# VCU
Virginia Commonwealth University

# Introduction to Linkage and Association for Quantitative Traits

## Michael C Neale
## Boulder Colorado Workshop March 2 2009

# Overview

- A brief history of SEM
- Regression
- Maximum likelihood estimation
- Models
  - Twin data
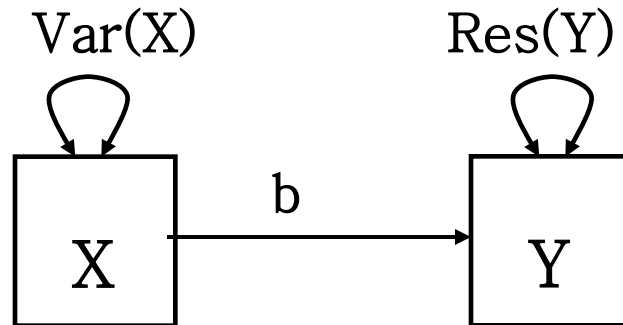  - Sib pair linkage analysis
  - Association analysis

VCU

# Origins of SEM

- Regression analysis
  - 'Reversion' Galton 1877: Biological phenomenon
  - Yule 1897 Pearson 1903: General Statistical Context
  - Initially Gaussian X and Y; Fisher 1922 Y|X

- Path Analysis
  - Sewall Wright 1918; 1921
  - Path Diagrams of regression and covariance relationships
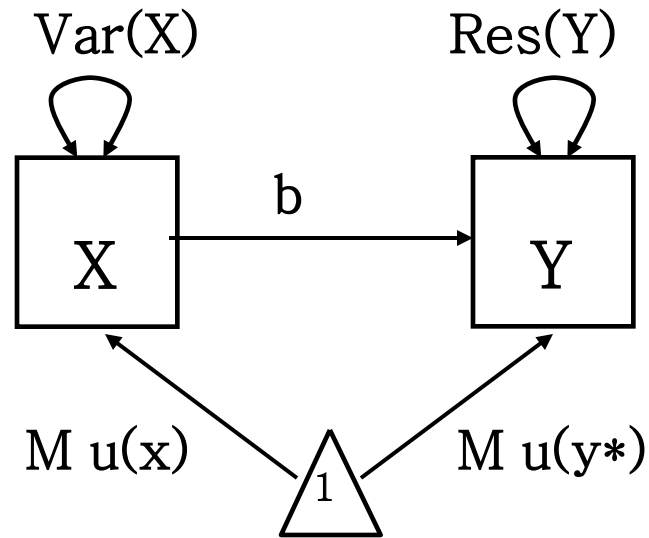
**VCU**

# Structural Equation Modeling Basics

- Two kinds of relationships
  - Linear *regression* X -> Y    single-headed
  - Unspecified *covariance* X<->Y   double-headed

- Four kinds of variable
  - Squares: observed variables
  - Circles: latent, not observed variables
  - Triangles: constant (zero variance) for specifying means
  - Diamonds: observed variables used as moderators (on paths)

**VCU**

# Linear Regression Covariance SEM



Var(X)    Res(Y)

X ——b——> Y

Models *covariances* only
Of historical interest

VCU

# Linear Regression SEM with means



Var(X)　　　　Res(Y)

X　　b　　Y

M u(x)　　1　　M u(y*)
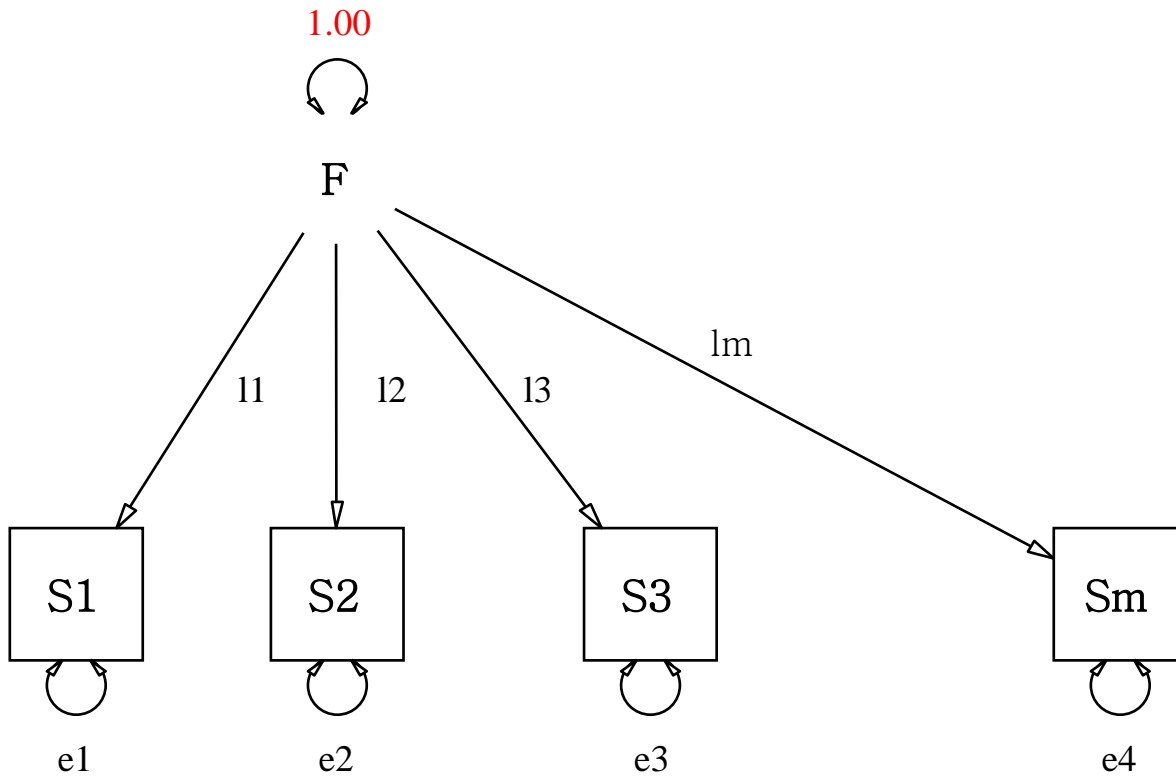
Models Means and Covariances

VCU

# Linear Regression SEM: Individual-level
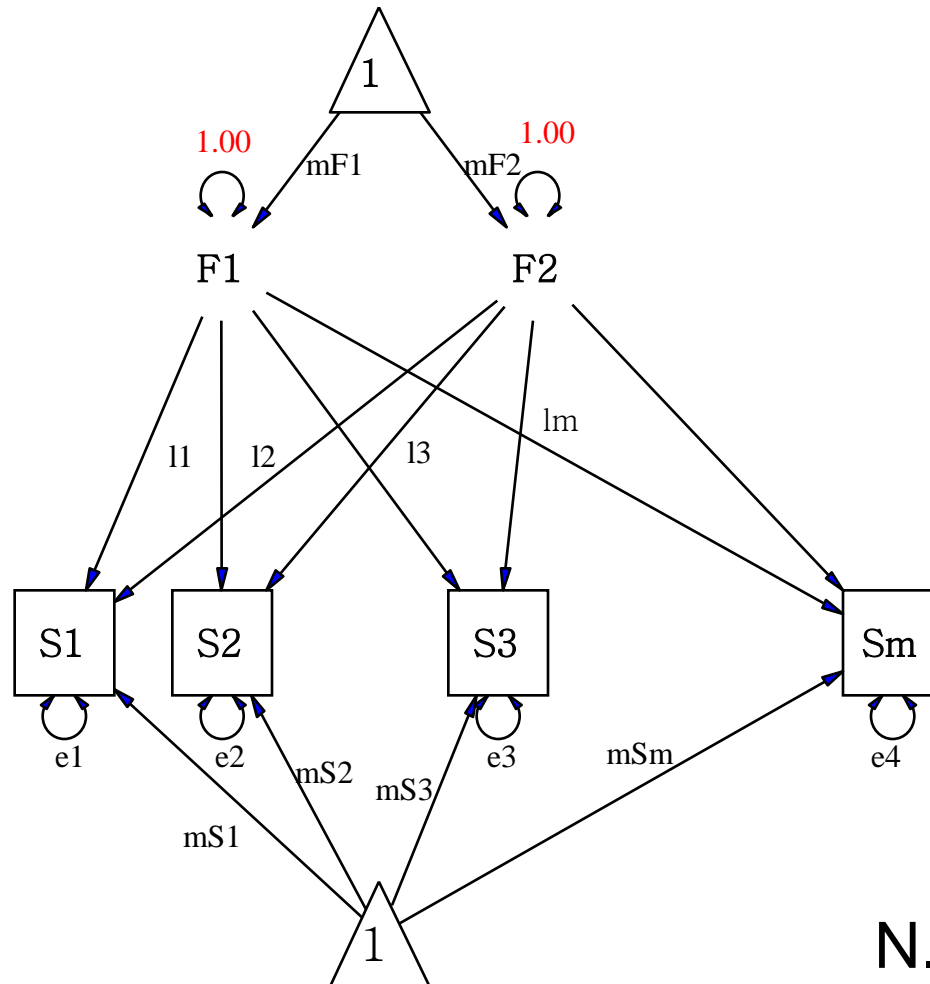
$Y_i = a + bX_i$



Models Mean and Covariance of Y *only*
Must have raw (individual level) data
$X_i$ is a *definition* variable
Mean of Y different for every observation

