



## 4. Dealing with epistatic interactions and non-linearities

gene x gene

gene x gene x gene

gene x gene x gene x gene

.....

(Alice in Wonderland)



## Statistical Interaction (fixed effects models)

$$y_{ijk} = \mu + A_i + B_j + AB_{ij} + e_{ijk}$$

$$E(y_{ijk}|A_i, B_j, AB_{ij}) = \mu + A_i + B_j + AB_{ij}$$

$$\begin{aligned} E(y_{ijk} - y_{ij'k'}|A_i, B_j, AB_{ij}, A_{i'}, B_j, AB_{i'j}) &= \mu + A_i + B_j + AB_{ij} \\ &\quad - (\mu + A_{i'} + B_j + AB_{i'j}) \\ &= A_i - A_{i'} + AB_{ij} - AB_{i'j} \end{aligned}$$

Difference between levels of factor A depends on level of B

If factor **A** has **a** levels and factor **B** has **b** levels, the degrees of freedom are:


- (a-1)
- (b-1)
- (a-1)(b-1) [assuming no-empty cells]

## Multi-SNP Fixed effects models?

(unraveling “physiological epistasis” a la Cheverud)

- Lots of “main effects”
- Splendid non-orthogonality
- Lots of 2-factor interactions
- Lots of 3-factor interactions
- Lots of non-estimability
- Lots of uninterpretable high-order interactions
- Run out of “degrees of freedom”

## Analysis of SNPs with random effects models?

MEUWISSEN et al. (2001)  Will talk about this later

GIANOLA et al. (2003) }  
XU (2003) } “Ridge regression-type”

--Use all SNP markers in statistical models  
--Mechanistic basis to mixed effects linear model  
(genetic effects treated as random variables)  
--Highly parametric models  
--Strong assumptions made

## What are ridge and Bayesian regression? (given some variance components or tuning parameters)

$$\hat{\beta}_{OLS} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Large values of  $\lambda$  “shrink” regressions towards 0 (induces bias, but higher precision than OLS)

$$\hat{\beta}_{RIDGE} = (\mathbf{X}'\mathbf{X} + \mathbf{I}\lambda)^{-1}\mathbf{X}'\mathbf{y}$$

Special case of Bayesian linear regression

$$\hat{\beta}_{BAYES} = \left( \mathbf{X}'\mathbf{X} + \mathbf{B}^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} \right)^{-1} \left( \mathbf{X}'\mathbf{y} + \frac{\sigma_e^2}{\sigma_\beta^2} \mathbf{B}^{-1} \beta_0 \right)$$

Bayes model assumes, a priori

$$\beta \sim N(\beta_0, \mathbf{B}\sigma_\beta^2)$$

Typically assumed 0

Typically identity matrix. However, can be given structure

## ORDINARY LEAST-SQUARES

“Full model”  $\Rightarrow y = X\beta + e$   
 $= X_1\beta_1 + X_2\beta_2 + e$

“OLS” estimator  $\Rightarrow \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{bmatrix}^{-1} \begin{bmatrix} X_1'y \\ X_2'y \end{bmatrix}$   
 $= [X'X]^{-1}X'y$   
 $E(\hat{\beta}|X) = [X'X]^{-1}X'E(y)$   
 $= [X'X]^{-1}X'X\beta = \beta$

“OLS” is biased If full model holds and one fits “smaller” model (e.g., single marker Regressions)

$\Rightarrow y = X_1\beta_1 + e$   
 $E(\tilde{\beta}_1|X_1) = (X_1'X_1)^{-1}E(y)$   
 $= (X_1'X_1)^{-1}[X_1\beta_1 + X_2\beta_2]$   
 $= \beta_1 + (X_1'X_1)^{-1}X_1'X_2\beta_2$

## RIDGE REGRESSION

Can assess by cross-validation

$$\begin{aligned}
 \hat{\beta}_{Ridge} &= [X'X + I\lambda]^{-1} X'y \\
 &= [I + (X'X)^{-1}\lambda]^{-1} (X'X)^{-1} X'y \\
 &= [I + (X'X)^{-1}\lambda]^{-1} \hat{\beta}_{OLS} \quad \text{Shrinkage towards 0} \\
 E(\hat{\beta}_{Ridge}|X) &= [I + (X'X)^{-1}\lambda]^{-1} E(\hat{\beta}_{OLS}) \\
 &= [I + (X'X)^{-1}\lambda]^{-1} \beta
 \end{aligned}$$

