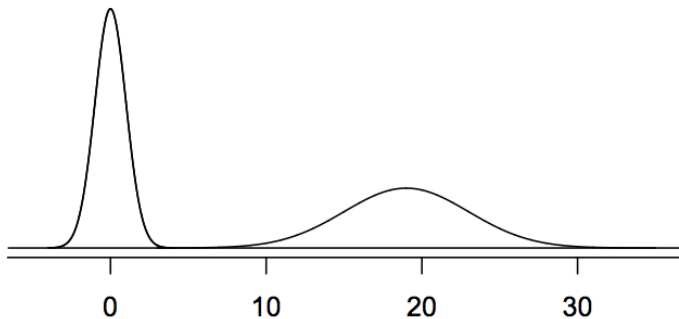
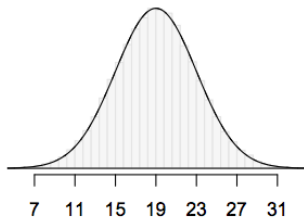
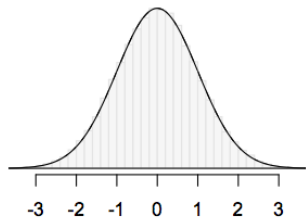


Normal distributions with different parameters

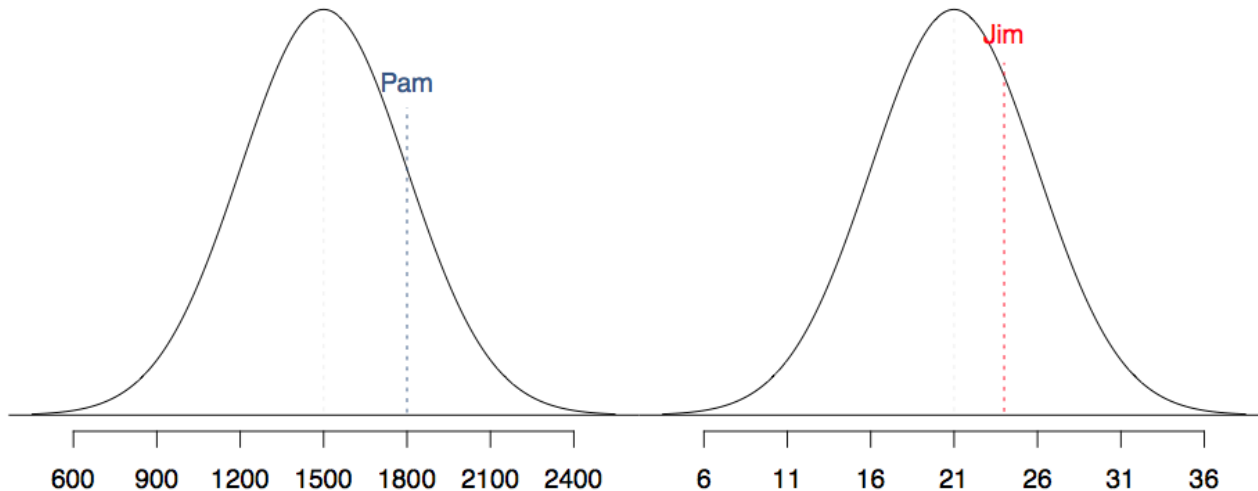
μ : mean, σ : standard deviation

$$N(\mu = 0, \sigma = 1)$$

$$N(\mu = 19, \sigma = 4)$$



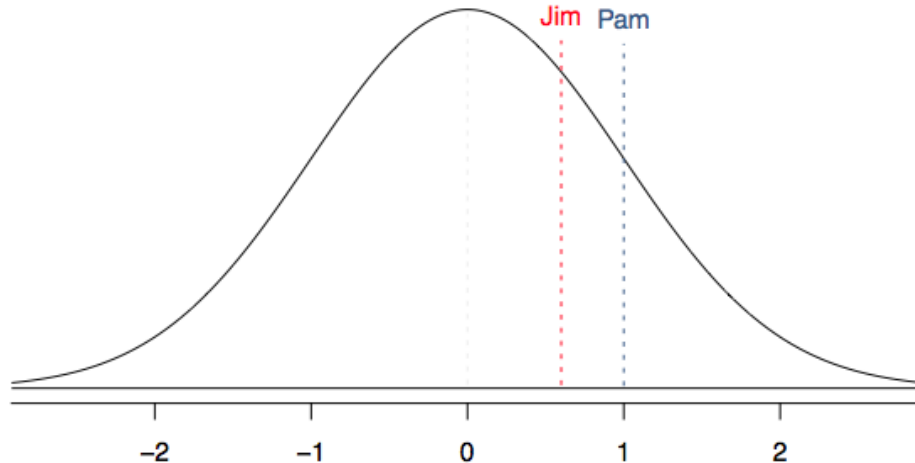
SAT scores are distributed nearly normally with mean 1500 and standard deviation 300. ACT scores are distributed nearly normally with mean 21 and standard deviation 5. A college admissions officer wants to determine which of the two applicants scored better on their standardized test with respect to the other test takers: Pam, who earned an 1800 on her SAT, or Jim, who scored a 24 on his ACT?



Standardizing with Z scores

Since we cannot just compare these two raw scores, we instead compare how many standard deviations beyond the mean each observation is.

