

# W05 - Maximum Likelihood - least squares - linear regression

Sivia ch 3, 5 ..

Lander and Botstein 1988

biol problem: Ding et al. 2018

## Maximum Likelihood (ML)

D - data

H - hypothesis

X - parameters of H

$$P(D|X, H)$$

$$P(X|D, H) = \frac{P(D|X, H) P(X|H)}{P(D)}$$

posterior:

$$P(X|D, H) \propto P(D|X, H) \cdot P(X|H)$$

if we assume an uninformative prior  $P(X|H) = \text{constant}$

$$\begin{aligned} \operatorname{argmax}_X P(X|D, H) &= \operatorname{argmax}_X P(D|X, H) \\ &= \operatorname{argmax}_X \log P(D|X, H) \end{aligned}$$

$X^*$  ML fit data is also  $X$  at which

posterior of  $X$  given  $D$  is also maximal.

(uninformative prior).

$$L(X) = \log P(D|X, H), \quad \left. \frac{\delta L}{\delta X} \right|_{X^*} = 0$$

$$L(x) = \log P(D|X H)$$

$$\left. \frac{\delta L(x)}{\delta x} \right|_{x^*} = 0 \quad \left. \frac{\delta^2 L}{\delta x^2} \right|_{x^*} = - \frac{1}{\sigma^2}$$

Distribution	$P(D X H)$	$x^*$	$\sigma^2$
exponential $D = \{t_1, \dots, t_N\}$	$\frac{1}{\lambda} e^{-t/\lambda}$ $\frac{e^{-\sum_i t_i/\lambda}}{\lambda^N}$	$x^* = \frac{\sum_i t_i}{N}$	$\lambda^2 / N$
binomial (P) $n, N$	$P^{\sum_i t_i} (1-P)^{N-\sum_i t_i}$ $\binom{N}{n} P^n (1-P)^{N-n}$	$P^* = \frac{n}{N}$	$\frac{P^*(1-P^*)}{N}$
Gaussian $\mu, \sigma$ $D = \{t_1, \dots, t_N\}$	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}}$ $\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^N e^{-\frac{\sum_i (t_i - \mu)^2}{2\sigma^2}}$	$\mu^* = \frac{\sum_i t_i}{N}$	$\sigma_r^2 = \frac{\sigma^2}{N}$ $\sigma_r^2 = \frac{\sigma^2}{2N}$

$$P(D|\mu, \sigma) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^N e^{-\frac{\sum_i (t_i - \mu)^2}{2\sigma^2}}$$

$$\log P(D|\mu, \sigma) = -N \log \sqrt{2\pi} - N \log \sigma - \frac{\sum_i (t_i - \mu)^2}{2\sigma^2}$$

$$\frac{\delta L(\mu, \sigma)}{\delta \mu} = -\frac{1}{2\sigma^2} \cdot 2 \sum_i (t_i - \mu) \cdot (-1) = + \frac{1}{\sigma^2} \sum_i (t_i - \mu)$$

$$\frac{\delta L(\mu, \sigma)}{\delta \sigma} = -\frac{N}{\sigma} - (-2) \frac{\sum_i (t_i - \mu)^2}{2\sigma^3} = -\frac{1}{\sigma} \left[ N - \frac{\sum_i (t_i - \mu)^2}{\sigma^2} \right]$$

$$\left. \frac{\delta L}{\delta \mu} \right| = 0 \quad \left\{ \mu^* = \frac{1}{N} \sum_i t_i \right.$$

$$\left. \frac{\delta L}{\delta \sigma} \right| = 0 \quad \left\{ \sigma^{*2} = \frac{\sum_i (t_i - \mu)^2}{N} \right.$$

$$\frac{\delta^2 L}{\delta \mu^2} = -\frac{N}{\sigma^2} \quad \left. \right|_{\mu^*, \sigma^*} = -\frac{N}{\sigma^{*2}} \Rightarrow \left\{ \sigma_{\mu}^{*2} = \frac{\sigma_r^2}{N} \right.$$

$$\begin{aligned} \left. \frac{\delta^2 L}{\delta \sigma^2} = +\frac{N}{\sigma^2} - 3 \frac{\sum_i (t_i - \mu)^2}{\sigma^4} \right|_{\mu^*, \sigma^*} &= \frac{N}{\sigma_r^2} - 3 \frac{N \sigma_r^2}{\sigma_r^4} \\ &= -2 \frac{N}{\sigma_r^2} \Rightarrow \left\{ \sigma_{\sigma}^{*2} = \frac{\sigma_r^2}{2N} \right. \end{aligned}$$