Biased estimator but more precise

## BAYESIAN REGRESSION (ASSUMING KNOWN VARIANCE COMPONENTS)

Prior



$$\beta \sim N(0, B\sigma_\beta^2)$$

$$\begin{aligned}
 \hat{\beta}_{Bayes} &= \left[ X'X + B^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} \right]^{-1} \left( X'y + \frac{\sigma_e^2}{\sigma_\beta^2} B^{-1} \beta_0 \right) \\
 E(\hat{\beta}_{Bayes}|\beta) &= \left[ X'X + B^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} \right]^{-1} \left( X'X\beta + \frac{\sigma_e^2}{\sigma_\beta^2} B^{-1} \beta_0 \right) \\
 &= \left[ I + (BX'X)^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} \right]^{-1} (X'X)^{-1} \left( X'X\beta + \frac{\sigma_e^2}{\sigma_\beta^2} B^{-1} \beta_0 \right) \\
 &= \left[ I + (BX'X)^{-1} \frac{\sigma_e^2}{\sigma_\beta^2} \right]^{-1} \left( \beta + \frac{\sigma_e^2}{\sigma_\beta^2} (BX'X)^{-1} \beta_0 \right)
 \end{aligned}$$

Conditionally biased

## ILLUSTRATION OF SOMEPOINTS

Standard analysis (fixed  $\mathbf{X}$ ) but random  $\beta$

Genotype

$$y = f + e = X\beta + e$$

$$\beta \sim N(0, I\sigma_\beta^2)$$

$$E(y|X) = X\beta$$

$$\begin{aligned} \text{Var}(y|X) &= \text{Var}(f) + \text{Var}(e) \\ &= XX'\sigma_\beta^2 + I\sigma_e^2 \end{aligned}$$

Prediction of marker effects: BLUP  
(iid marker effects)

$$\begin{aligned} \left[ X'X + \frac{\sigma_e^2}{\sigma_\beta^2} I \right] \hat{\beta} &= X'y \\ \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} (X'X)^{-1} \right] \hat{\beta} &= (X'X)^{-1} X'y \\ \hat{\beta} &= \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} (X'X)^{-1} \right]^{-1} \tilde{\beta}_{\text{OLS}} \Rightarrow \text{SHRINKAGE} \end{aligned}$$

Prediction of signal ( $X\beta$ ) to phenotype

$$\begin{aligned} \text{Var}(X\beta|y) &= X\text{Var}(\beta|y)X' \\ &= X \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} (X'X)^{-1} \right]^{-1} X' \sigma_e^2 \end{aligned}$$

### Prediction of future record

$$y^* = X^* \beta + e^*$$

$$\begin{aligned} E(X^* \beta + e^* | y, X, X^*) &= X^* E(\beta | y, X) \\ &= X^* \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} (X'X)^{-1} \right]^{-1} \tilde{\beta}_{OLS} \end{aligned}$$

$$\text{Var}(X^* \beta + e^* | y, X, X^*) = X^* \text{Var}(\beta | y, X) X^* + I^* \sigma_e^2$$

### GAUSSIAN PROCESS ANALYSIS (IID MARKER EFFECTS)

$$y = f + e = X\beta + e$$

$$\beta \sim N(0, I\sigma_\beta^2)$$

$$X \sim F$$

} Assume  $\mathbf{X}$  and  $\beta$  are independent

$$E(y|X, \beta) = X\beta$$

$$E(y|\beta) = E_X E(y|X, \beta) = E(X)\beta$$

$$E(y) = E_\beta[E(X)\beta] = E(X)E(\beta) = 0$$

$$\text{Var}(y) = \text{Var}(f) + \text{Var}(e) = \text{Var}(f) + I\sigma_e^2$$

$$\text{Var}(f) = \text{Var}(X\beta)$$

$$= E_X(\text{Var}(X\beta|X) + \text{Var}_X[E(X\beta|X)])$$

$$= E_X[X\text{Var}(\beta)X'] + \text{Var}_X[XE(\beta)]$$

$$= E_X[XX'\sigma_\beta^2] + \text{Var}_X(0)$$

$$= \sigma_\beta^2 E_X[XX'],$$

BP= "best predictor"  
(MULVN assumed)

$$\hat{f} = \text{BP}(f)$$

$$\left[ \frac{1}{\sigma_e^2} I + \text{Var}^{-1}(f) \right] \hat{f} = \frac{1}{\sigma_e^2} y$$

$$\left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} E_X^{-1}[XX'] \right] \hat{f} = y$$

$$E_X^{-1}[XX'] \left[ E_X[XX'] + \frac{\sigma_e^2}{\sigma_\beta^2} I \right] \hat{f} = y$$

$$\left[ E_X[XX'] + \frac{\sigma_e^2}{\sigma_\beta^2} I \right] \hat{f} = E_X[XX']y$$

Under multivariate normality

$$\text{Var}(f|y) = \text{Var}(f) - \text{Cov}(f, y)\text{Var}^{-1}(y)\text{Cov}'(f, y)$$

$$= \text{Var}(f) - \text{Var}(f)[\text{Var}(f) + I\sigma_e^2]^{-1}\text{Var}(f)$$

$$= \sigma_\beta^2 E_X[XX'] - \sigma_\beta^2 E_X[XX'] [\sigma_\beta^2 E_X[XX'] + I\sigma_e^2]^{-1} \sigma_\beta^2 E_X[XX']$$

$$= \sigma_\beta^2 E_X[XX'] - \sigma_\beta^2 E_X[XX'] \frac{E_X^{-1}[XX']}{\sigma_\beta^2} \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} E_X^{-1}[XX'] \right]^{-1} \sigma_\beta^2 E_X[XX']$$

$$= \left\{ I - \left[ I + \frac{\sigma_e^2}{\sigma_\beta^2} E_X^{-1}[XX'] \right]^{-1} \right\} \sigma_\beta^2 E_X[XX'].$$