- Pam's score is $(1800 - 1500) / 300 = 1$ standard deviation above the mean.
- Jim's score is $(24 - 21) / 5 = 0.6$ standard deviations above the mean.



Standardizing with Z scores (cont.)

These are called **standardized** scores, or **Z scores**.

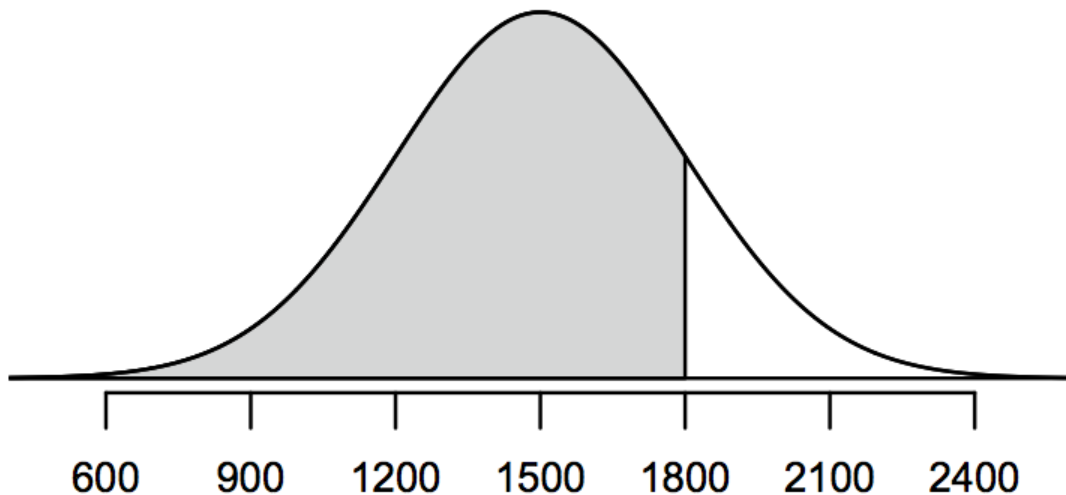
- Z score of an observation is the number of standard deviations it falls above or below the mean.

$$Z = (\text{observation} - \text{mean}) / \text{SD}$$

- Z scores are defined for distributions of any shape, but only **when the distribution is normal can we use Z scores to calculate percentiles**.
- Observations that are more than 2 SD away from the mean ($|Z| > 2$) are usually considered unusual.

Percentiles

- **Percentile** is the percentage of observations that fall below a given data point.
- Graphically, percentile is the area below the probability distribution curve to the left of that observation.