## Least-squares fit

$$D = \{d_i\}_{i=1}^N$$

$$P(D|XH) = \prod_{i=1}^N P(d_i|XH)$$

Assume we know how ideal data behaves except for a noise

$$d_i = f_i(x)$$

$$d_i - f_i(x) \sim \mathcal{W}(0, \sigma_i) \quad \begin{array}{l} \text{Gaussian noise} \\ \cdot \text{symmetric} \\ \cdot \text{not too large} \end{array}$$

$$\sigma_i = \sigma$$

$$P(D|XH) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^N e^{-\frac{\sum_i (d_i - f_i(x))^2}{2\sigma^2}}$$

$$\log P(D|XH) \propto -N \log \sigma - \frac{\sum_i (d_i - f_i(x))^2}{2\sigma^2}$$

$$\chi^2 = \sum_i (d_i - f_i(x))^2 \quad \text{least}$$

$$\left\{ \arg \max_x \log P(D|XH) = \arg \max_x \chi^2 \right.$$

least square fit  
 $L^2$ -norm.

Particular case fit to a line

$$d_i = \{y_i, x_i\}_{i=1}^N$$

$$y_i - (a + bx_i) \sim N(0, \sigma)$$

$$P(D|a, b) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^N e^{-\frac{\sum_i (y_i - a - bx_i)^2}{2\sigma^2}}$$

$$\operatorname{argmax}_{a, b} P(D|a, b) = \operatorname{argmax}_{a, b} \sum_i (y_i - (a + bx_i))^2$$

linear regression

# QTL analysis

crossable  $S_1$   $S_2$  } have a phenotypic difference  
 which genomic region is responsible for it?

noisy: many loci  
 large environment contribution.

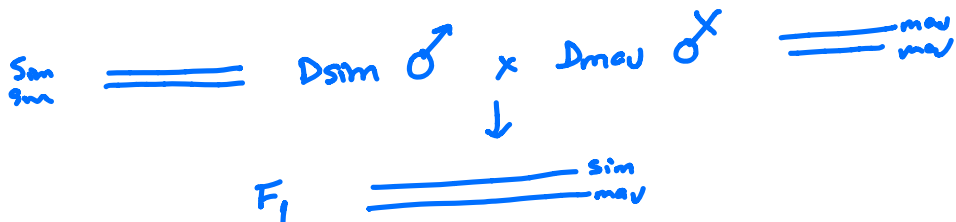
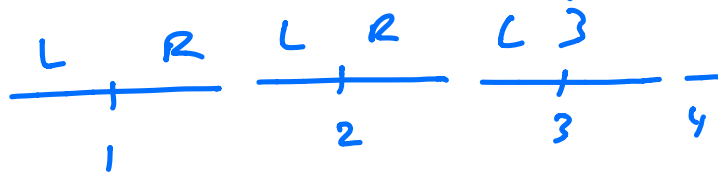
To identify the loci responsible for a particular trait.

