

# Ordinal Data

Mike Neale

Brad Verhulst

Sarah Medland

et al

Boulder Workshop Tuesday March 3 2020

Special thanks to Frühling Rijdsijk and those who came before

# Analysis of ordinal variables

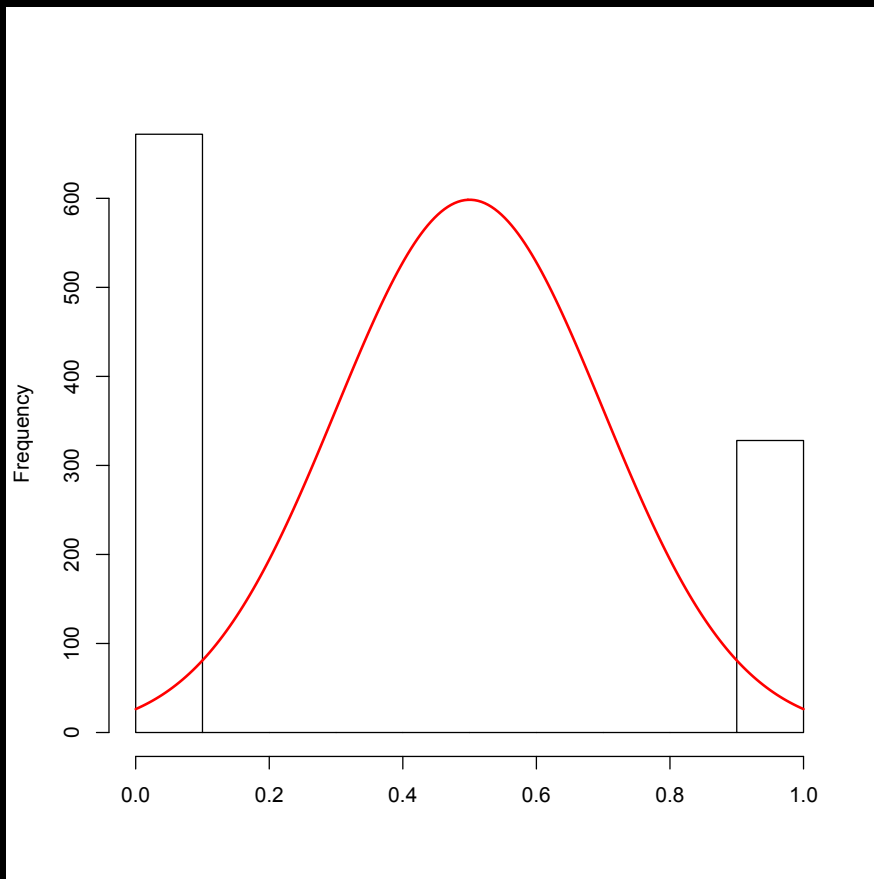
- Provide intuitive sense of how we estimate correlations from ordinal data
- Introduce the concept of liability threshold models
- Provide a mathematical description of the model
- Example exercise

# Ordinal data

- Measuring instrument discriminates between two or a few ordered categories, e.g.,
  - Absence (0) or presence (1) of a disorder
  - Severity of a disorder
  - Score on a single Likert item 'none/some/lots'
  - Number of symptoms (far from ideal)
- In such cases the data take the form of counts, i.e. the number of individuals within each category of response

# Problems with the treating ordinal variables as continuous

- Normality – Ordinal variables are not distributed normally, *obviously*.



- This means that the error terms cannot be normally distributed

# Two Ways of Thinking about Binary Dependent Variables

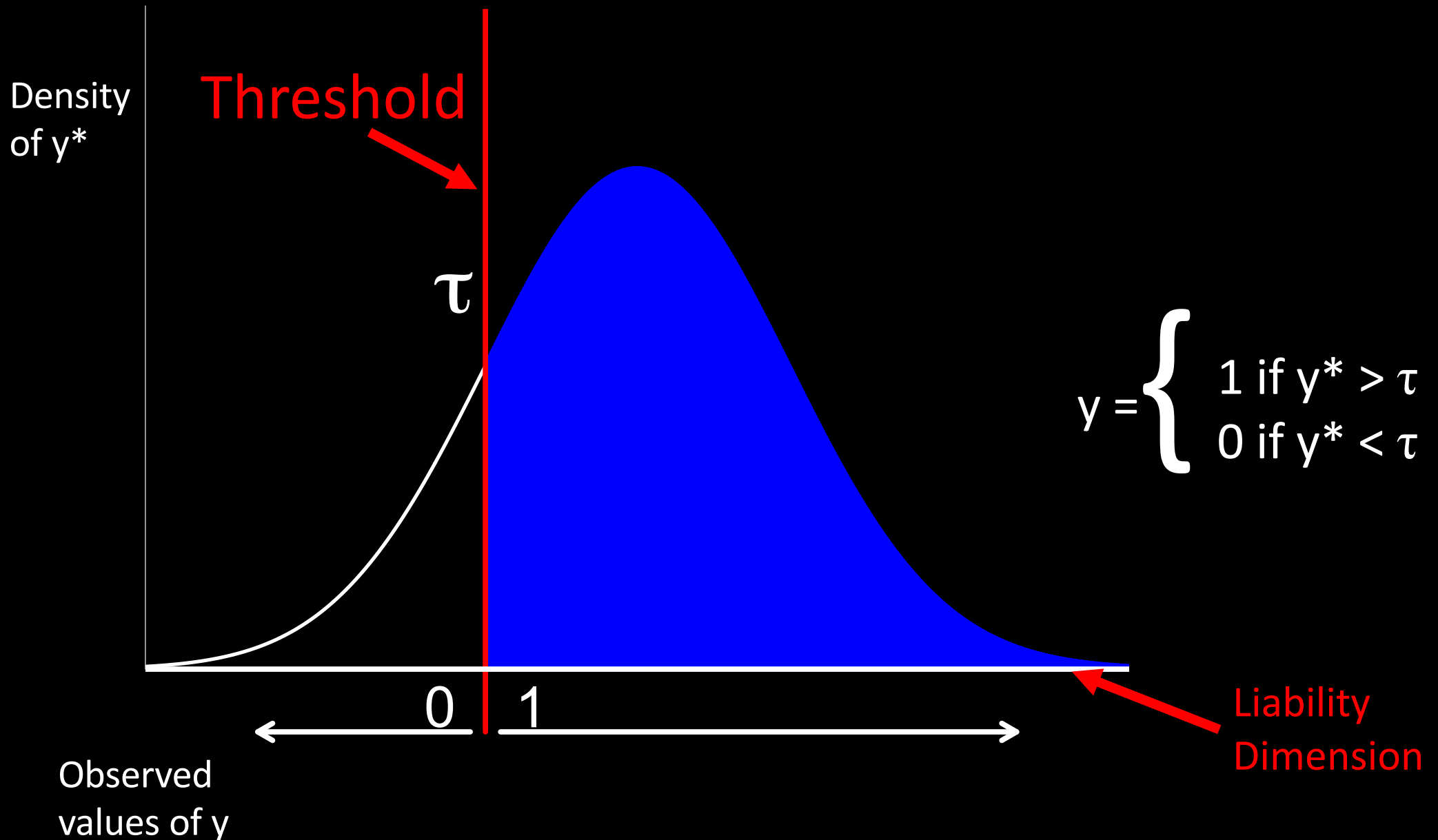
1. Assume that the observed binary variable is indicative of an underlying, latent (unobserved) continuous, normally distributed variable.
  - We call the unobserved variable a **Liability**
2. Assume the Binary Variable as a random draw from a Binomial (or Bernoulli) Distribution (Non-Linear Probability Model). Genuinely categorical responses, no underlying continuous distribution.

