

# More Linear Regression

March 7, 2019

Data Science CSCI 1951A

Brown University

Instructor: Ellie Pavlick

HTAs: Wennie Zhang, Maulik Dang, Gurnaaz Kaur

# Announcements

- Blog post and MR due soon (3/15)
- Stats assignment coming out today
- Anything? Questions? Comments? Concerns?

# Today

- Finish last lecture—derivation of slope/intercept
- Linear Regression (Part 2) — interpretation, controls, dummy variables

# Background

Sample Average:  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

Sample Variance:  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$

Sample Co-Variance:  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$

# Linear Regression

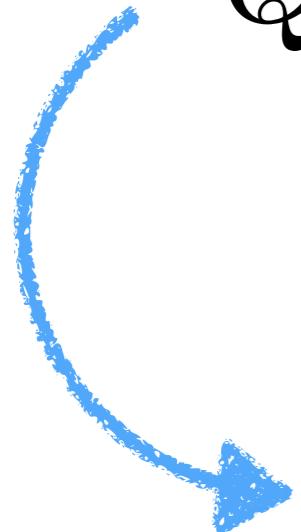
$$Q = \sum_{i=1}^n (Y_i - \hat{Y})^2$$

# Linear Regression

$$Q = \sum_{i=1}^n (Y_i - (mX_i + b))^2$$

# Linear Regression

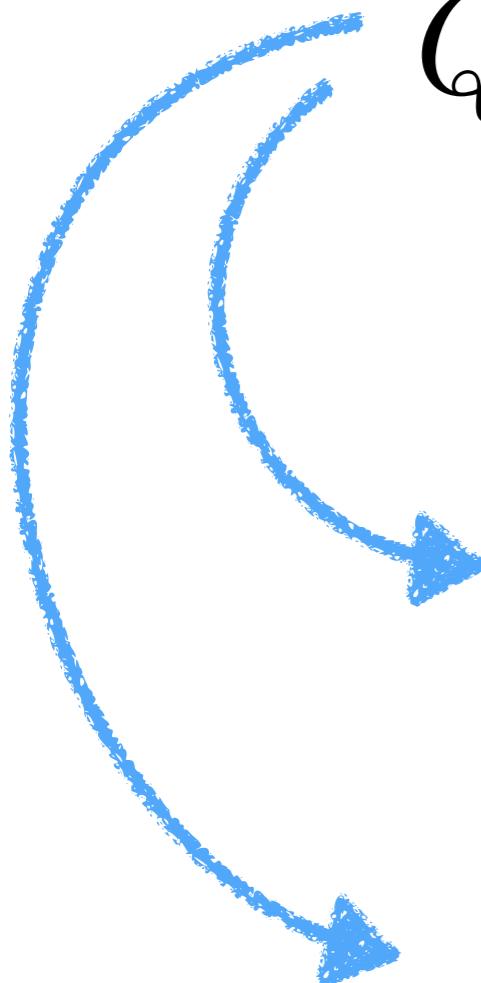
$$Q = \sum_{i=1}^n (Y_i - (mX_i + b))^2$$



$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

intercept at minimum

# Linear Regression

$$Q = \sum_{i=1}^n (Y_i - (mX_i + b))^2$$

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

intercept at minimum

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

slope at minimum

# Linear Regression

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

# Linear Regression

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

$$2\left(nb + m \sum_{i=1}^n X_i - \sum_{i=1}^n Y_i\right) = 0$$

# Linear Regression

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

$$2\left(nb + m \sum_{i=1}^n X_i - \sum_{i=1}^n Y_i\right) = 0$$

$$b = \frac{1}{n} \left( \sum_{i=1}^n Y_i - m \sum_{i=1}^n X_i \right)$$

# Linear Regression

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

$$2\left(nb + m \sum_{i=1}^n X_i - \sum_{i=1}^n Y_i\right) = 0$$

$$b = \frac{1}{n} \left( \sum_{i=1}^n Y_i - m \sum_{i=1}^n X_i \right)$$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

# Linear Regression

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

$$2\left(nb + m \sum_{i=1}^n X_i - \sum_{i=1}^n Y_i\right) = 0$$

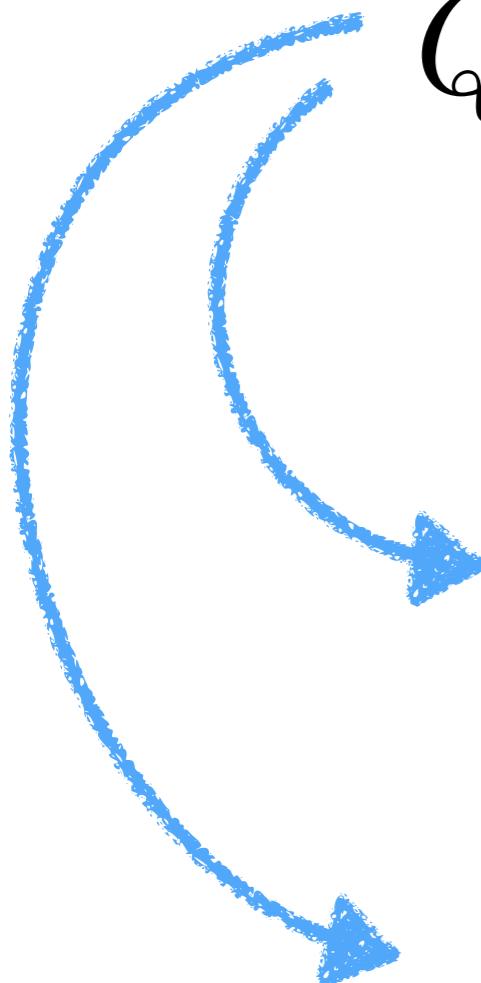
$$b = \frac{1}{n} \left( \sum_{i=1}^n Y_i - m \sum_{i=1}^n X_i \right)$$

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

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

# Linear Regression

$$Q = \sum_{i=1}^n (Y_i - (mX_i + b))^2$$


*intercept at minimum*

$$\frac{\partial Q}{\partial b} = \sum_{i=1}^n -2(Y_i - mX_i - b) = 0$$

*slope at minimum*

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

# Linear Regression

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

# Linear Regression

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

$$\sum_{i=1}^n -2(Y_iX_i - bX_i - mX_i^2) = 0$$

# Linear Regression

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

$$\sum_{i=1}^n -2(Y_iX_i - bX_i - mX_i^2) = 0$$

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

# Linear Regression

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

$$\sum_{i=1}^n -2(Y_iX_i - bX_i - mX_i^2) = 0$$

$$b = \bar{Y} - m\bar{X}$$
$$\sum_{i=1}^n -2(Y_iX_i - \bar{Y}X_i + m\bar{X}X_i - mX_i^2) = 0$$

Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

# Linear Regression

$$\frac{\partial Q}{\partial m} = \sum_{i=1}^n -2X_i(Y_i - b - mX_i) = 0$$

$$b = \bar{Y} - m\bar{X}$$
$$\sum_{i=1}^n -2(Y_iX_i - \bar{Y}X_i + m\bar{X}X_i - mX_i^2) = 0$$

$$\sum_{i=1}^n (Y_iX_i - \bar{Y}X_i) - m \sum_{i=1}^n X_i^2 - \bar{X}X_i = 0$$

Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

# Linear Regression

$$\sum_{i=1}^n (Y_i X_i - \bar{Y} X_i) - m \sum_{i=1}^n X_i^2 - \bar{X} X_i = 0$$

# Linear Regression

$$\sum_{i=1}^n (Y_i X_i - \bar{Y} X_i) - m \sum_{i=1}^n X_i^2 - \bar{X} X_i = 0$$

$$m = \frac{\sum_{i=1}^n (Y_i X_i - \bar{Y} X_i)}{\sum_{i=1}^n X_i^2 - \bar{X} X_i}$$

# Linear Regression

$$\sum_{i=1}^n (Y_i X_i - \bar{Y} X_i) - m \sum_{i=1}^n X_i^2 - \bar{X} X_i = 0$$

$$m = \frac{\sum_{i=1}^n (Y_i X_i - \bar{Y} X_i)}{\sum_{i=1}^n X_i^2 - \bar{X} X_i}$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

# Linear Regression

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$\sum_{i=1}^n \bar{X}\bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

$$m = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

$$m = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

$$m = \frac{Cov(X, Y)}{Var(X)}$$

Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

$$m = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

$$m = \frac{Cov(X, Y)}{Var(X)}$$

Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

Sums/Averages  
of Random  
Variables

$$m = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

$$m = \frac{Cov(X, Y)}{Var(X)}$$

Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

# Linear Regression

$$\sum_{i=1}^n \bar{X}^2 - X_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i) - n \bar{X} \bar{Y}}{\sum_{i=1}^n (X_i^2) - n \bar{X}^2}$$

$$\sum_{i=1}^n \bar{X} \bar{Y} - Y_i \bar{X} = 0$$

$$m = \frac{\sum_{i=1}^n (X_i Y_i - X_i \bar{Y}) + \sum_{i=1}^n (\bar{X} \bar{Y} - Y_i \bar{X})}{\sum_{i=1}^n (X_i^2 - X_i \bar{X}) + \sum_{i=1}^n (\bar{X}^2 - X_i \bar{X})}$$