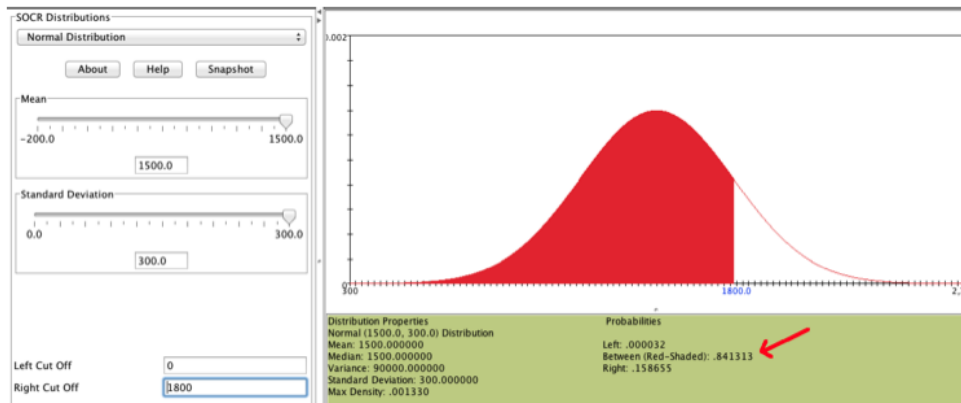


Calculating percentiles -- using computation

There are many ways to compute percentiles/areas under the curve. R:

```
> pnorm(1800, mean = 1500, sd = 300)
[1] 0.8413447
```

Applet: www.socr.ucla.edu/htmls/SOCR_Distributions.html



Calculating percentiles -- using tables

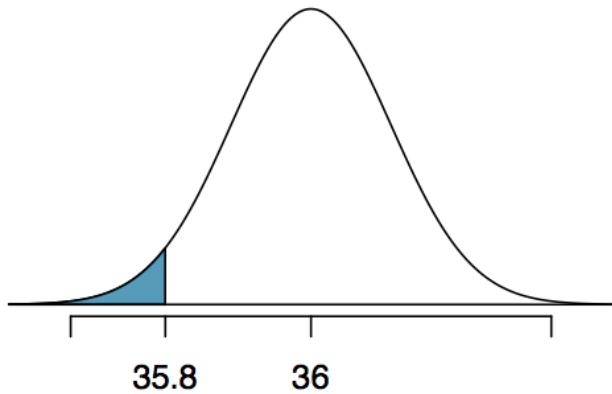
Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015

Quality control

At Heinz ketchup factory the amounts which go into bottles of ketchup are supposed to be normally distributed with mean 36 oz. and standard deviation 0.11 oz. Once every 30 minutes a bottle is selected from the production line, and its contents are noted precisely. If the amount of ketchup in the bottle is below 35.8 oz. or above 36.2 oz., then the bottle fails the quality control inspection. What percent of bottles have less than 35.8 ounces of ketchup?

- Let $X = \text{amount of ketchup in a bottle}$: $X \sim N(\mu = 36, \sigma = 0.11)$