# Single Factor Covariance Model
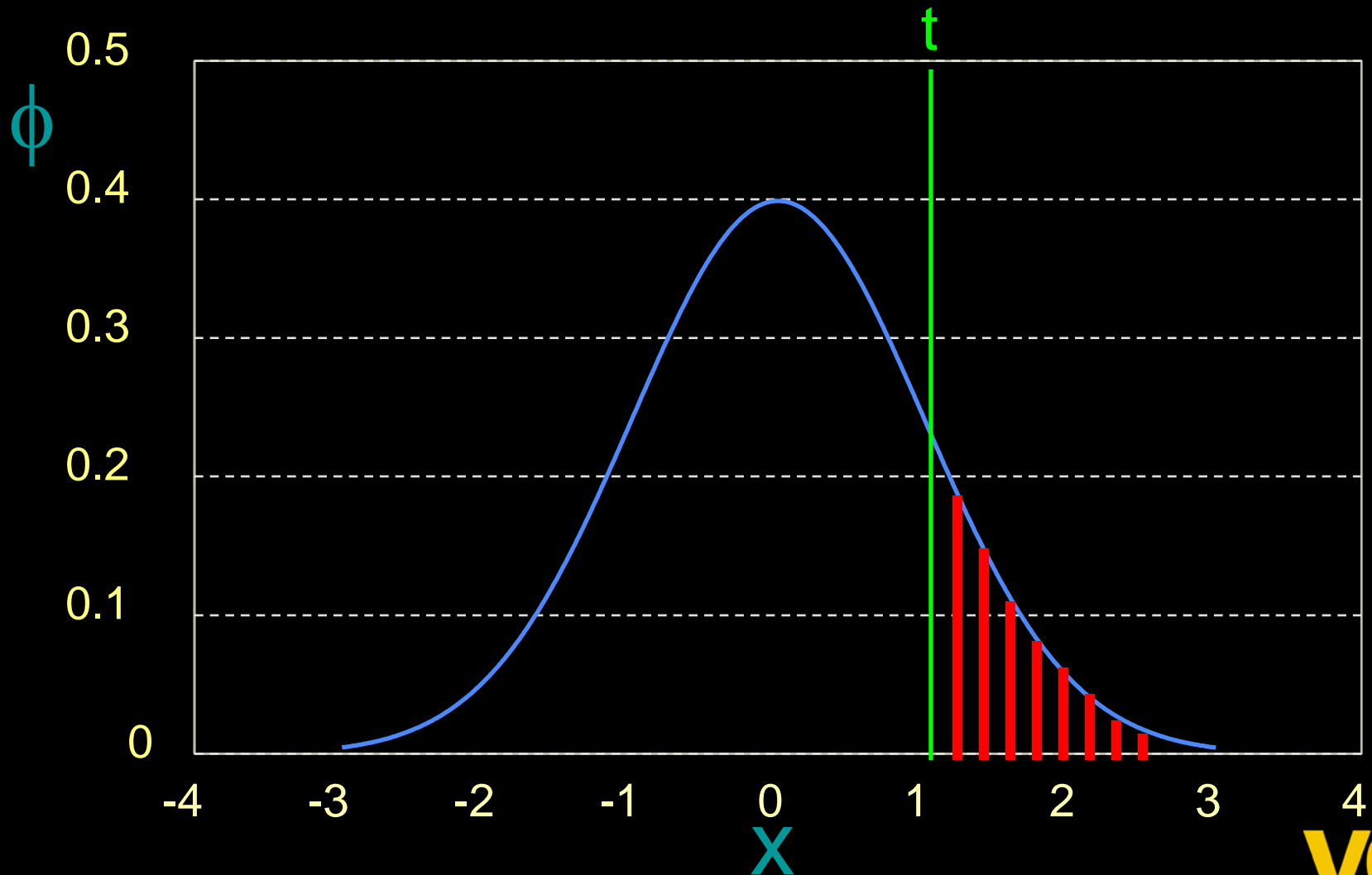
# Two Factor Model with Covs & Means



N.B. Not identified

# Factor model essentials

- In SEM the factors are typically assumed to be normally distributed

- May have more than one latent factor

- The error variance is typically assumed to be normal as well

- May be applied to binary or ordinal data
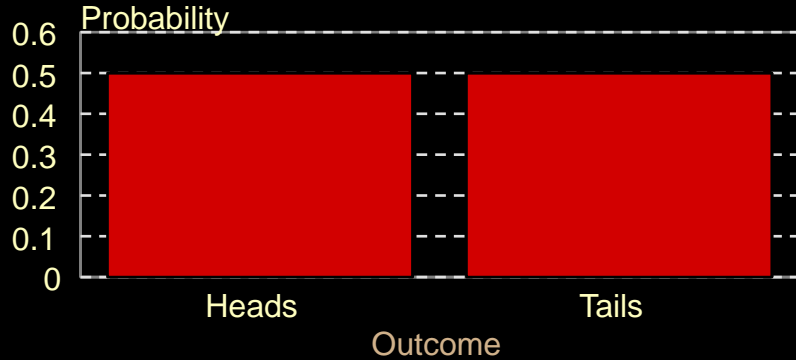  - Threshold model

**VCU**

# Measuring Variation

- Distribution
  - Population
  - Sample
  - Observed measures

- Probability density function 'pdf'
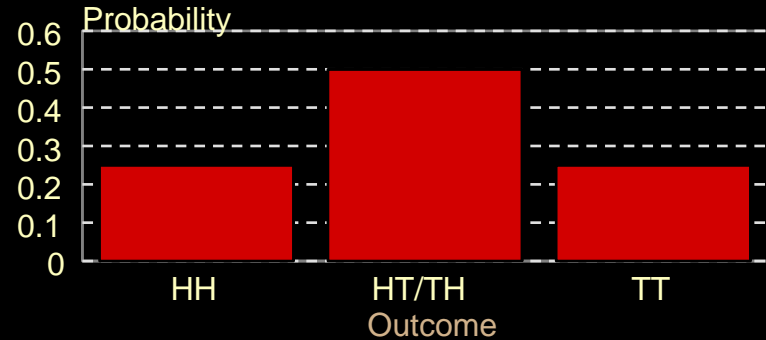  - Smoothed out histogram
  - f(x) >= 0  for all x

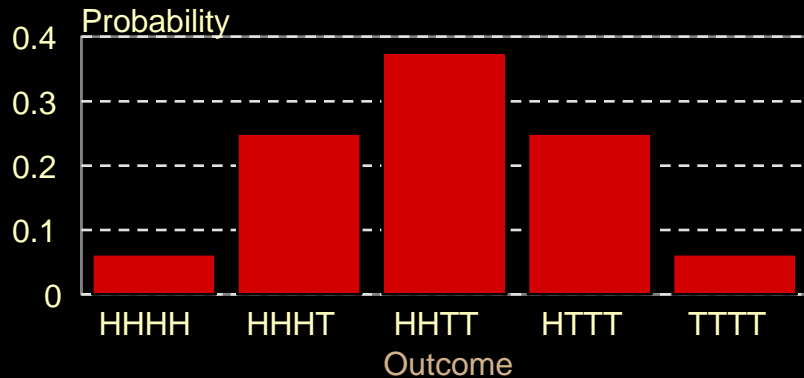$$\int_{-\infty}^{\infty} f(x)\, dx = 1.$$

**VCU**

# Flipping Coins

1 coin: 2 outcomes

Probability

| | |
|---|---|
| 0.6 | |
| 0.5 | |
| 0.4 | |
| 0.3 | |
| 0.2 | |
| 0.1 | |
| 0 | |

Heads          Tails

Outcome

2 coins: 3 outcomes

Probability

| | |
|---|---|
| 0.6 | |
| 0.5 | |
| 0.4 | |
| 0.3 | |
| 0.2 | |
| 0.1 | |
| 0 | |

