

# Correlating data

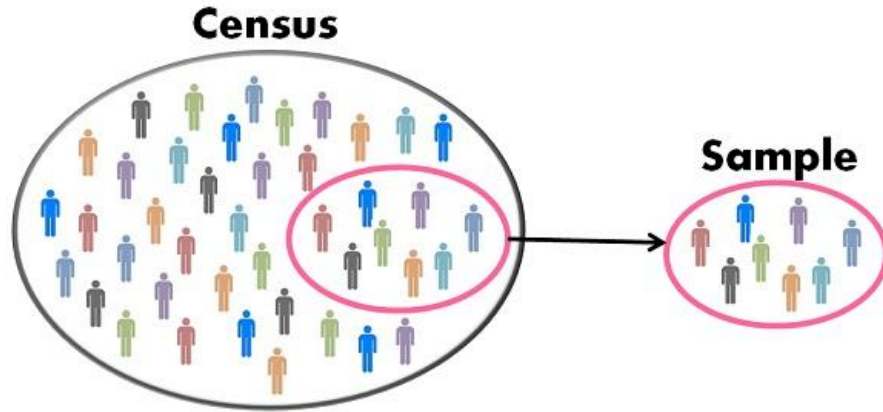
# Introduction

## Topics

- Surveys
- Mean/Average
- Bell/Normal/Gaussian distribution curve
- Standard Deviation
- Z-Scores
- Scatterplot graph
- Equation of a line
- Line of best fit
- Pearson's  $r$  correlation coefficient
- Regression line

# Surveys

You have probably heard of a census. This is a measure of information about the population. A full census is expensive, so it is common to survey a subset of the population using a sample.



# Surveys

Population distribution symbols are usually written with greek letters, and sample distribution symbols with roman letters, sampling distribution symbols are a combination of greek letters with a subscript of roman letters

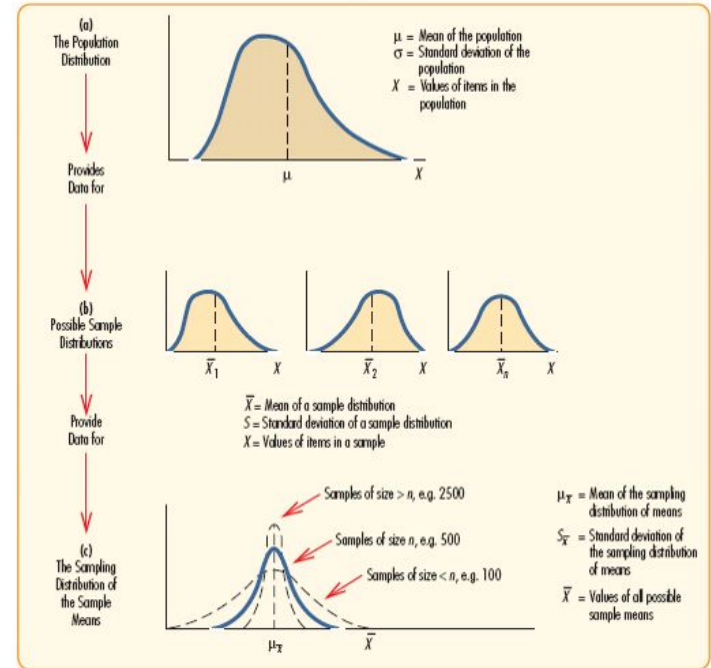
	Population	Sample	Sampling
Mean/Average	$\mu$ (mu)	$\bar{x}$ (x-bar)	$\mu_{\bar{x}}$
Standard deviation	$\sigma$ (sigma)	$s$ (s)	$\sigma_{\bar{x}}$
Number/Size	$N$ (nu)	$n$ (n)	

# Fundamental types of distributions

- Population distribution
- Sample distribution(s)
- Sampling distribution

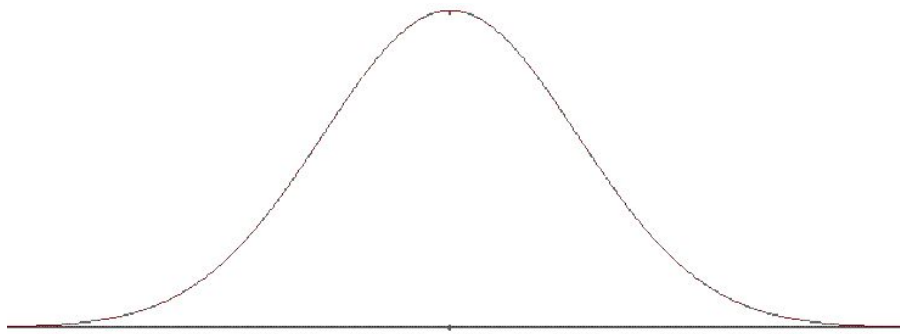
Even if the population or sample distributions are not normally distributed, the sampling distribution will be normally distributed, thanks to the central limit theorem.