Future record:

$$\begin{aligned}
 f^* &= X^* \beta + e^* \\
 E(f^*|f) &= E(f^*) + \text{Cov}(X^* \beta, \beta X') \text{Var}^{-1}(f) f \\
 E(f^*|y) &= E_{f|y} E(f^*|f, y) = E_{f|y} E(f^*|f) \\
 &= E_{f|y} [\text{Cov}(X^* \beta, \beta X') \text{Var}^{-1}(f) f] \\
 &= \text{Cov}(X^* \beta, \beta X') \text{Var}^{-1}(f) \hat{f}
 \end{aligned}$$

$$\begin{aligned}
 \text{Cov}(X^* \beta, \beta X') &= E_{X, X^*} [\text{Cov}(X^* \beta, \beta X') | X, X^*] \\
 &\quad + \text{Cov}_{X, X^*} [E(X^* \beta), E(\beta X') | X, X^*] \\
 &= \sigma_\beta^2 E_{X, X^*} [X^* X'] + \text{Cov}_{X, X^*} (0, 0) \\
 &= \sigma_\beta^2 E_{X, X^*} [X^* X']
 \end{aligned}$$

Dealing with interactions (“statistical epistasis”): much of this took place in inspiring lowan landscapes...



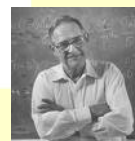
Bayesians,  
keep out!

SOME CORN



PIGS AGAIN

MORE  
PIGS HERE



C. C. C

$$\sum_i \sum_j \sum_k \sum_l \text{pig}_{ijkl}^2 - (\sum_i \sum_j \sum_k \sum_l \text{pig}_{ijkl})^2 / \text{as many pigs as you got}$$



**RANDOM EFFECTS MODELS**  
**FOR ASSESSING EPISTASIS REST ON:**  
**Cockerham (1954) and Kempthorne (1954)**

--Orthogonal partition of genetic variance into additive, dominance, additive x additive, etc. **ONLY** if



- ☐ No selection
- ☐ No inbreeding
- ☐ No assortative mating
- ☐ No mutation
- ☐ No migration
- ☐ Linkage equilibrium



A standard decomposition of phenotypic value in quantitative genetics (Falconer & Mackay, 1996) is

$$y = \mu + a + d + i + e,$$

where  $a$ ,  $d$  and  $i$  are additive, dominance and epistatic effects, respectively, and  $e$  is a residual, reflecting environmental (residual) variability. This linear de



The  $i$  effect can be decomposed into additive  $\times$  additive, additive  $\times$  dominance, dominance  $\times$  dominance, etc., deviates. In what has been termed 'statistical epistasis' (Cheverud & Routman, 1995), these deviates are assumed to be random draws from some distributions

The degrees of freedom of the distribution are NOT GIVEN by the number of levels.

There is now 1 df for each type of genetic effect.

$$\begin{aligned}
 &N(0, \sigma_a^2) \\
 &N(0, \sigma_d^2) \\
 &N(0, \sigma_{aa}^2) \\
 &N(0, \sigma_{ad}^2) \\
 &N(0, \sigma_{dd}^2) \\
 &\dots \\
 &N(0, \sigma_{ddd\dots d}^2)
 \end{aligned}$$

Matrix representation

$$\begin{aligned}
 \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + (\mathbf{a} + \mathbf{d} + \mathbf{i}_{aa} + \mathbf{i}_{ad} + \mathbf{i}_{dd}) + \mathbf{e} \\
 &= \mathbf{X}\boldsymbol{\beta} + \mathbf{g} + \mathbf{e},
 \end{aligned} \tag{1}$$

where  $\boldsymbol{\beta}$  is some nuisance location vector (equal to  $\mu$  if it contains a single element);  $\mathbf{X}$  is a known incidence matrix;  $\mathbf{a}$  and  $\mathbf{d}$  are vectors of additive and dominance effects, respectively;  $\mathbf{i}_{aa}$ ,  $\mathbf{i}_{ad}$  and  $\mathbf{i}_{dd}$  are epistatic effects, and  $\mathbf{g} = \mathbf{a} + \mathbf{d} + \mathbf{i}_{aa} + \mathbf{i}_{ad} + \mathbf{i}_{dd}$  is the 'total' genetic value. Assuming that  $\mathbf{g}$  and  $\mathbf{e}$  are uncorrelated, the variance-covariance decomposition is

Variance-covariance

$$\mathbf{V}_y = \mathbf{V}_g + \mathbf{V}_e, \tag{2}$$

where  $\mathbf{V}_y$ ,  $\mathbf{V}_g$  and  $\mathbf{V}_e$  are the phenotypic, genetic and residual variance-covariance matrices, respectively. Further,