HH          HT/TH          TT

Outcome

4 coins: 5 outcomes

Probability

| | |
|---|---|
| 0.4 | |
| 0.3 | |
| 0.2 | |
| 0.1 | |
| 0 | |

HHHH    HHHT    HHTT    HTTT    TTTT

Outcome

8 coins: 9 outcomes

Probability

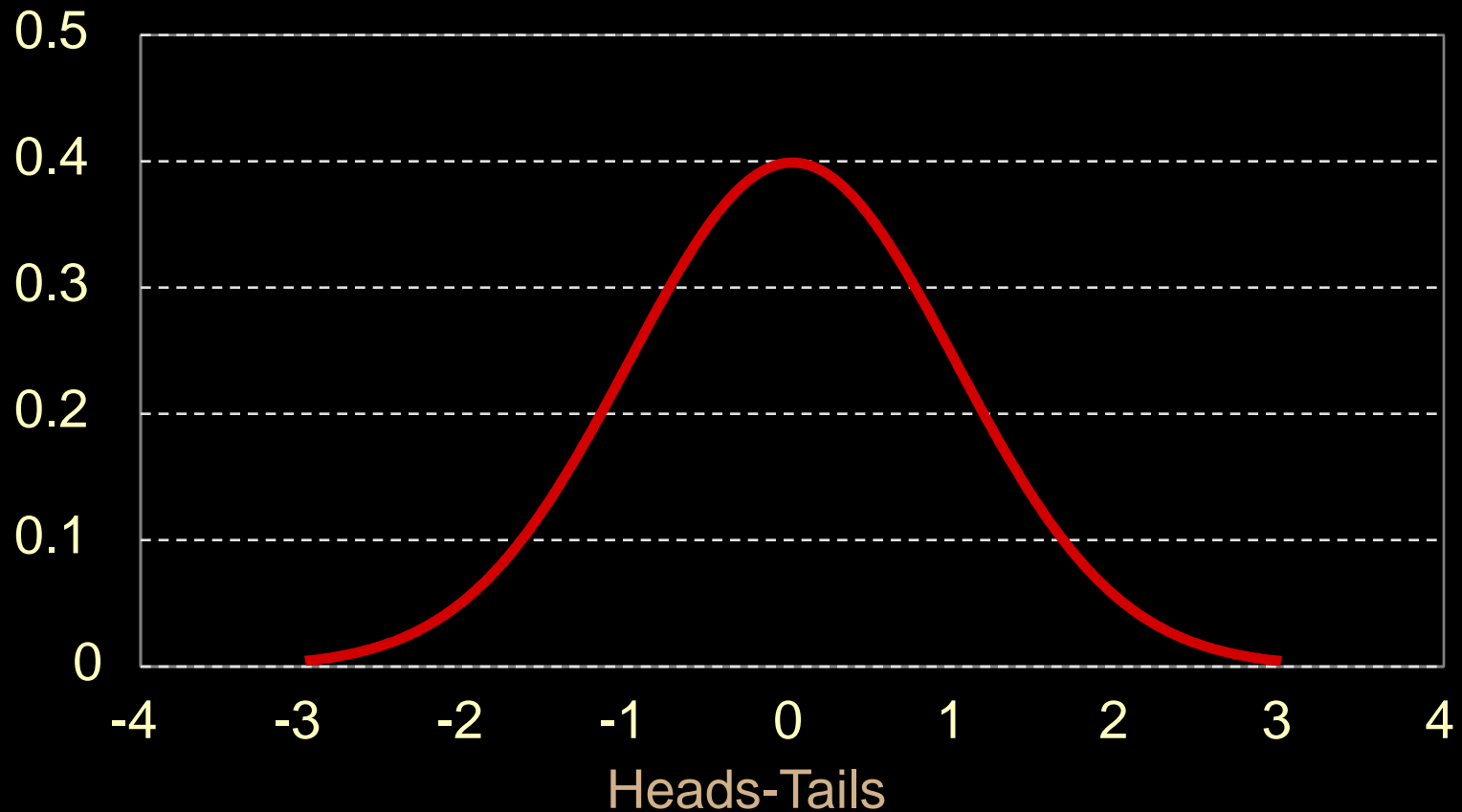| | |
|---|---|
| 0.3 | |
| 0.25 | |
| 0.2 | |
| 0.15 | |
| 0.1 | |
| 0.05 | |
| 0 | |

