



# Multivariable Analysis: Confounding or Interaction? How Much Complexity?

Giota Touloumi and Nikos Pantazis

Department of Hygiene, Epidemiology and Medical Statistics, Medical School, National and Kapodistrian University of Athens, Athens, Greece





# The components of GLMs

- The random part: the distribution the components of  $Y$  have
- The systematic component: covariates  $x_1, x_2, \dots, x_p$

produce a linear predictor  $\eta$  given by

$$\eta = \sum_{j=1}^p x_{ij} \beta_j$$

- The link between the random and the systematic components:  $g(\mu) = \eta$

○ Normal

$$\eta = \mu$$

*Identity*

○ Poisson

$$\eta = \log(\mu) = \log(\lambda)$$

*Log of rate*

○ Binomial

$$\eta = \log[\pi / (1 - \pi)]$$

*Logit*





# Multivariable Models

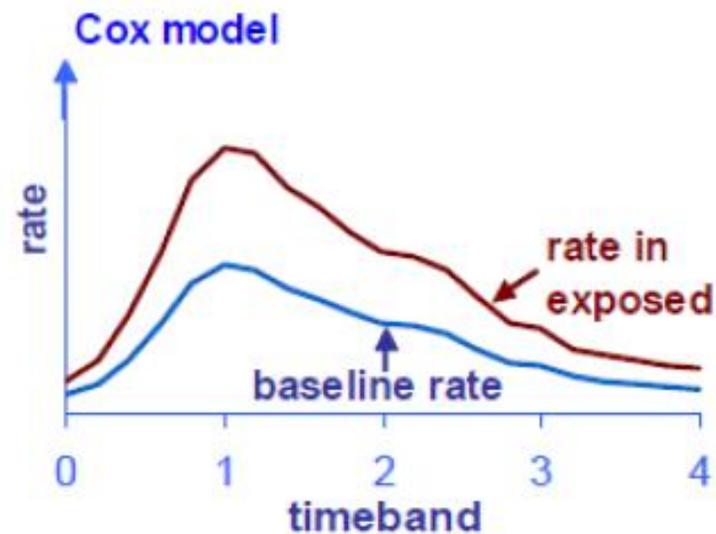
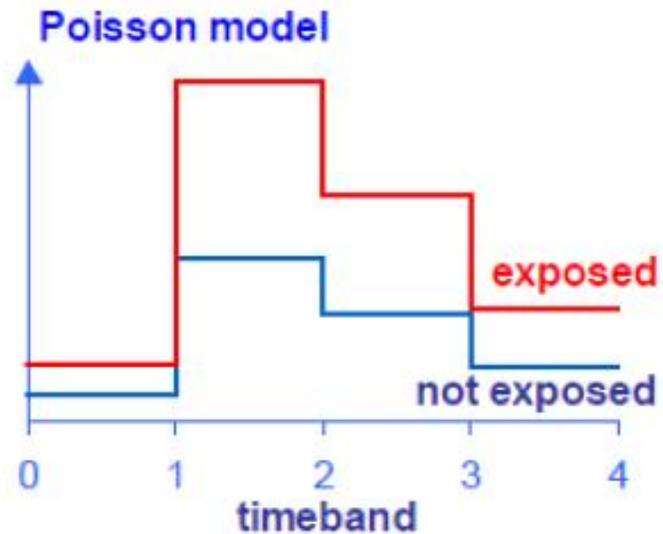
$$g(\mu) = b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p = \sum_0^p b_iX_i$$

- $g(\mu)=\mu$  linear regression
  - Mean change in response variable
- $g(\mu)=\log(\text{odds})=\log(\pi/1-\pi)$  Logistic regression
  - Probability of an event (odds ratio)
- $g(\mu)=\log(\mu)=\log(\text{rate})$  Poisson regression
  - rate of a new event (rate ratio)
  
- $\ln(h_t) = \ln(h_0(t)) + \sum_1^p b_iX_i$  Cox proportional hazards model
  - Time to event (Hazard ratio)





# Multivariable Models: Poisson vs Cox model





# Multivariable Models: How much complexity?

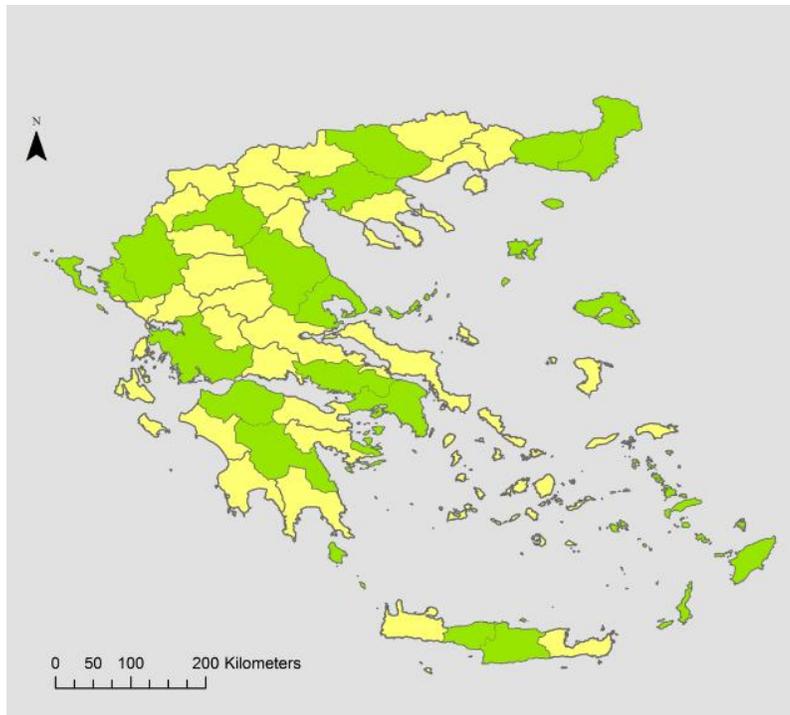
- **Linearity**
  - How to check, how to present/interpret
- **Confounding**
  - What is it? How we control for?
- **Interactions**
  - What is it? How we present/interpret results?
  - Confounding or Interactions?



# National Survey of Morbidity and Risk Factors (EMENO)



Stage 1: Prefectures/urbanization



Stage 2: City blocks

Stage 3: Households

Stage 4: Individual  
(with most recent  
birth date)

Sample selection: Multistage  
Stratified Random sampling

o Interviews "door-to-door"



Ιατρική Σχολή Πανεπιστημίου Αθηνών  
Εργ. Υγιεινής, Επιδημιολογίας & Ιατρικής Στατιστικής

6006 participants

# EMENO: Health Examination Survey

Questionnaire

Physical examination

Blood sample exams

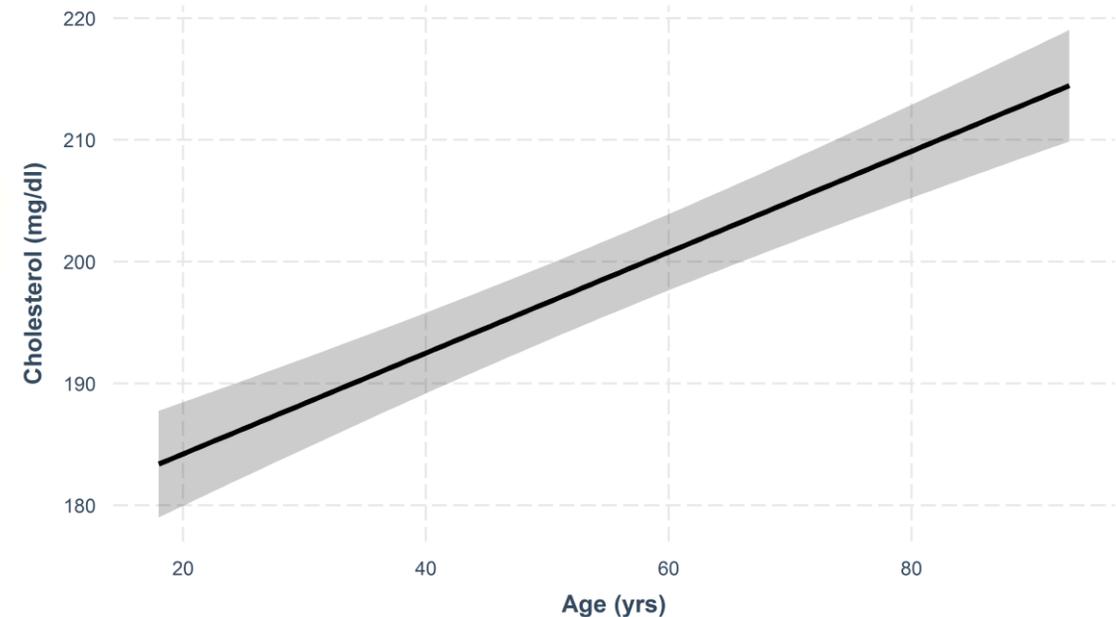
Height  
Weight  
Blood pressure  
Total, HDL, LDL Cholesterol  
Lipids  
Glucose



# Linearity $g(\mu) = b_0 + b_1X_1 + b_2X_2 + \dots + b_pX_p = \sum_0^p b_iX_i$

$$E(\text{Total Chol}) = b_0 + b_1\text{Age} + b_2\text{Female} + b_3\text{Diabetes} + b_4\text{Alcohol} + b_5 > 30\text{Walking}$$

Predicted cholesterol levels by age (cont. vars at mean, factors at base level)

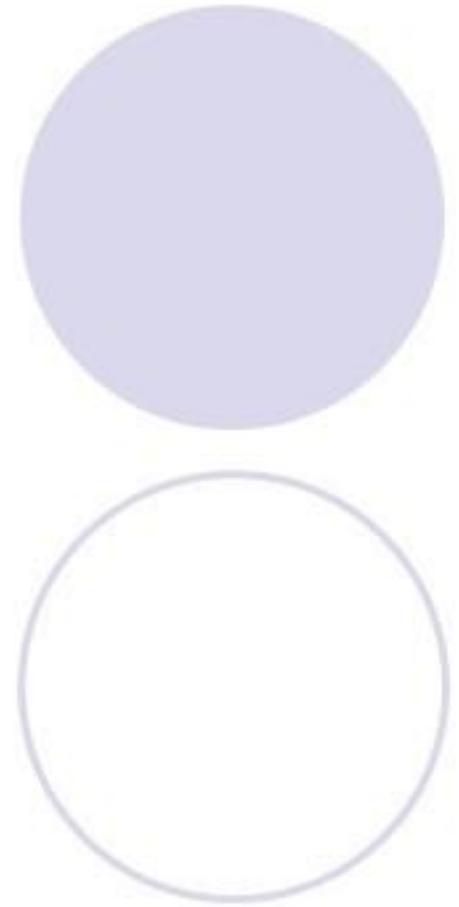
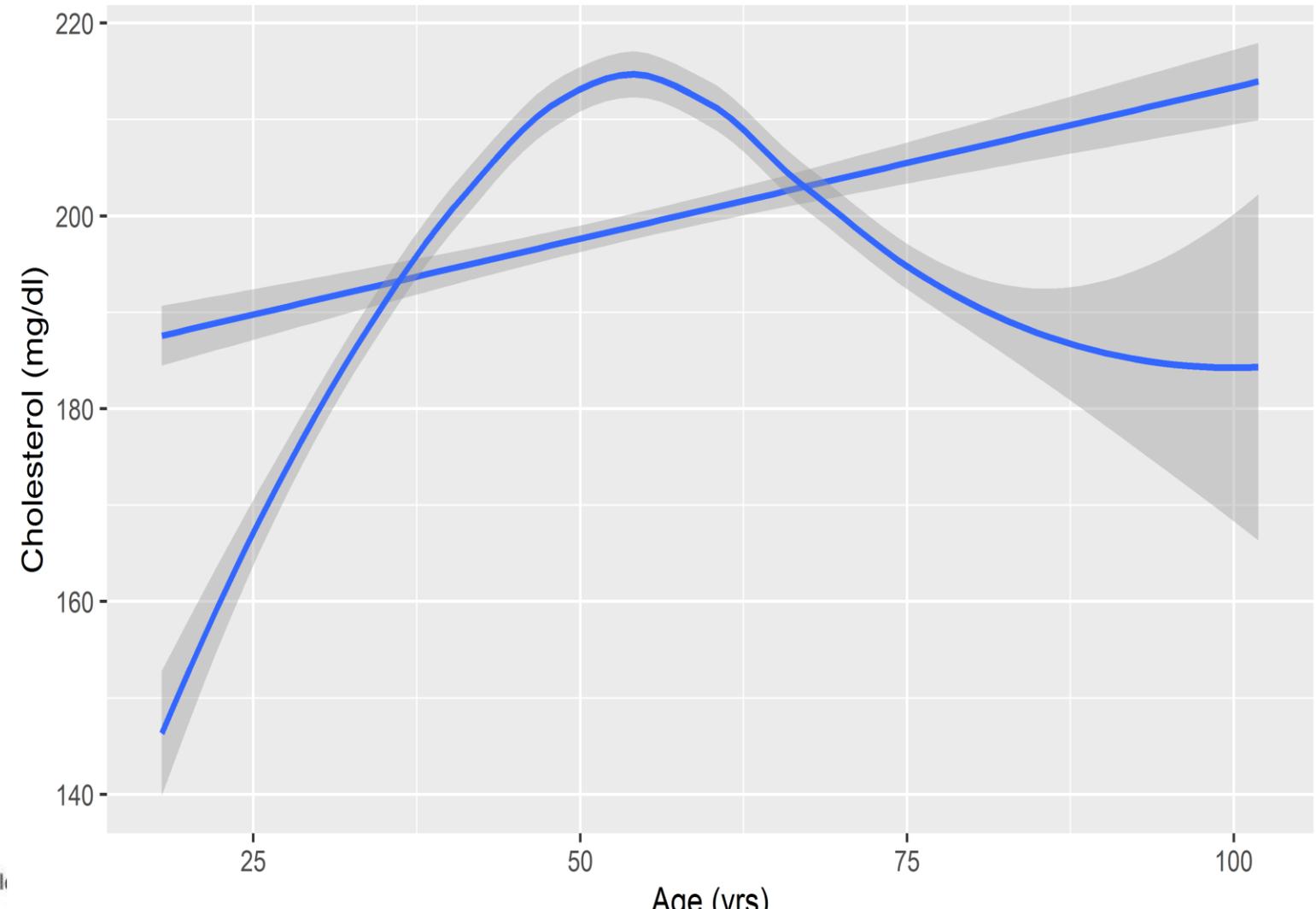