## Ding et al 2016

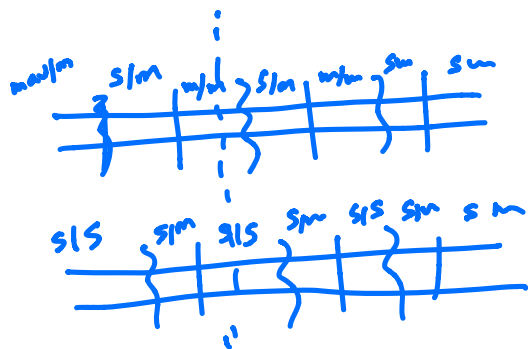
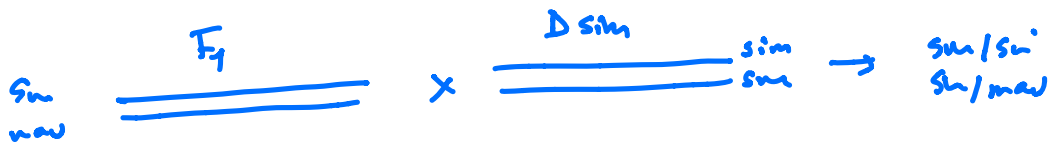
$\times$  (*Drosophila Simulans* (Dsim)  $\mu_{sim} \approx 176$  Hz  
*Drosophila mauritiana* (Dmau)  $\mu_{mau} = 186$  Hz)

male courtship sine song  
 where does the difference come from?

Dsim, Dmau ~ 98% similarity of genomes



## Back crosses



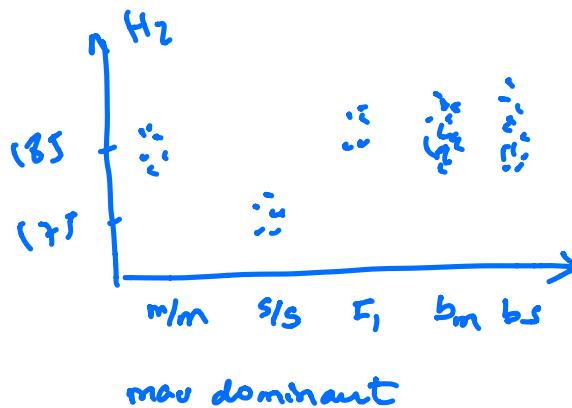
~ 1 breakpoint across

From back crosses collect in final population  $i$

+ phenotype  $\{p_i\}_{i=1}^L$

+ genotype  $\{g_i\}_{i=1}^L$

$$g_i = \begin{cases} 1 & \text{if } m/m \\ & m/s \\ 0 & \text{if } s/s \end{cases}$$



notice: how to get the haplotype assignment is another computational question in itself that we will study in wof/wof using hidden Markov models (HMMs)

The models assume are loci

+ QTL model  $f_i = a + b g_i + \epsilon_i$   
 $\uparrow f_{ly}$   $\uparrow A_{ly}$

$$P(D|a,b, QTL) = \left( \frac{1}{\sqrt{2\pi} \sigma} \right)^N e^{-\frac{\sum_i (f_i - a - b g_i)^2}{2\sigma^2}}$$

+ NQTL model  $f_i = c + \epsilon_i$

$$P(D|c, NQTL) = \left( \frac{1}{\sqrt{2\pi} \sigma} \right)^N e^{-\frac{\sum_i (f_i - c)^2}{2\sigma^2}}$$

QTL  $\rightarrow$  2 parameters  $a, b$

NQTL  $\rightarrow$  1 "  $c$  (if  $a=c$  and  $b=0$  identical)

$$\frac{P(QTL|D)}{P(NQTL|D)} = \frac{P(D|QTL) \cdot P(QTL)}{P(D|NQTL) \cdot P(QTL)} = \frac{P(D|QTL)}{P(D|NQTL)}$$

equal words

$$P(D|QTL) = \int_{-\infty}^{+\infty} da \int_{-\infty}^{+\infty} db P(D|a,b) P(a,b|QTL)$$

$$P(D|NQTL) = \int_{-\infty}^{+\infty} dc P(D|c) P(c|NQTL)$$



priors for parameters

$$P(a, b | QTL) = P(a | QTL) \cdot P(b | QTL)$$

$$P(a | QTL) = \begin{cases} \frac{1}{a^+ - a^-} = \frac{1}{\sigma_a} & \bar{a} < a < a^+ \\ 0 & \text{otherwise} \end{cases}$$

$$P(b | QTL) = \begin{cases} \frac{1}{b^+ - b^-} = \frac{1}{\sigma_b} & \bar{b} < b < b^+ \\ 0 & \text{otherwise} \end{cases}$$

$$P(c | NQTL) = \begin{cases} \frac{1}{c^+ - c^-} = \frac{1}{\sigma_c} & \bar{c} < c < c^+ \\ 0 & \text{otherwise} \end{cases}$$