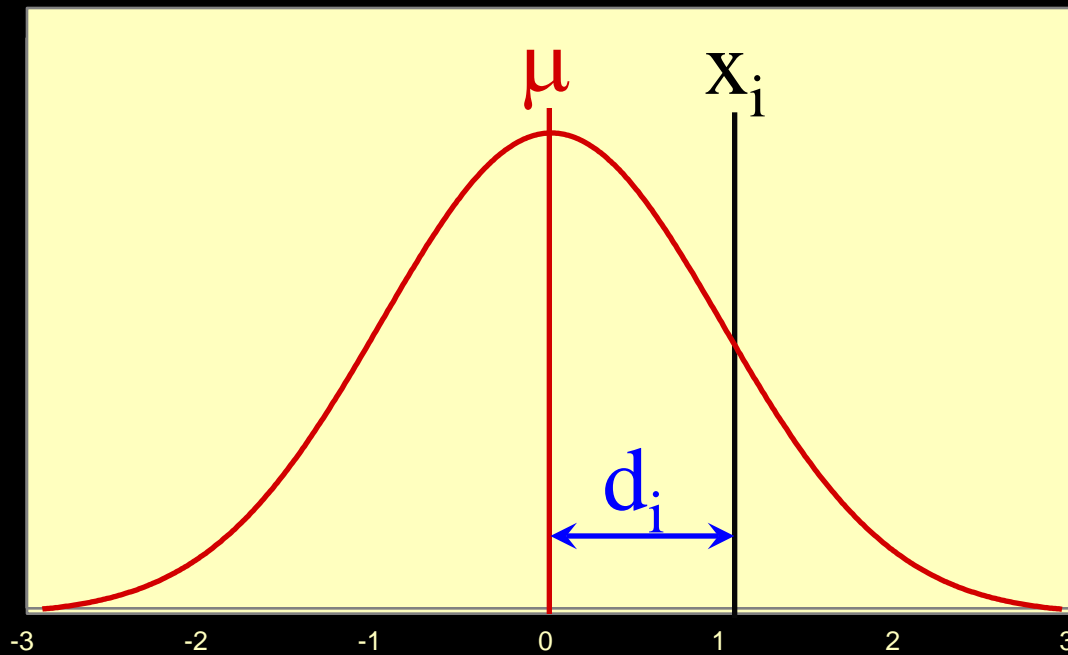Outcome

VCU

# Bank of China Coin Toss

Infinite outcomes



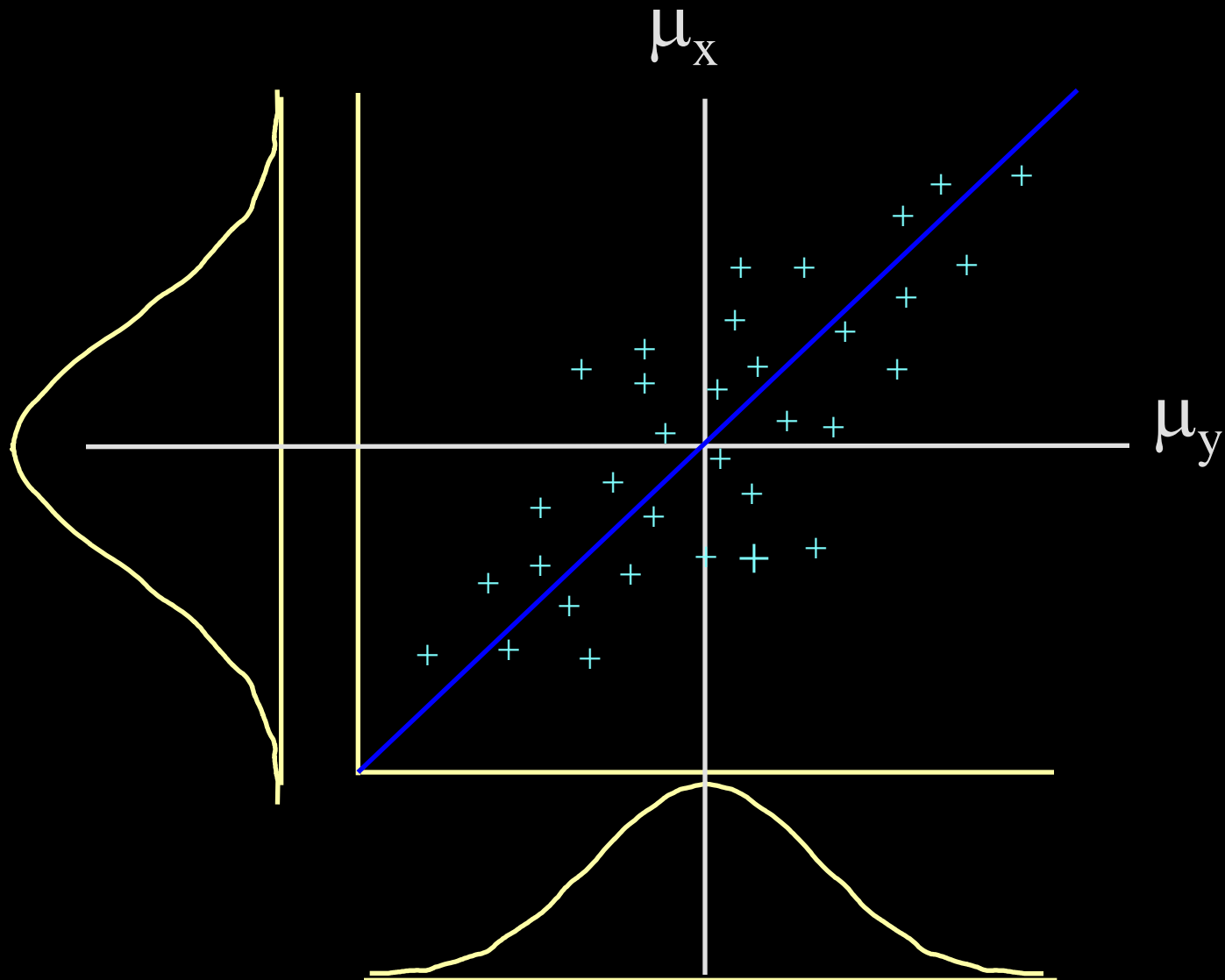De Moivre 1733 Gauss 1827

**VCU**

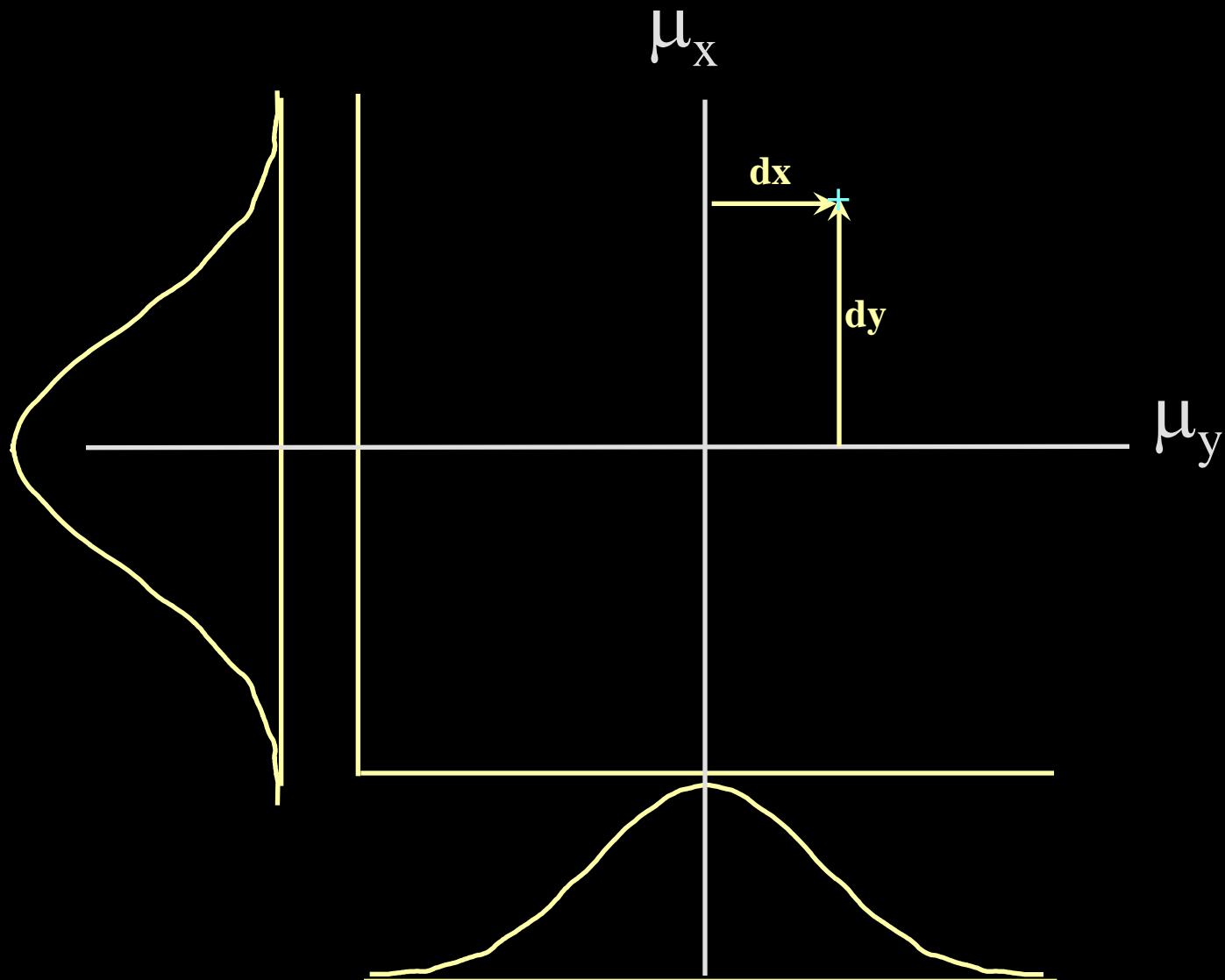# Variance: Average squared deviation

## Normal distribution



$$\text{Variance} = \Sigma\, d_i^2/N$$

VCU

# Deviations in two dimensions

# Deviations in two dimensions: dx x dy

# Covariance

- Measure of association between two variables

- Closely related to variance

- Useful to partition variance
  - "Analysis of Variance" term coined by Fisher

**VCU**

# Variance covariance matrix

### Univariate Twin/Sib Data

$$
\begin{bmatrix}
\text{Var(Twin1)} & \text{Cov(Twin1,Twin2)} \\
\text{Cov(Twin2,Twin1)} & \text{Var(Twin2)}
\end{bmatrix}
$$

Suitable for modeling when no missing data
Good conceptual perspective

VCU

# Maximum Likelihood Estimates: Nice Properties

1.  **Asymptotically unbiased**

    - Large sample estimate of p -> population value

2.  **Minimum variance "Efficient"**

    - Smallest variance of all estimates with property 1

3.  **Functionally invariant**

    - If g(a) is one-to-one function of parameter a

    - and MLE (a) = a*

    - then MLE g(a) = g(a*)

- See http://wikipedia.org

**VCU**

# Full Information Maximum Likelihood (FIML)

Calculate height of curve for each raw data vector

- Univariate - height of normal pdf
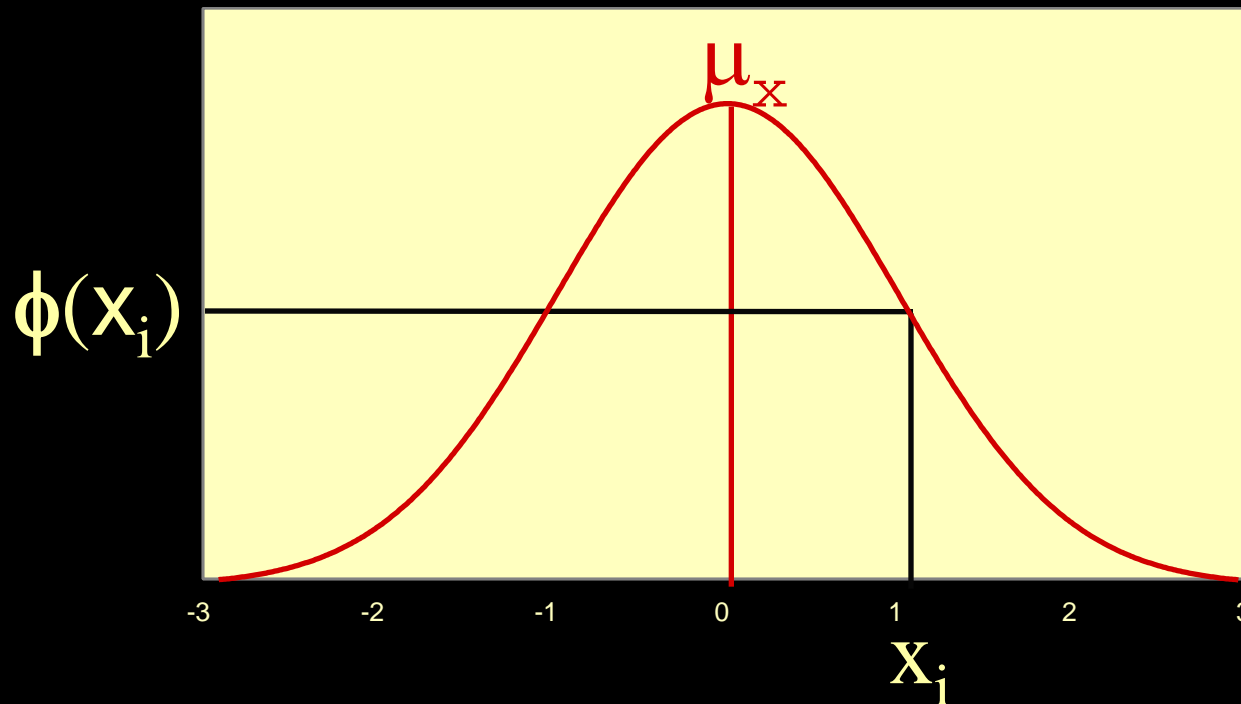
  - $\phi(x) =$