Decomposition

$$\mathbf{V}_g = \mathbf{A}\sigma_a^2 + \mathbf{D}\sigma_d^2 + (\mathbf{A}\#\mathbf{A})\sigma_{aa}^2 + (\mathbf{A}\#\mathbf{D})\sigma_{ad}^2 + (\mathbf{D}\#\mathbf{D})\sigma_{dd}^2. \tag{3}$$

Here,  $\mathbf{A}$  is the numerator relationship matrix;  $\mathbf{D}$  is a matrix due to dominance relationships which can be computed from entries in  $\mathbf{A}$ , and the remaining matrices involve Hadamard (element by element) products of matrices  $\mathbf{A}$  or  $\mathbf{D}$ . Thus, under CK, all



DO THESE ASSUMPTIONS HOLD?

**RANDOM EFFECTS MODELS**  
**FOR ASSESSING EPISTASIS REST ON:**  
**Cockerham (1954) and Kempthorne (1954)**

--Orthogonal partition of genetic variance into additive, dominance, additive x additive, etc. **ONLY** if

- ☐ No selection
- ☐ No inbreeding
- ☐ No assortative mating
- ☐ No mutation
- ☐ No migration
- ☐ Linkage equilibrium

ALL  
 ASSUMPTIONS  
 VIOLATED!

Just consider  
 Linkage disequilibrium



## Digression: linkage disequilibrium

Let the genotypes at the first locus be  $A_1A_1$ ,  $A_1A_2$  and  $A_2A_2$  and  $B_1B_1$ ,  $B_1B_2$  and  $B_2B_2$  at the second locus. Let the frequency of the  $A_1$  allele be  $p_{A_1}$ , of the  $A_2$  allele be  $p_{A_2}$ , ( $p_{A_1} + p_{A_2} = 1$ ) and at the  $B$  locus the equivalent frequencies are  $p_{B_1}$  and  $p_{B_2}$ , ( $p_{B_1} + p_{B_2} = 1$ ). The four possible gametes are  $A_1B_1$ ,  $A_1B_2$ ,  $A_2B_1$ ,  $A_2B_2$ , with respective frequencies  $p_{A_1B_1}$ ,  $p_{A_1B_2}$ ,  $p_{A_2B_1}$ ,  $p_{A_2B_2}$  and  $p_{A_1B_1} + p_{A_1B_2} + p_{A_2B_1} + p_{A_2B_2} = 1$ . Notice that

$$p_{A_1} = p_{A_1B_1} + p_{A_1B_2},$$

$$p_{A_2} = p_{A_2B_1} + p_{A_2B_2},$$

$$p_{B_1} = p_{A_1B_1} + p_{A_2B_1},$$

$$p_{B_2} = p_{A_1B_2} + p_{A_2B_2}.$$

If the allelic state at locus  $A$  is independent of that at locus  $B$ , one expects  $p_{A_1B_1} = p_{A_1}p_{B_1}$ ,  $p_{A_1B_2} = p_{A_1}p_{B_2}$ , and so on. The system is said to be in *linkage equilibrium*: the alleles at loci  $A$  and  $B$  are independent and their joint frequency is given by the product of their

marginal frequencies. If this is not the case, the dependence between alleles at loci  $A$  and  $B$  is measured by their covariance, known as *linkage disequilibrium* and symbolised by  $D$ . Define the random variable  $X$  which takes the value 1 if in gametes,  $A_i$  is present at locus  $A$  and zero otherwise, and the random variable  $Y$  which takes the value 1 if  $B_j$  is present at locus  $B$  zero otherwise. The expected value of  $X$  is  $p_{A_i}$  and that of  $Y$  is  $p_{B_j}$ . The expected value of  $(XY)$  is  $p_{A_i B_j}$  and the covariance between  $X$  and  $Y$  is by definition

$$\begin{aligned} D &= \text{Cov}(X, Y) \\ &= E(X, Y) - E(X) E(Y) \\ &= p_{A_i B_j} - p_{A_i} p_{B_j}. \end{aligned} \quad (1)$$