CTL applies,  
can compute  
expected value,  
p-values

$$m = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

$$m = \frac{Cov(X, Y)}{Var(X)}$$

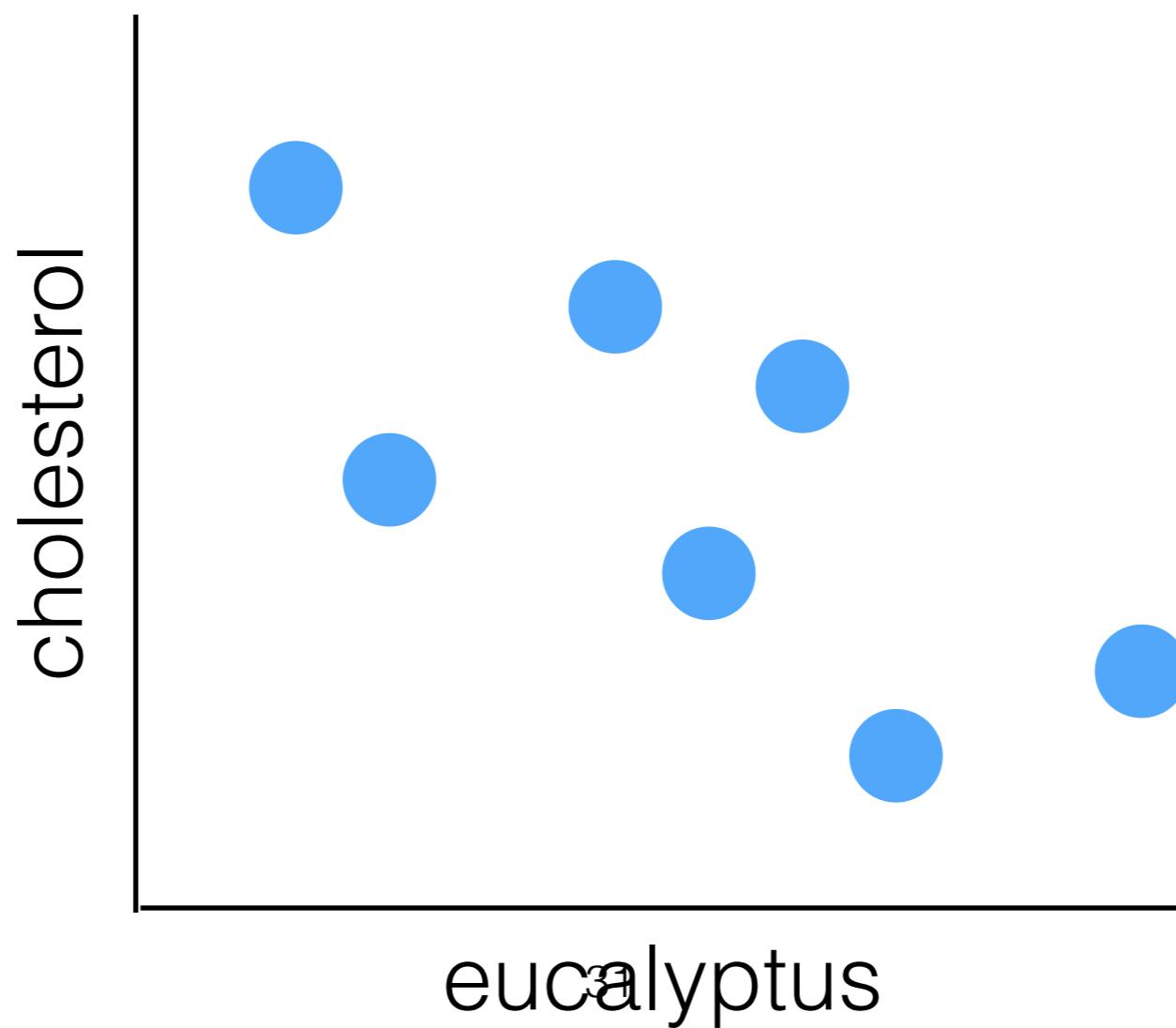
Data Analysis Toolkit #10: Simple linear regression

Copyright © 1996, 2001 Prof. James Kirchner

[http://seismo.berkeley.edu/~kirchner/eps\\_120/Toolkits/Toolkit\\_10.pdf](http://seismo.berkeley.edu/~kirchner/eps_120/Toolkits/Toolkit_10.pdf)