Characteristic	Beta	95% CI	p-value
Age (yrs)	0.414	0.328, 0.501	<0.001
sex			
Male			
Female	5.54	2.54, 8.54	<0.001
diab			
No			
Yes	-16.9	-21.2, -12.6	<0.001
Alcohol (categories)			
0-6			
7+	5.00	0.777, 9.23	0.020
Walking			
<30 min/day			
>=30 min/day	-2.81	-5.68, 0.069	0.056



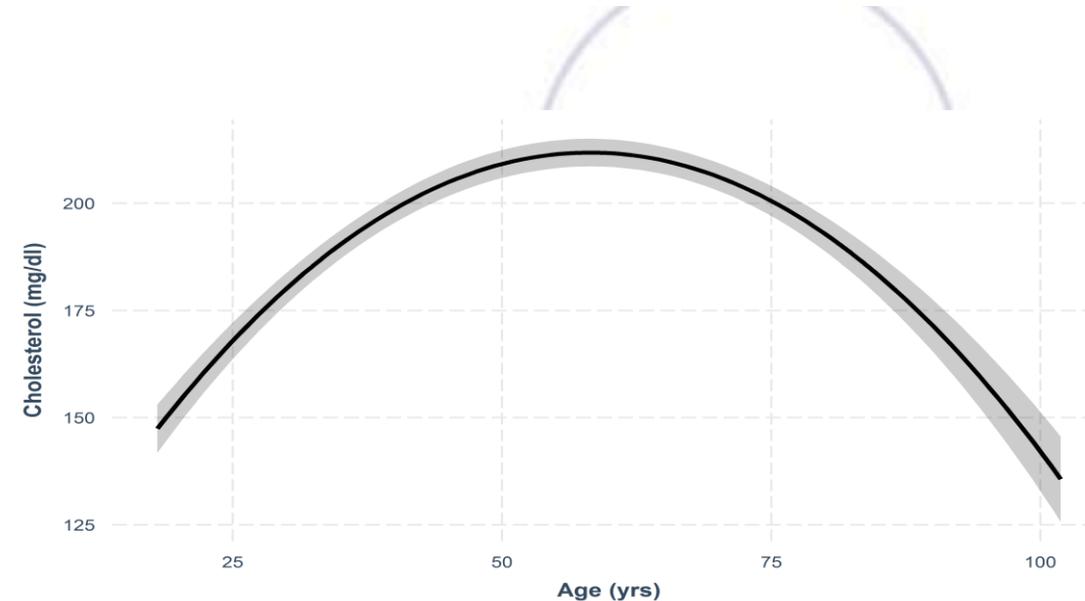
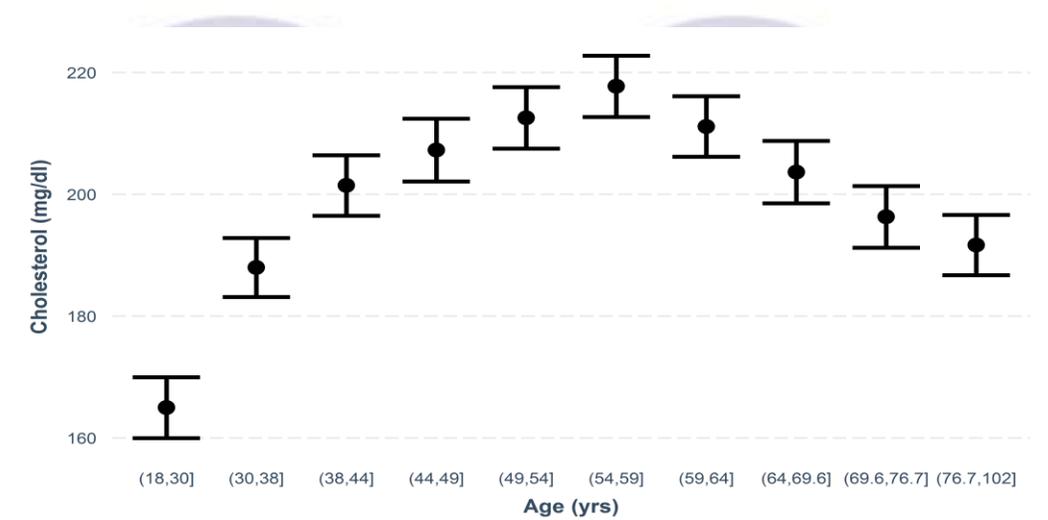
# Cholesterol age relationship: initial investigation

Scatterplot with smooth (loess) line



# Cholesterol age relationship: model checking (1)

Characteristic	Beta	95% CI	p-value
Age (categories - yrs)			
(18,30]			
(30,38]	23.0	17.1, 28.9	<0.001
(38,44]	36.5	30.5, 42.4	<0.001
(44,49]	42.3	36.2, 48.4	<0.001
(49,54]	47.6	41.6, 53.5	<0.001
(54,59]	52.8	46.8, 58.7	<0.001
(59,64]	46.2	40.1, 52.2	<0.001
(64,69.6]	38.7	32.6, 44.7	<0.001
(69.6,76.7]	31.3	25.2, 37.4	<0.001
(76.7,102]	26.7	20.6, 32.8	<0.001

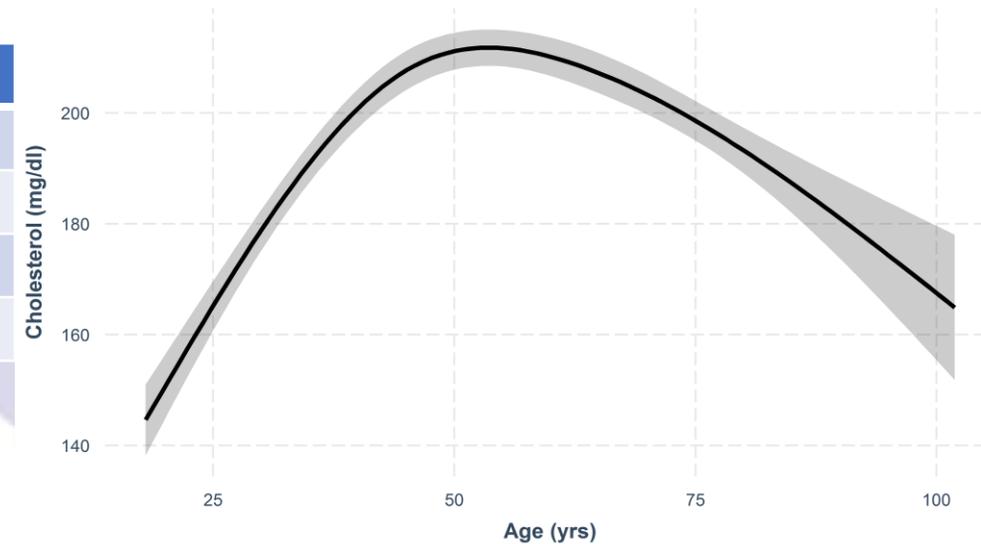


Characteristic	Beta	95% CI	p-value
Age			
Age	4.64	4.19, 5.09	<0.001
Age <sup>2</sup>	-0.040	-0.044, -0.036	<0.001



# Cholesterol age relationship: model checking (2)

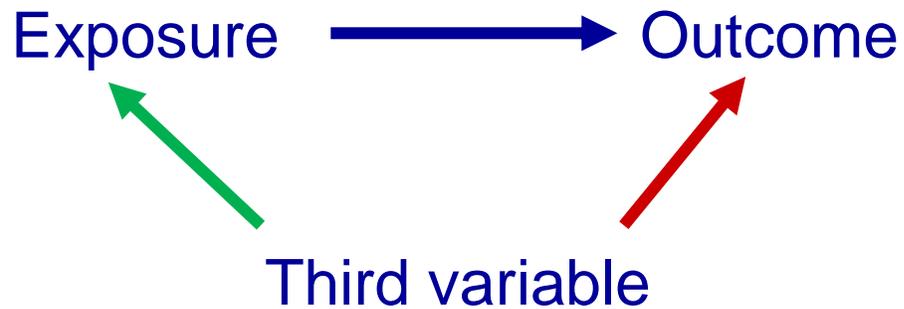
Characteristic	Beta	95% CI	p-value
Natural Spline of Age			
ns(Age, df = 3)1	38.9	32.3, 45.4	<0.001
ns(Age, df = 3)2	100	84.4, 116	<0.001
ns(Age, df = 3)3	-16.2	-29.1, -3.37	0.013



	df	AIC	$\Delta$ (AIC)	LogLikelihood
Linear	7	39341.32	351.94	-19663.64
Quadratic	8	39006.60	17.22	-19495.28
<b>Spline</b>	<b>9</b>	<b>38989.38</b>	<b>0.00</b>	<b>-19485.67</b>
Categorical	15	39001.58	12.20	-19485.73