Laplace approximation: for  $P(D|a,b,QTL)$   
 $P(D|c,NQTL)$

$$L^Q(a,b) =: \log P(D|a,b,QTL)$$

$$L^N(c) =: \log P(D|c,NQTL)$$

Taylor expansion around ML values  $a^*, b^*, c^*$

$$L^Q(a,b) \approx L^Q(a^*, b^*) + \frac{1}{2} \left. \frac{\partial^2 L^Q}{\partial a^2} \right|_{a^*} (a - a^*)^2$$

$$\left. \frac{\partial L^Q}{\partial a} \right|_{a^*, b^*} = 0 \quad + \frac{1}{2} \left. \frac{\partial^2 L^Q}{\partial b^2} \right|_{a^*} (b - b^*)^2$$

$$\left. \frac{\partial L^Q}{\partial b} \right|_{a^*, b^*} = 0 \quad + \left. \frac{\partial^2 L^Q}{\partial a \partial b} \right|_{a^*} (a - a^*)(b - b^*)$$

$$L^N(c) \approx L^N(c^*) + \frac{1}{2} \left. \frac{\partial^2 L^N}{\partial c^2} \right|_{c^*} (c - c^*)^2$$

$$\left. \frac{\partial L^N}{\partial c} \right|_{c^*} = 0$$

## ML rates of QTL/NQTL parameters

$$L^Q_{OC} = \frac{\sum_i (f_i - a - b g_i)^2}{2\sigma^2}$$

$$L^N_{OC} = \frac{\sum_i (f_i - c)^2}{2\sigma^2}$$

$$\frac{\delta L^Q}{\delta a} = - \frac{2}{2\sigma^2} \sum_i (f_i - a - b g_i) (-1) = \frac{1}{\sigma^2} \sum_i (f_i - a - b g_i)$$

$$\frac{\delta L^Q}{\delta b} = - \frac{2}{2\sigma^2} \sum_i (f_i - a - b g_i) (-g_i) = \frac{1}{\sigma^2} \sum_i (f_i - a - b g_i) g_i$$

$$= \frac{1}{\sigma^2} \left( \sum_i f_i g_i - a \sum_i g_i - b \sum_i g_i^2 \right)$$

$$\bar{f} = \frac{1}{N} \sum_i f_i \quad \bar{f}^2 = \frac{1}{N} \sum_i f_i^2$$

$$\bar{g} = \frac{1}{N} \sum_i g_i \quad \bar{g}^2 = \overline{g^2} = \frac{1}{N} \sum_i g_i^2$$

$$\overline{fg} = \overline{f g} = \frac{1}{N} \sum_i f_i g_i$$

$$\frac{\delta L^Q}{\delta a} = \frac{N}{\sigma^2} (\bar{f} - a - b \bar{g}) \Big|_{=0} \quad \bar{f} = a + b \bar{g}$$

$$\frac{\delta L^Q}{\delta b} = \frac{N}{\sigma^2} (\overline{fg} - a \bar{g} - b \overline{g^2}) \Big|_{=0} \quad \overline{fg} = a \bar{g} + b \overline{g^2}$$

$$a = \bar{f} - b \bar{g}$$

$$\frac{\delta L^N}{\delta c} = \frac{N}{\sigma^2} (\bar{f} - c) \Big|_{=0} \Rightarrow \boxed{c = \bar{f}}$$