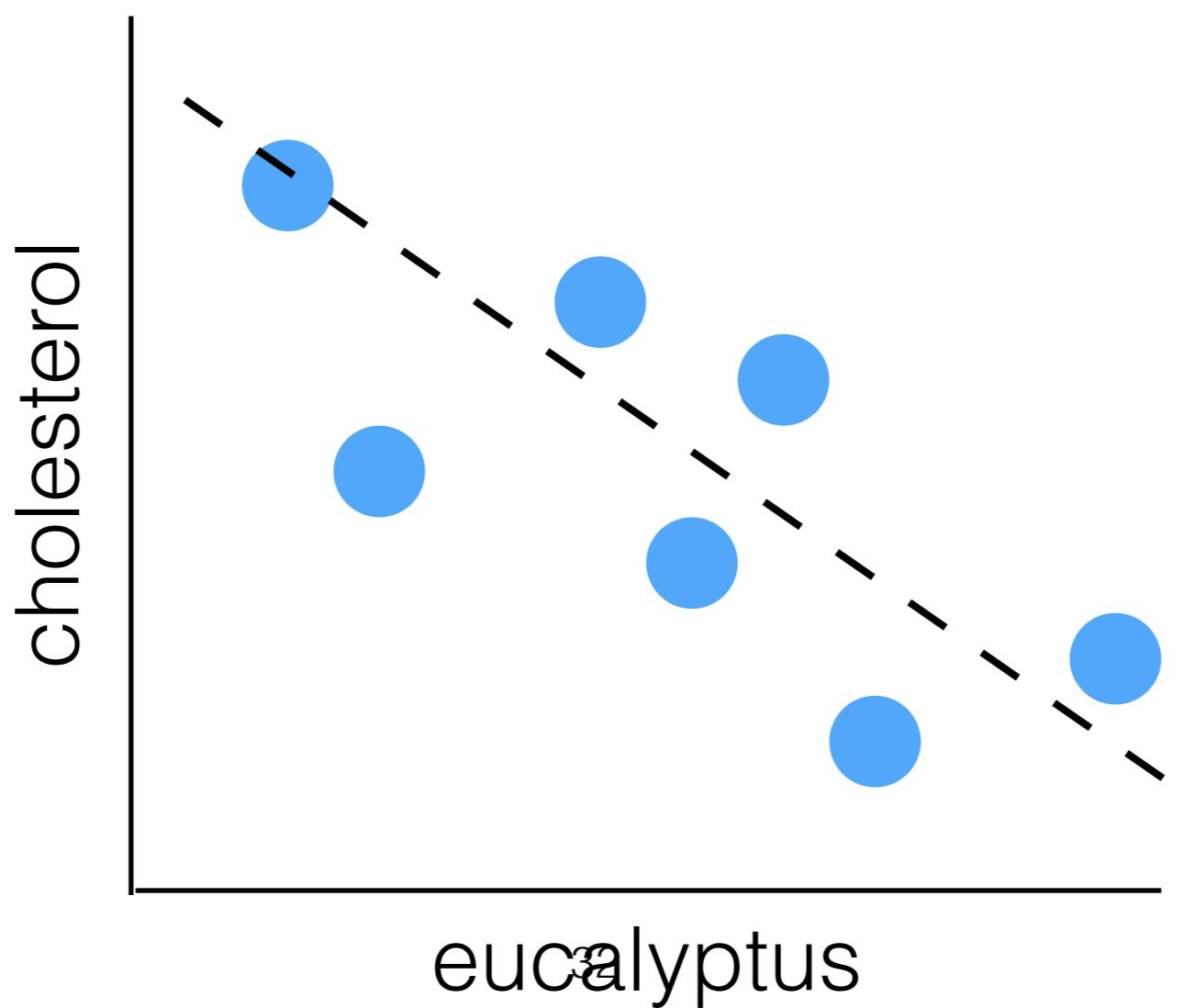
# Linear Regression

cholesterol = m(eucalyptus) + b



# Linear Regression

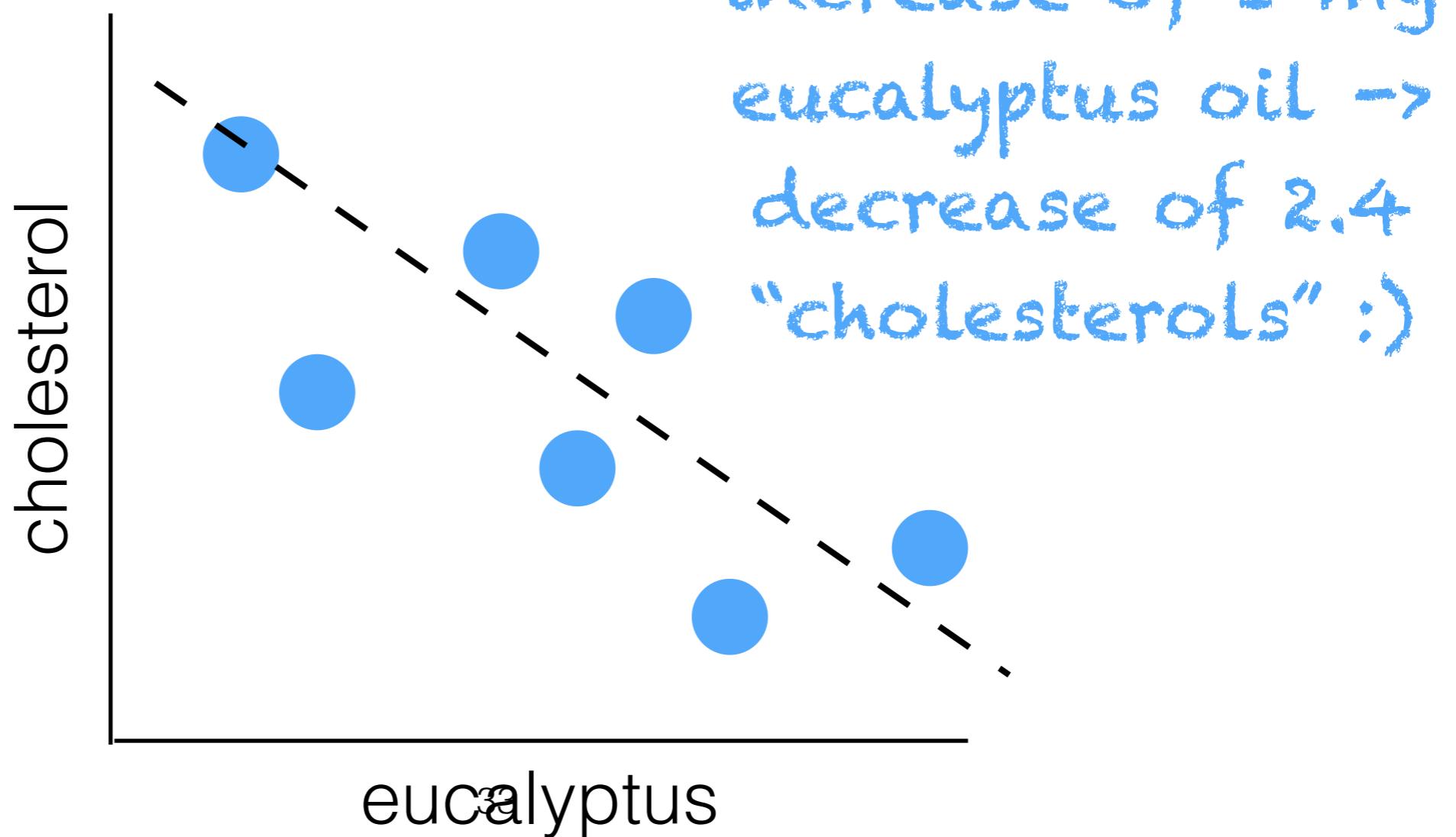
$$\text{cholesterol} = m(\text{eucalyptus}) + b$$
$$m = -2.4$$



# Linear Regression

$$\text{cholesterol} = m(\text{eucalyptus}) + b$$

$$m = -2.4$$

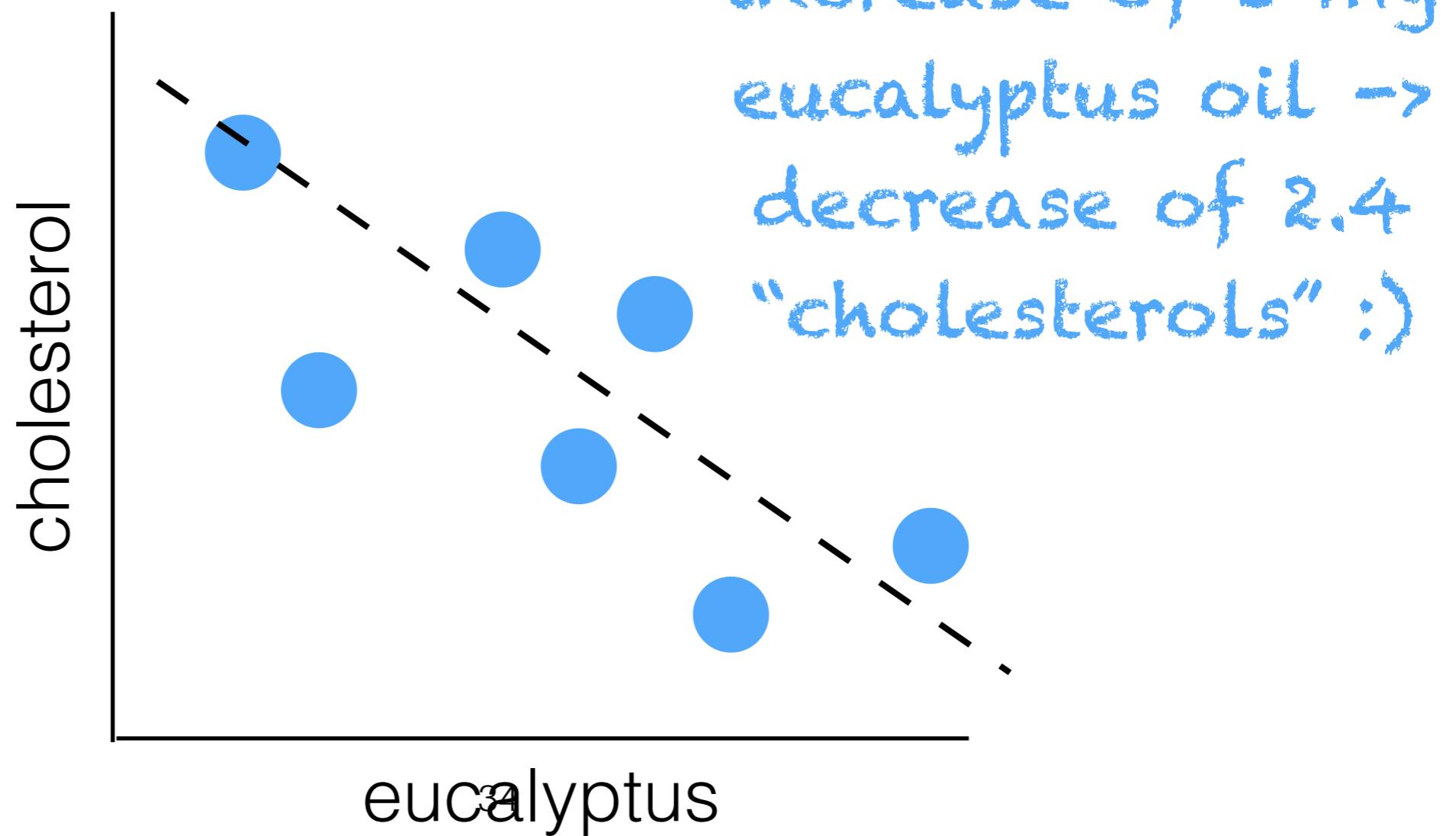


# Linear Regression

$$\text{cholesterol} = m(\text{eucalyptus}) + b$$

$$m = -2.4$$

¿que  
pasa?



# Discussion-question-thinly-veiled-as-a-clicker Question!

# Discussion-question-thinly-veiled-as-a-clicker Question!

What should we make of the observed relationship between use of eucalyptus oil and cholesterol levels?

- (a) There probably actually is a relationship. Linear regression is a legitimate method, so we should trust the result.
- (b) There is probably no actual relationship. We are confusing correlation with causation.
- (c) There is probably no actual relationship. We are measuring eucalyptus oil in the wrong units, so it just appears correlated.
- (d) There is probably no actual relationship. We are failing to capture other relevant variables.
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

# Discussion-question-thinly-veiled-as-a-clicker Question!

What should we make of the observed relationship between use of eucalyptus oil and cholesterol levels?

- (a) ~~There probably actually is a relationship. Linear regression is a legitimate method, so we should trust the result.~~
- (b) There is probably no actual relationship. We are confusing correlation
- (c) There is probably a relationship. Eucalyptus oil could be the case, but we want to do due diligence before concluding this... Measuring cholesterol levels is correlated.
- (d) There is probably a relationship. We have not captured all the variables involved in cholesterol levels before concluding this...
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

## Difference-in-difference-thinly-voiled-as-a-

Yes and no. We **\*are\*** confusing correlation with causation, but linear regression does this by construction (even when we are looking at a real relationship).

- (a) There is probably a causal relationship. Linear regression is a legitimate method, so we should trust the result.
- (b) There is probably no actual relationship. We are confusing correlation with causation.
- (c) There is probably no actual relationship. We are measuring eucalyptus oil in the wrong units, so it just appears correlated.
- (d) There is probably no actual relationship. We are failing to capture other relevant variables.
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

# Discussion-question-thinly-veiled-as-a-clicker Question!

What should we make of the observed relationship between use of eucalyptus oil and

Units should not matter, since differences in units are usually equivalent up to linear transformation

- (a) There probably is a legitimate relationship.
- (b) There is probably a correlation.
- (c) There is probably no actual relationship. We are measuring eucalyptus oil in the wrong units, so it just appears correlated.
- (d) There is probably no actual relationship. We are failing to capture other relevant variables.
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

# Discussion-question-thinly-veiled-as-a-clicker Question!

What should we make of the observed relationship between use of eucalyptus oil and cholesterol levels?

- (a) There probably actually is a relationship. Linear regression is a legitimate method.
- (b) There is probably no correlation with car
- (c) There is probably no eucalyptus oil in the wrong units, so it just appears correlated.
- (d) There is probably no actual relationship. We are failing to capture other relevant variables.
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

This is a good answer!

Let's spend multiple more slides on it.

# Discussion-question-thinly-veiled-as-a-clicker Question!

What should we make of the observed relationship between use of eucalyptus oil and cholesterol levels?