# Confounding

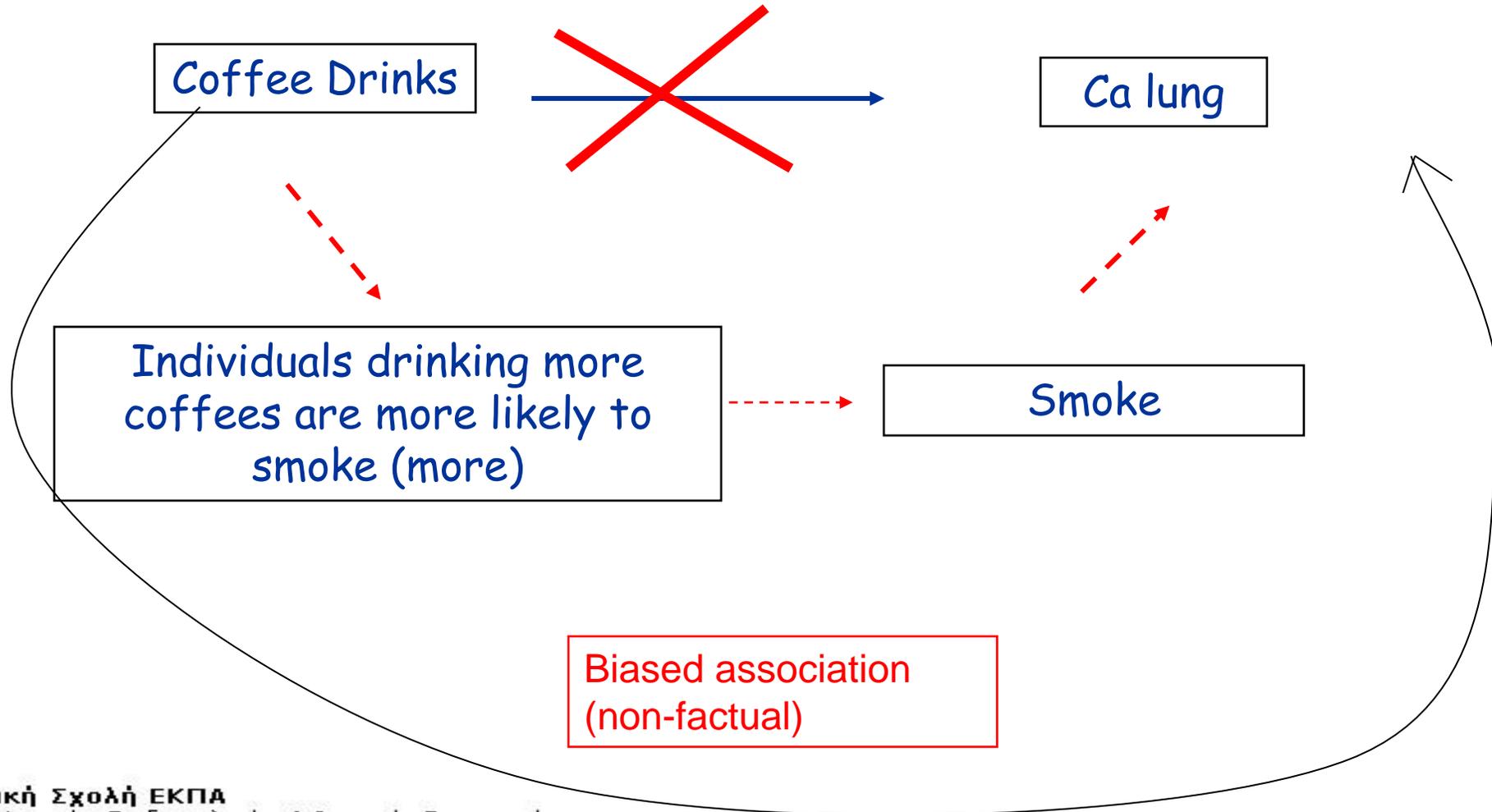


1. Prognostic factor of outcome
2. Associated with the third variable
3. It is not in the causal pathway

HDL in the causal pathway  
HDL: Not a confounder;  
HDL: No need for adjustment

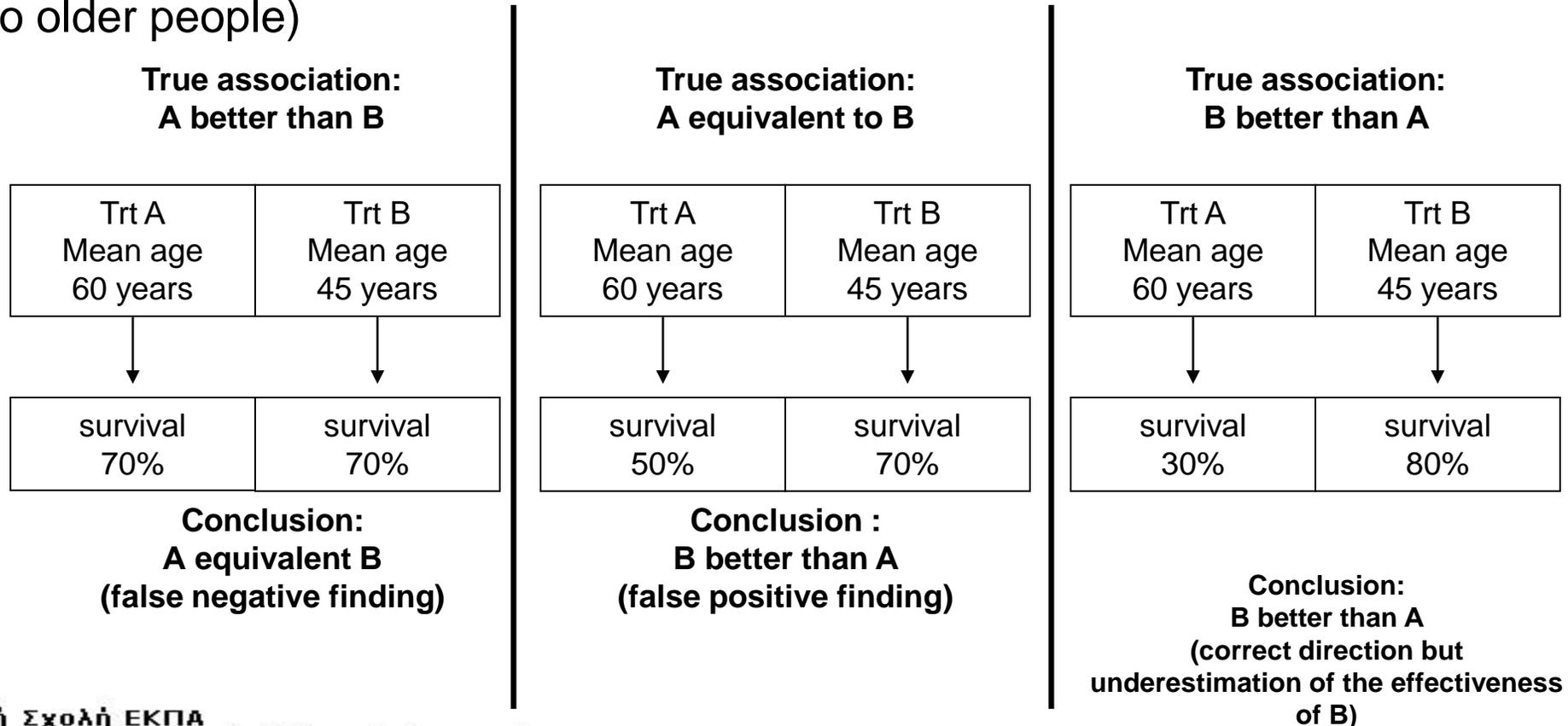


# Confounding: example



# Direction of bias: Confounding

- Compare treatment A with treatment B (exposure)
- Event: 5-years survival
- Confounder: age (lower survival at older ages; treatment A tended to be given to older people)



# Channeling effect

Cohort studies:

Hormone Replacement  
Therapy

**Association:**

Lower rates of CHD in  
women who used HRT

**Lower** rates of Coronary  
Heart Disease



# Channeling effect

Randomized Controlled  
Trials:

Hormone Replacement  
Therapy

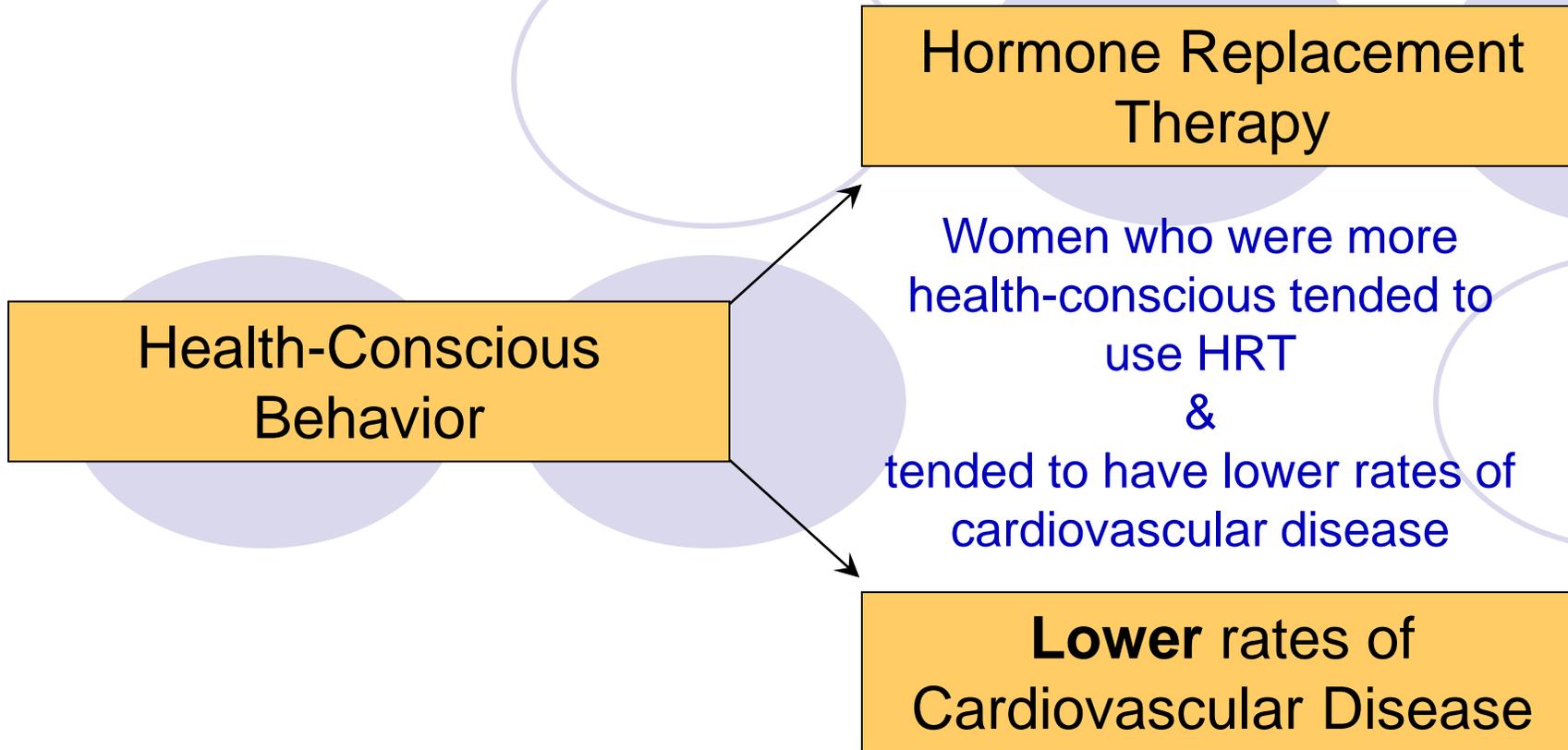
showed **HIGHER** rates  
of Cardiovascular  
Disease