$$Z = \frac{35.8 - 36}{0.11} = -1.82$$



Finding the exact probability -- using the Z table

Second decimal place of Z										Z
0.09	0.08	0.07	0.06	0.05	0.04	0.03	0.02	0.01	0.00	
0.0014	0.0014	0.0015	0.0015	0.0016	0.0016	0.0017	0.0018	0.0018	0.0019	-2.9
0.0019	0.0020	0.0021	0.0021	0.0022	0.0023	0.0023	0.0024	0.0025	0.0026	-2.8
0.0026	0.0027	0.0028	0.0029	0.0030	0.0031	0.0032	0.0033	0.0034	0.0035	-2.7
0.0036	0.0037	0.0038	0.0039	0.0040	0.0041	0.0043	0.0044	0.0045	0.0047	-2.6
0.0048	0.0049	0.0051	0.0052	0.0054	0.0055	0.0057	0.0059	0.0060	0.0062	-2.5
0.0064	0.0066	0.0068	0.0069	0.0071	0.0073	0.0075	0.0078	0.0080	0.0082	-2.4
0.0084	0.0087	0.0089	0.0091	0.0094	0.0096	0.0099	0.0102	0.0104	0.0107	-2.3
0.0110	0.0113	0.0116	0.0119	0.0122	0.0125	0.0129	0.0132	0.0136	0.0139	-2.2
0.0143	0.0146	0.0150	0.0154	0.0158	0.0162	0.0166	0.0170	0.0174	0.0179	-2.1
0.0183	0.0188	0.0192	0.0197	0.0202	0.0207	0.0212	0.0217	0.0222	0.0228	-2.0
0.0233	0.0239	0.0244	0.0250	0.0256	0.0262	0.0268	0.0274	0.0281	0.0287	-1.9
0.0294	0.0301	0.0307	0.0314	0.0322	0.0329	0.0336	0.0344	0.0351	0.0359	-1.8
0.0367	0.0375	0.0384	0.0392	0.0401	0.0409	0.0418	0.0427	0.0436	0.0446	-1.7
0.0455	0.0465	0.0475	0.0485	0.0495	0.0505	0.0516	0.0526	0.0537	0.0548	-1.6
0.0559	0.0571	0.0582	0.0594	0.0606	0.0618	0.0630	0.0643	0.0655	0.0668	-1.5

Finding the exact probability -- using the Z table

Second decimal place of Z							0.02			Z
0.09	0.08	0.07	0.06	0.05	0.04	0.03		0.01	0.00	
0.0014	0.0014	0.0015	0.0015	0.0016	0.0016	0.0017	0.0018	0.0018	0.0019	-2.9
0.0019	0.0020	0.0021	0.0021	0.0022	0.0023	0.0023	0.0024	0.0025	0.0026	-2.8
0.0026	0.0027	0.0028	0.0029	0.0030	0.0031	0.0032	0.0033	0.0034	0.0035	-2.7
0.0036	0.0037	0.0038	0.0039	0.0040	0.0041	0.0043	0.0044	0.0045	0.0047	-2.6
0.0048	0.0049	0.0051	0.0052	0.0054	0.0055	0.0057	0.0059	0.0060	0.0062	-2.5
0.0064	0.0066	0.0068	0.0069	0.0071	0.0073	0.0075	0.0078	0.0080	0.0082	-2.4
0.0084	0.0087	0.0089	0.0091	0.0094	0.0096	0.0099	0.0102	0.0104	0.0107	-2.3
0.0110	0.0113	0.0116	0.0119	0.0122	0.0125	0.0129	0.0132	0.0136	0.0139	-2.2
0.0143	0.0146	0.0150	0.0154	0.0158	0.0162	0.0166	0.0170	0.0174	0.0179	-2.1
0.0183	0.0188	0.0192	0.0197	0.0202	0.0207	0.0212	0.0217	0.0222	0.0228	-2.0
0.0233	0.0239	0.0244	0.0250	0.0256	0.0262	0.0268	0.0274	0.0281	0.0287	-1.9
0.0294	0.0301	0.0307	0.0314	0.0322	0.0329	0.0336	0.0344	0.0351	0.0359	-1.8
0.0367	0.0375	0.0384	0.0392	0.0401	0.0409	0.0418	0.0427	0.0436	0.0446	-1.7
0.0455	0.0465	0.0475	0.0485	0.0495	0.0505	0.0516	0.0526	0.0537	0.0548	-1.6
0.0559	0.0571	0.0582	0.0594	0.0606	0.0618	0.0630	0.0643	0.0655	0.0668	-1.5

Practice

What percent of bottles pass the quality control inspection?

(a) 1.82%

(d) 93.12%

(b) 3.44%

(e) 96.56%

(c) 6.88%

Practice

What percent of bottles pass the quality control inspection?

