Examples of continuous probability distributions:

The normal and standard normal

The Normal Distribution



The Normal Distribution: as mathematical function (pdf)



The Normal PDF

It's a probability function, so no matter what the values of μ and σ , must integrate to 1!

$$\int_{-\infty}^{+\infty} \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx = 1$$

Normal distribution is defined by its mean and standard dev.

$$\mathsf{E}(\mathsf{X}) = \mu = \int_{-\infty}^{+\infty} x \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx$$

$$\operatorname{Var}(\mathsf{X}) = \sigma^2 = \int_{-\infty}^{+\infty} x^2 \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx - \mu^2$$

Standard Deviation(X)= σ

**The beauty of the normal curve:

No matter what μ and σ are, the area between μ - σ and μ + σ is about 68%; the area between μ - 2σ and μ + 2σ is about 95%; and the area between μ - 3σ and μ + 3σ is about 99.7%. Almost all values fall within 3 standard deviations.



68-95-99.7 Rule in Math terms...



How good is rule for real data?

Check some example data: The mean of the weight of the women = 127.8The standard deviation (SD) = 15.5 68% of 120 = .68x120 = ~ 82 runners

In fact, 79 runners fall within 1-SD (15.5 lbs) of the mean.



95% of 120 = .95 x 120 = ~ 114 runnersIn fact, 115 runners fall within 2-SD's of the mean.



99.7% of 120 = .997 x 120 = 119.6 runners In fact, all 120 runners fall within 3-SD's of the mean.



Example

- Suppose SAT scores roughly follows a normal distribution in the U.S. population of college-bound students (with range restricted to 200-800), and the average math SAT is 500 with a standard deviation of 50, then:
 - 68% of students will have scores between 450 and 550
 - 95% will be between 400 and 600
 - 99.7% will be between 350 and 650

BUT...

What if you wanted to know the math SAT score corresponding to the 90th percentile (=90% of students are lower)?

$P(X \leq Q) = .90 \rightarrow$

$$\int_{200}^{Q} \frac{1}{(50)\sqrt{2\pi}} \bullet e^{-\frac{1}{2}(\frac{x-500}{50})^2} dx = .90$$

Solve for Q?....Yikes!

The Standard Normal (Z): "Universal Currency"

The formula for the standardized normal probability density function is

 $-\frac{1}{2}(\frac{Z-0}{1})^2$ $\frac{1}{\sqrt{2}} \cdot e$ $\frac{1}{(1)\sqrt{2\pi}} \cdot e$

The Standard Normal Distribution (Z)

All normal distributions can be converted into the standard normal curve by subtracting the mean and dividing by the standard deviation:

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

Somebody calculated all the integrals for the standard normal and put them in a table! So we never have to integrate!

Even better, computers now do all the integration.



Example

• For example: What's the probability of getting a math SAT score of 575 or less, μ =500 and σ =50?

$$Z = \frac{575 - 500}{50} = 1.5$$

•i.e., A score of 575 is 1.5 standard deviations above the mean

$$\therefore P(X \le 575) = \int_{200}^{575} \frac{1}{(50)\sqrt{2\pi}} \cdot e^{-\frac{1}{2}(\frac{x-500}{50})^2} dx \longrightarrow \int_{-\infty}^{1.5} \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2}Z^2} dz$$

Yikes!

But to look up Z= 1.5 in standard normal chart (or enter into SAS) \rightarrow no problem! = .9332

Practice problem

If birth weights in a population are normally distributed with a mean of 109 oz and a standard deviation of 13 oz,

- a. What is the chance of obtaining a birth weight of 141 oz *or heavier* when sampling birth records at random?
- b. What is the chance of obtaining a birth weight of 120 *or lighter*?

Answer

a. What is the chance of obtaining a birth weight of 141 oz *or heavier* when sampling birth records at random?

$$Z = \frac{141 - 109}{13} = 2.46$$

From the chart or SAS \rightarrow Z of 2.46 corresponds to a right tail (greater than) area of: P(Z \ge 2.46) = 1-(.9931)= .0069 or .69 %

Answer

b. What is the chance of obtaining a birth weight of 120 *or lighter*?

$$Z = \frac{120 - 109}{13} = .85$$

From the chart or SAS \rightarrow Z of .85 corresponds to a left tail area of: P(Z \le .85) = .8023 = 80.23%

Looking up probabilities in the standard normal table



Normal probabilities in SAS

data _null_;
 theArea=probnorm(1.5);
 put theArea;
run;

0.9331927987

The "probnorm(Z)" function gives you the probability from negative infinity to Z (here 1.5) in a standard normal curve.

And if you wanted to go the other direction (i.e., from the area to the Z score (called the so-called "Probit" function \rightarrow

data _null_;
 theZValue=probit(.93);
 put theZValue;
run;
1.4757910282

The "probit(p)" function gives you the Z-value that corresponds to a left-tail area of p (here .93) from a standard normal curve. The probit function is also known as the inverse standard normal function.

Probit function: the inverse

 $\phi(area)$ = Z: gives the Z-value that goes with the probability you want For example, recall SAT math scores example. What's the score that corresponds to the 90th percentile?