**Lower** rates of Coronary  
Heart Disease



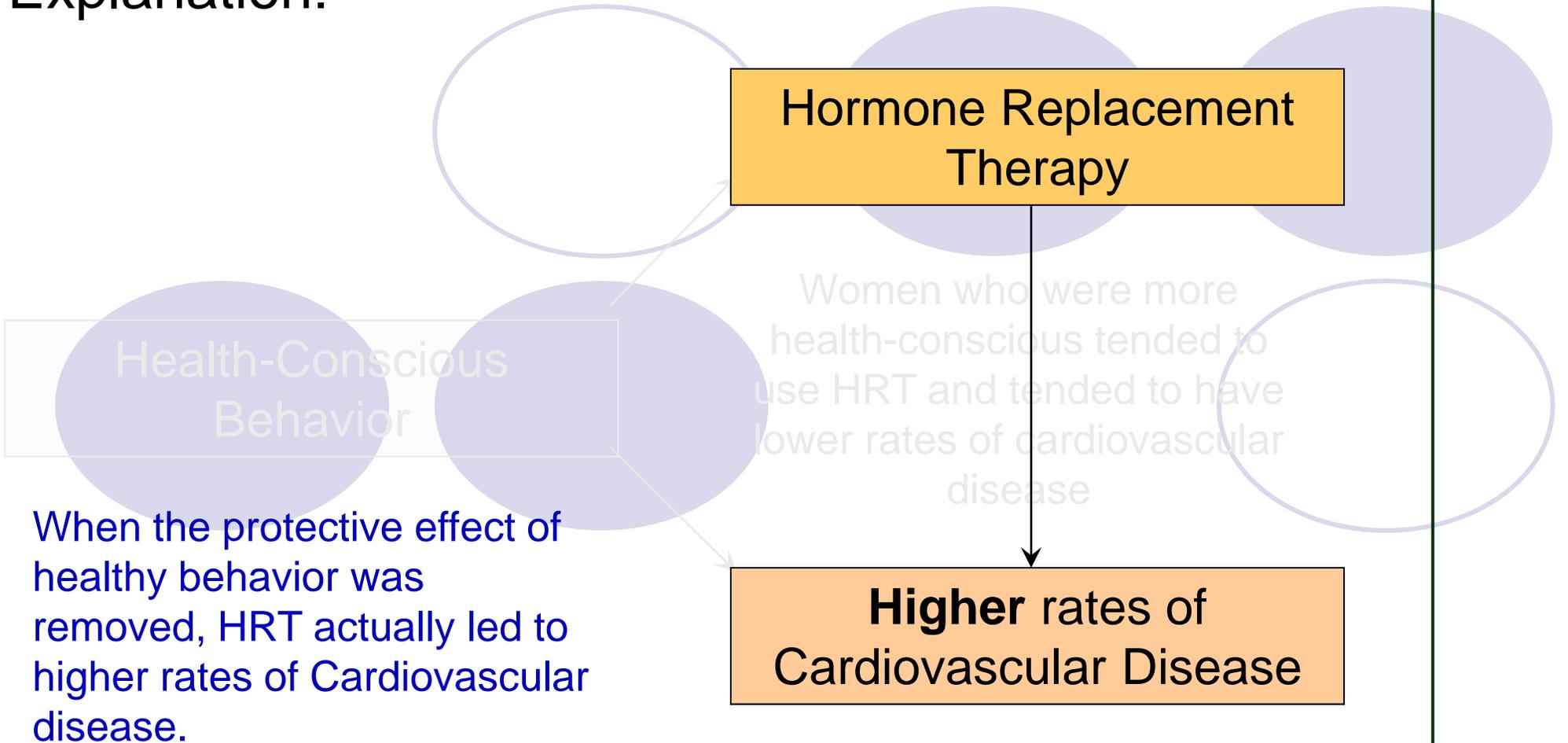
# Channeling effect

Explanation:



# Channeling effect

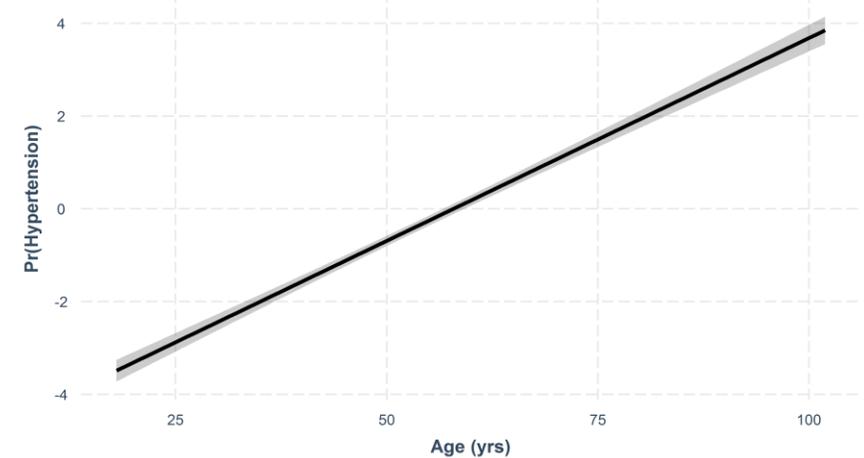
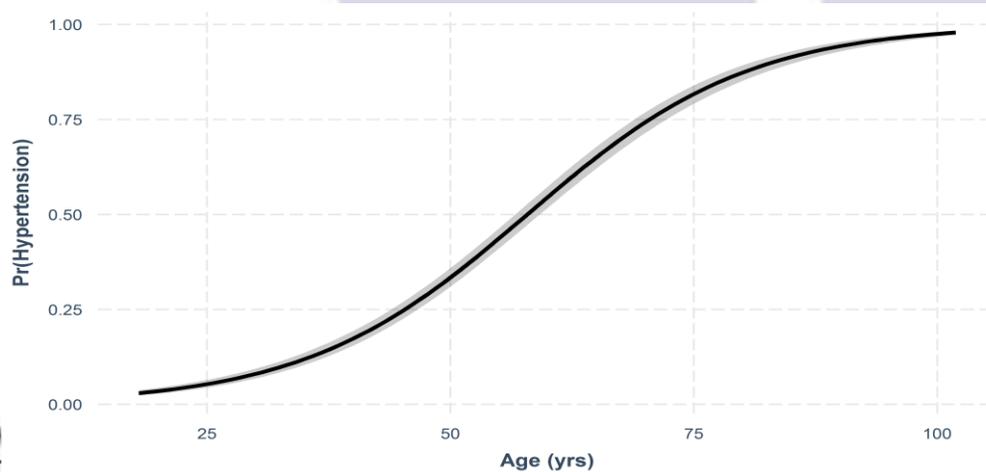
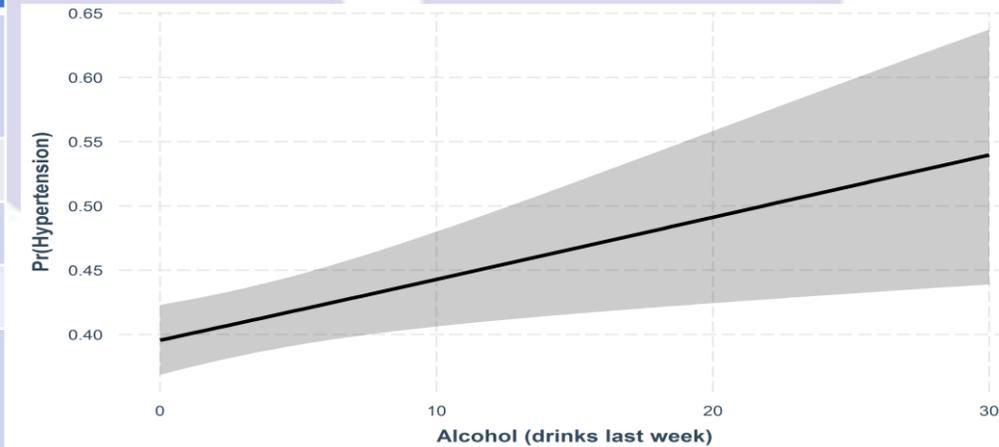
Explanation:



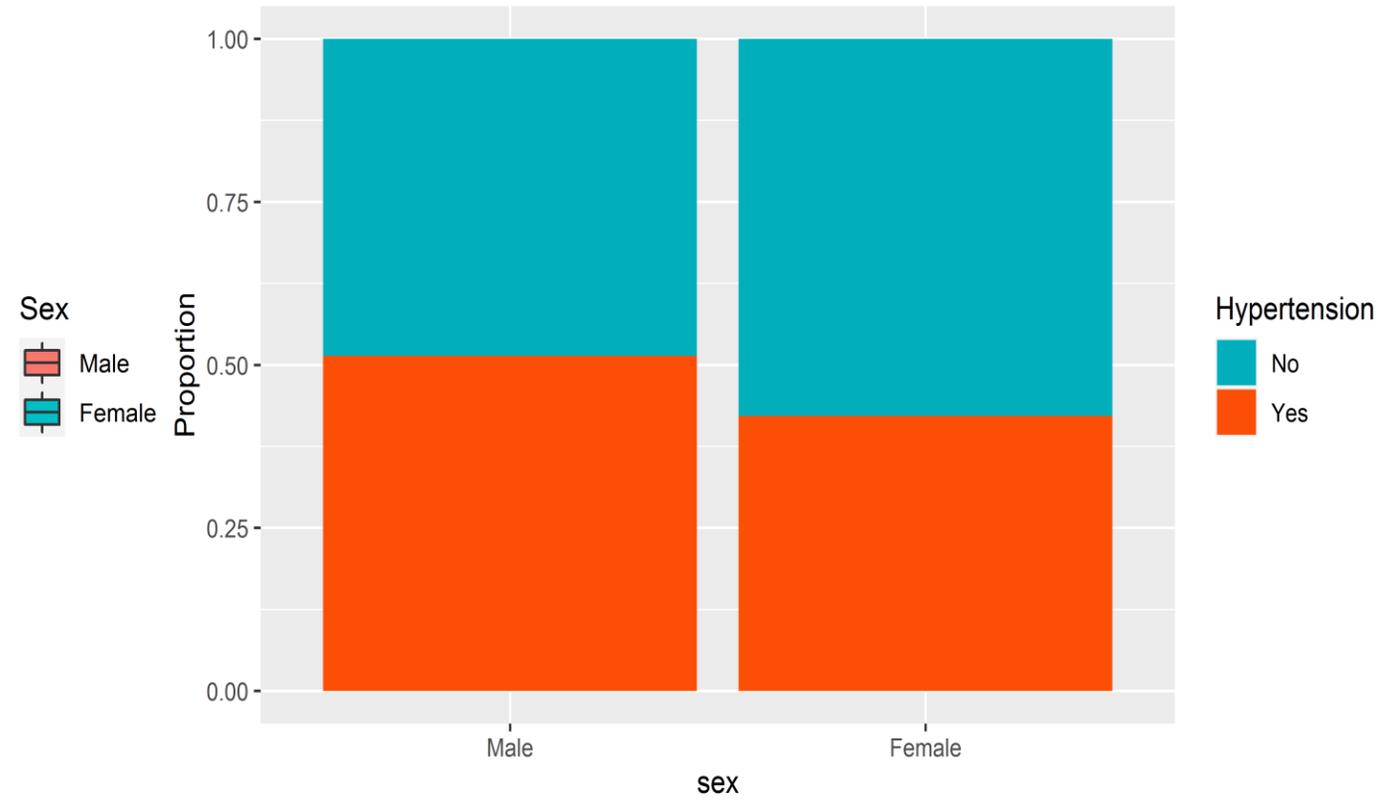
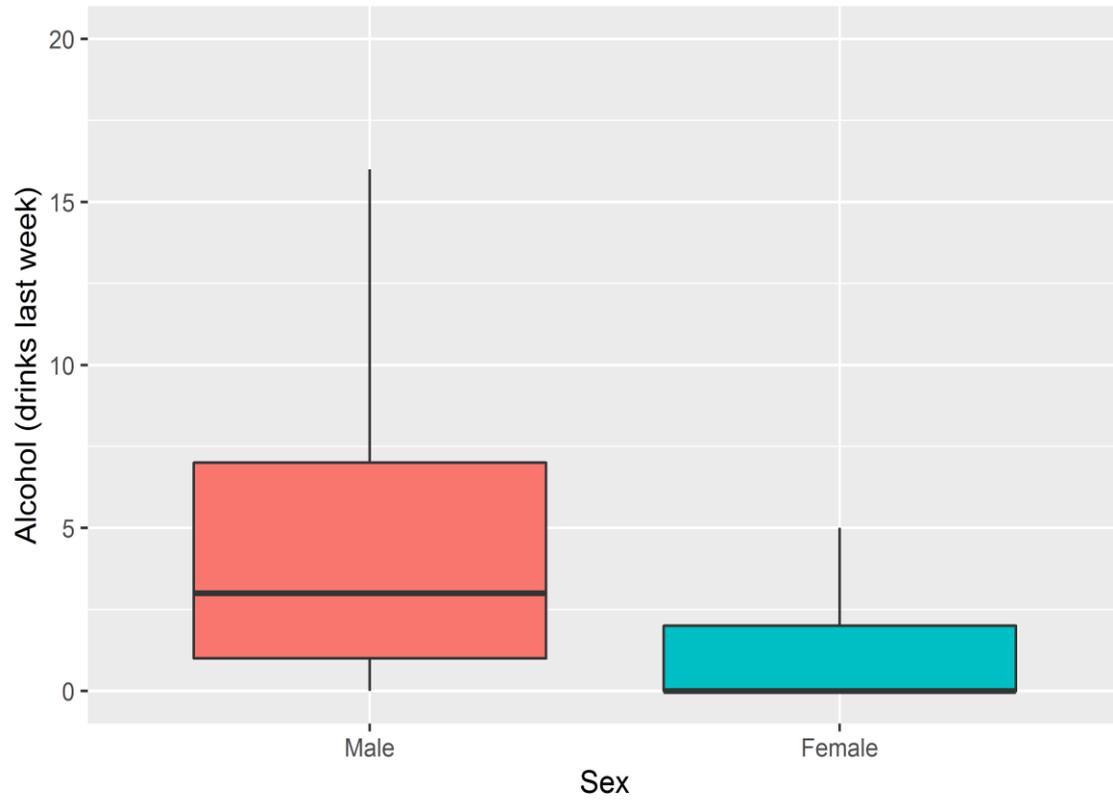
# Confounding: Pr (Hypertension)

$$\ln\left(\frac{\pi}{1-\pi}\right) = \log(\text{odds}) = b_0 + b_1 \text{Alcohol} + b_2 \text{Age} + b_3 \text{Semi-urban} + b_4 \text{Rural}$$