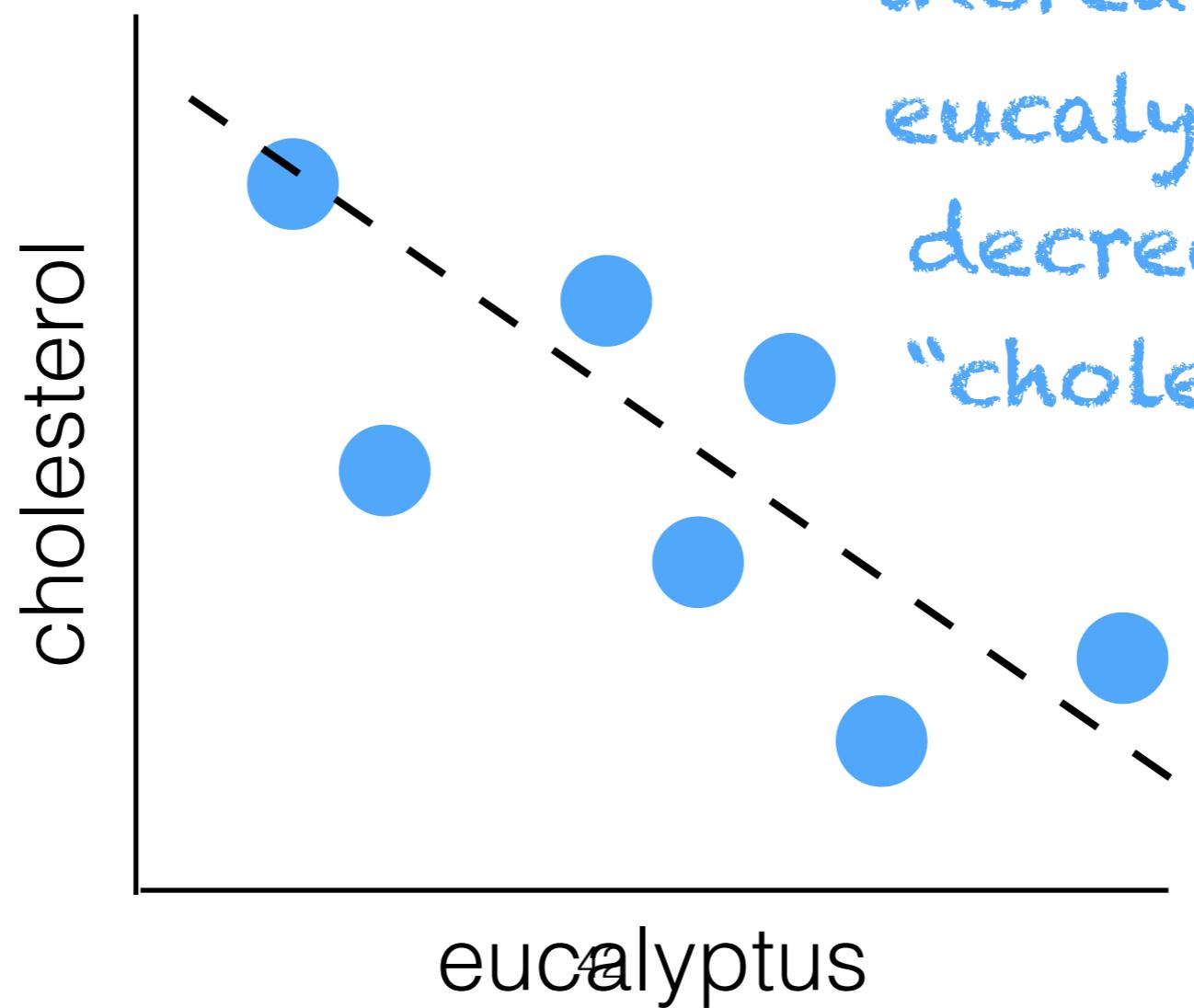
- (a) There probably actually is a relationship. Linear regression is a legitimate method, so we should trust the result.
- (b) There is probably no correlation with eucalyptus oil. Not mad, just disappointed
- (c) There is probably no eucalyptus oil in the wrong units, so it just appears correlated.
- (d) There is probably no actual relationship. We are failing to capture other relevant variables.
- (e) We should click on the obviously sharky answer and see if Ellie gets mad.

# Linear Regression

$$\text{cholesterol} = m(\text{eucalyptus}) + b$$

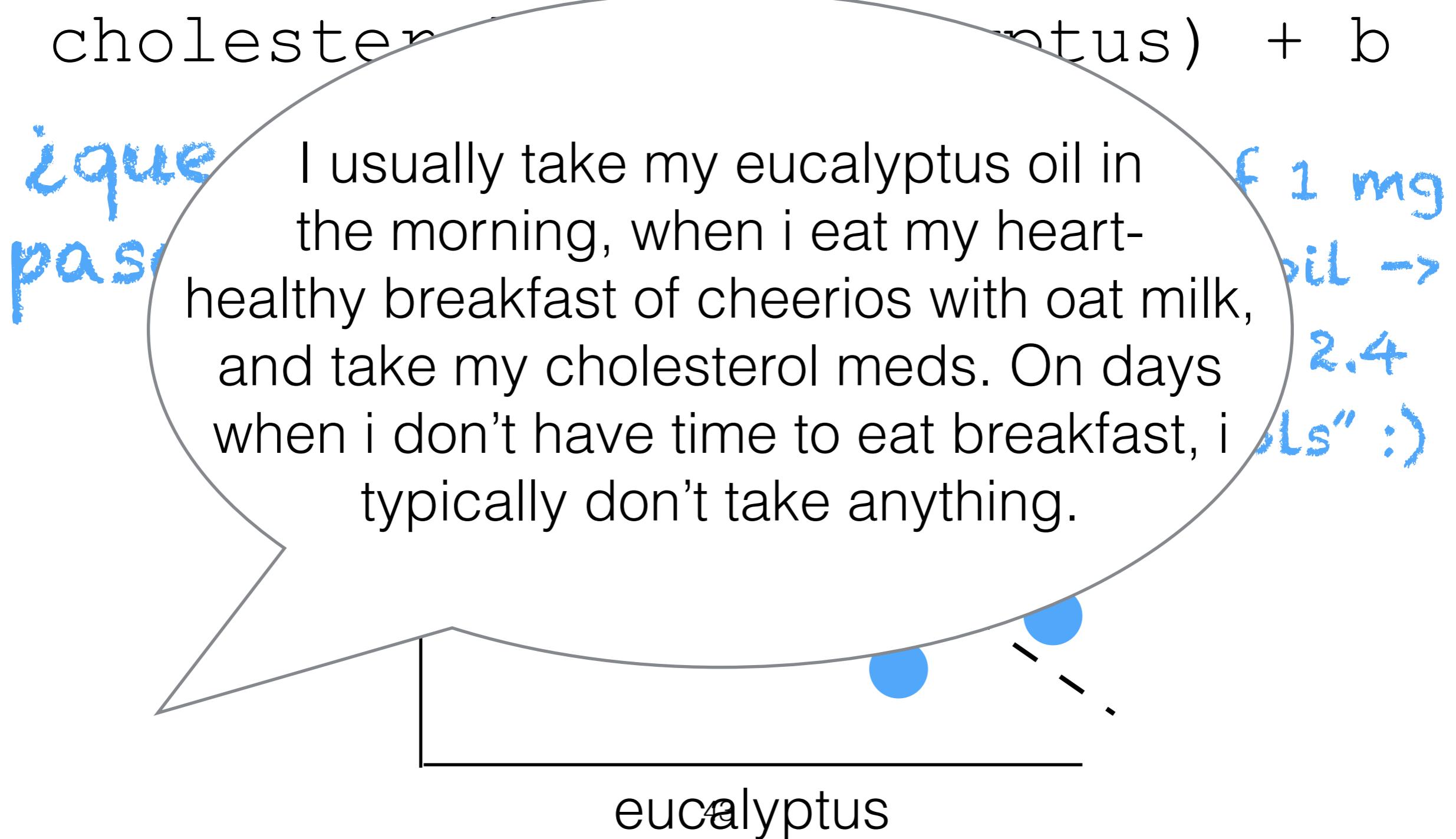
$$m = -2.4$$

¿que  
pasa?