EXHIBIT 12.6 Schematic of the three fundamental types of distributions<sup>8</sup>



# Normal distribution curve

The normal distribution is useful because of the central limit theorem. In its most general form, under some conditions, the central limit theorem states that average of samples of observations independently drawn from separate independent distributions become normally distributed when the number of observations is sufficiently large.



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

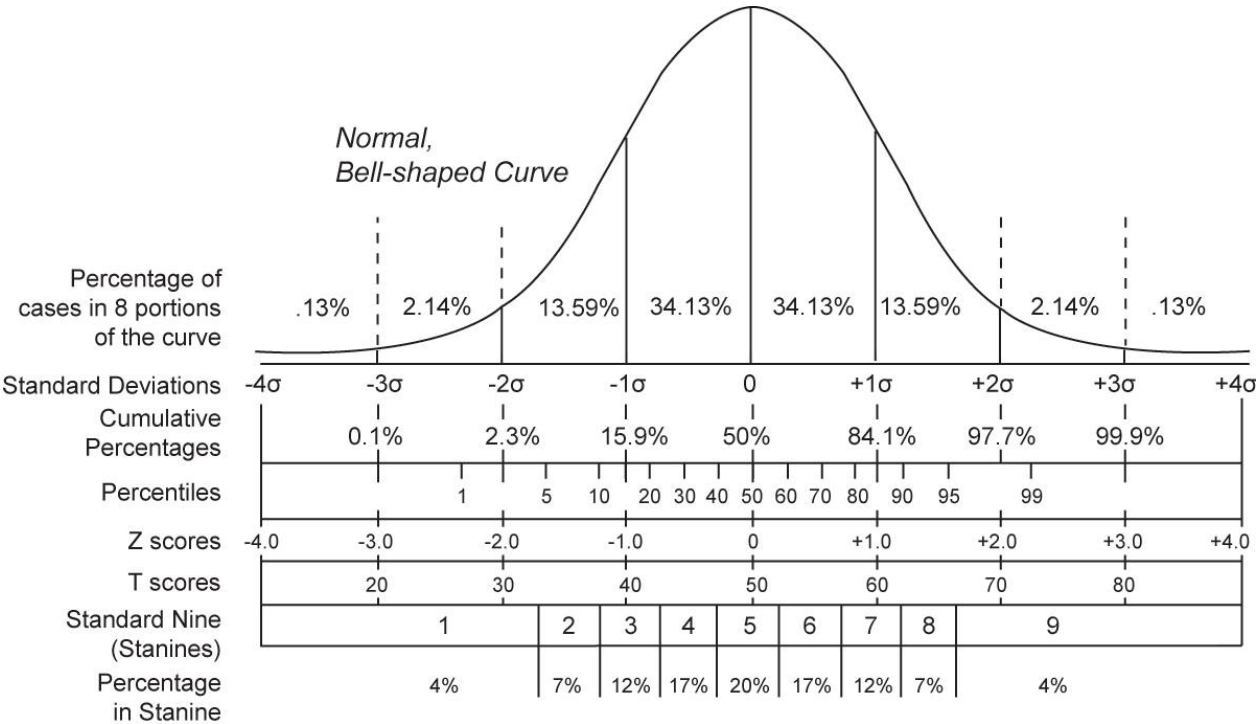
$\mu$  = Mean

$\sigma$  = Standard Deviation

$\pi \approx 3.14159 \dots$

$e \approx 2.71828 \dots$

# Normal distribution curve



# Central limit theorem

Imagine you roll a single dice, the probability of predicting the outcome is equal for all values.  $\frac{1}{6}$  chance

However consider rolling 2 dice simultaneously and adding the values. Select a number between 2 and 12.

The numbers 2 and 12 can only occur with the combinations 1,1 and 6,6.