  - $(2\pi\sigma^2)^{-.5} e^{-.5((x_i - \mu)^2)/\sigma^2}$

- Multivariate - height of multinormal pdf

  - $|2\pi\Sigma|^{-n/2} e^{-.5((\mathbf{x}_i - \mu)\Sigma^{-1}(\mathbf{x}_i - \mu)')}$

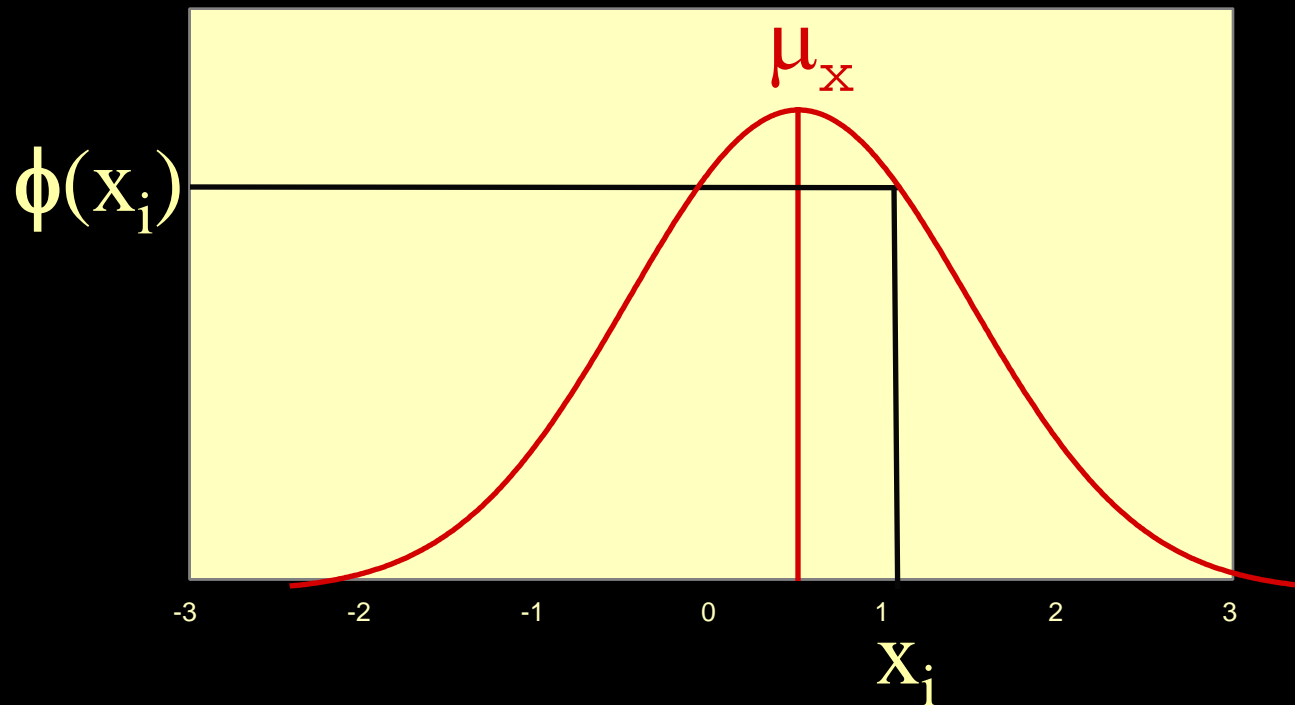# Height of normal curve: $\mu_x = 0$

## Probability density function

$\phi(x_i)$ is the likelihood of data point $x_i$ for particular mean & variance estimates

**VCU**

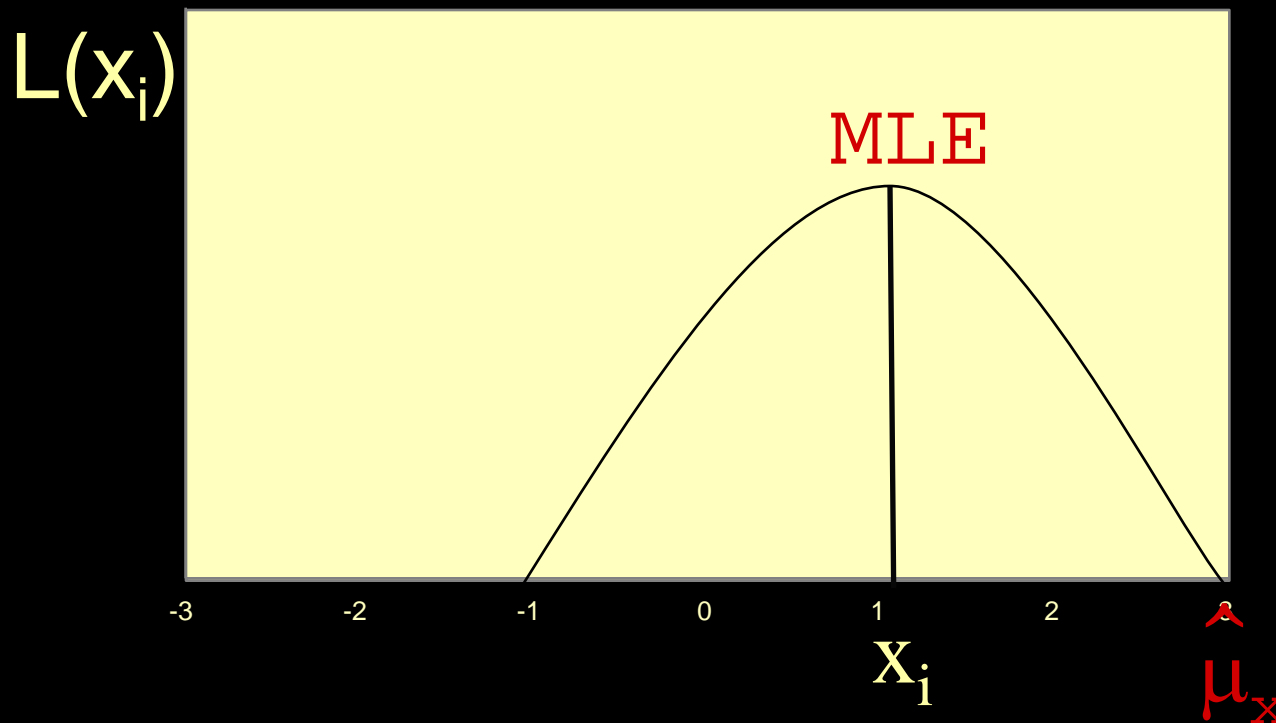# Height of normal curve at xi: $\mu_x$ = .5

Function of *mean*



Likelihood of data point $x_i$ *increases* as $\mu_x$ approaches $x_i$

**VCU**

# Likelihood of $x_i$ as a function of $\mu$

Likelihood function

$L(x_i)$

MLE

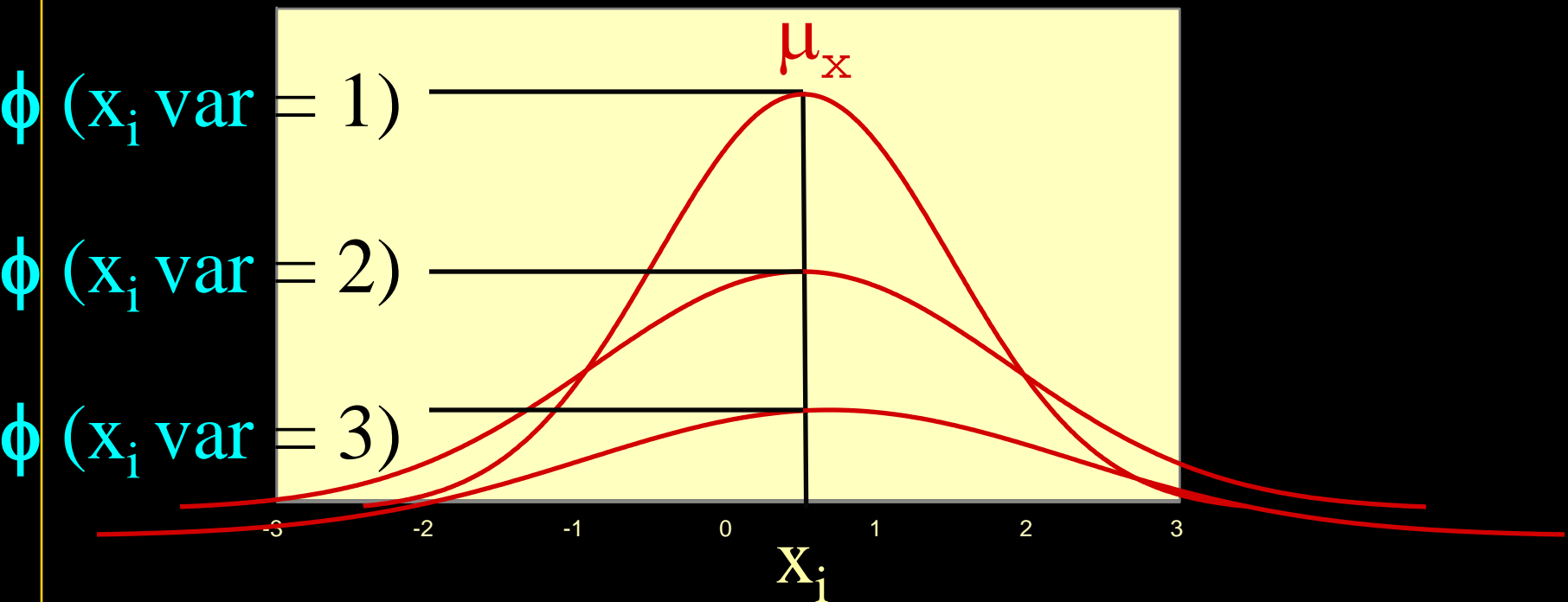-3    -2    -1    0    1    2    $\hat{}$