For example, if we set arbitrarily  $i = 1, j = 1$ , then

$$p_{A_1 B_1} = p_{A_1} p_{B_1} + D, \quad (2)$$

and

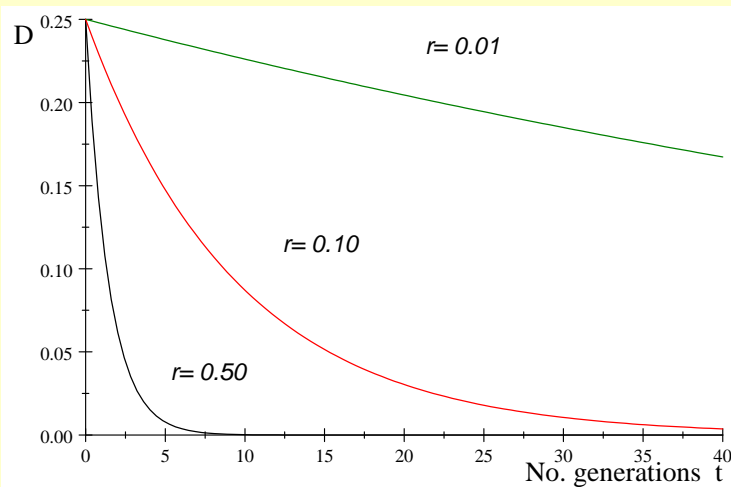
$$\begin{aligned} D &= p_{A_1 B_1} - p_{A_1} p_{B_1} \\ &= p_{A_1 B_1} (p_{A_1 B_1} + p_{A_1 B_2} + p_{A_2 B_1} + p_{A_2 B_2}) - (p_{A_1 B_1} + p_{A_1 B_2}) (p_{A_1 B_1} + p_{A_2 B_1}) \\ &= p_{A_1 B_1} p_{A_2 B_2} - p_{A_1 B_2} p_{A_2 B_1}, \end{aligned} \quad (3)$$

the difference between the product of the frequencies of the *coupling* and *repulsion* gametic phases. Choosing  $i = 1, j = 1$  resulted in (2) and in

$$\begin{aligned} p_{A_1 B_2} &= p_{A_1} p_{B_2} - D, \\ p_{A_2 B_1} &= p_{A_2} p_{B_1} - D, \\ p_{A_2 B_2} &= p_{A_2} p_{B_2} + D. \end{aligned}$$

### Evolution of linkage disequilibrium as a function of recombination rate

$$D_t = (1 - r)^t D_0$$



# The pattern of linkage disequilibrium in German Holstein cattle

S. Qanbari\*, E. C. G. Pimentel\*, J. Tetens†, G. Thaller†, P. Lichtner†, A. R. Sharifi\* and H. Simianer\*

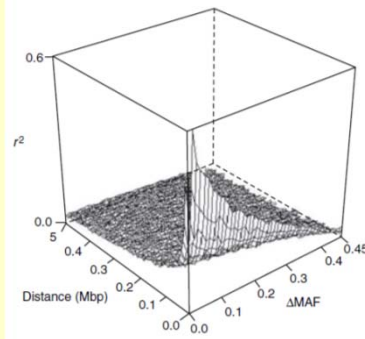
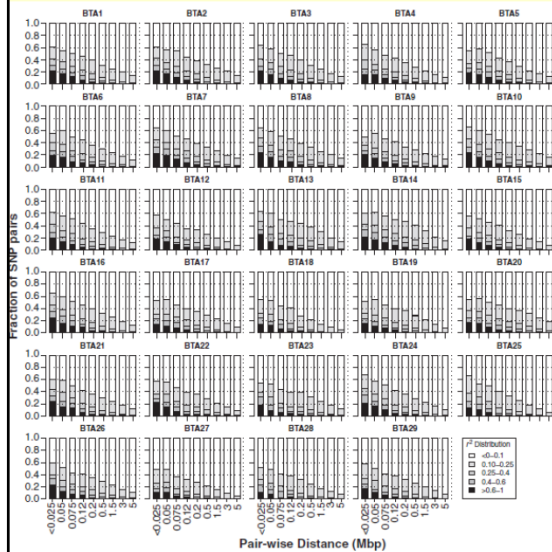


Figure 5 Three-dimensional surface plot depicting the decay of linkage disequilibrium vs. inter-marker distance and minor allelic frequency interval.

## A VIEW OF LINEAR MODELS (as employed in q. genetics)

Mathematically, can be viewed as a “local” approximation of a complex process

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f^{(3)}(a)}{3!}(x-a)^3 + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + \dots$$