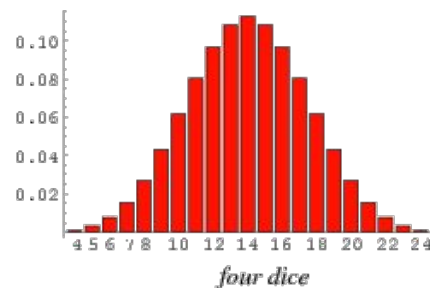
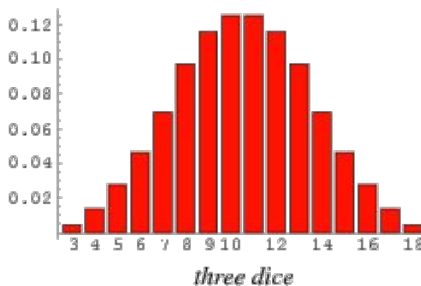
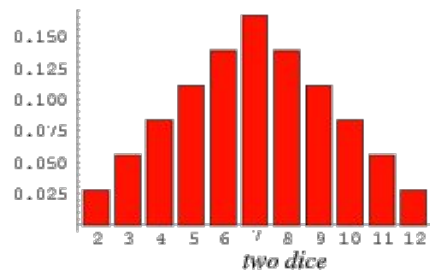
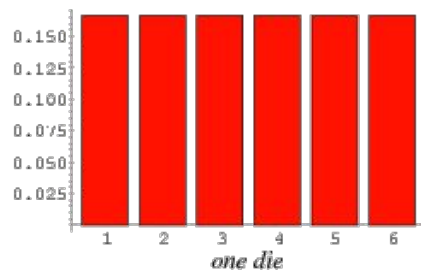
The number 10 can occur with 5,5 and 5,5 and 4,6 and 6,4.



# Central limit theorem

Consider adding more dice.

As more dice are added  
the average sum of the values  
approaches a normal distribution.



# Average

The average is defined as the sum of the values divided by the total number of values. The full sigma notation is shown below.

Sample Mean	Population Mean
$\bar{x} = \frac{\sum x}{n}$	$\mu = \frac{\sum x}{N}$

where  $\sum x$  is sum of all data values

$N$  is number of data items in population

$n$  is number of data items in sample

$$\text{average} = \frac{\text{sum of values}}{\text{number of values}}$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

# Standard deviation

The variance is (average difference)<sup>2</sup>

The standard deviation is the sqrt( (average difference)<sup>2</sup> )

You may notice that the sample standard deviation uses n-1 and the population standard deviation uses N. This is to account for the degrees of freedom.

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

$$s_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

$n$  = The number of data points

$\bar{x}$  = The mean of the  $x_i$

$x_i$  = Each of the values of the data

# Z-score

The standard deviation is an overall measure of the variability of the data and applies to the entire data set.

The Z-score is the number of standard deviations an individual data point lies from the mean.

To calculate Z-score, simply subtract the mean from the data point and divide the result by the standard deviation.

$$z = \frac{x - \mu}{\sigma}$$

$\mu$  = Mean

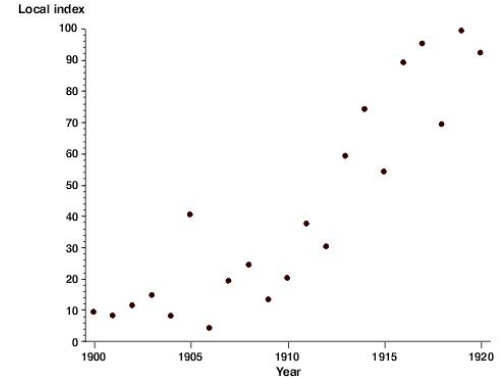
$\sigma$  = Standard Deviation

# Scatter plot graphs

A scatter plot graph is a type of diagram using (x,y) coordinates to display values for typically two variables of a data set.

A scatter plot can be used either when one variable is under the control of the experimenter (customarily plotted along the horizontal axis) and the other depends on it (customarily plotted along the vertical axis).

When both variables are independent either can be plotted on either axis and a scatter plot will illustrate only the degree of correlation (not causation) between two variables.



# Line equation

Given 2 points  $(x_1, y_1)$  and  $(x_2, y_2)$ , or  $(-3, 3)$  and  $(3, -1)$

A line can be defined by the equation of a line  $y=mx+b$  and the slope can be calculated using the slope formula.