In Table, find the Z-value that corresponds to area of $.90 \rightarrow Z= 1.28$ <u>Or use SAS</u>

```
data _null_;
```

```
theZValue=probit(.90);
```

```
put theZValue;
```

run;

1.2815515655

If Z=1.28, convert back to raw SAT score \rightarrow 1.28 = $\frac{X-500}{50}$ X - 500 =1.28 (50) X=1.28(50) + 500 = 564 (1.28 standard deviations above the mean!)

Are my data "normal"?

- Not all continuous random variables are normally distributed!!
- It is important to evaluate how well the data are approximated by a normal distribution

Are my data normally distributed?

- Look at the histogram! Does it appear bell shaped?
- 2. Compute descriptive summary measures—are mean, median, and mode similar?
- Do 2/3 of observations lie within 1 std dev of the mean? Do 95% of observations lie within 2 std dev of the mean?
- 4. Look at a normal probability plot—is it approximately linear?
- 5. Run tests of normality (such as Kolmogorov-Smirnov). But, be cautious, highly influenced by sample size!





































The Normal Probability Plot

- Normal probability plot
 - Order the data.
 - Find corresponding standardized normal quantile values: i^{th} quantile = $\phi(\frac{i}{n+1})$

where ϕ is the probit function, which gives the Z value that corresponds to a particular left - tail area

- Plot the observed data values against normal quantile values.
- Evaluate the plot for evidence of linearity.

Normal probability plot coffee...

Normal Probability Plot, Coffee



Normal probability plot love of writing...

Normal Probability Plot, Love of Writing



Neither right-skewed or left-skewed, but big gap at 6.

Norm prob. plot Exercise...

Normal Probability Plot, Exercise



Norm prob. plot Wake up time

Normal Probability Plot, Wake up times



Closest to a straight line...

Formal tests for normality

- Results:
- Coffee: Strong evidence of non-normality (p<.01)
- Writing love: Moderate evidence of nonnormality (p=.01)
- Exercise: Weak to no evidence of nonnormality (p>.10)
- Wakeup time: No evidence of non-normality (p>.25)

Normal approximation to the binomial

When you have a binomial distribution where *n* is large and *p* is middle-of-the road (not too small, not too big, closer to .5), then the binomial starts to look like a normal distribution \rightarrow in fact, this doesn't even take a particularly large $n \rightarrow$

<u>Recall:</u> What is the probability of being a smoker among a group of cases with lung cancer is .6, what's the probability that in a group of 8 cases you have less than 2 smokers?

Normal approximation to the binomial

When you have a binomial distribution where n is large and p isn't too small (rule of thumb: mean>5), then the binomial starts to look like a normal distribution->



Starting to have a normal shape even with fairly small n. You can imagine that if n got larger, the bars would get thinner and thinner and this would look more and more like a continuous function, with a bell curve shape. Here np=4.8.



What is the probability of fewer than 2 smokers? <u>Exact</u> binomial probability (from before) = .00065 + .008 = .00865

Normal approximation probability: $\mu=4.8$ $\sigma=1.39$ $Z \approx \frac{2-(4.8)}{1.39} = \frac{-2.8}{1.39} = -2$

$$P(Z < 2) = .022$$

A little off, but in the right ballpark... we could also use the value to the left of 1.5 (as we really wanted to know less than but not including 2; called the "continuity correction")...

$$Z \approx \frac{1.5 - (4.8)}{1.39} = \frac{-3.3}{1.39} = -2.37$$

 $P(Z \le -2.37) = .0069$

A fairly good approximation of the exact probability, .00865.

Practice problem

 You are performing a cohort study. If the probability of developing disease in the exposed group is .25 for the study duration, then if you sample (randomly) 500 exposed people, What's the probability that <u>at</u> <u>most</u> 120 people develop the disease?

Answer

By hand (yikes!):

 $P(X \le 120) = P(X=0) + P(X=1) + P(X=2) + P(X=3) + P(X=4) + \ldots + P(X=120) =$

$$\begin{pmatrix} 500 \\ 120 \end{pmatrix} (.25)^{120} (.75)^{380} + \begin{pmatrix} 500 \\ 2 \end{pmatrix} (.25)^{2} (.75)^{498} + \begin{pmatrix} 500 \\ 1 \end{pmatrix} (.25)^{1} (.75)^{499} + \begin{pmatrix} 500 \\ 0 \end{pmatrix} (.25)^{0} (.75)^{500} \cdots$$

OR Use SAS:

data _null_;

Cohort=cdf('binomial', 120, .25, 500);

put Cohort;

run;

0.323504227

OR use, normal approximation:

 $\mu = np = 500(.25) = 125$ and $\sigma^2 = np(1-p) = 93.75$; $\sigma = 9.68$

$$Z = \frac{120 - 125}{9.68} = -.52$$

P(Z<-.52)= .3015

Proportions...

- The binomial distribution forms the basis of statistics for proportions.
- A proportion is just a binomial count divided by n.
 - For example, if we sample 200 cases and find 60 smokers, X=60 but the observed proportion=.30.
- Statistics for proportions are similar to binomial counts, but differ by a factor of n.



It all comes back to Z...

 Statistics for proportions are based on a normal distribution, because the binomial can be approximated as normal if np>5