Characteristic	OR	95% CI	p-value
Alcohol (drinks last week)	1.02	1.01, 1.03	0.008
Age (yrs)	1.09	1.09, 1.10	<0.001
Urbanity			
Urban	—	—	
Semi-urban	1.20	0.98, 1.47	0.077
Rural	1.32	1.11, 1.57	0.002

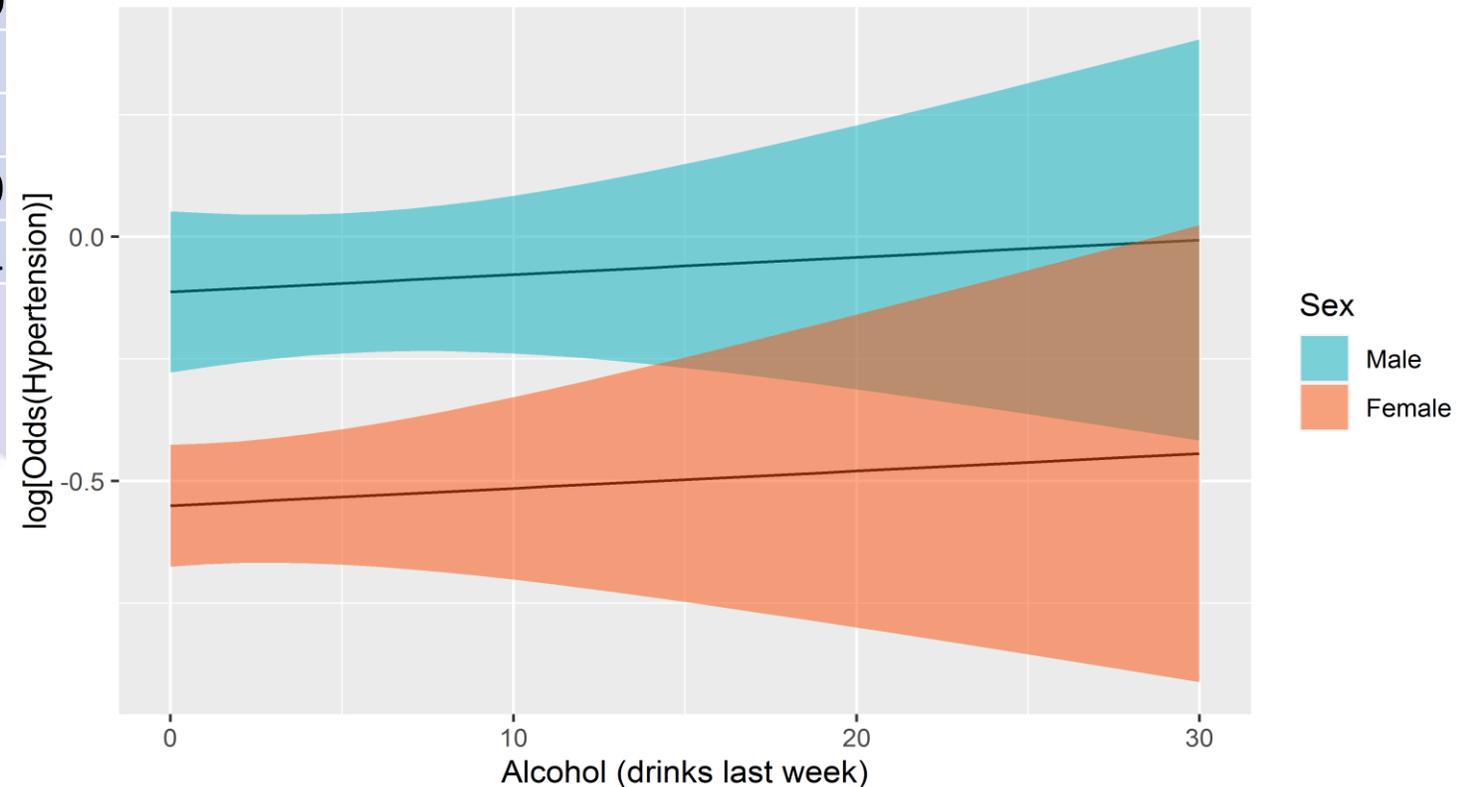


# Sex: Confounder?



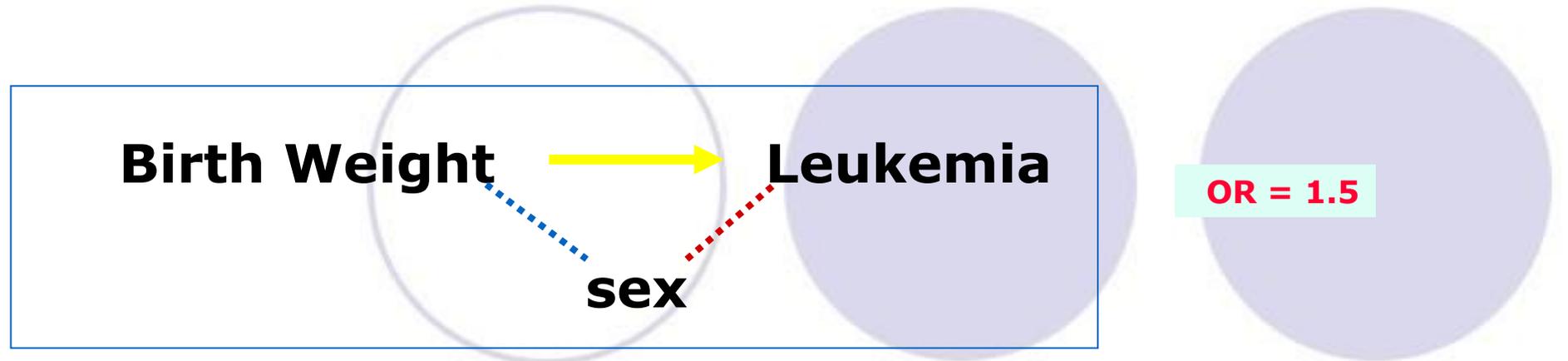
# Adjusting for sex

Characteristic	OR	95% CI	p-value
Alcohol (drinks last week)	<b>1.00</b>	<b>0.99, 1.02</b>	<b>0.655</b>
Sex			
Male			
<b>Female</b>	<b>0.646</b>	<b>0.545, 0.764</b>	<b>&lt;0.001</b>
Age (yrs)	1.09	1.00, 1.19	<0.001
Urbanity			
Urban			
Semi-urban	1.20	0.9, 1.5	0.001
Rural	1.31	1.1, 1.5	<0.001

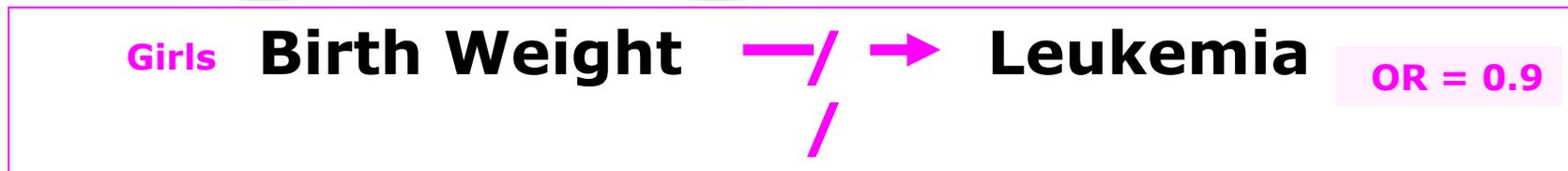




# Effect modifier (Interactions)



**Is the weight effect differentiated by sex?**



# Effect Modifier



---

- ❑ If the estimated effect of a risk factor on the outcome differs across the levels of a third variable
- ❑ Statistically: Interaction
- ❑ There is no point/intention to adjust for an effect modifier
- ❑ The effect of the risk factor should be presented separately at the different levels of the effect modifier

# Effect Modifier vs confounder



---

## Effect Modifier

**Belong to «nature»!**

**Different effect at different levels of the effect modifier**

**Useful information...**

**It improves our knowledge for the underlying mechanisms**

**Application to public health interventions and to personalized medicine**

## Confounder

**Belong to the study!**

**Same effect across the levels of the third variable**

**We need to adjust for the third variable (crude and adjusted estimates)**

**It caused confusion in the data and the results**

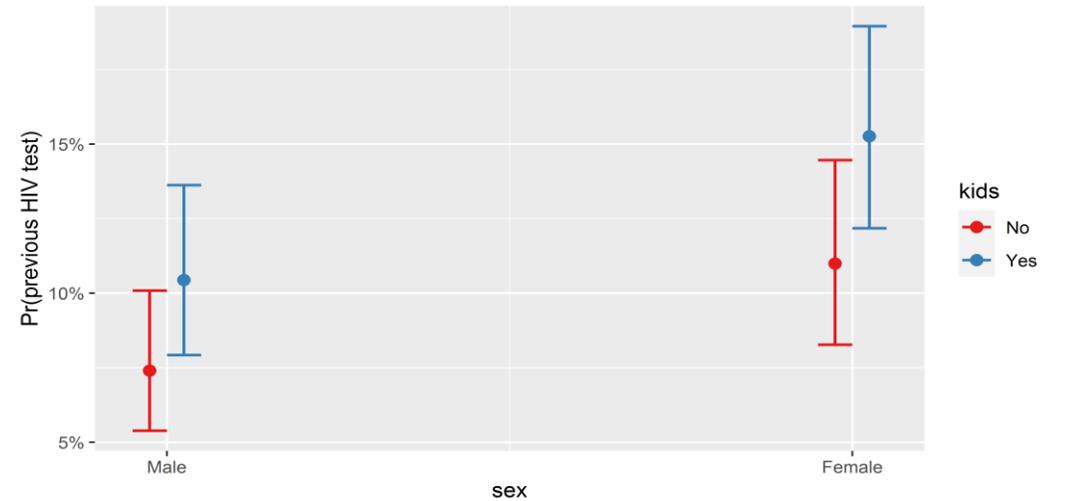
**We can deal with a confounder during the data analysis**

# Interactions: Binary x Binary

$$\text{Log}(\text{Odds}) = b_0 + b_1 \text{Fem} + b_2 \text{Kids} + b_3 \text{age} + b_4 \text{Ins.} + b_5 \text{Sex}(6 - 10) + b_6 \text{sex}(11+)$$

$\pi$  = Probability of past testing for HIV

Characteristic	OR	95% CI	p-value
<b>Sex</b>			
Male			
Female	1.55	1.28, 1.87	<0.001
<b>Kids</b>			
No			
Yes	1.46	1.18, 1.81	<0.001
<b>Age (yrs)</b>			
	0.97	0.96, 0.97	<0.001
<b>Insurance</b>			
No			
Yes	1.33	1.03, 1.72	0.033
<b>Number of Sexual partners</b>			
0-5			
6-10	2.02	1.61, 2.52	<0.001
11+	3.84	3.04, 4.86	<0.001



# Interactions: Binary x Binary

$$\text{Log(Odds)} = b_0 + b_1 \text{Fem} + b_2 \text{Kids} + b_3 \text{age} + b_4 \text{Ins.} + b_5 \text{Sex}(6 - 10) + b_6 \text{sex}(11 +) + b_7 (\text{Fem} * \text{Kids})$$