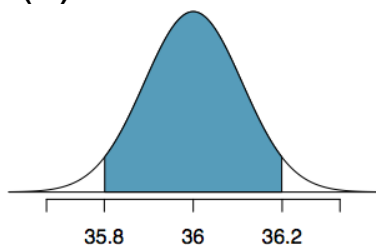
(a) 1.82%

(b) 3.44%

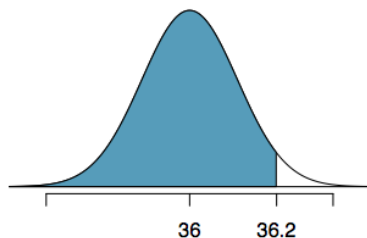
(c) 6.88%

(d) 93.12%

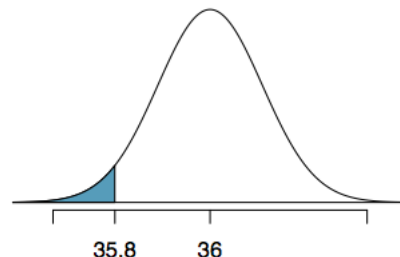
(e) 96.56%



=



-



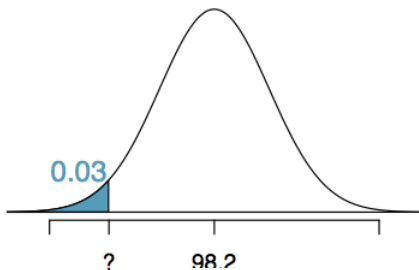
$$Z_{35.8} = \frac{35.8 - 36}{0.11} = -1.82$$

$$Z_{36.2} = \frac{36.2 - 36}{0.11} = 1.82$$

$$P(35.8 < X < 36.2) = P(-1.82 < Z < 1.82) = 0.9656 - 0.0344 = 0.9312$$

Finding cutoff points

Body temperatures of healthy humans are distributed nearly normally with mean 98.2°F and standard deviation 0.73°F. What is the cutoff for the lowest 3% of human body temperatures?



0.09	0.08	0.07	0.06	0.05	Z
0.0233	0.0239	0.0244	0.0250	0.0256	-1.9
0.0294	0.0301	0.0307	0.0314	0.0322	-1.8
0.0367	0.0375	0.0384	0.0392	0.0401	-1.7

$$P(X < x) = 0.03 \rightarrow P(Z < -1.88) = 0.03$$

$$Z = \frac{obs - mean}{SD} \rightarrow \frac{x - 98.2}{0.73} = -1.88$$

$$x = (-1.88 \times 0.73) + 98.2 = 96.8^{\circ}F$$

Mackowiak, Wasserman, and Levine (1992), A Critical Appraisal of 98.6 Degrees F, the Upper Limit of the Normal Body Temperature, and Other Legacies of Carl Reinhold August Wunderlick.

Practice

Body temperatures of healthy humans are distributed nearly normally with mean 98.2°F and standard deviation 0.73°F . What is the cutoff for the highest 10% of human body temperatures?

(a) 97.3°F

(b) 99.1°F

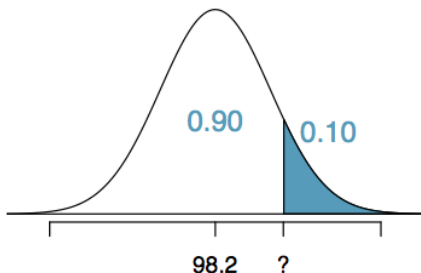
(c) 99.4°F

(d) 99.6°F

Practice

Body temperatures of healthy humans are distributed nearly normally with mean 98.2°F and standard deviation 0.73°F. What is the cutoff for the highest 10% of human body temperatures?