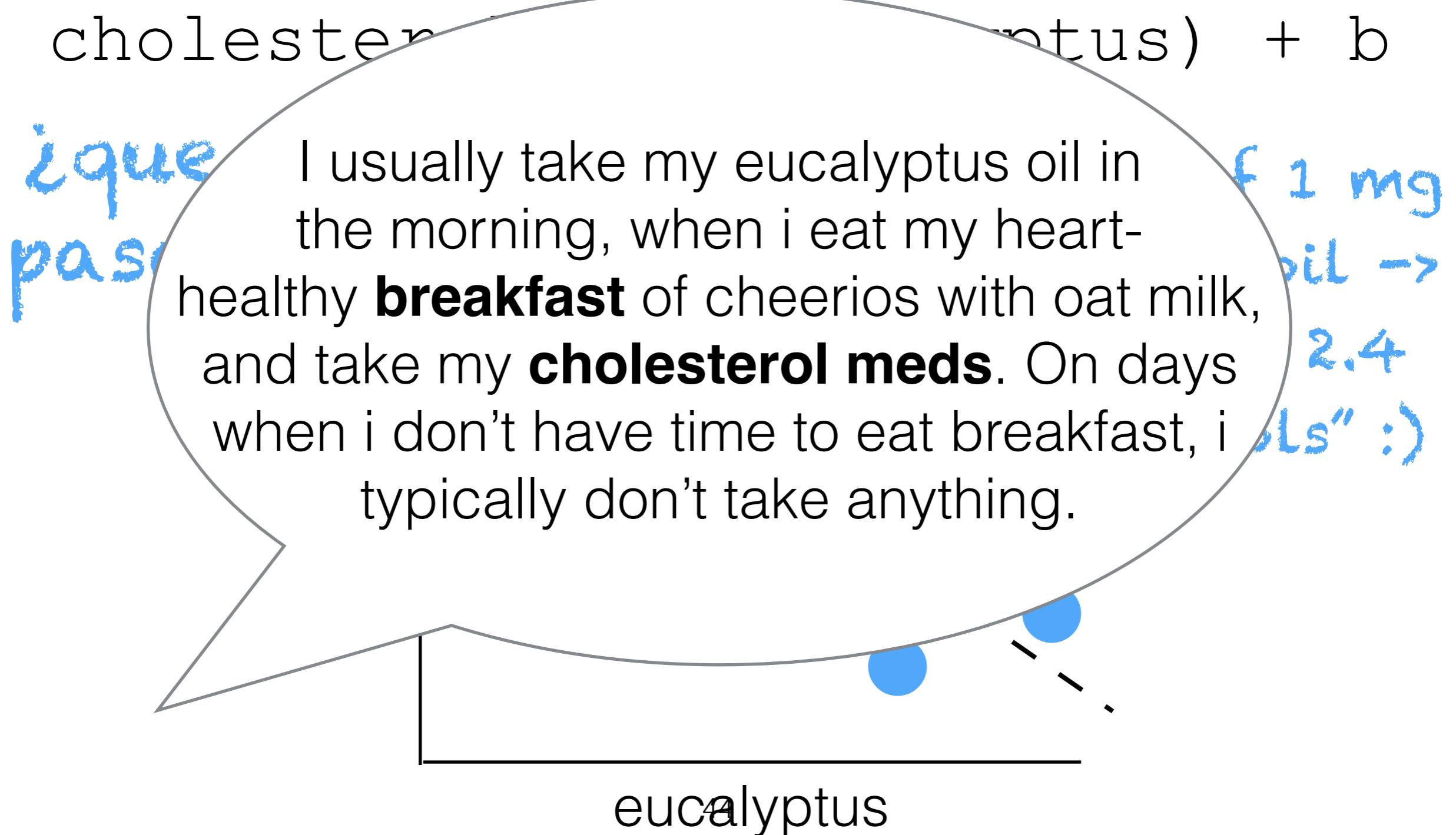


increase of 1 mg  
eucalyptus oil →  
decrease of 2.4  
“cholesterols” :)

# Linear Regression



# Linear Regression



# Omitted Variable Bias

# Omitted Variable Bias

- By construction, we assume that the dependent variable can be predicted from the explanatory variables only

# Omitted Variable Bias

- By construction, we assume that the dependent variable can be predicted from the explanatory variables only
- We assume changes in the dependent variable that are correlated with the explanatory variable are *because of* the explanatory variable

# Omitted Variable Bias

- By construction, we assume that the dependent variable can be predicted from the explanatory variables only
- We assume changes in the dependent variable that are correlated with the explanatory variable are *because of* the explanatory variable
- We assume that changes in the dependent variable that are *not* explained by the explanatory variables is “noise”

# Multiple Linear Regression

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4$$

Y: cholesterol level

X1: eucalyptus

X2: cholesterol meds

X3: breakfast

X4: constant term

# Multiple Linear Regression

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4$$

Y: cholesterol level

intercept

X1: eucalyptus

X2: cholesterol meds

X3: breakfast

X4: constant term

# Multiple Linear Regression

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4$$

Y: cholesterol level  
X1: eucalyptus  
X2: cholesterol meds  
X3: breakfast  
X4: constant term

slopes/  
coefficients/  
effects

# Multiple Linear Regression

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Q = \sum_{i=1}^n (Y_i - (m_1 X_{1i} + m_2 X_{2i} + m_3 X_{3i} + m_4 X_{4i}))^2$$

# Multiple Linear Regression

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Q = \sum_{i=1}^n (Y_i - (m_1 X_{1i} + m_2 X_{2i} + m_3 X_{3i} + m_4 X_{4i}))^2$$

$$\frac{\partial Q}{\partial m_1} = f(X_1, X_2, X_3, X_4, Y)$$

# Multiple Linear Regression

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Q = \sum_{i=1}^n (Y_i - (m_1 X_{1i} + m_2 X_{2i} + m_3 X_{3i} + m_4 X_{4i}))^2$$

depends on other  
explanatory variables

$$\frac{\partial Q}{\partial m_1} = f(X_1, X_2, X_3, X_4, Y)$$

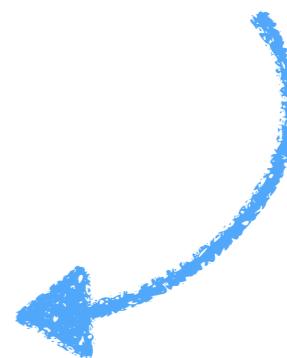
# Multiple Linear Regression

$$Y = m_1 X_1$$

change in cholesterol

$$Q = \sum_{i=1}^n (Y_i - \text{mg eucalyptus oil, holding other variables constant})$$

$$\frac{\partial Q}{\partial m_1} = f(X_1, X_2, X_3, X_4, Y)$$



# Multiple Linear Regression

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

# LinAlg Detour

$$Y = m \text{ Matrices of observations } m \times 4$$
$$Y = X\beta$$

$$\hat{\beta} = (X'X)^{-1}X'Y$$

# LinAlg Detour

$$Y = m_1 X_1$$

Vector of coefficients

$$Y = X \beta$$

$$\hat{\beta} = (X'X)^{-1} X' Y$$

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Y = X\beta$$

$$X = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$\hat{\beta} = (X'X)^{-1}X'Y$$

$$X' = \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

X Transpose  
60

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Y = X\beta$$

$$\hat{\beta} = (X'X)^{-1}X'Y$$

Inverse  
61

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$Y = X\beta$$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

doesn't  
always  
exist...

Inverse  
62

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 2 & 6 \\ 4 & 9 & 12 \end{bmatrix}$$

linearly  
dependent/  
co-linear

$$\hat{\beta} = (X'X)^{-1} X'Y$$

Inverse  
63

# LinAlg Detour

$$Y = m_1 X_1 + m_2 X_2 + m_3 X_3 + m_4 X_4$$

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 2 & 6 \\ 4 & 9 & 12 \end{bmatrix}$$

linearly  
dependent/  
co-linear

$$\hat{\beta} = (X'X)^{-1} X'Y$$

"Pseudo-Inverse"

# Dummy Variables

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4$$

Y: cholesterol level    ???  
X1: eucalyptus  
X2: cholesterol meds  
X3: breakfast  
X4: constant term