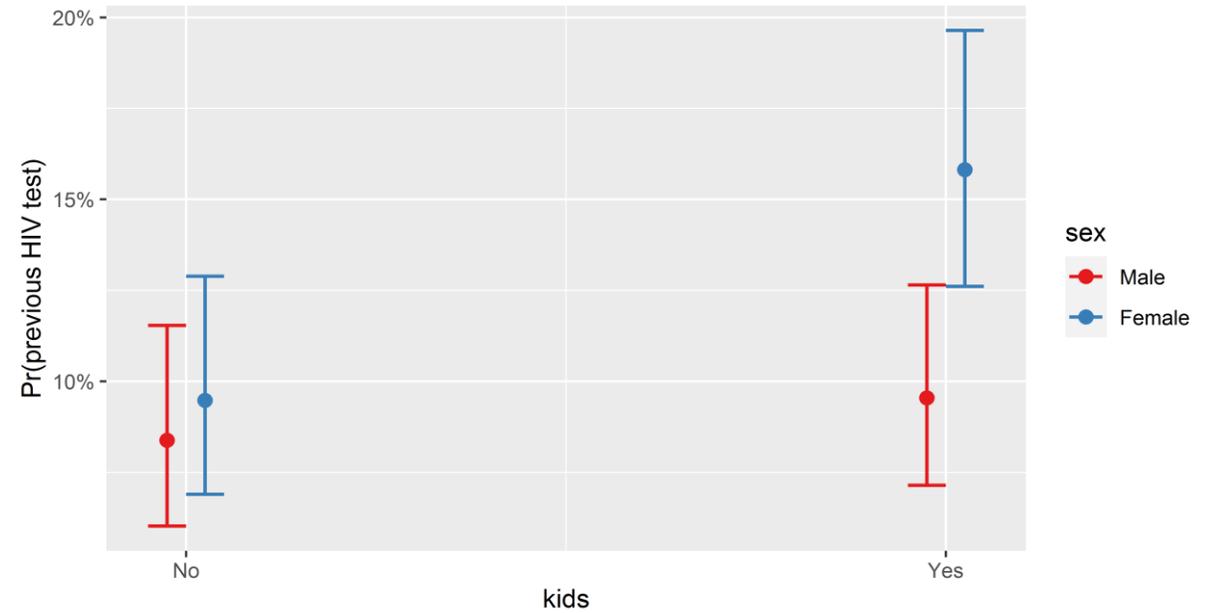
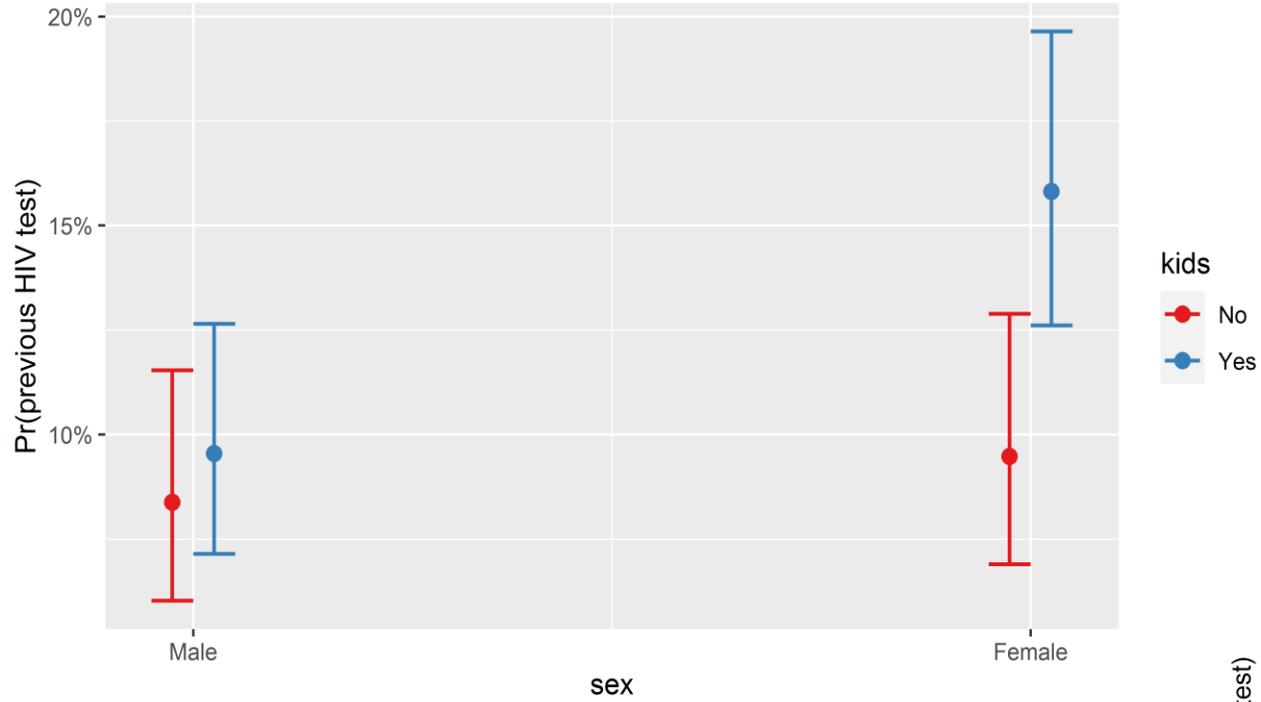
Characteristic	log(OR)	p-value
Sex		
Male		
Female	0.135	0.391
Kids		
No		
Yes	0.143	0.326
Sex * Kids		
Female * Yes	0.441	<b>0.016</b>

Characteristic	OR	95% CI	p-value
Sex			
Male			
Female	1.14	0.840, 1.56	0.391
Kids			
No			
Yes	1.15	0.868, 1.54	0.326
Female * Yes	1.55	1.09, 2.23	0.016

	No Kids	Yes Kids
Male	1	$e^{0.143} = 1.15$
Female	$e^{0.135} = 1.14$	$e^{0.135+0.143+0.441} = 2.05$

	Female/Male
No Kids	$e^{0.135} = 1.14 (0.84-1.56)$
Yes Kids	$e^{0.135+0.441} = 1.78 (1.43 - 2.22)$
	Yes Kids/No Kids
Male	$e^{0.143} = 1.15 (0.87-1.54)$
Female	$e^{0.143+0.441} = 1.79 (1.36 - 2.36)$

# Effect Modifier: Graphical representation



# Interactions: How to estimate combinations of bs

$$\text{Log(Odds)} = b_0 + b_1 \text{Fem} + b_2 \text{Kids} + b_3 \text{age} + b_4 \text{Ins.} + b_5 \text{Sex}(6 - 10) + b_6 \text{sex}(11 +) + b_7 (\text{Fem} * \text{Kids})$$

---

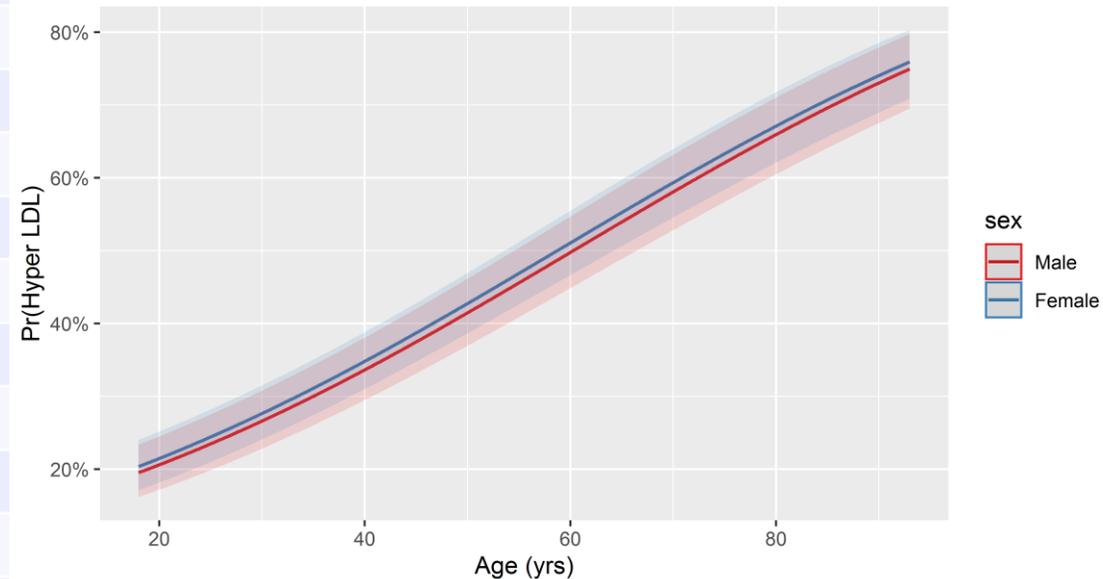
1. Define Constrain  $C = (0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 1)'$
2. Compute  $C*b'$ ,  $b$ : 1xp vector of bs
3. Compute  $C*V(b)*C'$

# Interactions: Binary x Continuous

$$\text{Log(Odds)} = b_0 + b_1 \text{Fem} + b_2 \text{age} + b_3 \text{overw} + b_4 \text{obese} + b_5 \text{walking30} + b_6 \text{walking30+}$$

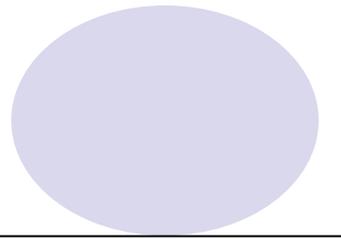
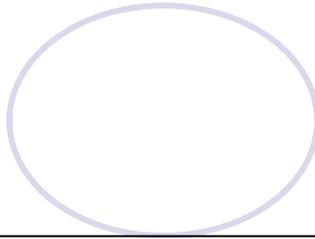
$\pi$  = Probability of elevated LDL

Characteristic	OR	95% CI	p-value
Sex			
Male			
Female	1.05	0.918, 1.21	0.457
Age (yrs)	1.03	1.03, 1.04	<0.001
BMI Categories			
Normal			
Overweight	1.94	1.62, 2.31	<0.001
Obese	1.85	1.55, 2.22	<0.001
Walking			
<30 min/day			
>=30 min/day	0.811	0.706, 0.931	0.003



# Interactions: Binary x Continuous

$$\text{Log(Odds)} = b_0 + b_1 \text{Fem} + b_2 \text{age} + b_3 \text{overw} + b_4 \text{obese} + b_5 \text{walking30} + b_6 \text{walking30} + b_7 \text{Fem} * \text{Age}$$



Characteristic	log(OR)	p-value
Sex		
Male		
Female	-0.667	0.005
Age (yrs)	0.026	<0.001
Sex * Age (yrs)		
Female * Age (yrs)	0.013	0.002

Characteristic	OR	95% CI	p-value
Sex			
Male			
Female	0.513	0.321, 0.820	0.005
Age (yrs)	1.03	1.02, 1.03	<0.001
Sex * Age (yrs)			
Female * Age (yrs)	1.01	1.01, 1.02	0.002

Male:  $0.026 * \text{age}$

OR (per 10 years): 1.30 (1.23-1.38)