(a) 97.3°F



(c) 99.4°F

Z	0.05	0.06	0.07	0.08	0.09
1.0	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9115	0.9131	0.9147	0.9162	0.9177

$$P(X > x) = 0.10 \rightarrow P(Z < 1.28) = 0.90$$

$$Z = \frac{\text{obs} - \text{mean}}{SD} \rightarrow \frac{x - 98.2}{0.73} = 1.28$$

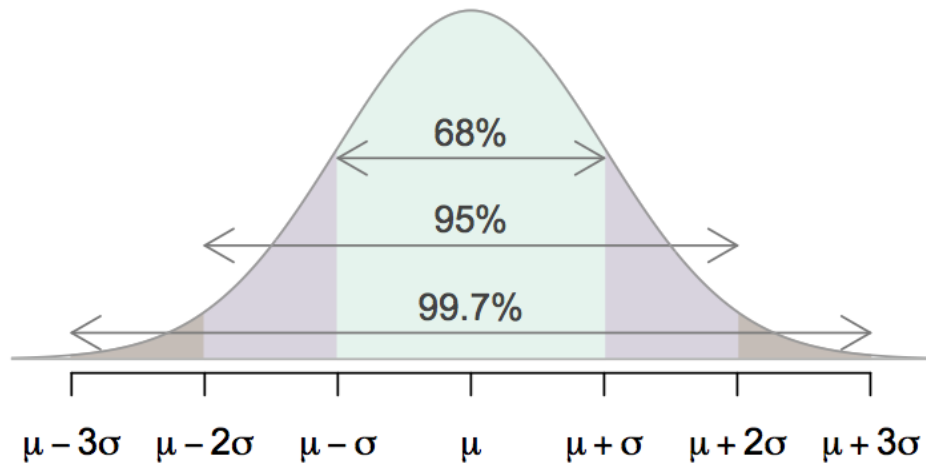
$$x = (1.28 \times 0.73) + 98.2 = 99.1$$