# Dummy Variables

- Used to encode qualitative features

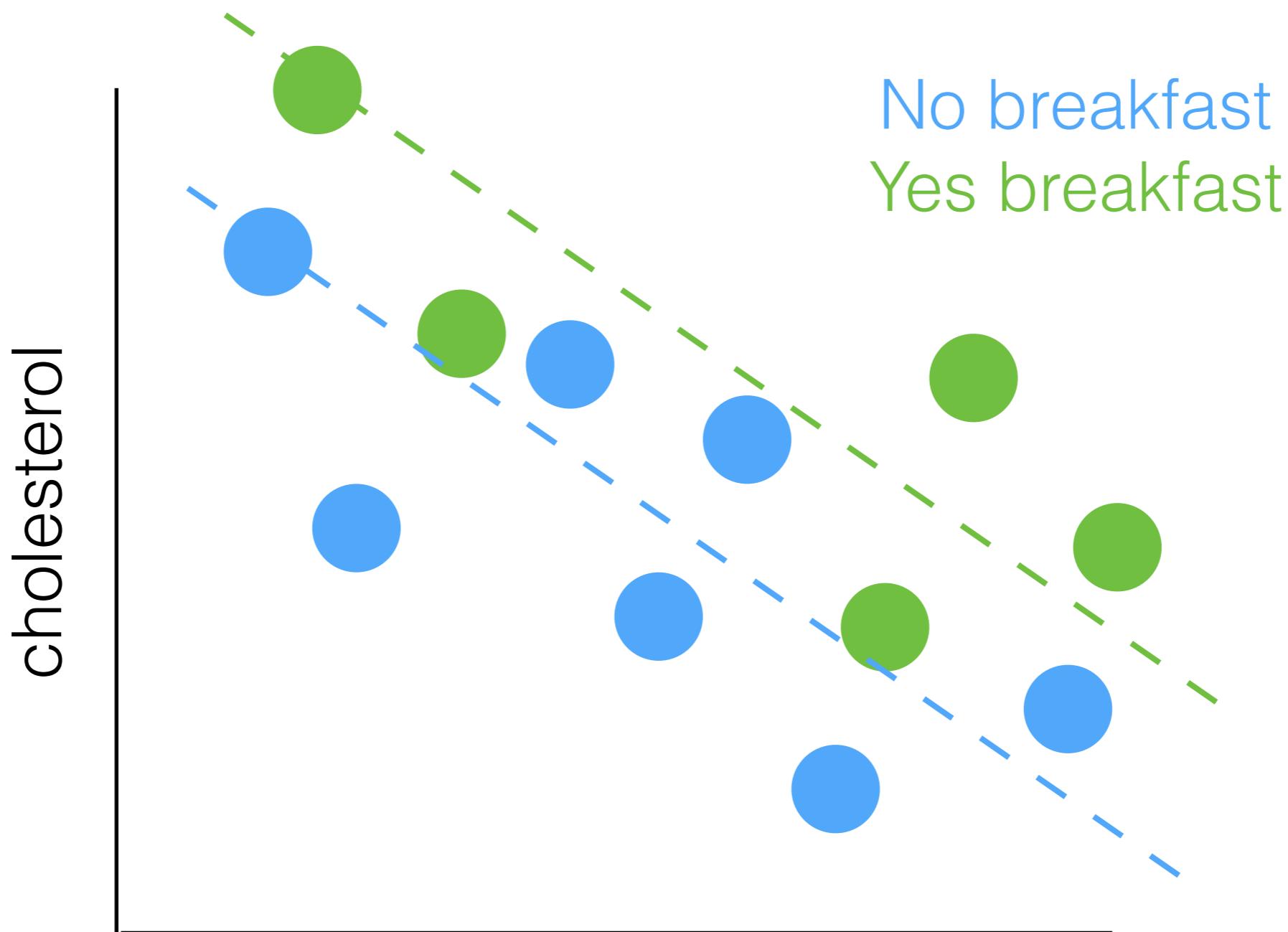
# Dummy Variables

- Used to encode qualitative features
- AKA indicator variables, Boolean variables, one-hot variables, sparse variables...

# Dummy Variables

- Used to encode qualitative features
- AKA indicator variables, Boolean variables, one-hot variables, sparse variables...
- Interpretable as shift in intercept for different groups

# Dummy Variables



# Dummy Variables

$X =$

cholesterol meds	20	31	0	1	1
	20	5	0	1	1
	20	40	0	1	1
	25	18	1	0	1

yes breakfast

eucalyptus

no breakfast

constant

# Dummy Variables

$X =$

cholesterol meds	20	31	0	1	1
	20	5	0	1	1
	20	40	0	1	1
	25	18	1	0	1

yes breakfast

eucalyptus

no breakfast

constant

Qualms?

# Dummy Variables

$X =$

cholesterol meds	eucalyptus	yes breakfast	no breakfast	constant
20	31	0	1	1
20	5	0	1	1
20	40	0	1	1
25	18	1	0	1

#!@\*\$!

linearly  
dependent

# Dummy Variables

X =

cholesterol meds	eucalyptus	yes breakfast	no breakfast	constant
20	31	0	1	1
20	5	0	1	1
20	40	0	1	1
25	18	1	0	1

"dummy  
variable  
trap"

# Dummy Variables

$X =$

cholesterol meds	yes breakfast	constant
20	31	0
20	5	0
20	40	0
25	18	1

$n-1$   
dummies  
(usually done  
for you)

eucalyptus

no breakfast

# Clicker Question!

# Clicker Question!

For the below model, how many parameters (coefficients) do we need to estimate?

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4 + m_5X_5$$

Y: happiness

X<sub>1</sub>: day of week (dummies M, T, W, Th, F, S, Su)

X<sub>2</sub>: bank account balance (real value)

X<sub>3</sub>: breakfast (dummies yes, no)

X<sub>4</sub>: whether you have found your inner peace  
(dummies yes, no, unclear)

**(a) 5**

**(b) 10**

**(c) 11**

**(d) infinite**

# Clicker Question!

For the below model, how many parameters (coefficients) do we need to estimate?

$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4 + m_5X_5$$

Y: happiness

X1: day of week (dummies M, T, W, Th, F, S, Su) **6**

X2: bank account balance (real value) **1**

X3: breakfast (dummies yes, no) **1**

X4: whether you have found your inner peace **2**  
(dummies yes, no, unclear)