Female:  $(0.026 + 0.013) * \text{age}$

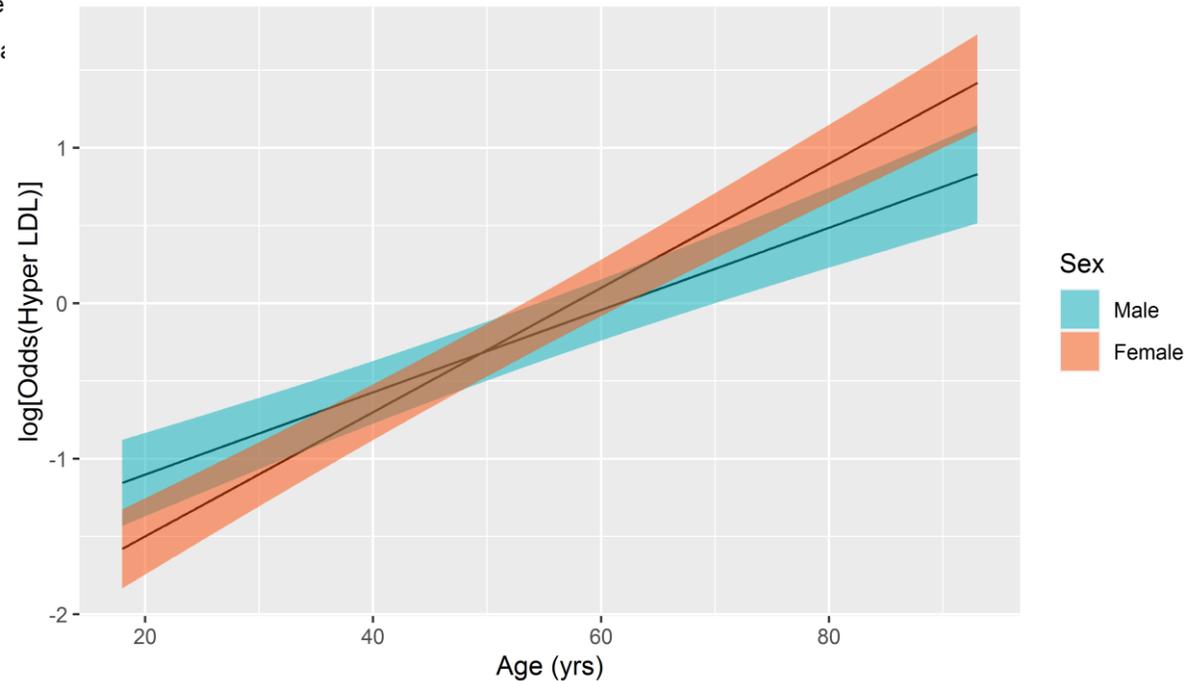
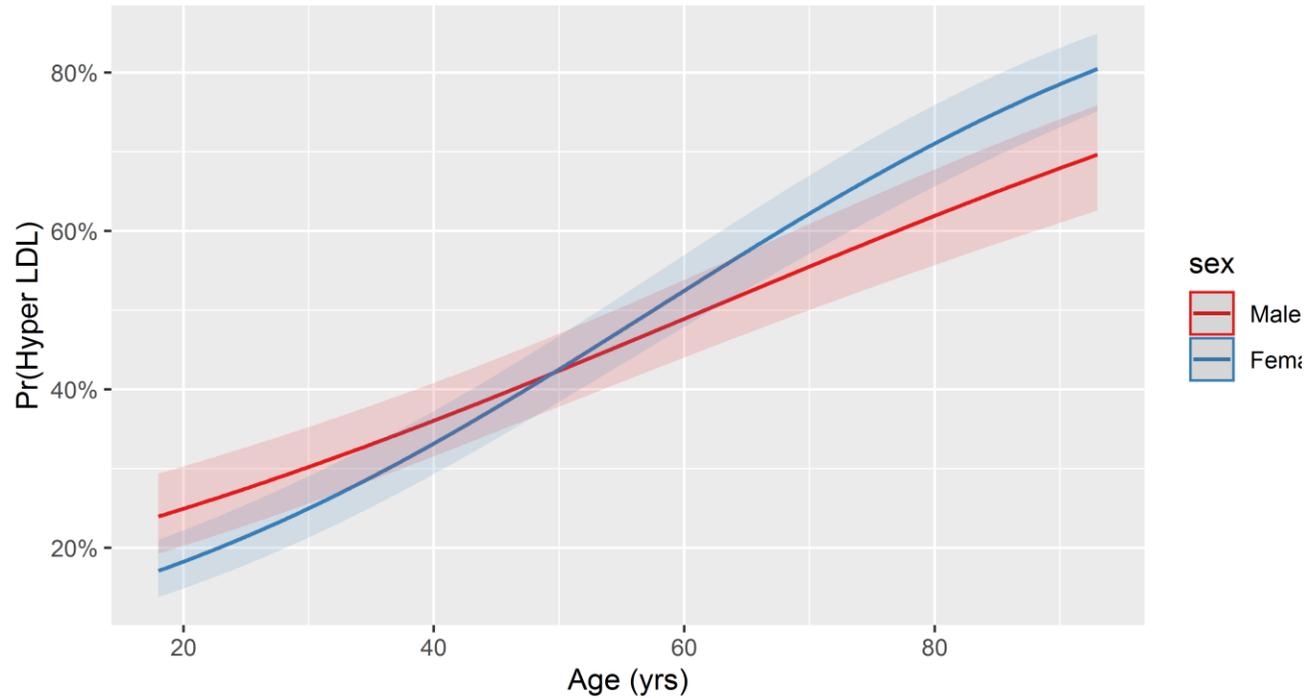
OR (per 10 years): 1.49 (1.40-1.58)

$$\frac{\text{Female}}{\text{Male}} : b_1 + b_6 * \text{age}$$

Age (years)	Female/Male: OR (95%CI)
20	0.67 (0.49-0.91)
50	1.01 (0.87-1.15)
80	1.51 (1.16-1.96)

# Interactions: Binary x Continuous; Graphically

$$\text{Log(Odds)} = b_0 + b_1 \text{Fem} + b_2 \text{age} + b_3 \text{overw} + b_4 \text{obese} + b_5 \text{walking30} + b_6 \text{walking30} + b_7 \text{Fem} * \text{Age}$$

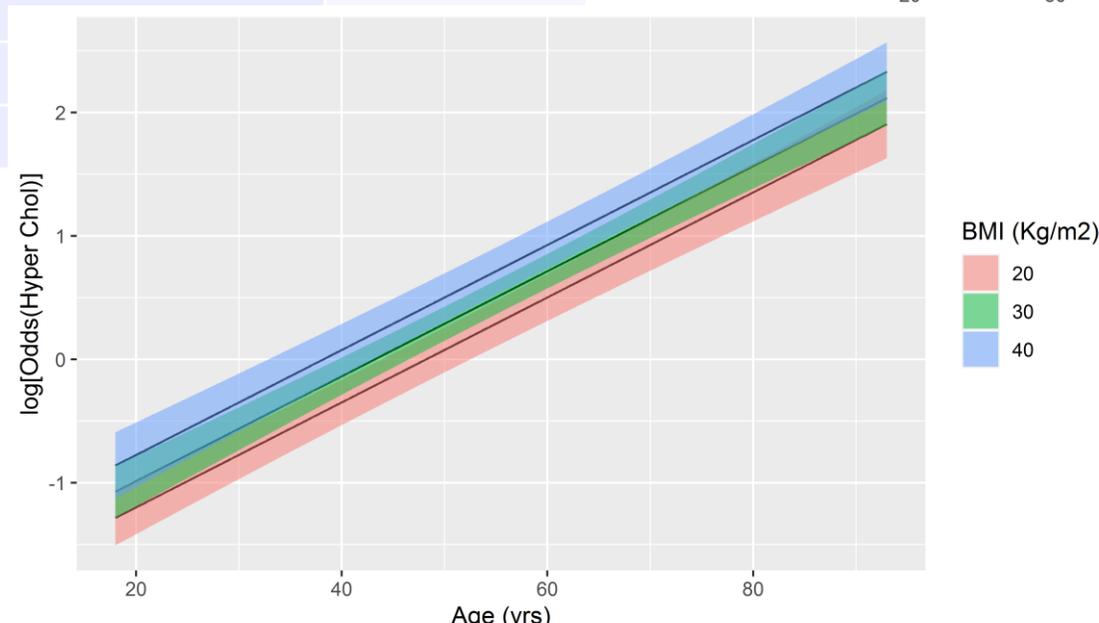
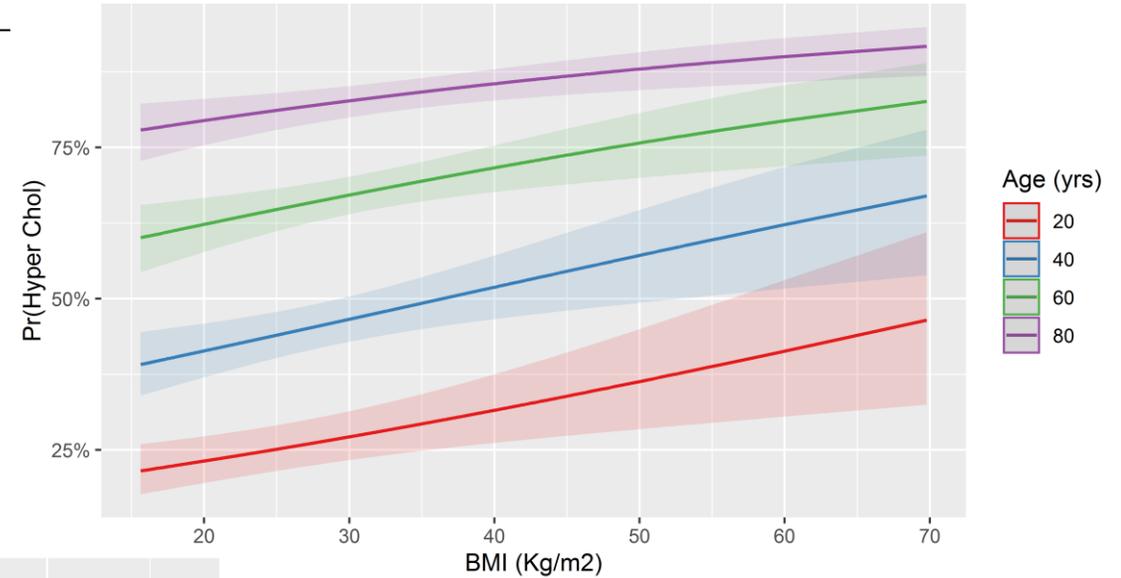


# Interactions: Continuous x Continuous

$$\text{Log}(\text{Odds}) = b_0 + b_1 \text{BMI} + b_2 \text{age} + b_3 \text{Fem} + b_4 \text{Walking30+}$$

$\pi$  = Probability of elevated Total Cholesterol

Characteristic	OR	95% CI	p-value
BMI (Kg/m <sup>2</sup> )	1.02	1.01, 1.03	0.001
Age (yrs)	1.04	1.04, 1.05	<0.001
Sex			
Male			
Female	1.17	1.02, 1.35	0.026
Walking			
<30 min/day			
>=30 min/day	0.848		



# Interactions: Continuous x Continuous

$$\text{Log(Odds)} = b_0 + b_1\text{BMI} + b_2\text{age} + b_3\text{Fem} + b_4\text{Walking30} + b_5\text{BMI} * \text{age}$$

Characteristic	log(OR)	p-value
BMI (Kg/m <sup>2</sup> )	0.170	<0.001
Age (yrs)	0.125	<0.001
Sex		
Male		
Female	0.197	0.006
Walking		
<30 min/day		
>=30 min/day	-0.177	0.015
BMI (Kg/m <sup>2</sup> ) *	-0.003	<0.001
Age (yrs)		

Characteristic	OR	95% CI	p-value
BMI (Kg/m <sup>2</sup> )	1.19	1.13, 1.24	<0.001
Age (yrs)	1.13	1.11, 1.16	<0.001
Sex			
Male			
Female	1.22	1.06, 1.40	0.006
Walking			
<30 min/day			
>=30 min/day	0.837	0.726, 0.97	0.015
BMI (Kg/m <sup>2</sup> ) *	1.00	1.00, 1.00	<0.001
Age (yrs)			