Linear approximation

Quadratic approximation

$n^{\text{th}}$  order approximation

**Example**

$$y = g(x) + e$$

Response variate

Model residual

Some function of a covariate  $x$

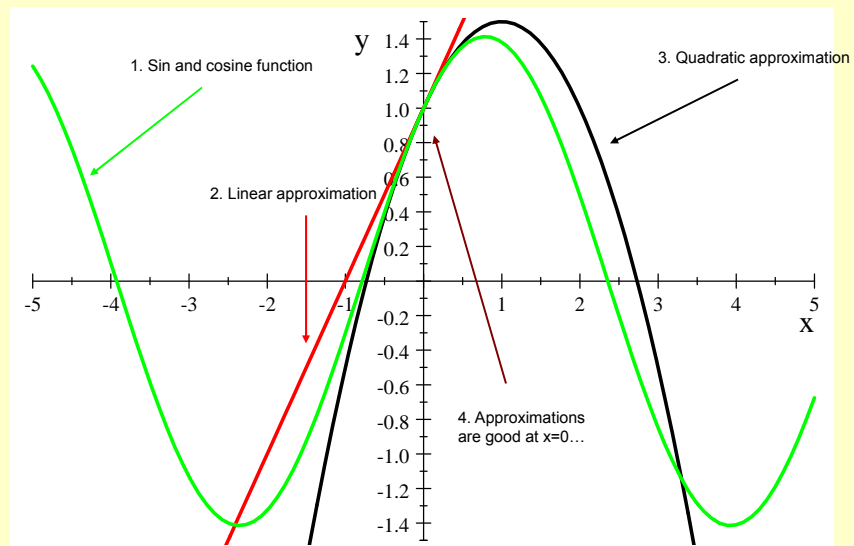
Suppose  $g(x) = \sin(x) + \cos(x)$

$$\begin{aligned}\frac{d}{dx} [\sin(x)] &= [\cos x] \\ \frac{d}{dx} [\cos(x)] &= [-\sin x] \\ \frac{d}{dx} [\sin(x) + \cos(x)] &= [\cos x - \sin x] \\ \frac{d}{dx} [\cos x - \sin x] &= [-\cos x - \sin x] \\ \frac{d^2}{(dx)^2} [\sin(x) + \cos(x)] &= [-\cos x - \sin x]\end{aligned}$$

Second-order Taylor series expansion about 0

$$\begin{aligned}[\sin(x) + \cos(x)] &\approx [\sin(0) + \cos(0)] + [\cos 0 - \sin 0](x - 0) + \frac{1}{2}[-\cos 0 - \sin 0](x - 0)^2 \\ &= 1 + x - \frac{x^2}{2}\end{aligned}$$

How good are the linear and quadratic approximations? Recall that a Taylor series provides a local approximation only...



## Finding structure from noisy data we have environmental noise...:

evaluate function  $\sin(x)+\cos(x)$  at  $x=0, 0.5$  and  $1$

True values are:

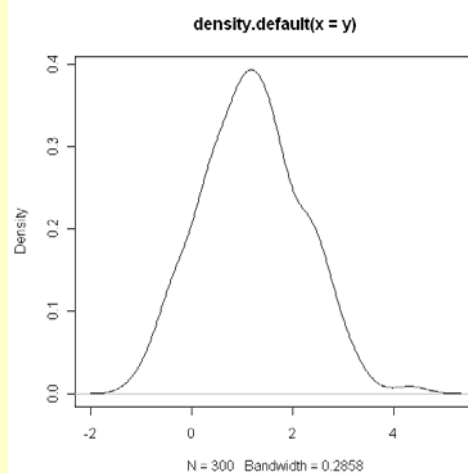
```
> sin(0)+cos(0)
[1] 1
> sin(0.5)+cos(0.5)
[1] 1.357008
> sin(1)+cos(1)
[1] 1.381773
```

VERY CLOSE TO EACH OTHER  
NOISE CAN MASK SIGNALS!

Create an R data set ( $N=300$ ) from adding 100  $N(0,1)$  residuals to each of the 3 values

```
> y0<-sin(0)+cos(0) +rnorm(100,0,1)
> y05<-sin(0.5)+cos(0.5)+rnorm(100,0,1)
> y1<-sin(1)+cos(1) +rnorm(100,0,1)
> y<-c(y0,y05,y1)
```

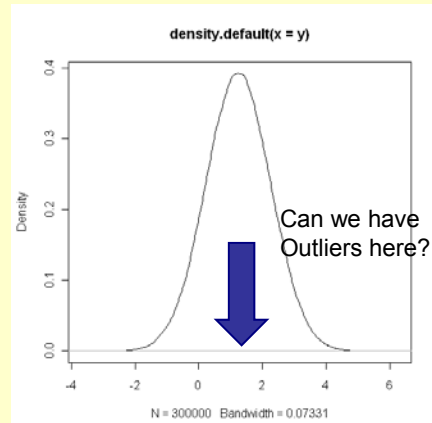
**MEASURING MACHINE 1**



Create a larger R data set (N=300000) by adding  
100000 N(0,1) residuals to each of the 3 values

```
> y0<-sin(0)+cos(0)      +rnorm(100000,0,1)
> y05<-sin(0.5)+cos(0.5) +rnorm(100000,0,1)
> y1<-sin(1)+cos(1)      +rnorm(100000,0,1)
> y<-c(y0,y05,y1)
```

CANNOT SEE UNDERLYING STRUCTURE.  
LARGE NOISE (ERROR VARIANCE)

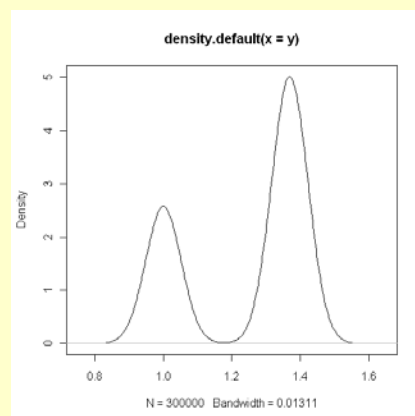


Now we get a more precise measuring instrument with variance 0.05

```
> y0<-sin(0)+cos(0)      +rnorm(100000,0,.05)
> y1<-sin(1)+cos(1)      +rnorm(100000,0,.05)
> y05<-sin(0.5)+cos(0.5) +rnorm(100000,0,.05)
```

### MEASURING MACHINE 2

STRUCTURE IS REVEALED BUT  
WE CANNOT DIFFERENTIATE  
BETWEEN TWO OF THE UNDERLYING  
VALUES





...SO WE BUY ANOTHER INSTRUMENT WITH VARIANCE 0.001!