# Binary Variables as indicators of Latent Continuous Variables

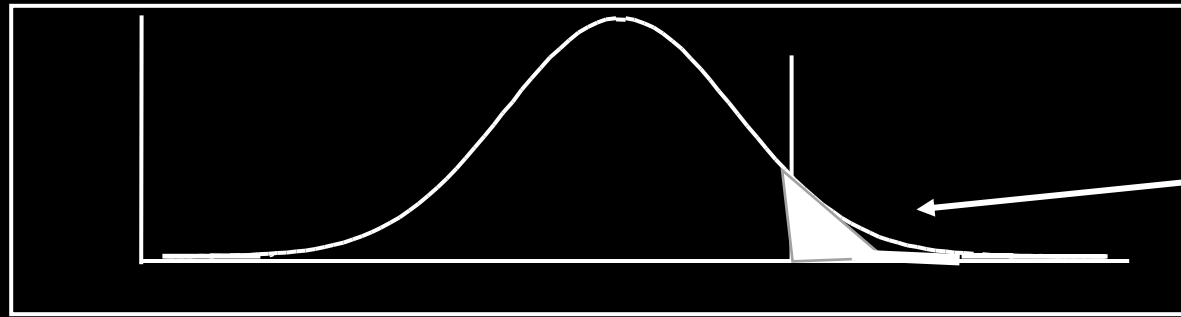
- Assume that the observed binary variable is indicative of an underlying, latent (unobserved) continuous, normally distributed variable.
- Assumptions:
  1. Categories reflect an imprecise measurement of an underlying *normal distribution* of liability. This liability is thought to be influenced by many many things, each of which does almost nothing. The Central Limit Theorem predicts that variation should be distributed according to the normal or Gaussian distribution.
  2. The liability distribution has 1 or more *thresholds*

# Fundamentals of the Threshold Model

Curnow & Smith 1975  
In mcn/2020/Papers



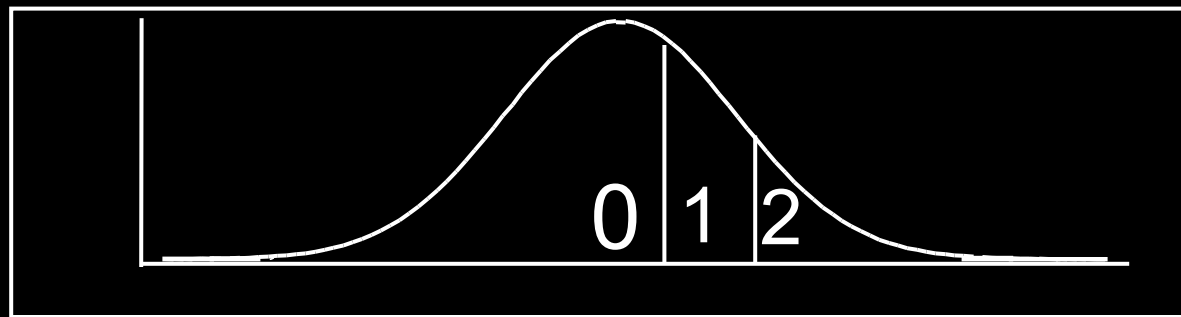
## For disorders:



Affected  
individuals

The **risk** or **liability** to a disorder is normally distributed  
When an individual exceeds a threshold they have the disorder.  
**Prevalence**: proportion of affected individuals.

## For a single questionnaire item score, e.g.,



0 = not at all  
1 = sometimes  
2 = always

Does not make sense to talk about prevalence: we simply count the endorsements of each response category



# Ideas behind the Liability Threshold Model (LTM)

- We can only observe binary outcomes, affected or unaffected, but people can be more or less affected.
- Since the variables are latent (and therefore not directly observed) we cannot estimate the means and variances we did for continuous variables.
- Thus, we have to make assumptions about them (pretend that they are some arbitrary value).

# Identifying Assumptions

## Mean Assumption

The intercept (mean) is 0

or

The threshold is 0 ( $\tau = 0$ )

- Either of these two assumptions provide equivalent model fit and the intercept is a transformation of  $\tau$ .

The traditional assumption

## Variance Assumption

$\text{Var}(\varepsilon|x) = 1$  in the normal-ogive model

$\text{Var}(\varepsilon|x) = \pi^2/3$  in the logit model.

The Probit Model

The Logit Model

## Assumption 3

The conditional mean of  $\varepsilon$  is 0.

- This is the same assumption as we make for continuous variables, and allows the parameters to be unbiased

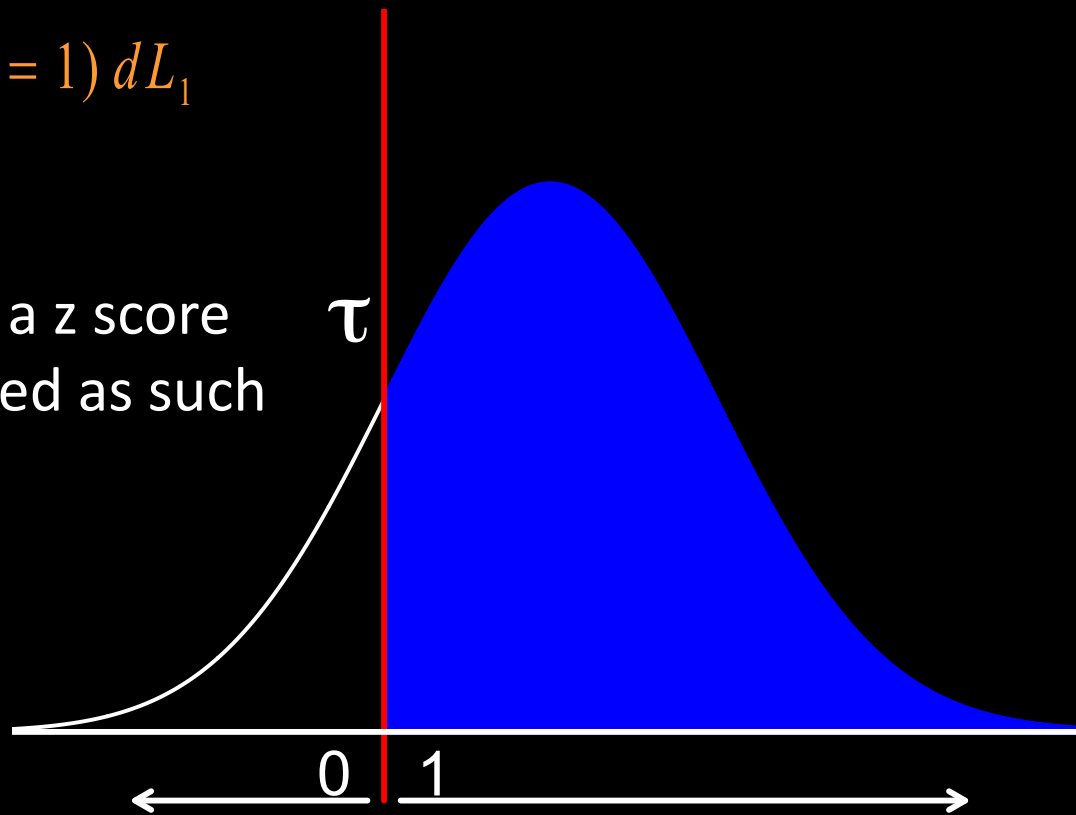
# Identifying Assumptions of Ordinal Associations