$$\overline{fg} = (\bar{f} - b \bar{g}) \bar{g} + b \overline{g^2}$$

$$b [\overline{g^2} - \bar{g}^2] = \overline{fg} - \bar{f} \bar{g}$$

$$\left. \begin{array}{l} \text{ML} \\ \text{QTL} \end{array} \right\} \begin{cases} b^* = \frac{\bar{f}\bar{g} - \bar{f}\bar{g}}{\bar{g}\bar{g} - \bar{g}^2} \\ a^* = \bar{f} - b^* \bar{g} \end{cases}$$

$$\text{NQTL} \left\{ \begin{array}{l} c^* = \bar{f} \end{array} \right.$$

$$\frac{\delta L^Q}{\delta a} = \frac{N}{\sigma^2} (\bar{f} - a - b\bar{g})$$

$$\frac{\delta L^Q}{\delta b} = \frac{N}{\sigma^2} (\bar{f}\bar{g} - a\bar{g} - b\bar{g}\bar{g})$$

$$\frac{\delta L^N}{\delta c} = \frac{N}{\sigma^2} (\bar{f} - c)$$

$$\frac{\delta^2 L^Q}{\delta a^2} = -\frac{N}{\sigma^2}$$

$$\frac{\delta^2 L^Q}{\delta a \delta b} = -\frac{N}{\sigma^2} \bar{g}$$

$$\frac{\delta^2 L^Q}{\delta b^2} = -\frac{N}{\sigma^2} \bar{g}\bar{g}$$

$$\frac{\delta^2 L^N}{\delta c^2} = -\frac{N}{\sigma^2}$$

$$P(D|a \ b \ Q_{TL}) \approx P(D|a^* \ b^* \ Q_{TL})$$

$$\times \exp \left[ -\frac{N}{2\sigma_a^2} (a-a^*)^2 - \frac{N}{2\sigma_b^2} (b-b^*)^2 - \frac{N}{\sigma_c^2} \bar{f} (a-a^*) (b-b^*) \right]$$

$$P(D|c \ N_{TL}) \approx P(D|c^* \ N_{TL})$$

$$\exp \left[ -\frac{1}{2} \frac{N}{\sigma_c^2} (c-c^*)^2 \right]$$

$$P(D|a \ b \ Q_{TL}) \approx P(D|a^* \ b^* \ Q_{TL}) \exp \left[ \bar{c} \begin{bmatrix} 1 & \bar{f} \\ -\bar{f} & 1 \end{bmatrix} \begin{bmatrix} a-a^* \\ b-b^* \end{bmatrix} \right] \frac{N}{\sigma_c^2}$$

$$P(D|Q_{TL}) = P(D|a^* \ b^* \ Q_{TL}) \cdot \int_{-\infty}^{\infty} da \int_{-\infty}^{\infty} db \ e^{-\frac{1}{2} (a-a^*, b-b^*) \mathcal{Q} \begin{pmatrix} a-a^* \\ b-b^* \end{pmatrix}} \times \frac{1}{\sigma_a} \frac{1}{\sigma_b}$$

$$\mathcal{Q} = \frac{N}{\sigma_c^2} \begin{bmatrix} 1 & \bar{f} \\ \bar{f} & 1 \end{bmatrix} \approx P(D|a^* \ b^* \ Q_{TL}) \frac{1}{\sigma_a} \frac{1}{\sigma_b} \frac{2\pi}{\sqrt{\det \mathcal{Q}}}$$

$$\det \mathcal{Q} = \left( \frac{N}{\sigma_c^2} \right)^2 (\bar{f}\bar{f} - \bar{f}^2) \approx P(D|a^* \ b^* \ Q_{TL}) \frac{2\pi\sigma_c^2}{\sigma_a\sigma_b} \frac{1/N}{\sqrt{\bar{f}\bar{f} - \bar{f}^2}}$$

$$\approx P(D|a^* \ b^* \ Q_{TL}) \frac{2\pi/N}{\sqrt{\bar{f}\bar{f} - \bar{f}^2}} \cdot \frac{\sigma_c}{\sigma_a} \frac{\sigma_c}{\sigma_b}$$