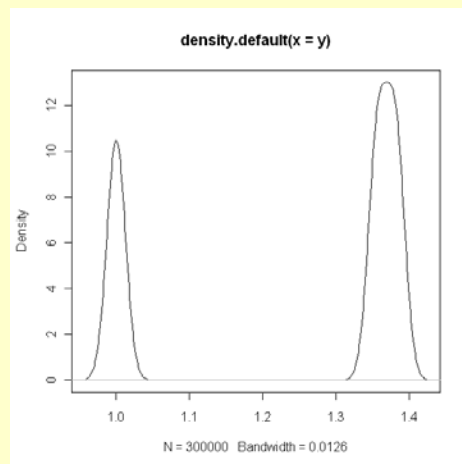
```
> y0<-sin(0)+cos(0)      +rnorm(100000,0,.001)
> y1<-sin(1)+cos(1)      +rnorm(100000,0,.001)
> y05<-sin(0.5)+cos(0.5) +rnorm(100000,0,.001)
> y<-c(y0,y05,y1)
```

### MEASUREMENT MACHINE 3

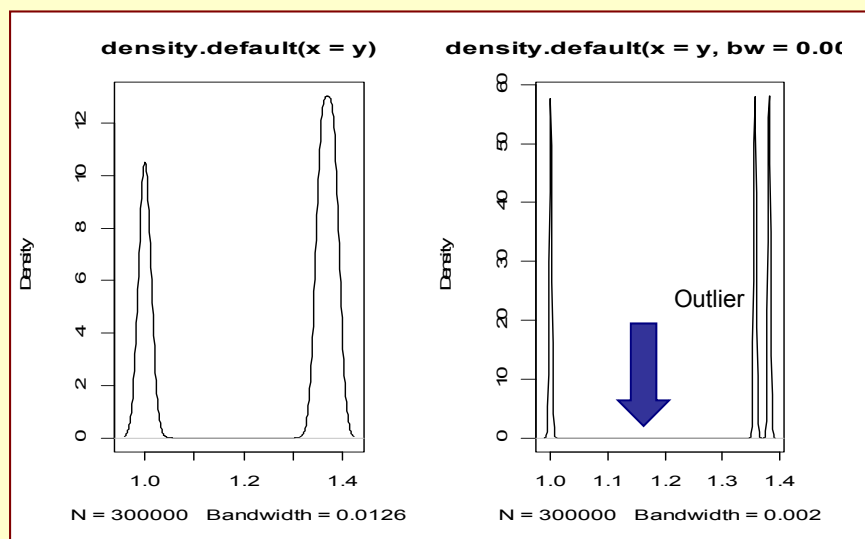
STILL CANNOT DIFFERENTIATE  
BETWEEN THE

```
> sin(0.5)+cos(0.5)
[1] 1.357008
```

```
> sin(1)+cos(1)
[1] 1.381773
```



HOWEVER, NON-PARAMETRIC DENSITY ESTIMATES DEPEND ON SOME  
BANDWIDTH PARAMETER. BY REDUCING IT, WE CAN SEE THE ENTIRE  
STRUCTURE OF THE PROBLEM...



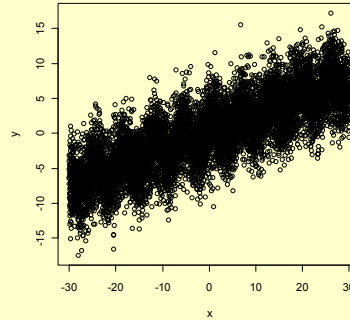
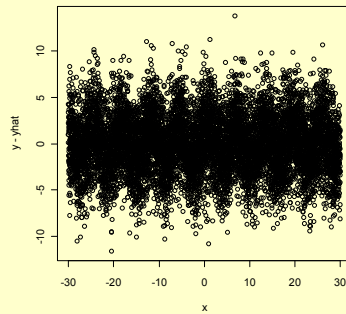
## FINDING “STRUCTURE” WITH A LINEAR MODEL

We are given (x,y) data (n=10,000). It looks like this and we run a linear regression

$$\hat{y} = 0.07936 + 0.24814x$$

```
> cor(x,y)
[1] 0.8064256
```

```
> cor(y,yhat)
[1] 0.8064256
```



RESIDUALS DISPLAY  
SINUSOIDAL BEHAVIOR

## TRUE MODEL

```
> e<-rnorm(10000,0,sqrt(9))
```

```
> x<-runif(10000,-30,30)
```

```
> a<-0.10
```

```
> b<-0.25
```