- **The assumptions are arbitrary**
  - The same model can be specified in different ways, but the parameters will estimate different things. The  $-2\ln L$  should be the same for models that are transformations of each other.
- **The assumptions are necessary.**
  - Because the latent dimension is only measured indirectly, by ordinal items, we have no direct information on its variance. The thresholds could expand or contract (think accordion) to completely compensate for a change in variance.

# Intuitive explanation of thresholds in the univariate normal distribution

$$\int_{z_T}^{\infty} \Phi(L_1; \mu = 0, \sigma^2 = 1) dL_1$$

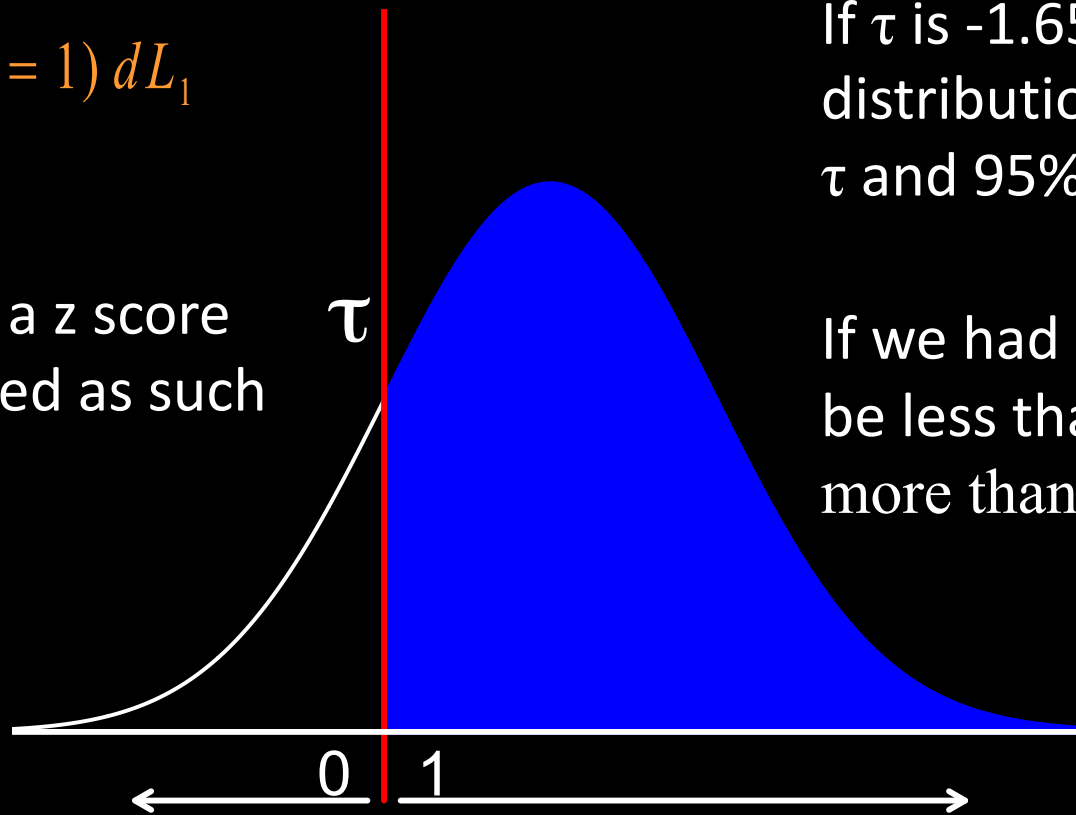
The threshold is just a z score and can be interpreted as such



# Intuitive explanation of thresholds in the univariate normal distribution

$$\int_{z_T}^{\infty} \Phi(L_1; \mu = 0, \sigma^2 = 1) dL_1$$

The threshold is just a z score and can be interpreted as such



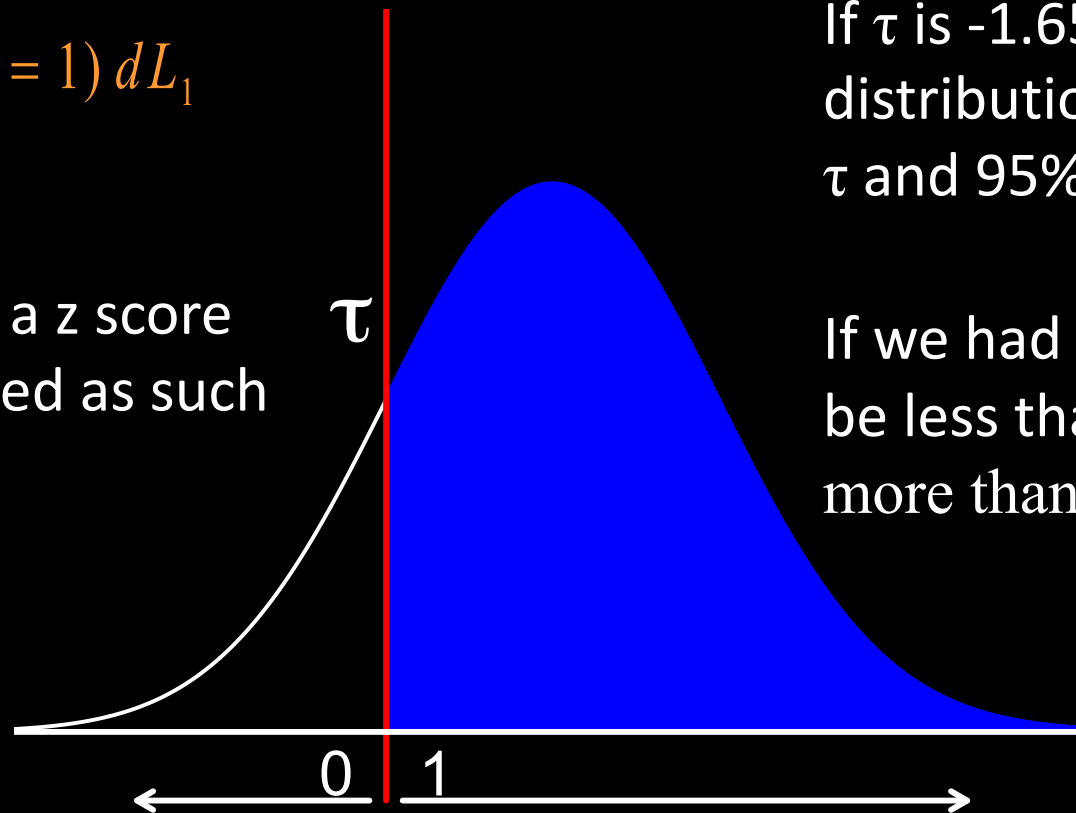
If  $\tau$  is -1.65 then 5% of the distribution will be to the left of  $\tau$  and 95% will be to the right