68-95-99.7 Rule

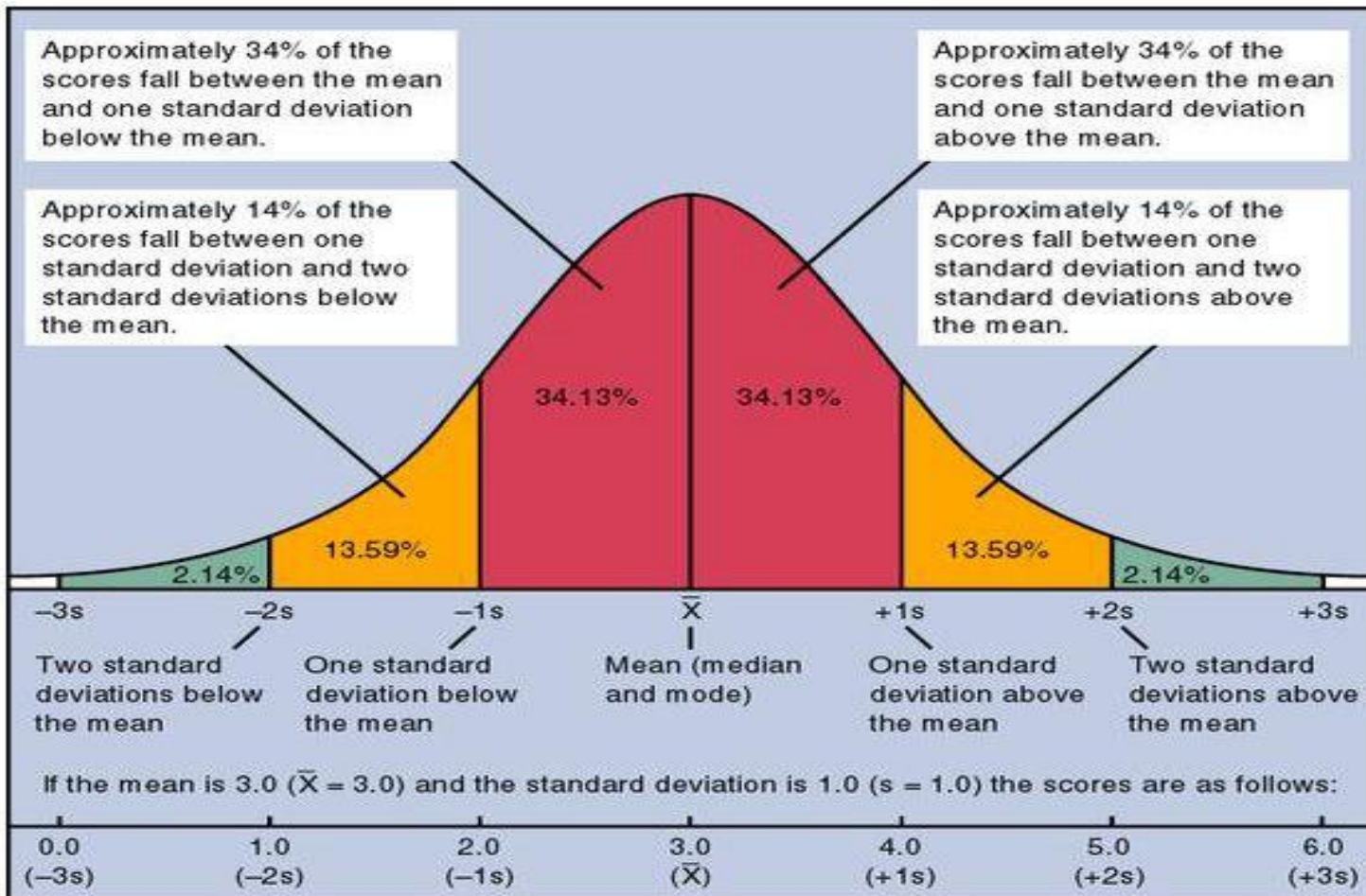
For nearly normally distributed data,

- about 68% falls within 1 SD of the mean,
- about 95% falls within 2 SD of the mean,
- about 99.7% falls within 3 SD of the mean.

It is possible for observations to fall 4, 5, or more standard deviations away from the mean, but these occurrences are very rare if the data are nearly normal.



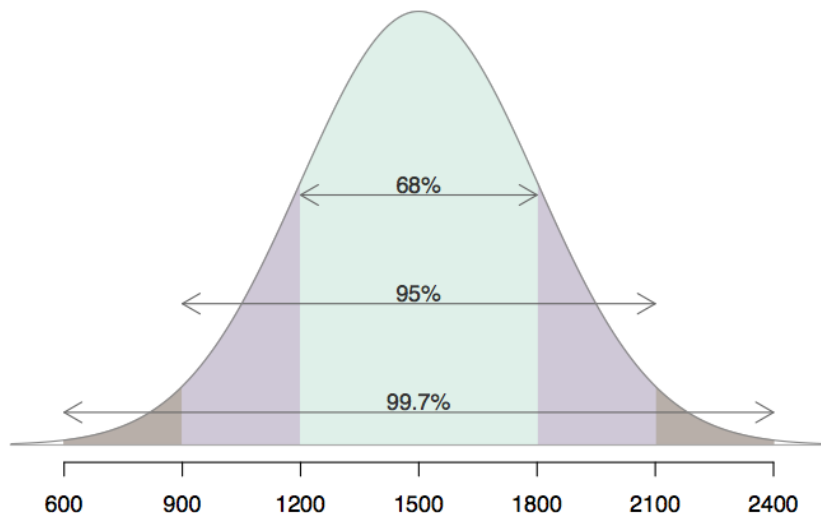
Normal Curve



Describing variability using the 68-95-99.7 Rule

SAT scores are distributed nearly normally with mean 1500 and standard deviation 300.

- ~68% of students score between 1200 and 1800 on the SAT.
- ~95% of students score between 900 and 2100 on the SAT.
- ~99.7% of students score between 600 and 2400 on the SAT.



Summary

- Probability
 - Independence
 - Independent events: $P(A \text{ or } B)$, $P(A \text{ and } B)$
- Probability Distributions
 - Normal Distribution
 - Z-Score (standardized scores)
 - Calculating probabilities