The constant  $b$  is the  $y$  intercept is where the line crosses the  $y$ -axis, where  $x=0$ . This can be done by substituting the values for the slope of the line and the coordinates of a single point  $(x_1, y_1)$  in the line equation, set  $x=0$  and then solve for  $b$ .

$$y = \underset{\substack{\uparrow \\ \text{slope}}}{m}x + \underset{\substack{\uparrow \\ \text{y-intercept}}}{b}$$

$$\text{Slope} = m = \frac{y_2 - y_1}{x_2 - x_1}$$

# Regression line equation

A regression line can be defined by its formula  $\hat{Y}=bX+a$

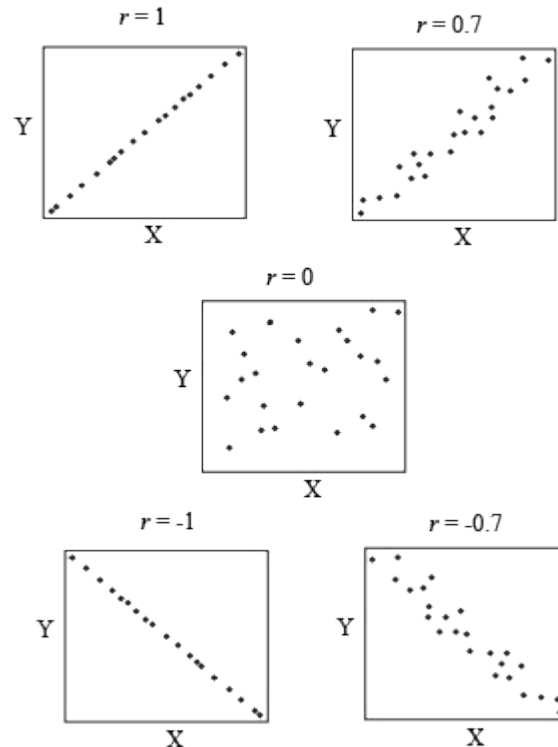
$\hat{Y}$  or y-hat is the predicted value for y based on the estimated variable values for the line of best fit. This is usually used with scatter plot data where the line does not fit perfectly along all the points.

# Pearson's correlation coefficient

Pearson's  $r$  correlation coefficient is a measure of the linear correlation between two variables  $X$  and  $Y$ .

The  $r$  value can range between 1.0 and -1.0

A value of 1 or -1 denotes a perfect positive or negative linear correlation, and a value of 0 denotes no linear correlation.





# Pearson's correlation coefficient

The pearson's correlation formula can be expressed as a relationship of z-scores.

$$\text{Sample: } r = \frac{\sum z_X z_Y}{n-1}$$

$$\text{Population: } \rho = \frac{\sum z_X z_Y}{N}$$

# Estimating regression line parameters

If the regression line is defined as  $\hat{Y}=bX+a$  then how can we estimate the variables  $a$  and  $b$ ?

$b$  can be estimated by multiplying the pearson coefficient value by the dividing the quotient of the standard deviation of the  $y$  values divided by the the standard deviation of the  $x$  values.

$a$  can be estimated using the estimated  $b$  value, the mean of  $y$  and the mean of  $x$ .

$$b = r \frac{S_Y}{S_X}$$

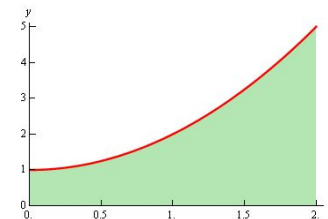
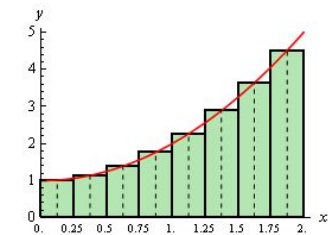
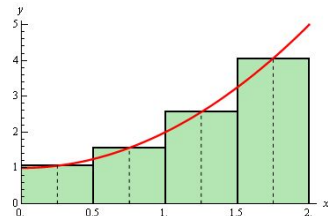
$$a = \bar{Y} - b\bar{X}$$

# Integrals for area under the curve

Integration is one of the two main operations of calculus, with its inverse, differentiation.

It is defined informally as the area of the region in the  $xy$ -plane that is bounded by the graph of  $f()$ , the  $x$ -axis and the vertical lines  $x = a$  and  $x = b$ .