If we had 1000 people, 50 would be less than  $\tau$  and 950 would be more than  $\tau$

# Intuitive explanation of thresholds in the univariate normal distribution

$$\int_{z_T}^{\infty} \Phi(L_1; \mu = 0, \sigma^2 = 1) dL_1$$

The threshold is just a z score and can be interpreted as such



If  $\tau$  is -1.65 then 5% of the distribution will be to the left of  $\tau$  and 95% will be to the right

If we had 1000 people, 50 would be less than  $\tau$  and 950 would be more than  $\tau$

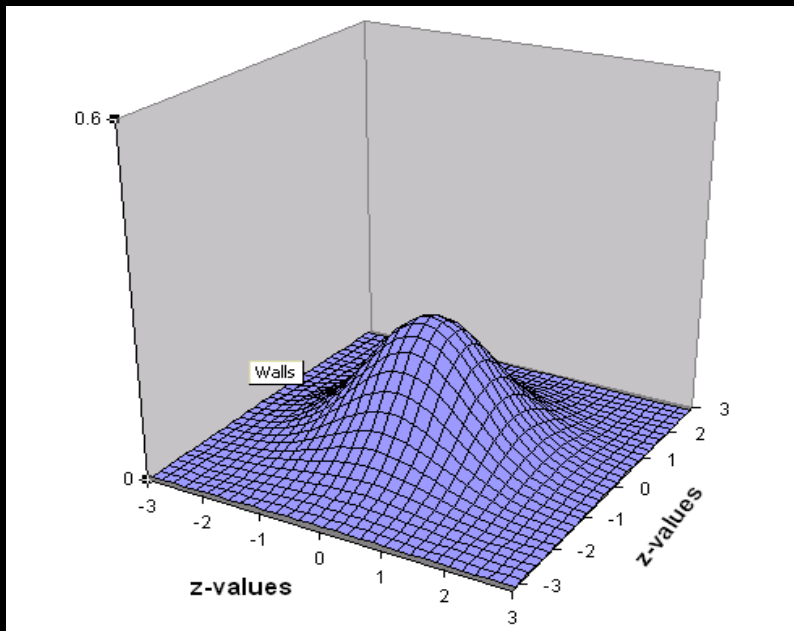
If  $\tau$  is 1.96 then 97.5% of the distribution will be to the left of  $\tau$  and .025% will be to the right

If we had 1000 people, 975 would be less than  $\tau$  and 25 would be more than  $\tau$

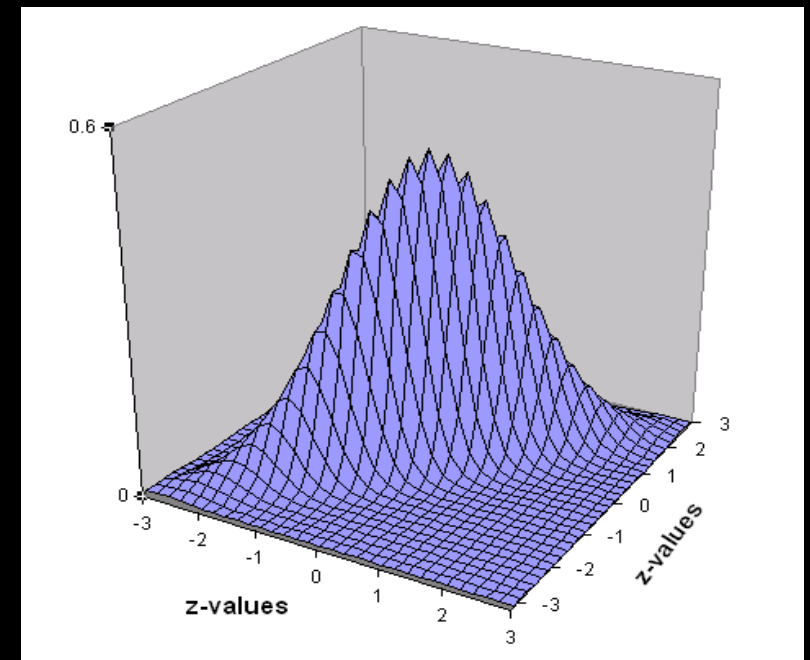
# Bivariate normal distribution

- Two variables:
  - Two Variances and Two Means
- What does Covariance do?

**$r = .00$**



**$r = .90$**



# Two binary traits (e.g., data from twins)

Contingency Table with 4 observed cells:

Cell a: pairs concordant for unaffected

Cells b&c: pairs discordant for the disorder

Cell d: pairs concordant for affected

Twin1 Twin2	0	1
0	a	b
1	c	d

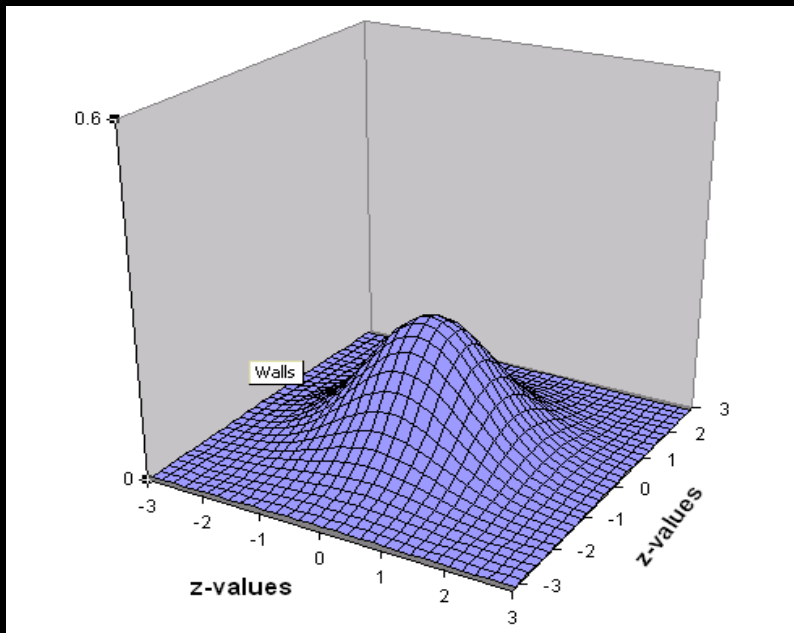
**0 = unaffected**  
**1 = affected**



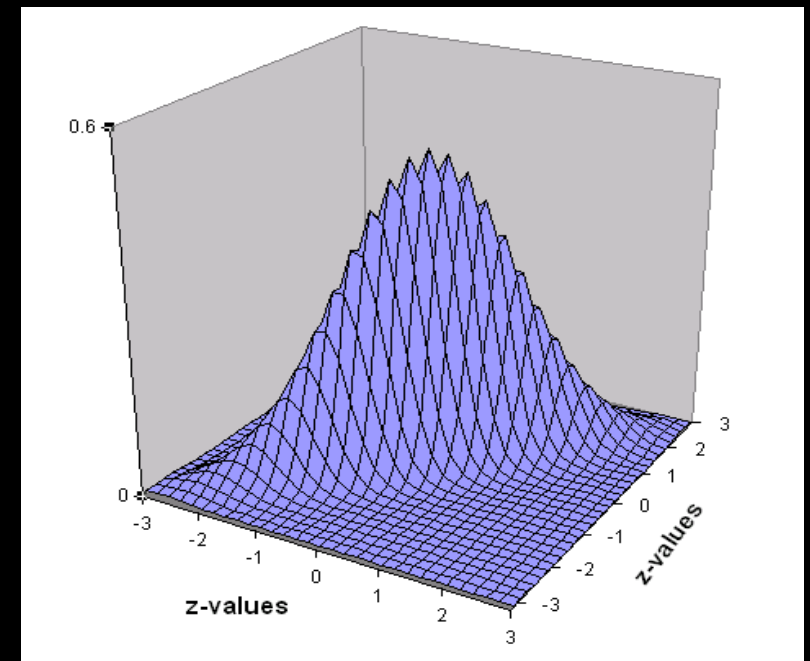
# Joint Liability Threshold Model for twin pairs

- Pairs are assumed to follow a **bivariate normal distribution**, where both traits have a mean of 0 and standard deviation of 1, and the **correlation** between them is what we want to know.
- The **shape** of a bivariate normal distribution is determined by the **correlation** between the traits

$r = .00$

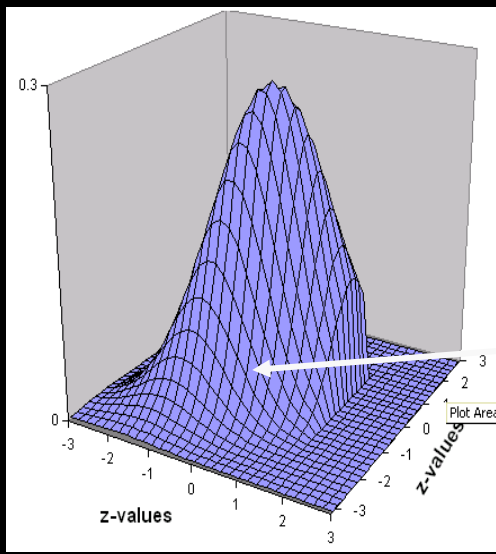


$r = .90$

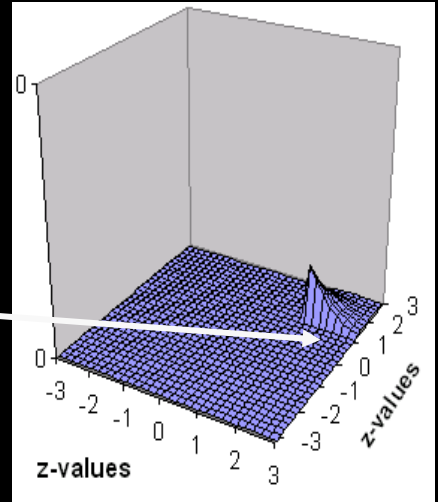
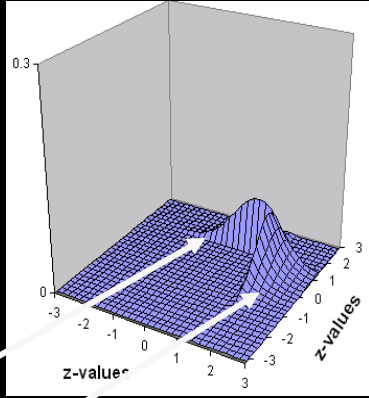


- The observed cell proportions relate to the proportions of the Bivariate Normal Distribution with a certain correlation between the latent variables ( $y_1$  and  $y_2$ ), each cut at a certain threshold

In other words, the joint probability of a certain response combination is the volume under the BND surface bounded by appropriate thresholds on each liability



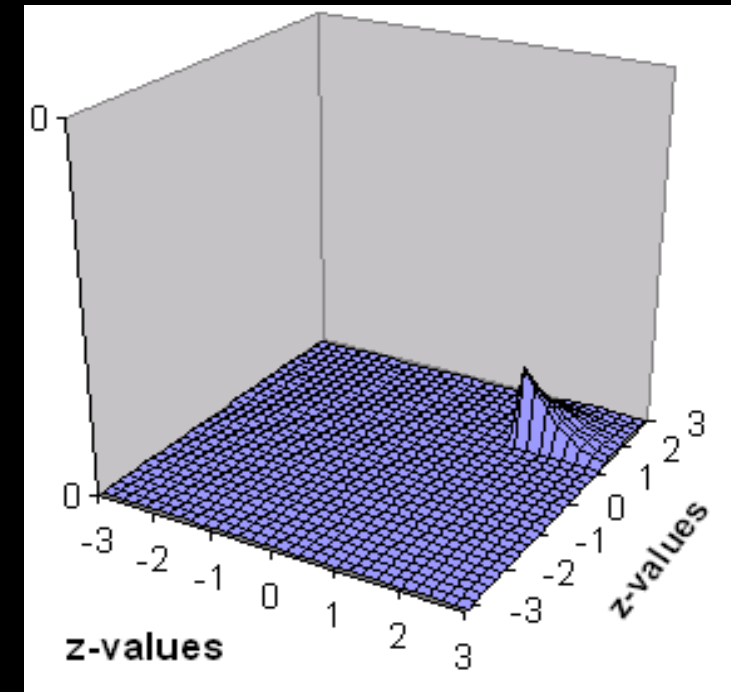
$y_1$	$y_2$	0	1
0	00	01	
1	10	11	



To calculate the cell proportions we rely on **Numerical Integration** of the Bivariate Normal Distribution over the two liabilities

e.g. the probability that both twins are above  $T_c$  :