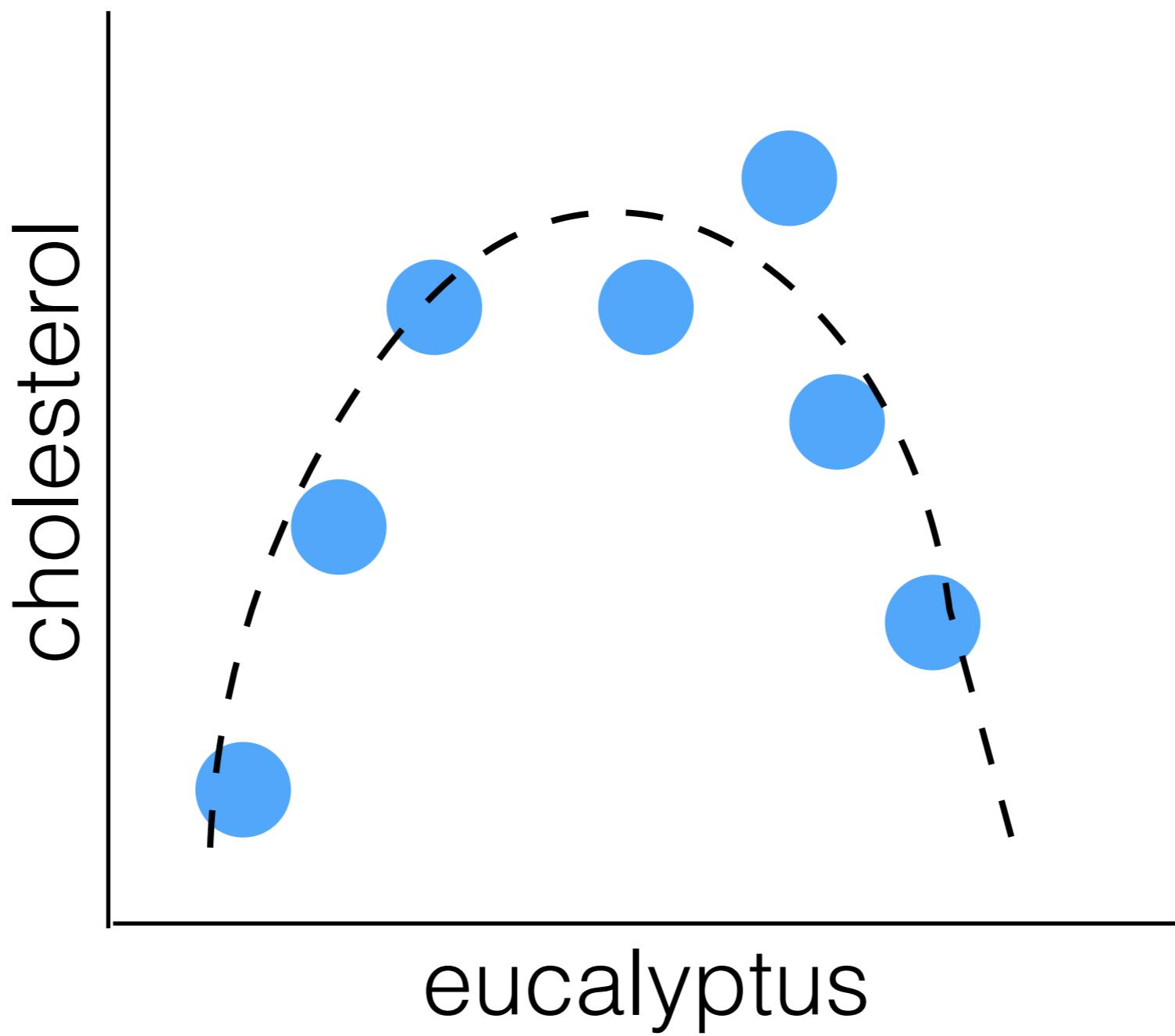
constant = 1

- (a) 5
- (b) 10

(c) 11

(d) infinite

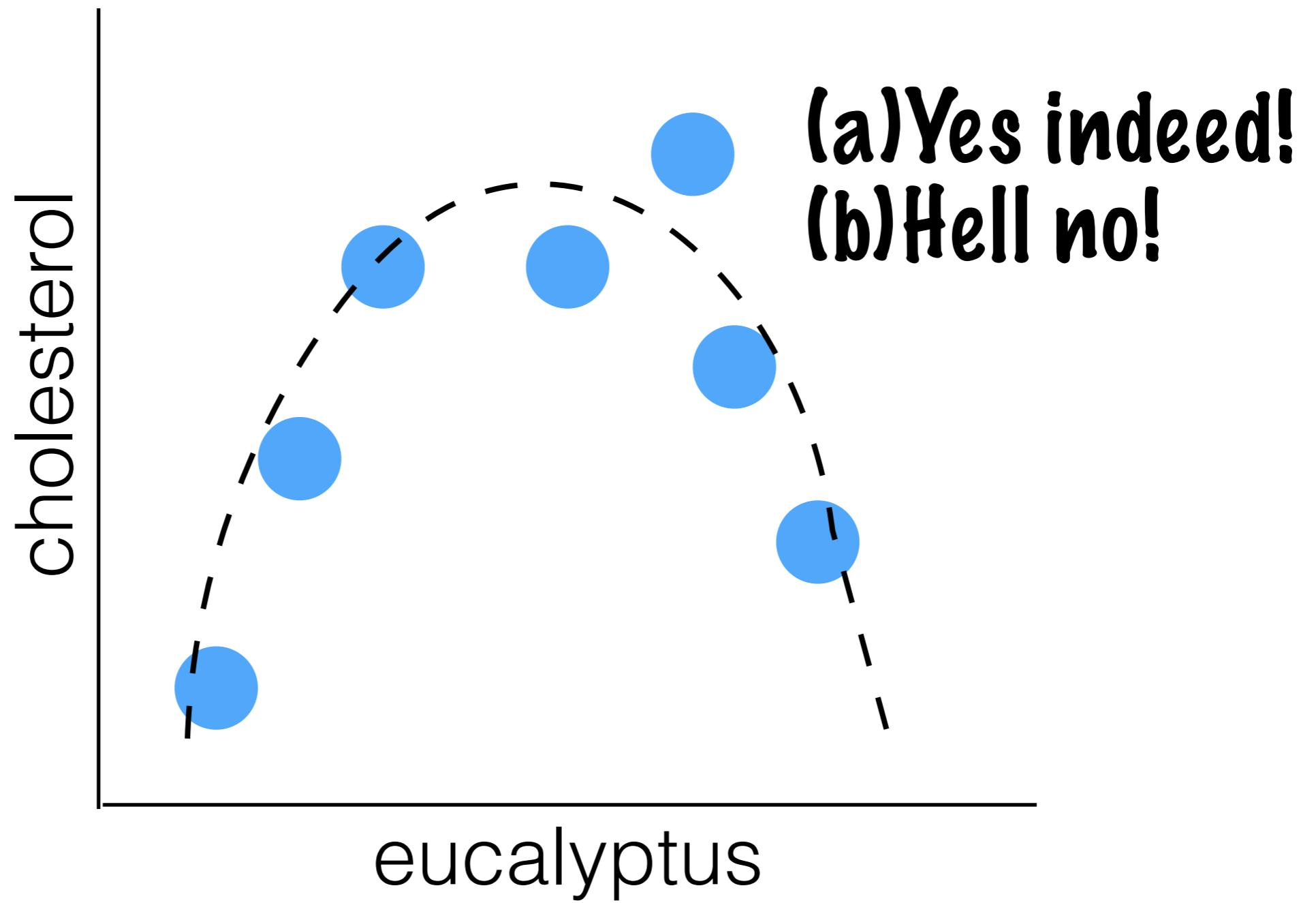
# Nonlinear Relationships



# Clicker Question!

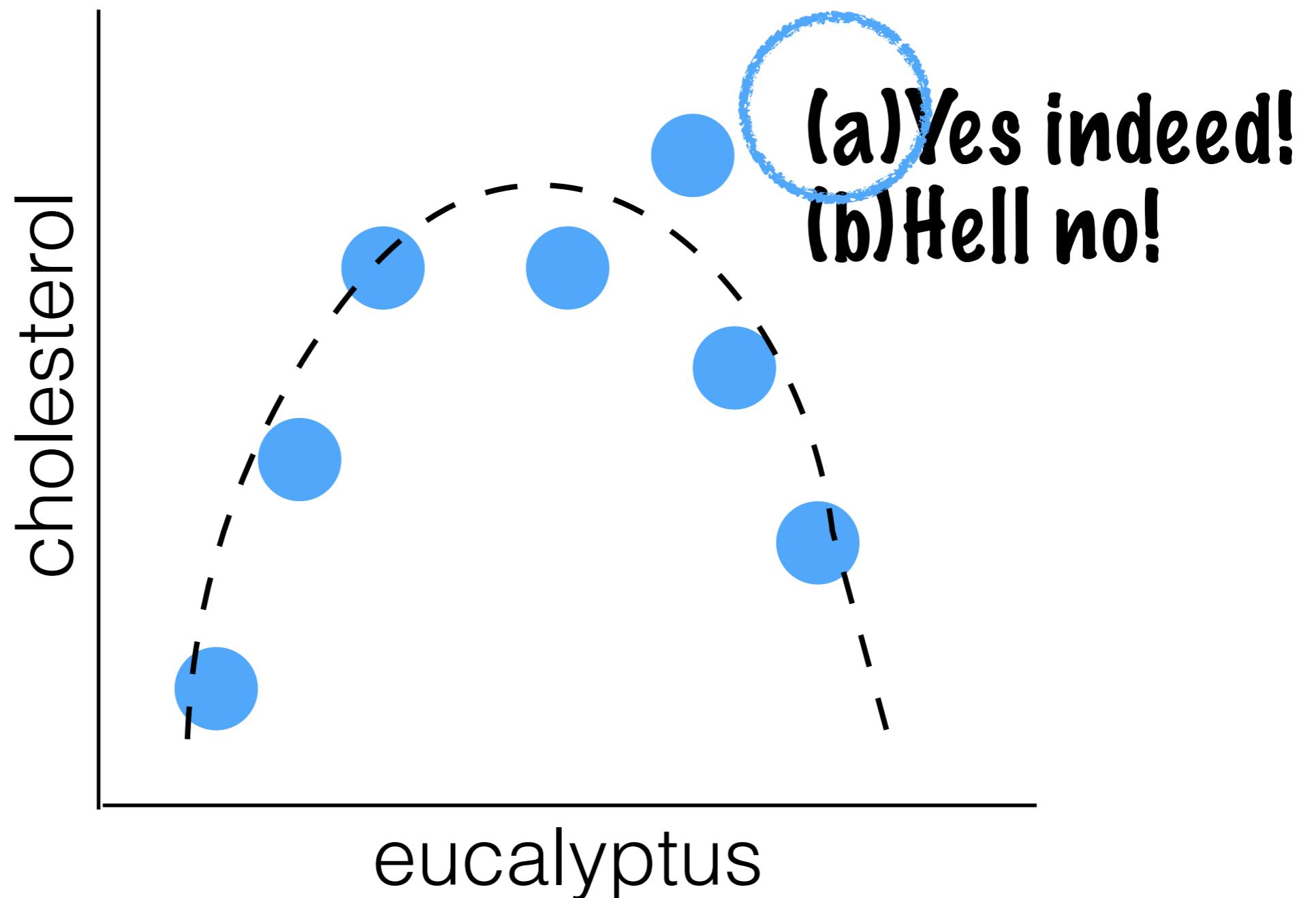
# Clicker Question!

Can we model this with linear regression?



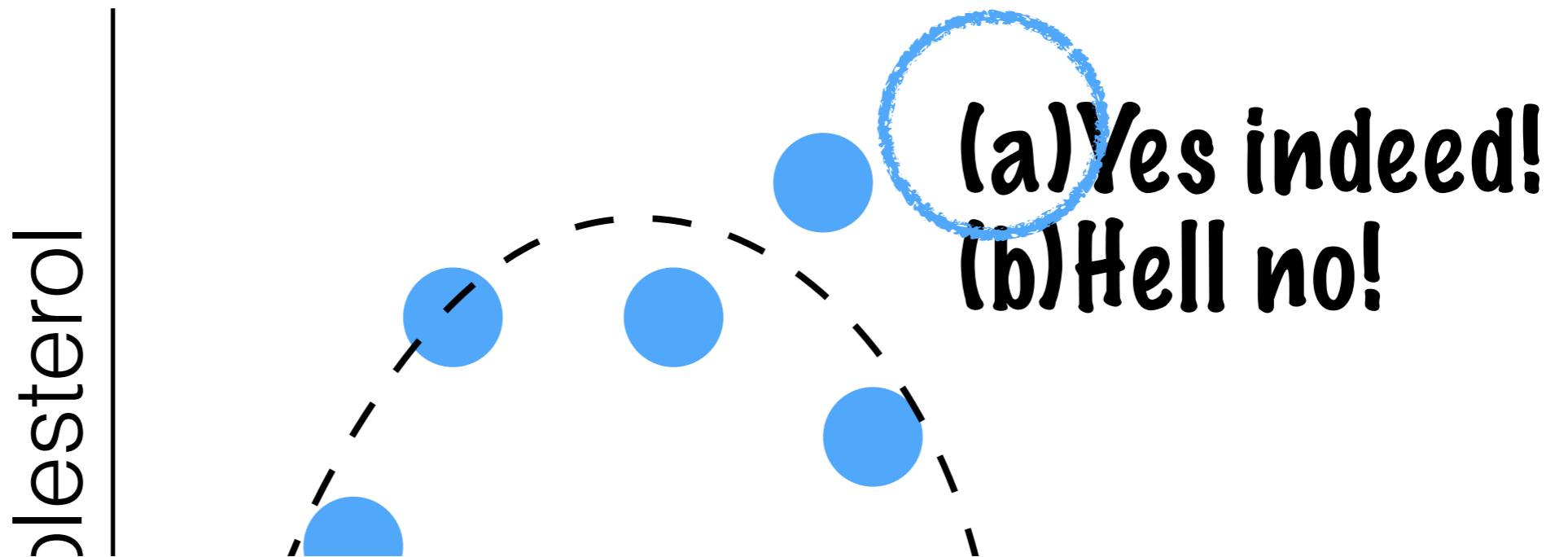
# Clicker Question!

Can we model this with linear regression?