$x_i$

$\hat{\mu}_x$

$L(x_i)$ is the likelihood of data point $x_i$ for particular mean & variance estimates

VCU

# Height of normal curve at x1

Function of *variance*

$\phi$ (x$_i$ var = 1)

$\phi$ (x$_i$ var = 2)

$\phi$ (x$_i$ var = 3)

$\mu_x$

x$_i$

Likelihood of data point x$_i$ *changes* as variance of distribution changes

**VCU**

# Height of normal curve at x1 and x2



$\phi (x_1 \, \text{var} = 1)$

$\phi (x_1 \, \text{var} = 2)$

$\mu_x$

$\phi (x_2 \, \text{var} = 2)$

$\phi (x_2 \, \text{var} = 1)$

$x_1$

$x_2$

$x_1$ has higher likelihood with var=1 whereas $x_2$ has higher likelihood with var=2

VCU

# Height of bivariate normal density function

## Likelihood varies as $f(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$
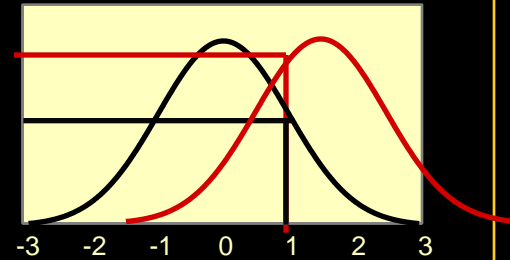
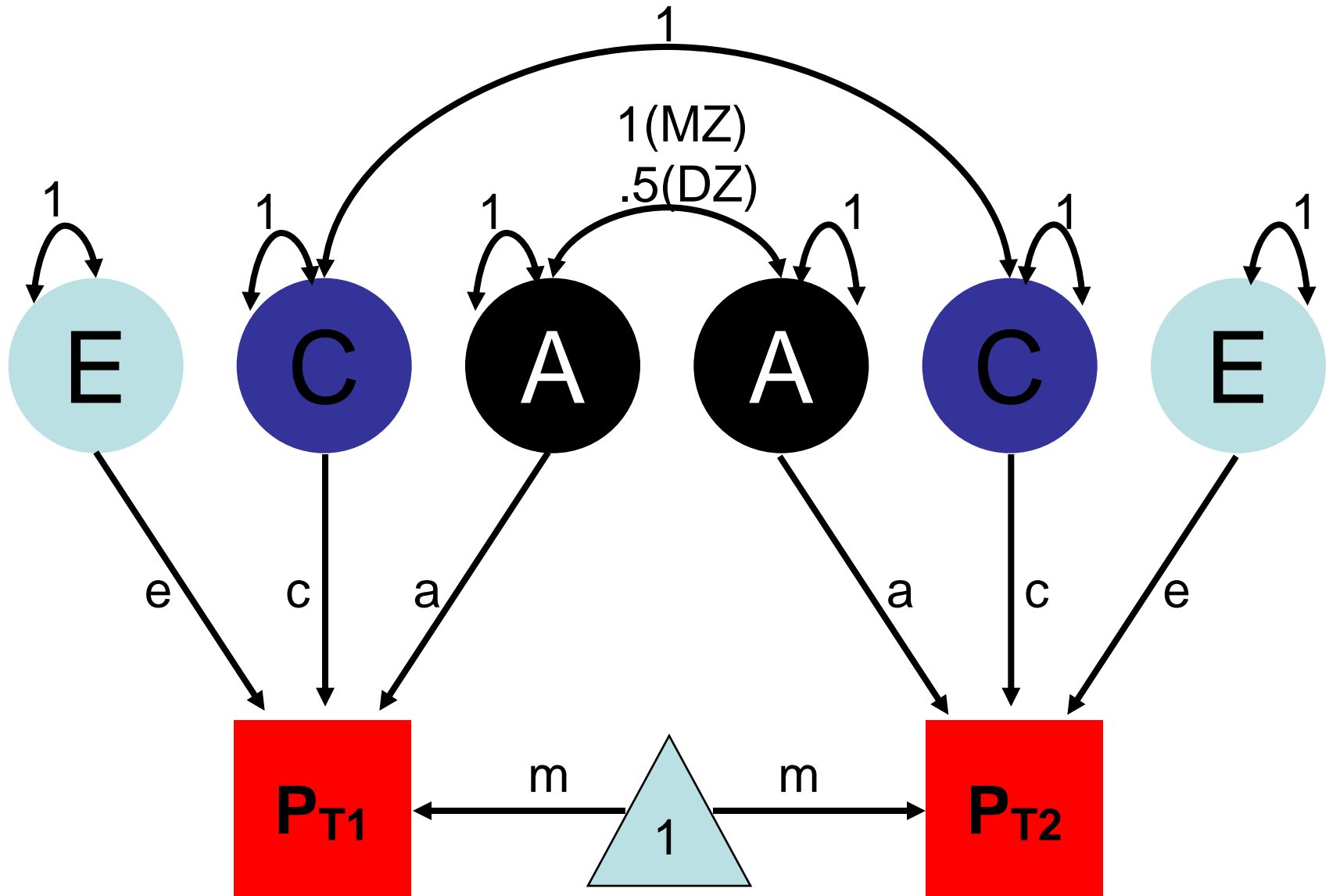# Likelihood of Independent Observations

- Chance of getting two heads
- $L(x_1 \ldots x_n) = \text{Product}(L(x_1), L(x_2), \ldots L(x_n))$
- $L(x_i)$ typically $< 1$

- Avoid vanishing $L(x_1 \ldots x_n)$
- Computationally convenient log-likelihood
- $\ln(a * b) = \ln(a) + \ln(b)$

- Minimization more manageable than maximization
  - Minimize $-2 \ln(L)$

**VCU**

# Likelihood Ratio Tests

- Comparison of likelihoods
- Consider *ratio*  L(data,model 1) / L(data, model 2)
- ln(a/b) = ln(a) - ln(b)
- Log-likelihood lnL(data, model 1) - ln L(data, model 2)

- Useful asymptotic feature when model 2 is a submodel of model 1
  $-2$ (lnL(data, model 1) - lnL(data, model 2)) $\sim \chi^2$
  df = # parameters of model 1 - # parameters of model 2

- BEWARE of gotchas!
  - Estimates of $a^2 q^2$ etc. have implicit bound of zero
  - Distributed as 50:50 mixture of 0 and $\chi_1^2$

**VCU**

# Linkage vs Association

## Linkage

1. Family-based

2. Matching/ethnicity generally unimportant

3. Few markers for genome coverage (300-400 STRs)

4. Can be weak design

5. Good for initial detection; poor for fine-mapping

6. Powerful for rare variants

## Association

1. Families or unrelated individuals

2. Matching/ethnicity crucial

3. Many markers req for genome coverage ($10^5 - 10^6$ SNPs)

4. Powerful design

5. Ok for initial detection; good for fine-mapping

6. Powerful for common variants; rare variants generally impossible

**VCU**

# Identity by Descent (IBD)

Number of alleles shared IBD at a locus,
parents AB and CD: Three subgroups of sibpairs