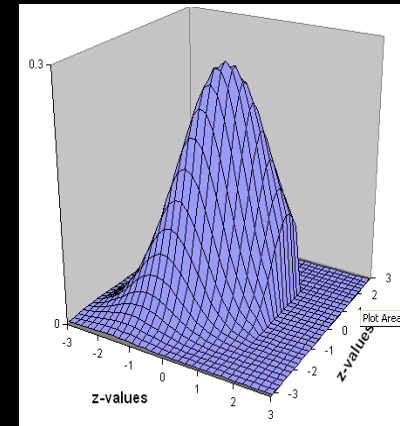
$$\int_{T_{c1}}^{\infty} \int_{T_{c2}}^{\infty} \Phi(y_1, y_2; \mu = 0, \Sigma) dy_1 dy_2$$



**$\Phi$**  is the bivariate normal probability density function,  
 **$y_1$**  and  **$y_2$**  are the liabilities of twin1 and twin2,  
with means of **0**, and  **$\Sigma$**  the correlation between the two liabilities  
 **$T_{c1}$**  is threshold (z-value) on  **$y_1$** ,  **$T_{c2}$**  is threshold (z-value) on  **$y_2$**

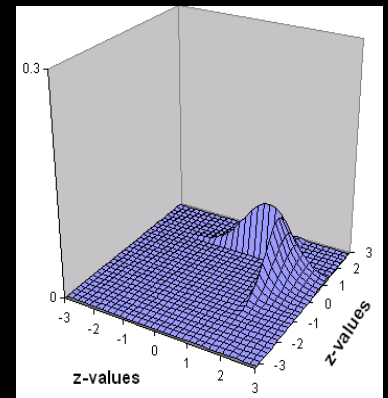
# Expected cell proportions

$$\int_{-\infty}^{T_{c1}} \int_{-\infty}^{T_{c2}} \Phi(y_1, y_2; \mu = 0, \Sigma) dy_1 dy_2$$



$$\int_{-\infty}^{T_{c1}} \int_{T_{c2}}^{\infty} \Phi(y_1, y_2; \mu = 0, \Sigma) dy_1 dy_2$$

$$\int_{T_{c1}}^{\infty} \int_{-\infty}^{T_{c2}} \Phi(y_1, y_2; \mu = 0, \Sigma) dy_1 dy_2$$



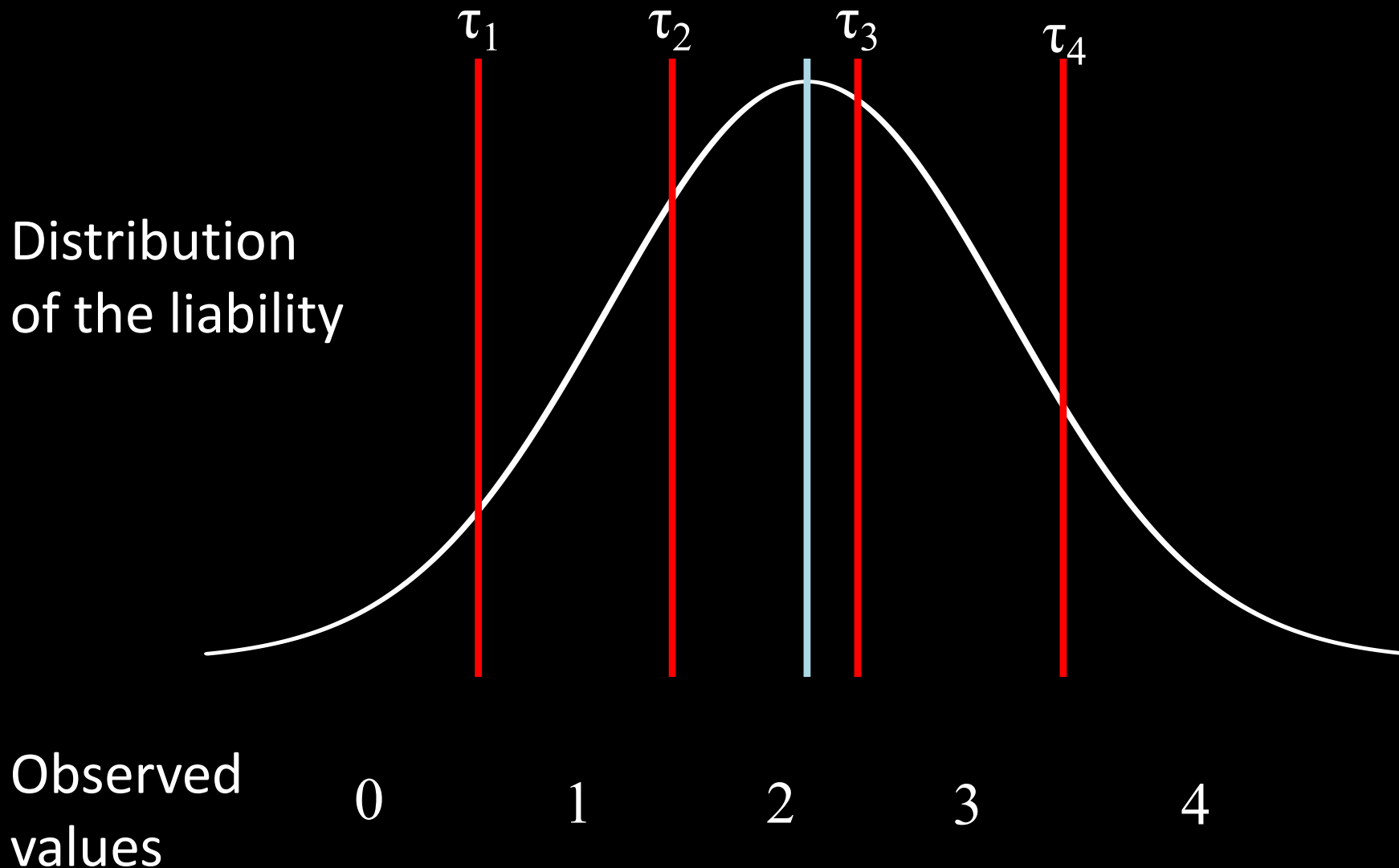
# Estimation of Correlations and Thresholds

- Since the Bivariate Normal distribution is a known mathematical distribution, for each correlation ( $\rho$ ) and any set of thresholds on the liabilities we know what the expected proportions are in each cell.
- Therefore, observed cell proportions of our data will inform on the most likely correlation and threshold on each liability.