$$P(D|c \ N_{TL}) \approx P(D|c^* \ N_{TL}) \frac{1}{\sigma_c} \int_{-\infty}^{\infty} e^{-\frac{1}{2} \frac{N}{\sigma_c^2} (c-c^*)^2} dc$$

$$\approx P(D|c^* \ N_{TL}) \frac{\sqrt{2\pi} \sigma_c / \sqrt{N}}{\sigma_c}$$

$$\frac{P(D|Q_{TL})}{P(D|N_{TL})} = \underbrace{\frac{P(D|a^* \ b^* \ Q_{TL})}{P(D|c^* \ N_{TL})}}_{n_L} \underbrace{\left[ \frac{\sqrt{2\pi/N}}{\sqrt{\bar{f}\bar{f} - \bar{f}^2}} \frac{\sigma_c/\sigma_a \ \sigma_c/\sigma_b}{\sigma_c} \right]}_{\text{Occam's razor}}$$

## LOD scores

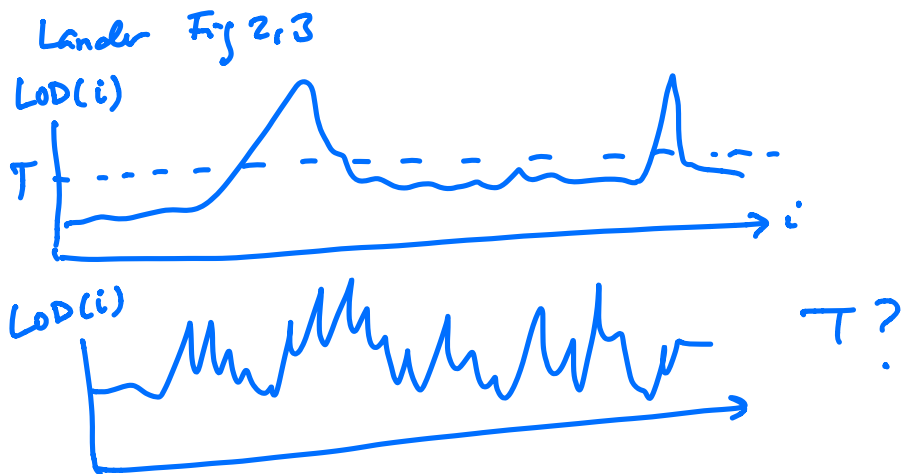
Dry et al were not bayesian's, but they used a method similar to Lander's.

$$\begin{aligned} \text{LOD} &= \log \frac{P(D|a^*b^* \& T_L)}{P(D|c^* N \& T_L)} \\ &= \log P(D|a^*b^* \& T_L) - \log P(D|c^* N \& T_L) \geq 0 \end{aligned}$$

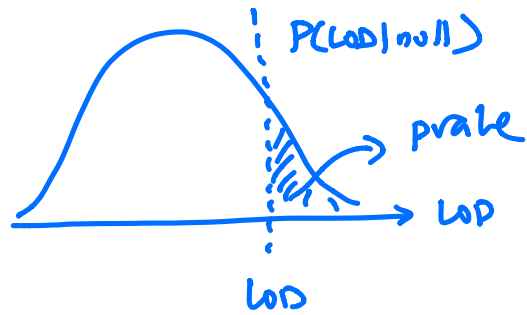
nested hypotheses

"Important issue"

"What LOD threshold  $T$  should be used in order to maintain an acceptable low positive rate?"



+ Take a collection of pairs  $(f_i, s_i)$  known to be non-associated -  $H_0$



$$\text{expected \#FP} \approx P(\text{LOD} > \text{LOD}(i) | \text{null}) \times N$$

↑  
number of loci tested

## Generalization

1 phenotype - 1 loci  $f_i = a + b g_i$

1- phenotype -  $K$  loci

$$f_i = a + \sum_{k=1}^K b^k g_i^k + \epsilon_i$$

## Assumptions

1- linear fit btw  $f_i \sim g_i^k$

2- Gaussian noise

3- all loci/ have same noise  
all individuals

4-  $\sigma$  is known.