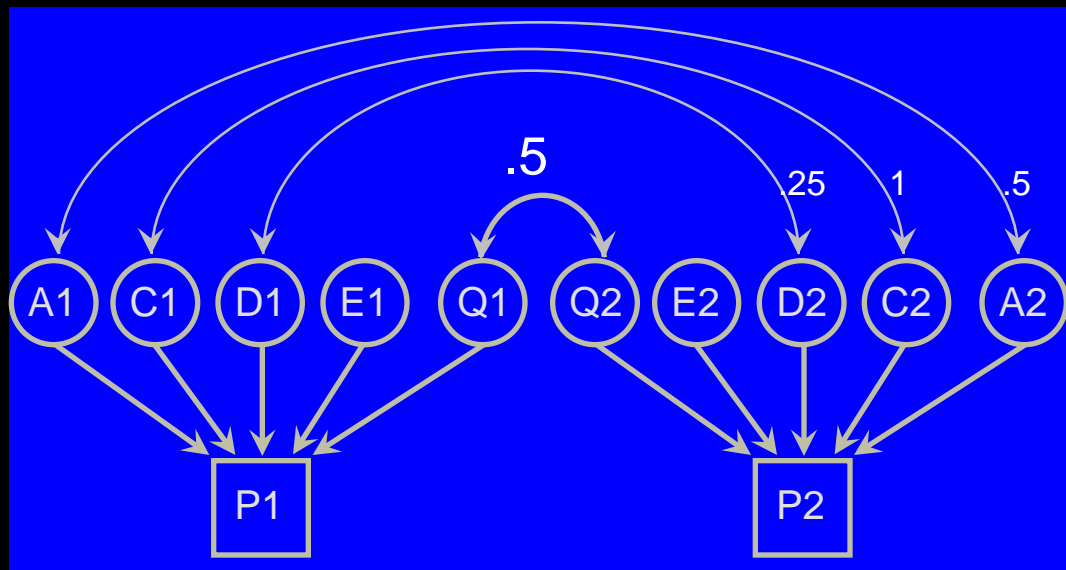|      | AC | AD | BC | BD |
|------|----|----|----|----|
| AC   | 2  | 1  | 1  | 0  |
| AD   | 1  | 2  | 0  | 1  |
| BC   | 1  | 0  | 2  | 1  |
| BD   | 0  | 1  | 1  | 2  |

vCU

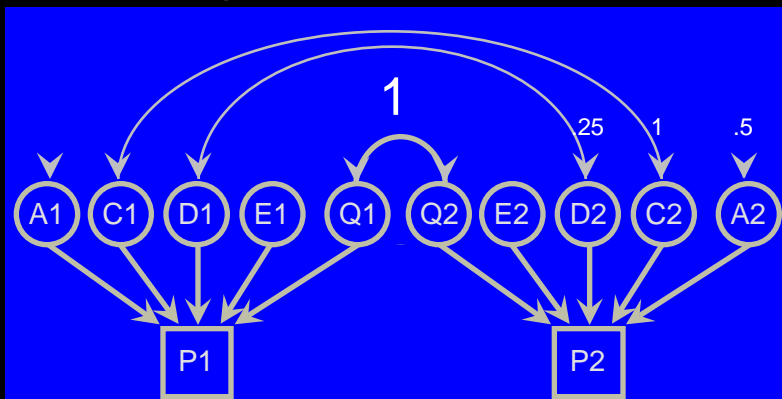# Partitioned Twin Analysis

- Nance & Neale (1989) Behav Genet 19:1

    - Separate DZ pairs into subgroups
        - IBD=0 IBD=1 IBD=2
    - Correlate Q with 0 .5 and 1 coefficients
    - Compute statistical power

**VCU**

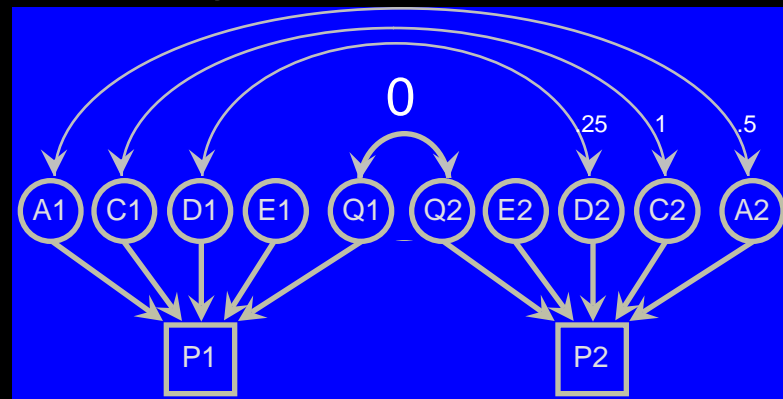# Partitioned Twin Analysis: Three DZ groups



IBD=1 group

.5    25    1    .5

A1 C1 D1 E1 Q1 Q2 E2 D2 C2 A2

P1    P2

IBD=2 group

1    25    1    .5

A1 C1 D1 E1 Q1 Q2 E2 D2 C2 A2

P1    P2

IBD=0 group

0    25    1    .5

A1 C1 D1 E1 Q1 Q2 E2 D2 C2 A2

P1    P2

# Problem 1 with Partitioned Twin analysis: Low Power

**Table II.** Twin Pairs Required to Reject False Hypotheses Under Two Research Designs: (A) MZ and DZ Twins; (B) Only DZ Twins

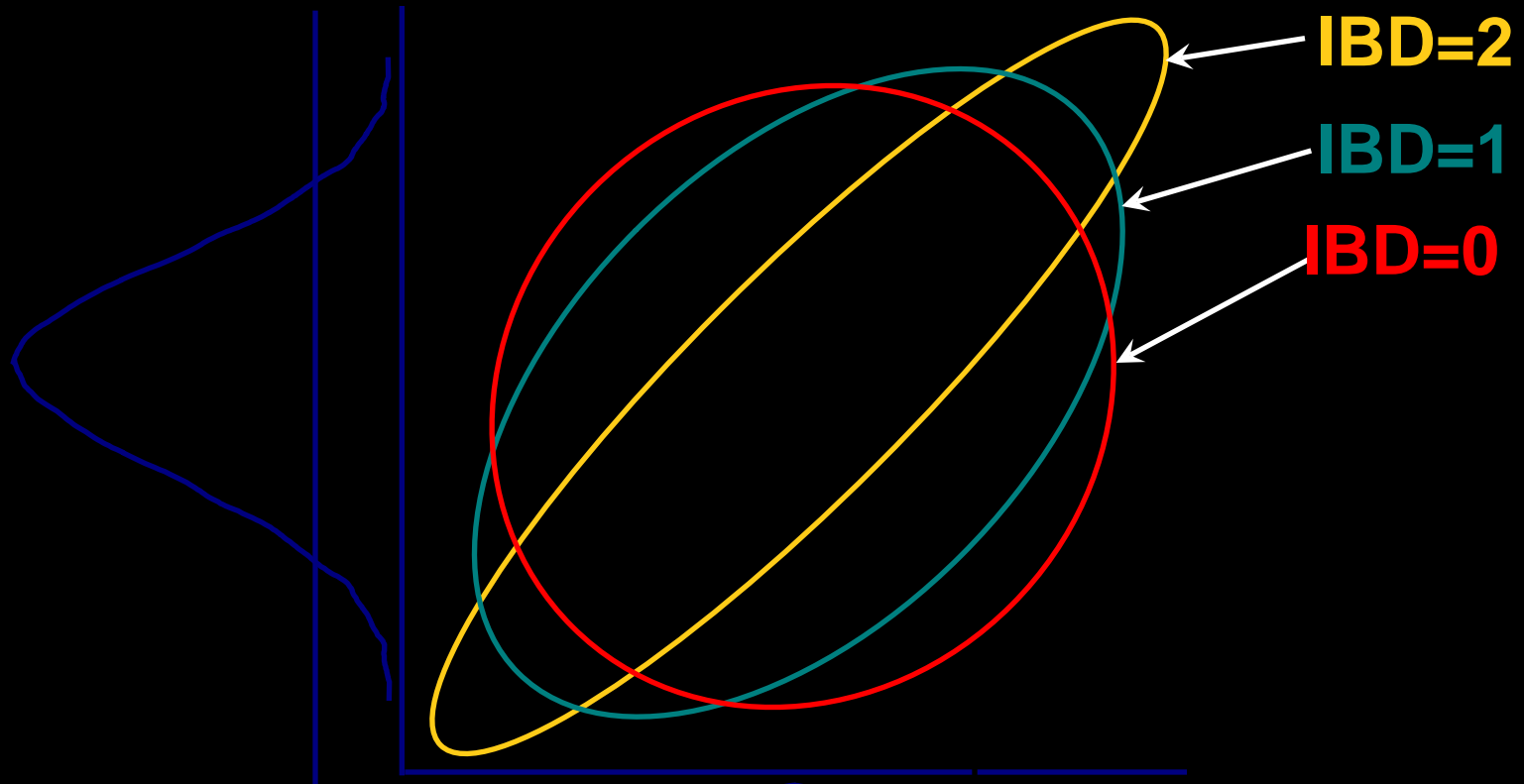| | True model (G, M, E) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Heritability | .9 | .6 | .3 | .9 | .6 | .3 | .9 | .6 | .3 |
| % marker effect | 75 | 75 | 75 | 50 | 50 | 50 | 25 | 25 | 25 |

(B) No. of pairs in three DZ subgroups in 1:2:1 ratio

# Problem 2: IBD is not known with certainty

- Markers may not be fully informative
  - Only so much heterozygosity in e.g., 20 allele microsatellite marker
  - Less in a SNP
  - Unlikely to have typed the exact locus we are looking for
  - Genome is big!