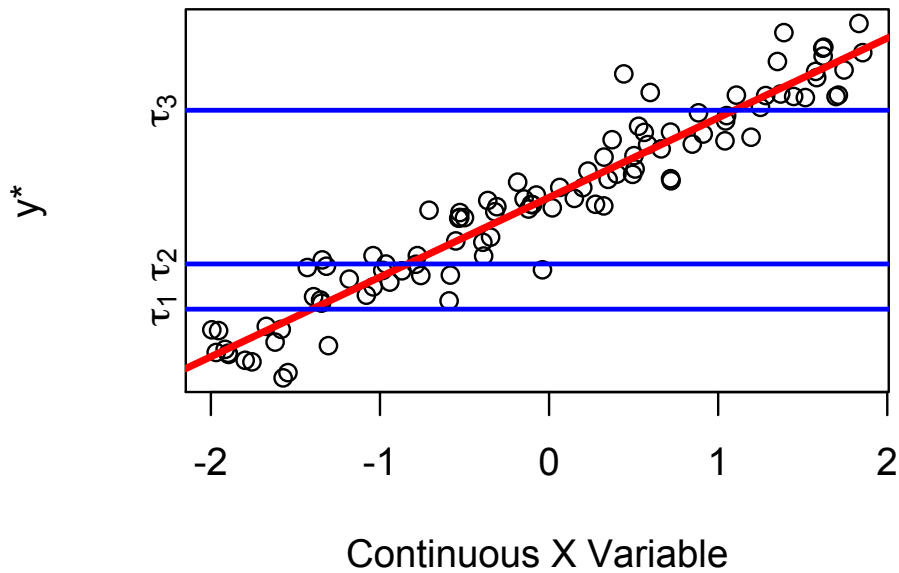
	y2	0	1
y1			
0		<b>.87</b>	<b>.05</b>
1		<b>.05</b>	<b>.03</b>

$$r = 0.60$$
$$T_{c1} = T_{c2} = 1.4 \text{ (z-value)}$$

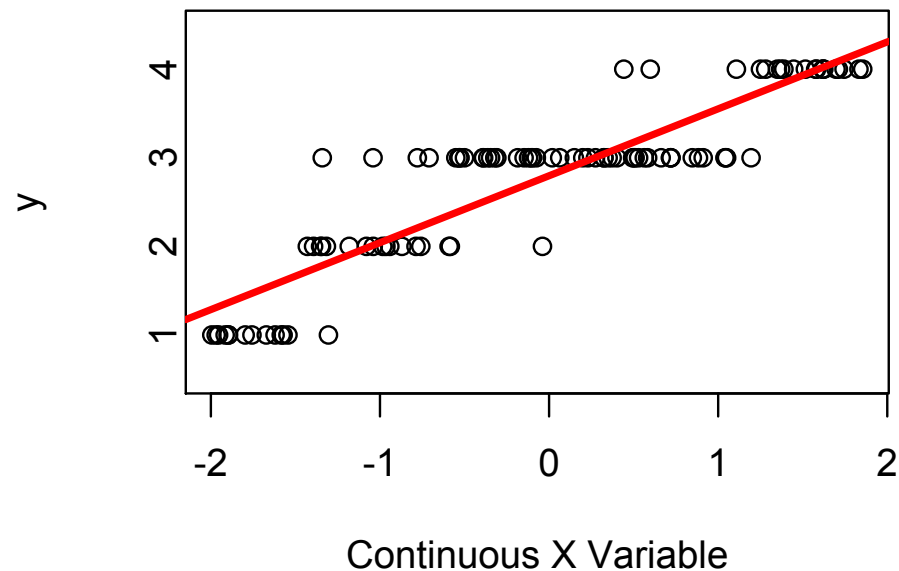
# Intuition behind the Multiple Threshold Liability model



Regression on the Latent Continuous Y Variable



Regression on the Observed Ordered Y Variable



Comparison between the regression of the latent  $y^*$  and the observed  $y$

It is important to keep in mind that the scale of the ordinal variable is arbitrary, and therefore it is virtually impossible to compare the slopes of the two graphs (even though they look pretty similar)

# Squeezing Interval Change From Ordinal Panel Data: Latent Growth Curves With Ordinal Outcomes

Paras D. Mehta  
University of Illinois at Chicago

Michael C. Neale  
Virginia Commonwealth University

Brian R. Flay  
University of Illinois at Chicago

What happens if we change the default assumptions?

## Mean Assumption

The intercept (mean) is 0

or

The threshold is 0 ( $\tau = 0$ )