The principles of integration were formulated independently by Isaac Newton and Gottfried Leibniz in the late 17th century, who thought of the integral as an infinite sum of rectangles of infinitesimal width.



# How to find normally distributed probabilities

The standard normal distribution table provides the probability that a normally distributed random variable  $Z$ , with mean equal to 0 and variance equal to 1, is less than or equal to  $z$ .

Each cell in the table represents the area  $P$  under the standard normal curve, below the respective  $z$ -statistic.

The label for rows contains the integer part and the first decimal place of  $Z$ .

The label for columns contains the second decimal place of  $Z$

# The standard normal table (z-table)

The z-table allows us to calculate probabilities of encountering values within range on a normally distributed variable.

Probabilities are calculated as a specified area under the standard normal curve.

As computing integrals is a lengthy process, tables were introduced to accelerate and simplify the procedure.

# How to find normally distributed probabilities

Since probability tables cannot be printed for every normal distribution, as there are an infinite variety of normal distributions, it is common practice to convert a normal to a standard normal and then use the standard normal table to find probabilities.

To convert any value to its standardised value, simply calculate its z-score. The collection of values is its distribution.

# T-distributions

The t-distribution can be used to say how confident you are that any given range of values would contain the true population mean, in situations where the sample size is small ( $<30$ ) and population standard deviation is unknown.

The t-distribution is similar to the standard normal distribution but it has thicker tails to account for more variance. The t-distribution distribution takes a parameter called df or degrees of freedom. This is simply the number of samples  $n-1$ . As the value for df increases the shape of the distribution approaches the shape of the standard normal distribution and converges at around  $df=30$ .

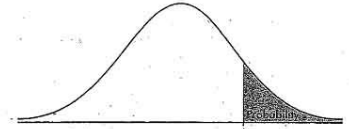
# T-tables

As computing integrals is a lengthy process, t-tables were introduced to accelerate and simplify the procedure.

Each row represents the degrees of freedom (n-1)

Each column represents the offset from a tail end (upper or lower)

Each cell represents the area under the curve between the tail end and the specified offset. The area represents the probability of capturing the mean within that range.