**VCU**

# IBD pairs vary in similarity



IBD=2

IBD=1

IBD=0

VCU

# Improving Power for Linkage

- Increase marker density (yaay SNP chips)
- Change design
  - Families
  - Larger Sibships
  - Selected samples
- Multivariate data
- More heritable traits with less error

**VCU**

# Problem 2: IBD is not known with certainty

- Markers may not be fully informative
  - Only so much heterozygosity in e.g., 20 allele microsatellite marker
  - Less in a SNP
  - Unlikely to have typed the locus that causes variation
  - Genome is big!
  - *The Universe is Big. Really big. It may seem like a long way to the corner chemist, but compared to the Universe, that's peanuts.* - D. Adams

**VCU**

# Center for STATISTICAL GENETICS

THE UNIVERSITY OF MICHIGAN · 1817

MERLIN

- Home
- Tutorial
- Download
- Register
- Reference
- FAQ

# MERLIN

Welcome!

MERLIN uses sparse trees to represent gene flow in pedigrees and is one of the fastest pedigree analysis packages around (Abecasis et al, 2002). Comments and suggestions are welcome, please e-mail goncalo@umich.edu.

Thanks to the Wizard of Draws for the cool cartoon!

©1998 Jeff Bucchino

University of Michigan | School of Public Health | Gonçalo Abecasis

# Using Merlin/Genehunter etc

- Several Faculty experts
  - Goncalo Abecasis
  - Sarah Medland
  - Stacey Cherny
- Possible to use Merlin via Mx GUI
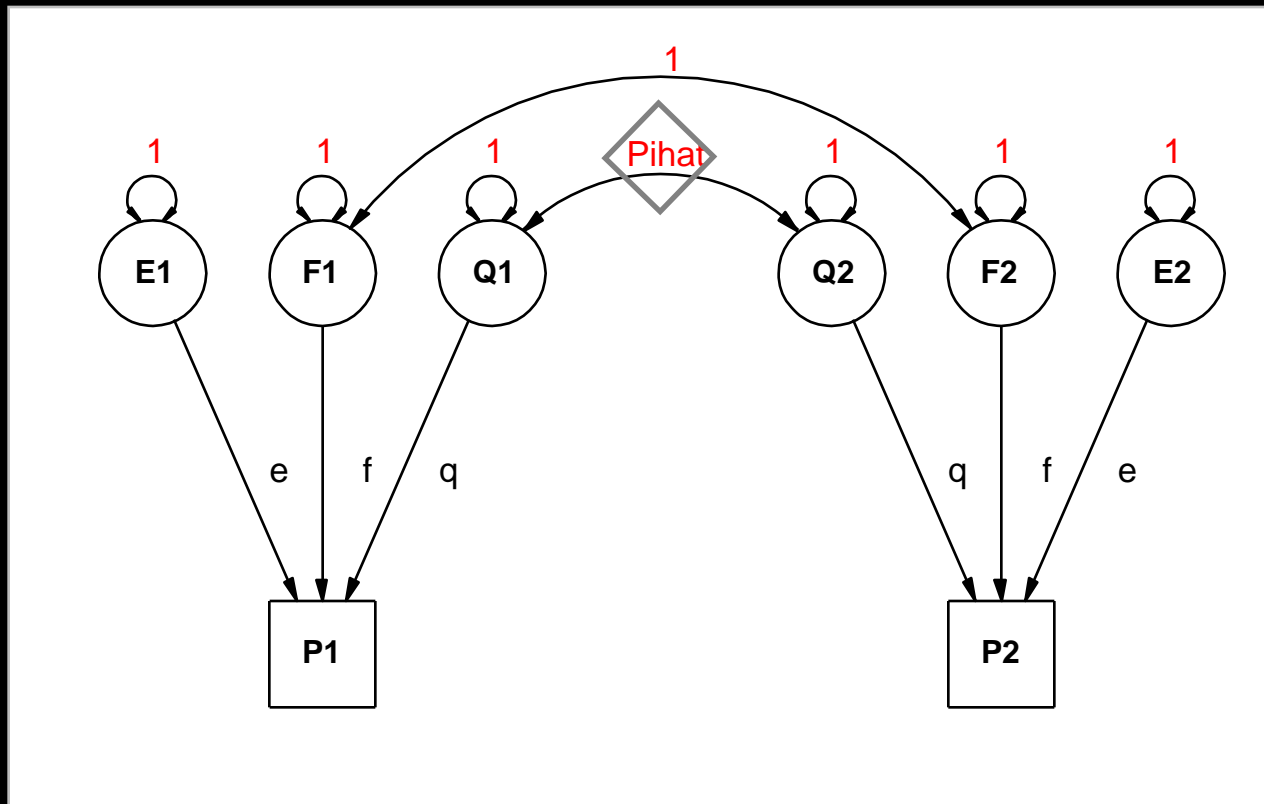
**VCU**

# "Pi-hat" approach

1  Pick a putative QTL location

2  Compute p(IBD=0) p(IBD=1) p(IBD=2) given
   marker data [Use Mapmaker/sibs or Merlin]

3  Compute $\hat{\pi}_i$ =  p(IBD=2) +  .5p(IBD=1)

4  Fit model

Repeat 1-4 as necessary for different locations across
   genome

Elston & Stewart

VCU

# Basic Linkage (QTL) Model

$$\pi_i = p(\widehat{IBD_i}=2) + .5\, p(IBD_i=1) \quad \textit{individual-level}$$



Q: QTL Additive Genetic     F: Family Environment     E: Random Environment
3 estimated parameters: q, f and e     Every sibship may have different model
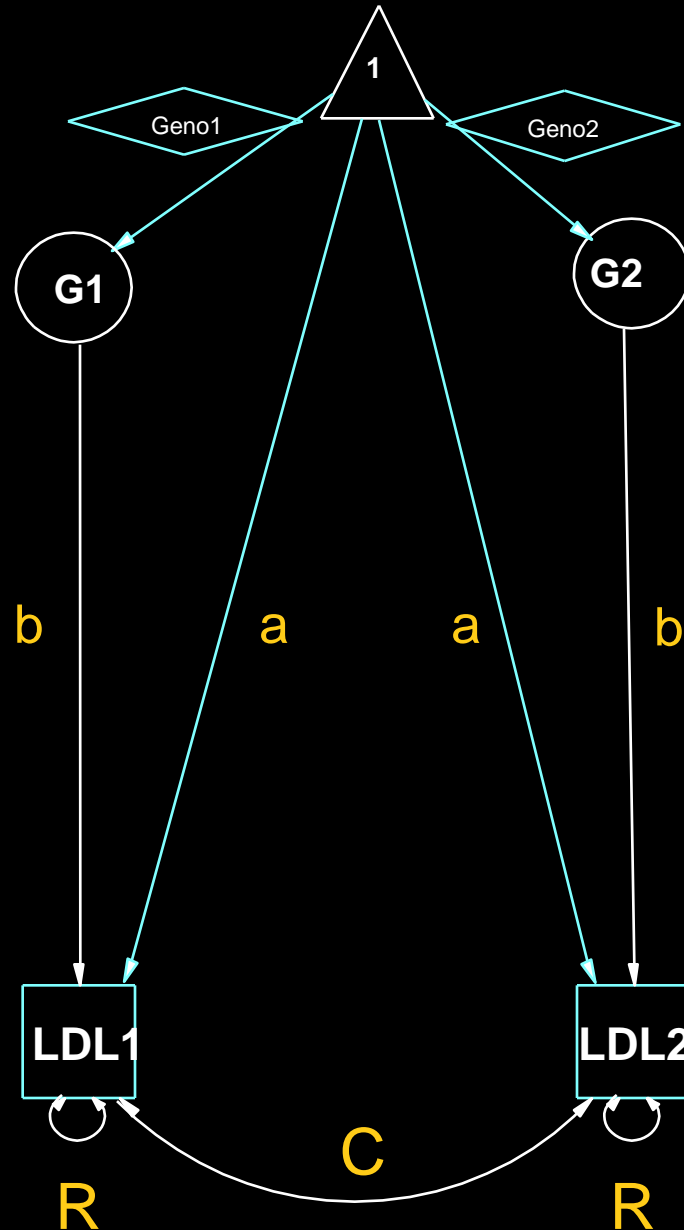
# Association Model

$$LDL1_i = a + b\ Geno1_i$$

$$Var(LDL_i) = R$$

$$Cov(LDL_1, LDL_2) = C$$

$C$ may be $f(\pi_i)$ in joint linkage & association

# Between/Within Fulker Association Model

Model for the means

$LDL1_i = .5bGeno1 + .5bGeno2 + .5wGeno1 - .5wGeno2$

$= .5( b(Geno1+Geno2) + w(Geno1-Geno2) )$