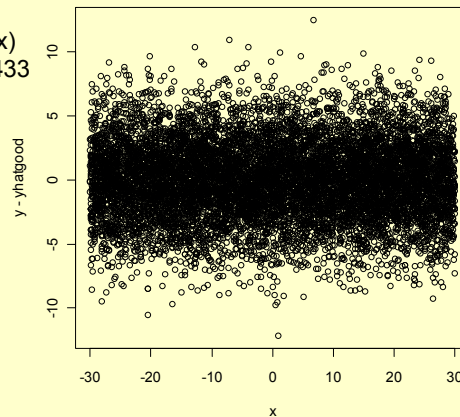
```
> y<-a+b*x+sin(x)+cos(x)+e
```

```
> model<-lm(y~x+sin(x)+cos(x))
```

>Coefficients:

```
>(Intercept)      x      sin(x)      cos(x)
> 0.1030      0.2489      0.9518      0.9433
```

RESIDUALS LOOK RANDOM



WE GENERATE A NEW SAMPLE AT THE SAME VALUES OF X

```
> enew<-rnorm(10000,0,sqrt(9))
> ynew<-a+b*x+sin(x)+cos(x)+enew
```

CALCULATE PREDICTIVE MEAN SQUARED ERROR

```
> msepredbadmodel<-sum((ynew-yhat)**2/10000)
> msepredbadmodel
[1] 9.725709
```

```
> msepredgoodmodel<-sum((ynew-yhatgood)**2/10000)
> msepredgoodmodel
[1] 8.729272
```

CALCULATE PREDICTIVE CORRELATIONS

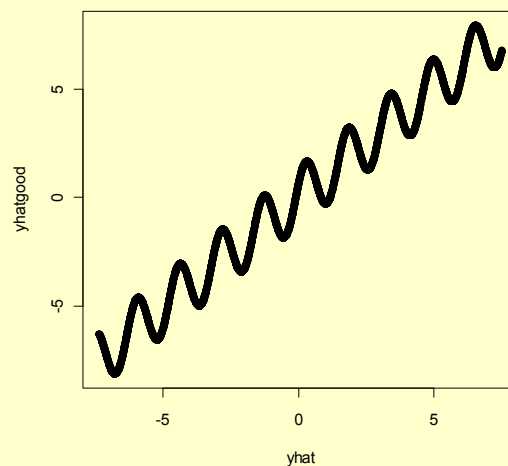
```
> cor(yhat,ynew)
[1] 0.8070097
> cor(yhatgood,ynew)
[1] 0.828854
```

MSE(Good)/MSE(Bad)=0.8975  
MSE(Bad)/MSE(Good)=1.1141

Cor(BAD)/Cor(GOOD)=0.9736

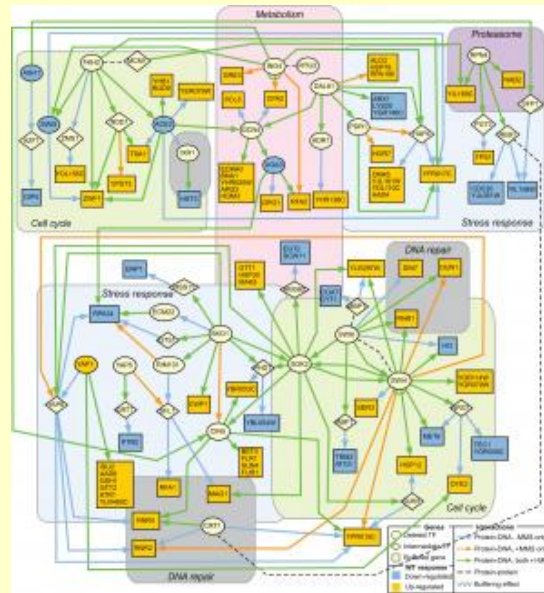
```
> lm(yhat~yhatgood)
Coefficients:
(Intercept)  yhatgood
  0.005653    0.953468
```

**DO NOT TRUST CORRELATIONS!**



## RECALLING COMPLEXITY...

How one  
Would model  
something like this?



**Heal Thyself: Systems Biology Model Reveals How Cells Avoid Becoming Cancerous.** ScienceDaily (May 21, 2006)

## What to do in genomic-assisted analysis of complex genetic signals?

- Include all markers, model all possible interactions? Unrealistic...
- Select sets of influential markers via model selection
  - ➔ Huge search space
  - ➔ Frequentist methods "err" probabilistically
  - ➔ Bayesian model selection (RJMC) difficult to tune
- Use LASSO (least absolute shrinkage and selection operator): Tibshirani (1996). What about interactions?
- Explore model-free techniques that have been used successfully in many domains
  - ➔ **semi-parametric regression**
  - ➔ **machine learning**: focus on prediction, learning mappings from inputs to outputs

# DEFINITION OF MACHINE LEARNING (Wikipedia)

**Machine learning:** subfield of [artificial intelligence](#) concerned with design and development of [algorithms](#) that allow [computers](#) (machines) to improve their performance over time (to [learn](#)) based on [data](#),

A major focus of machine learning research is to automatically produce (induce) [models](#), such as [rules](#) and [patterns](#), from data. Hence, machine learning is closely related to fields such as [data mining](#), [statistics](#), [inductive reasoning](#), [pattern recognition](#), and [theoretical computer science](#).