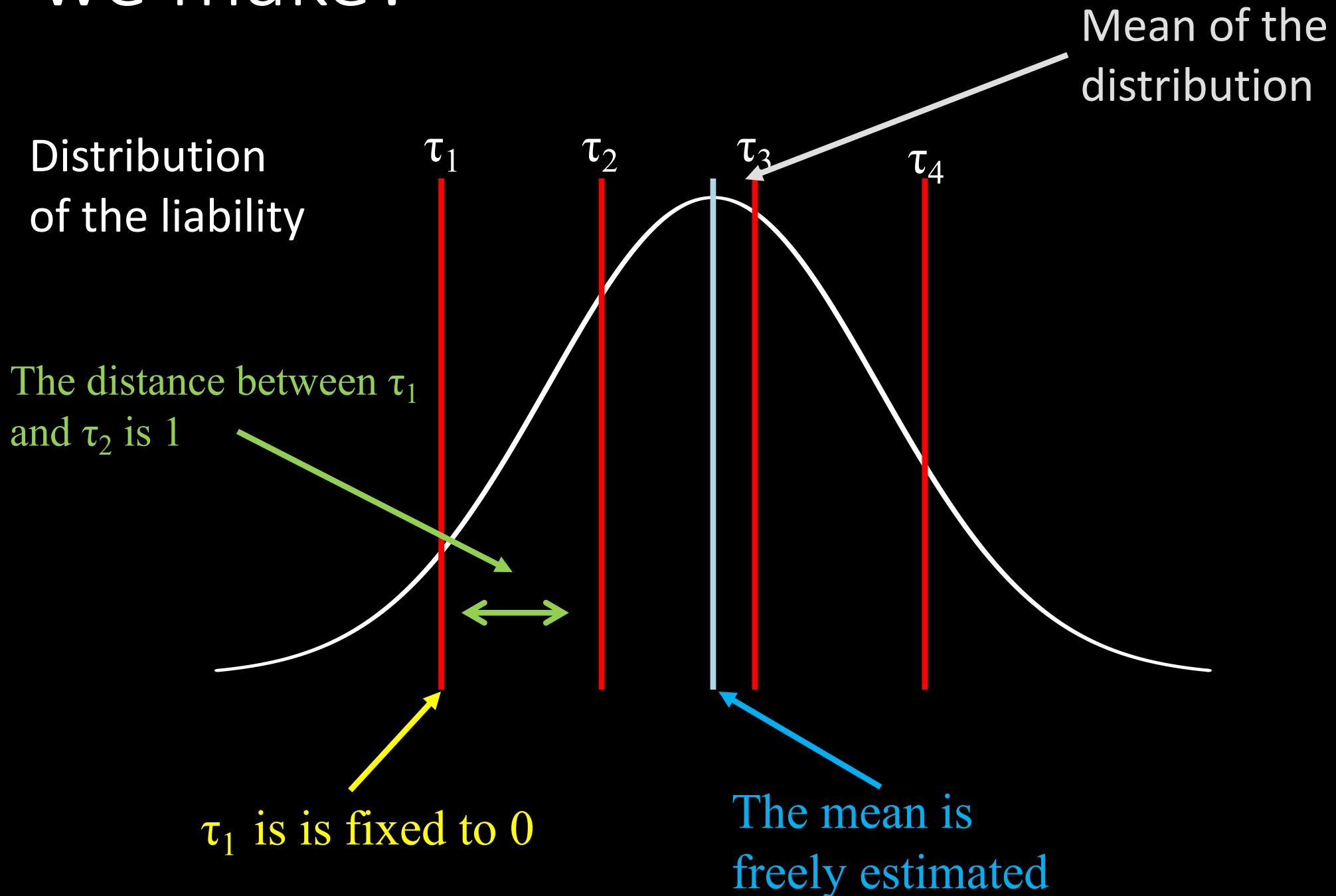
## Variance Assumption

$\text{Var}(\varepsilon | x) = 1$  in the normal-ogive model

Remember that we can make slightly different assumptions with equal model fit



# What alternative assumptions could we make?



# Squeezing Interval Change From Ordinal Panel Data: Latent Growth Curves With Ordinal Outcomes

Paras D. Mehta  
University of Illinois at Chicago

Michael C. Neale  
Virginia Commonwealth University

Brian R. Flay  
University of Illinois at Chicago

It is important to reiterate that the model fit is the same and that all the parameters can be transformed from one set of assumptions to another.

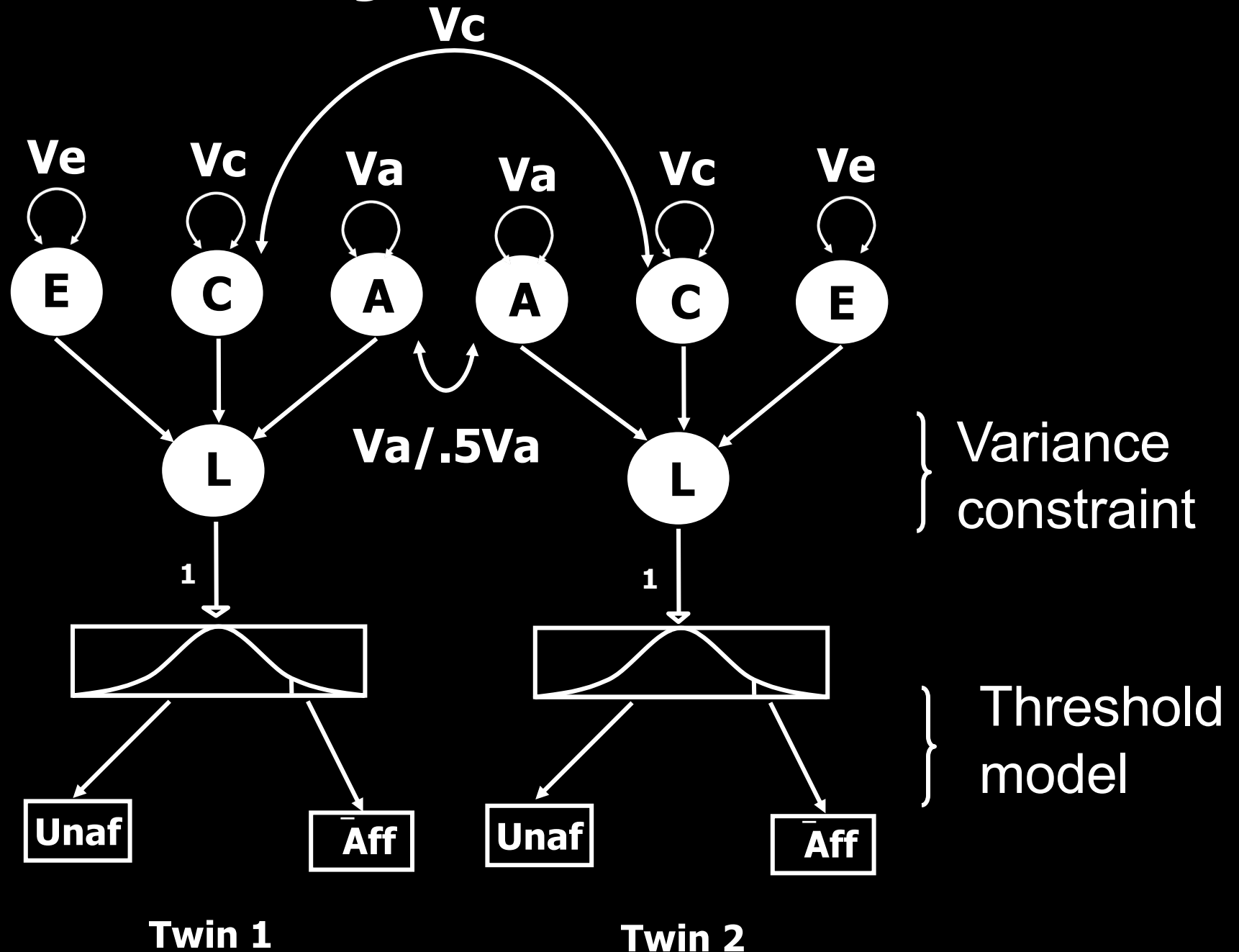
# Bivariate Ordinal Likelihood

- The likelihood for each observed ordinal response pattern is computed by the expected proportion in the corresponding cell of the bivariate normal distribution
- The maximum-likelihood equation for the whole sample is the sum of  $-2 \times \log$  of the likelihood of each row of data (e.g. twin pairs)
- This  $-2LL$  is minimized to obtain the maximum likelihood estimates of the correlation and thresholds
- Tetra-choric correlation if  $y_1$  and  $y_2$  reflect 2 categories (1 Threshold); Poly-choric when  $>2$  categories per liability

# Twin Models

- Estimate correlation in liabilities separately for MZ and DZ pairs from contingency table
- Variance decomposition (A, C, E) can be applied to the *liability* of the trait
- Correlations in liability are determined by path model
- Estimate of the heritability of the *liability*

# ACE Liability Model



# Summary

- OpenMx models ordinal data under a threshold model
- Assumptions about the (joint) distribution of the data (Standard Bivariate Normal)
- The relative proportions of observations in the cells of the Contingency Table are translated into proportions under the Multivariate Normal Distribution
- The most likely thresholds and correlations are estimated
- Genetic/Environmental variance components are estimated based on these correlations derived from MZ and DZ data
- See [Medland et al 2005](#) for GxE modelling of ordinal data

# Power issues

- Ordinal data / Liability Threshold Model: less power than analyses on continuous data

Neale, Eaves & Kendler 1994

(*alpha* = .05 should say .025)

- Solutions:
  1. Bigger samples
  2. Use more categories

Please do not categorize continuous variables