**TABLE B: *t*-DISTRIBUTION CRITICAL VALUES**

df	Tail probability $p$											
	.25	.20	.15	.10	.05	.025	.01	.005	.0025	.001	.0005	.0001
1	1.000	1.376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3	696.6
2	.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33	31.60
3	.727	.978	1.250	1.638	2.262	3.438	3.919	5.481	7.704	10.83	15.09	22.92
4	.704	.941	1.190	1.533	2.132	3.299	3.799	5.347	7.404	10.59	17.13	26.10
5	.727	.920	1.156	1.476	2.015	3.271	3.757	5.365	7.403	10.73	17.59	26.86
6	.718	.906	1.134	1.440	1.943	3.247	3.742	5.341	7.377	10.47	17.28	25.99
7	.711	.899	1.119	1.415	1.895	3.236	3.731	5.298	7.349	10.42	17.48	25.08
8	.706	.893	1.109	1.397	1.860	3.206	3.699	5.286	7.335	10.41	17.50	25.04
9	.703	.888	1.103	1.385	1.845	3.193	3.689	5.280	7.326	10.40	17.50	25.04
10	.700	.887	1.099	1.372	1.812	3.228	3.598	5.264	7.169	10.381	17.44	24.87
11	.697	.876	1.088	1.363	1.796	3.201	3.528	5.218	7.106	10.349	17.407	24.45
12	.695	.873	1.083	1.356	1.782	3.179	3.503	5.201	7.082	10.328	17.40	24.21
13	.694	.870	1.079	1.350	1.771	3.160	3.482	5.200	7.071	10.322	17.382	24.21
14	.693	.868	1.076	1.345	1.761	3.154	3.464	5.204	7.067	10.320	17.377	24.16
15	.691	.866	1.074	1.341	1.753	3.131	3.449	5.202	7.047	10.326	17.373	24.03
16	.690	.865	1.071	1.337	1.746	3.120	3.425	5.283	7.021	10.325	17.366	24.01
17	.689	.863	1.069	1.333	1.740	3.110	3.422	5.267	7.008	10.322	17.366	23.95
18	.688	.862	1.067	1.330	1.734	3.101	3.412	5.252	7.000	10.319	17.361	23.92
19	.687	.861	1.066	1.327	1.729	3.092	3.403	5.245	6.991	10.314	17.359	23.88
20	.687	.860	1.064	1.325	1.725	3.086	3.397	5.238	6.984	10.312	17.358	23.85
21	.686	.859	1.063	1.323	1.721	3.080	3.391	5.238	6.981	10.311	17.357	23.81
22	.686	.858	1.061	1.321	1.717	3.074	3.383	5.230	6.975	10.310	17.350	23.792
23	.685	.858	1.060	1.319	1.714	3.069	3.377	5.207	6.970	10.314	17.348	23.768
24	.685	.857	1.059	1.318	1.711	3.064	3.372	5.192	6.964	10.309	17.347	23.745
25	.684	.856	1.058	1.316	1.708	3.060	3.367	5.185	6.958	10.313	17.345	23.721
26	.684	.856	1.058	1.315	1.706	3.056	3.362	5.182	6.953	10.309	17.346	23.705
27	.684	.855	1.057	1.314	1.703	3.052	3.358	5.173	6.947	10.307	17.341	23.690
28	.683	.855	1.056	1.313	1.701	3.048	3.354	5.167	6.943	10.304	17.340	23.674
29	.683	.854	1.055	1.311	1.699	3.045	3.350	5.162	6.938	10.308	17.336	23.659
30	.683	.854	1.055	1.310	1.698	3.042	3.347	5.157	6.934	10.305	17.334	23.641
40	.681	.851	1.050	1.303	1.692	3.021	3.325	5.142	6.920	10.294	17.324	23.551
50	.679	.849	1.047	1.299	1.676	3.009	3.310	5.103	6.908	10.283	17.317	23.496
60	.679	.848	1.045	1.296	1.671	3.000	3.300	5.099	6.900	10.281	17.315	23.466
80	.678	.846	1.043	1.292	1.664	2.990	3.288	5.078	6.887	10.273	17.311	23.440
100	.677	.845	1.042	1.290	1.660	2.984	3.281	5.069	6.880	10.271	17.310	23.430
1000	.675	.842	1.037	1.282	1.646	2.956	3.250	5.031	6.841	10.258	17.298	23.388
∞	.674	.841	1.036	1.282	1.645	2.950	3.250	5.030	6.840	10.257	17.297	23.391
50% 60% 70% 80% 90% 95% 96% 98% 99% 99.5% 99.8% 99.9% 99.95%												
Confidence level $C$												



# T-values

To solve the problem of calculating a t-values with an unknown population standard deviation, the value is estimated using the sample standard deviation.

This leads to the following formula:  $\bar{X}$  plus and minus the z-score for the 95% confidence level times the estimated standard deviation of the sampling distribution of the sample mean, which equals the sample standard deviation divided by the square root of the sample size. We call this estimated standard deviation of the sampling distribution the standard error. But because we now estimate the standard deviation we add extra error in our computation. For that reason we employ another distribution than the standard normal distribution (also called the z-distribution) we employed previously. Because of the extra error we now use the T-distribution. That leads to this formula:  $\bar{X}$  plus and minus the t-score for the 95% confidence level times the estimated standard deviation of the

# Logistic Regression

Sum the values of the log likelihood for each value of x and y

$$\beta_0=0.0, \beta_1=0.0, -2 \sum_{i=1}^n y_i \log_e \left( \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right) + (1 - y_i) \log_e \left( 1 - \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \right)$$

To begin we give the parameters  $B_0$  and  $B_1$  some arbitrary values with the understanding that these starting value will be replaced by optimized values later on. We start with small value and the estimation process will increase them gradually. After these initial parameters are estimated, the search process repeats using newton's method until the Log Likelihood is maximised and does not change significantly.

# Maximum likelihood estimation

The regression coefficients are usually estimated using maximum likelihood estimation. Unlike linear regression with normally distributed values, it is not possible to find a closed-form expression for the coefficient values that maximize the likelihood function, so iterative processes must be used instead.

Typically, the log likelihood is maximized using a something like a gradient descent algorithm such as Generalized Reduced Gradient (GRG) or Newton's method.

These processes begins with a tentative solution, revises it slightly to see if it can be improved, and repeats this revision until no more improvement is made, at which point the process is said to have converged.

# Logistic Regression

To graph  $y$ , simply plug in the constant, coefficient and  $x$  value

$$P(y|x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$