# Clicker Question!

Can we model this with linear regression?



$$Y = m_1X_1 + m_2X_2 + m_3X_3$$

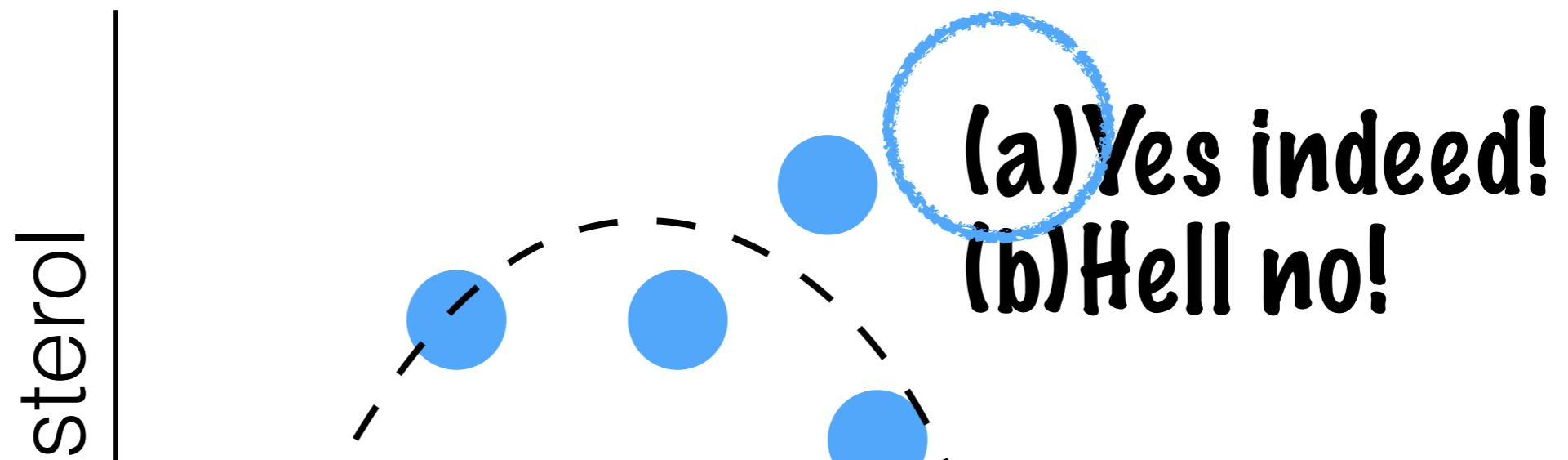
Y: cholesterol

X<sub>1</sub>: eucalyptus

X<sub>2</sub>: eucalyptus<sup>2</sup>

# Clicker Question!

Can we model this with linear regression?



$$Y = m_1X_1 + m_2X_2 + m_3X_3 + m_4X_4$$

Y: cholesterol

“interaction term”

X1: eucalyptus

X2: cholesterol meds

X3: X1  $\times$  X2

# statsmodels

```
import statsmodels.api as sm

y, X = read_data()
X = sm.add_constant(X)
model = sm.OLS(y, X)
results = model.fit()
print(results.summary())
```

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

```
import statsmodels.api as sm
import statsmodels.formula.api as smf
# M has column headers w/ names
M = read_data()
X = sm.add_constant(X)
eq = "chol ~ eucalyptus + meds + breakfast"
model = smf.ols(formula=eq, data=M)
results = model.fit()
print(results.summary())
```

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

```
import statsmodels.api as sm
import statsmodels.formula.api as smf
# M has column headers w/ names
M = read_data()
X = sm.add_constant(X)          interaction term
eq = "chol ~ eucalyptus + meds + breakfast
+ eucalyptus:meds"
model = smf.ols(formula=eq, data=M)
results = model.fit()
print(results.summary())
```

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

```
import statsmodels.api as sm
import statsmodels.formula.api as smf
# M has column headers w/ names
M = read_data()
X = sm.add_constant(X)      squared terms
eq = "chol ~ eucalyptus + meds + breakfast
+ eucalyptus^2"
model = smf.ols(formula=eq, data=M)
results = model.fit()
print(results.summary())
```

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

## OLS Regression Results

```
=====
Dep. Variable:                      y      R-squared:                 1.000
Model:                          OLS      Adj. R-squared:            1.000
Method:                         Least Squares      F-statistic:             4.020e+06
Date:                Tue, 26 Feb 2019      Prob (F-statistic):        2.83e-239
Time:                  04:42:47      Log-Likelihood:           -146.51
No. Observations:                  100      AIC:                     299.0
Df Residuals:                      97      BIC:                     306.8
Df Model:                           2
Covariance Type:                nonrobust
=====
```

	coef	std err	t	P> t	[ 0.025	0.975 ]
const	1.3423	0.313	4.292	0.000	0.722	1.963
x1	-0.0402	0.145	-0.278	0.781	-0.327	0.247
x2	10.0103	0.014	715.745	0.000	9.982	10.038

<=====

Omnibus:	2.042	Durbin-Watson:	2.274
Prob(Omnibus):	0.360	Jarque-Bera (JB):	1.875
Skew:	0.234	Prob(JB):	0.392
Kurtosis:	2.519	Cond. No.	144.

<=====

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

OLS Regression Results						
Dep. Variable:	y	R-squared:	1.000			
Model:	OLS	Adj. R-squared:	1.000			
Method:	Least Squares	F-stat:				0.020e+06
Date:	Tue, 26 Feb 2019	Prob > F:	8.3e-239			
Time:	04:42:47	Log-likelihood:	-146.51			
No. Observations:	100	AIC:	299.0			
Df Residuals:	97	BIC:	model (SSE)	306.8		
Df Model:	2					
Covariance Type:	nonrobust					
-----						
	coef	std err	t	P> t	[ 0.025	0.975 ]
-----						
const	1.3423	0.313	4.292	0.000	0.722	1.963
x1	-0.0402	0.145	-0.278	0.781	-0.327	0.247
x2	10.0103	0.014	715.745	0.000	9.982	10.038
-----						
Omnibus:	2.042	Durbin-Watson:	2.274			
Prob(Omnibus):	0.360	Jarque-Bera (JB):	1.875			
Skew:	0.234	Prob(JB):	0.392			
Kurtosis:	2.519	Cond. No.	144.			
-----						

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

## OLS Regression Results

Dep. Variable:	y	R-squared:	1.000
Model:	OLS	Adj. R-squared:	1.000
Method:	Least Squares	F-statistic:	4.020e+06
Date:	Tue, 26 Feb 2019	Prob (F-statistic):	2.83e-239
Time:	04:42:47	Log-Likelihood:	-146.51
	100	AIC:	299.0
	97	BIC:	306.8
	2		
	bust		

coefficients  
(i.e. effect sizes)

	coef	std err	t	P> t	[ 0.025	0.975 ]
const	1.3423	0.313	4.292	0.000	0.722	1.963
x1	-0.0402	0.145	-0.278	0.781	-0.327	0.247
x2	10.0103	0.014	715.745	0.000	9.982	10.038

Omnibus:	2.042	Durbin-Watson:	2.274
Prob(Omnibus):	0.360	Jarque-Bera (JB):	1.875
Skew:	0.234	Prob(JB):	0.392
Kurtosis:	2.519	Cond. No.	144.

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# statsmodels

## OLS Regression Results

Dep. Variable:	y	R-squared:	1.000
Model:	OLS	Adj. R-squared:	1.000
Method:	Least Squares	F-statistic:	4.020e+06
Date:	Tue, 26 Feb 2019	Prob (F-statistic):	2.83e-239
Time:	04:42:47	Log-Likelihood:	-146.51
No. Observations:	100	AIC:	292.0
Df Residuals:	97	BIC:	
Df Model:	2		
Covariance Type:	nonrobust		

p-values

	coef	std err	t	P> t	[ 0.025	0.975 ]
const	1.3423	0.313	4.292	0.000	0.722	1.963
x1	-0.0402	0.145	-0.278	0.781	-0.327	0.247
x2	10.0103	0.014	715.745	0.000	9.982	10.038

Omnibus:	2.042	Durbin-Watson:	2.274
Prob(Omnibus):	0.360	Jarque-Bera (JB):	1.875
Skew:	0.234	Prob(JB):	0.392
Kurtosis:	2.519	Cond. No.	144.

<https://www.statsmodels.org/dev/examples/notebooks/generated/ols.html>

[https://www.statsmodels.org/dev/generated/statsmodels.regression.linear\\_model.OLS.html](https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html)

# Discussion Question!

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

Going to college corresponds to a increase of \$20K  
in salary, assuming other variables are fixed.

# Clicker Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

```
income: salary ($)  
edu: 1=college  
gender: 1=F  
parent_edu: 1=col  
parent_income:  
salary($)
```

var	const	P> t
edu	20000	0.03
gender	-12000	0.06
parent_edu	15000	0.07
parent_income	1.8	0.01
edu:parent_income	2.3	0.02

Did you get it right

- (a) Yes, obv.
- (b) No, I cannot tell a lie.

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

Being female corresponds to a decrease of 12K in salary, holding all other things fixed.

# Clicker Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

Did you get it right

- (a) Yes, obv.
- (b) No, I cannot tell a lie.

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

# Discussion Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

How to we interpret this?

Conditioned on your having gone to college, an increase of \$1 in parents' salary corresponds to an increase of \$2.3 in your salary.

# Clicker Question!\*

```
income ~ education + gender + parent_edu +  
parent_income + education:parent_income
```

income: salary (\$)	var	const	P> t
edu: 1=college	edu	20000	0.03
gender: 1=F	gender	-12000	0.06
parent_edu: 1=col	parent_edu	15000	0.07
parent_income:	parent_income	1.8	0.01
salary(\$)	edu:parent_income	2.3	0.02

Did you get it right

- (a) Yes, obv.
- (b) No, I cannot tell a lie.

run along