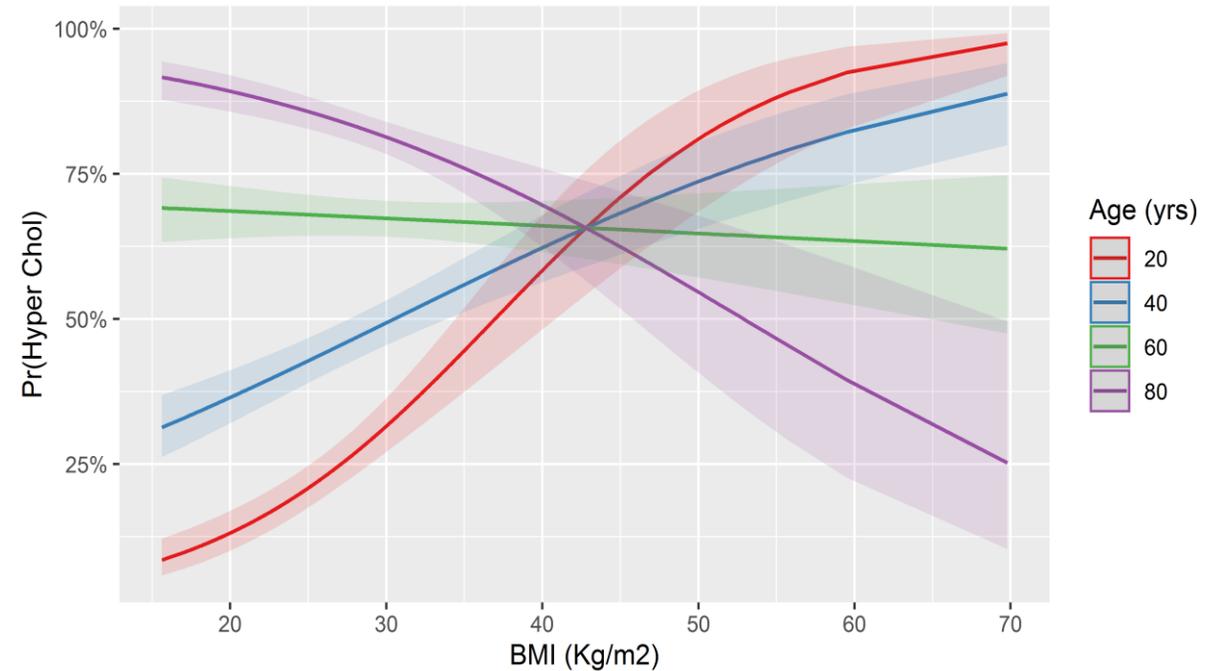
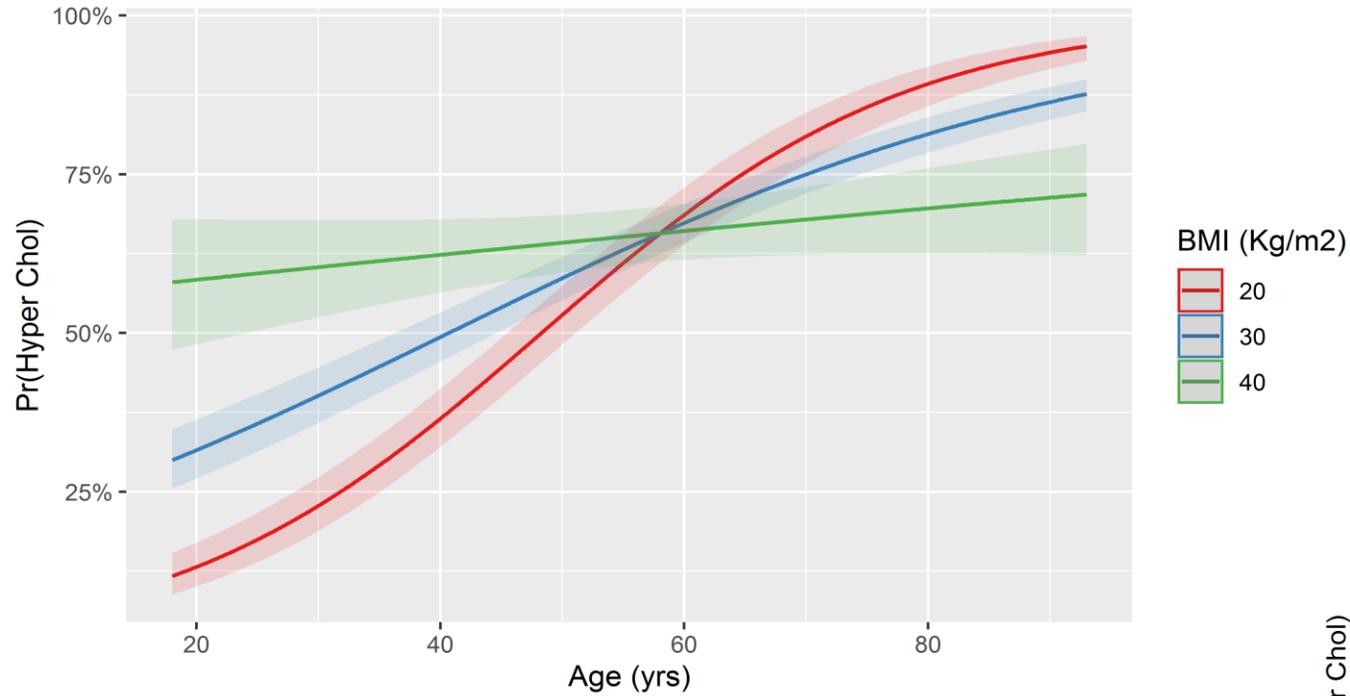
BMI (kg/m <sup>2</sup> )	Age (per 10yrs)
20	1.94 (1.79-2.12)
30	1.45 (1.39-1.52)
Age (yrs)	BMI (per 5 kg/m <sup>2</sup> )
30	1.51 (1.35-1.68)
60	0.73 (0.64-0.83)

BMI (per unit):  $b_1 + b_5 * \text{age}$

Age (per unit):  $b_2 + b_5 * \text{BMI}$

# Interactions: Continuous x Continuous graphically

$$\text{Log}(\text{Odds}) = b_0 + b_1 \text{BMI} + b_2 \text{age} + b_3 \text{Fem} + b_4 \text{Walking30} + b_5 \text{BMI} * \text{age}$$



## Prevalence of tobacco smoking and association with other unhealthy lifestyle risk factors in the general population of Greece: Results from the EMENO study

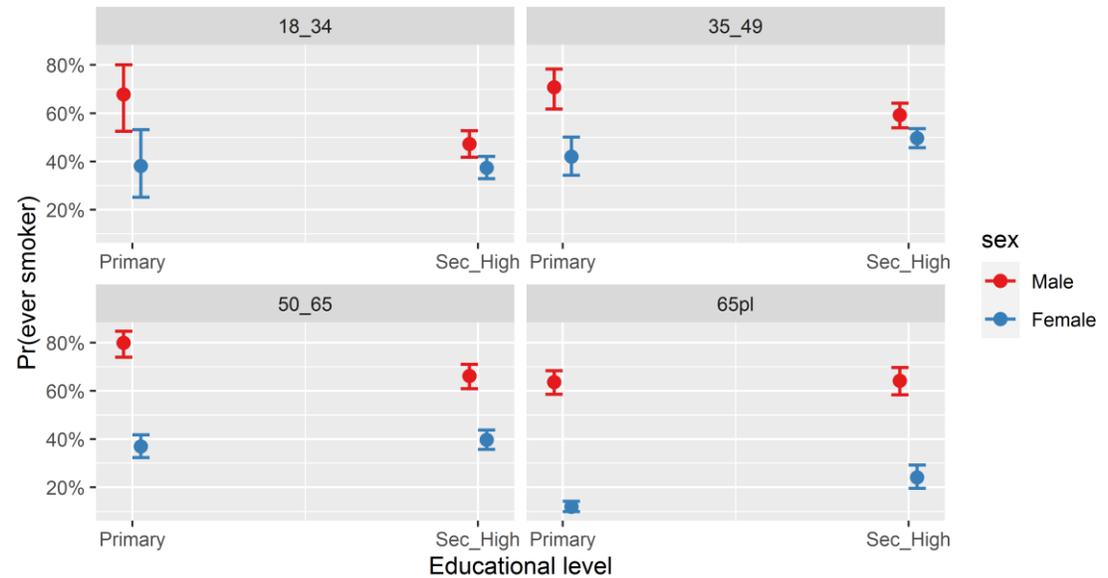
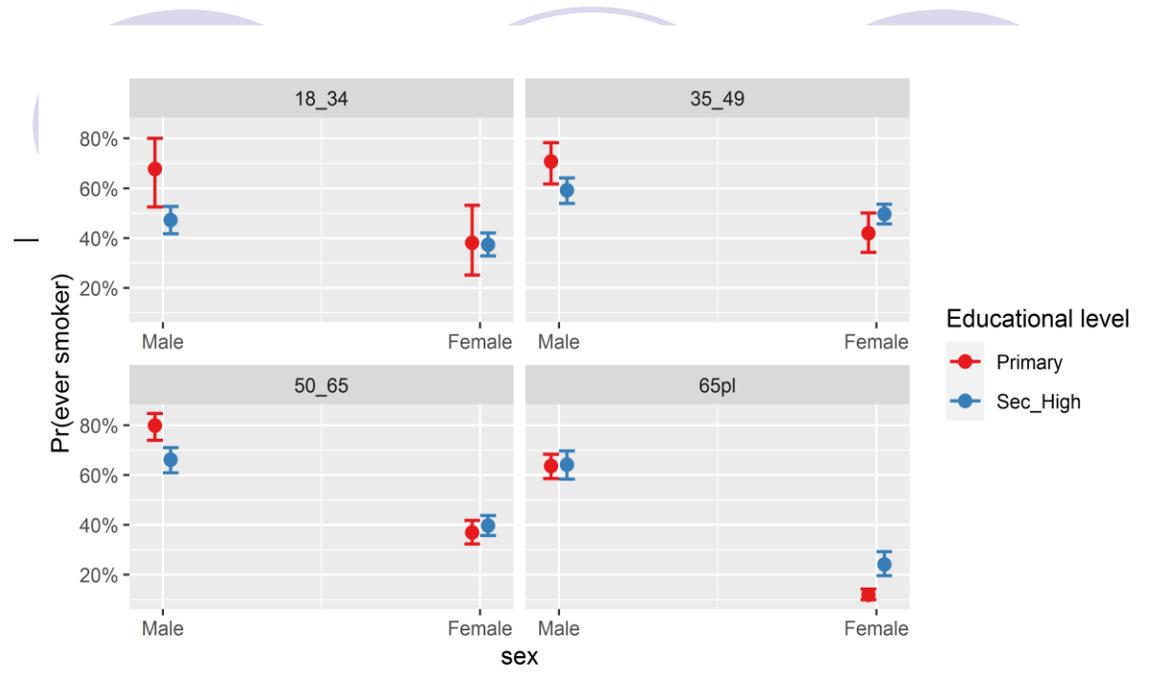
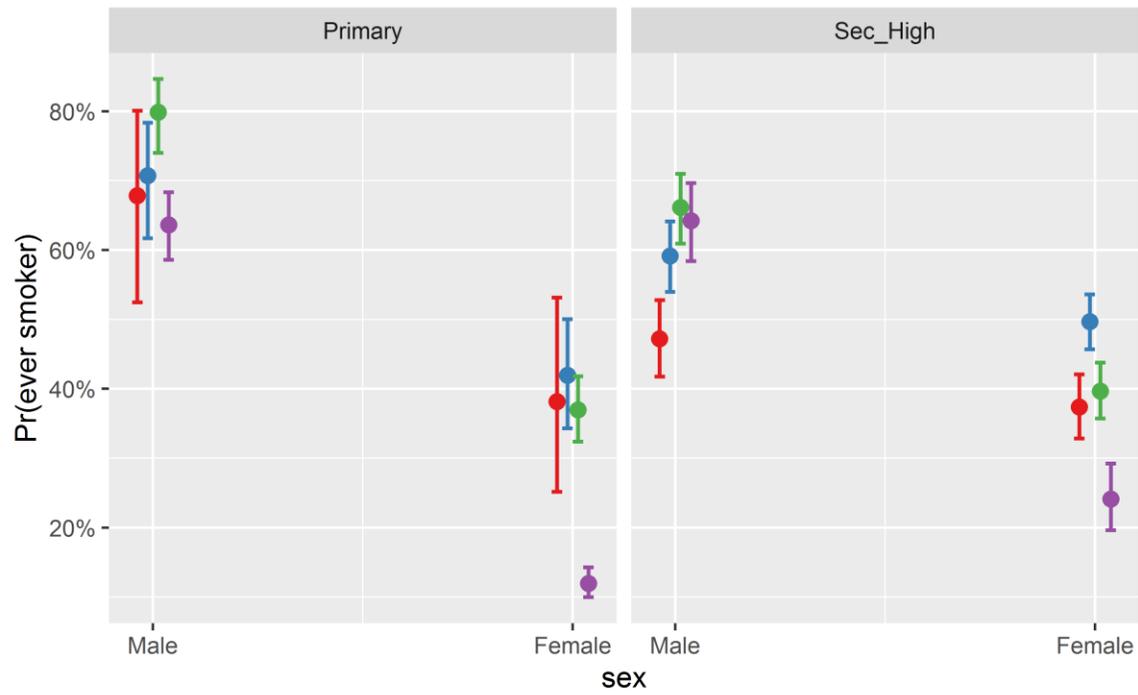
*Maria Gangadi<sup>1</sup>, Natasa Kalpourtzi<sup>2</sup>, Magda Gavana<sup>3</sup>, Apostolos Vantarakis<sup>4</sup>, Gregory Chlouverakis<sup>5</sup>, Christos Hadjichristodoulou<sup>6</sup>, Gregory Trypsianis<sup>7</sup>, Paraskevi V. Voulgari<sup>8</sup>, Yannis Alamanos<sup>9</sup>, Argiro Karakosta<sup>2</sup>, Giota Touloumi<sup>2\*</sup>, Anna Karakatsani<sup>1\*</sup>*

Tob. Prev. Cessation, 2021

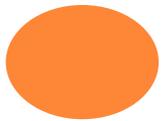
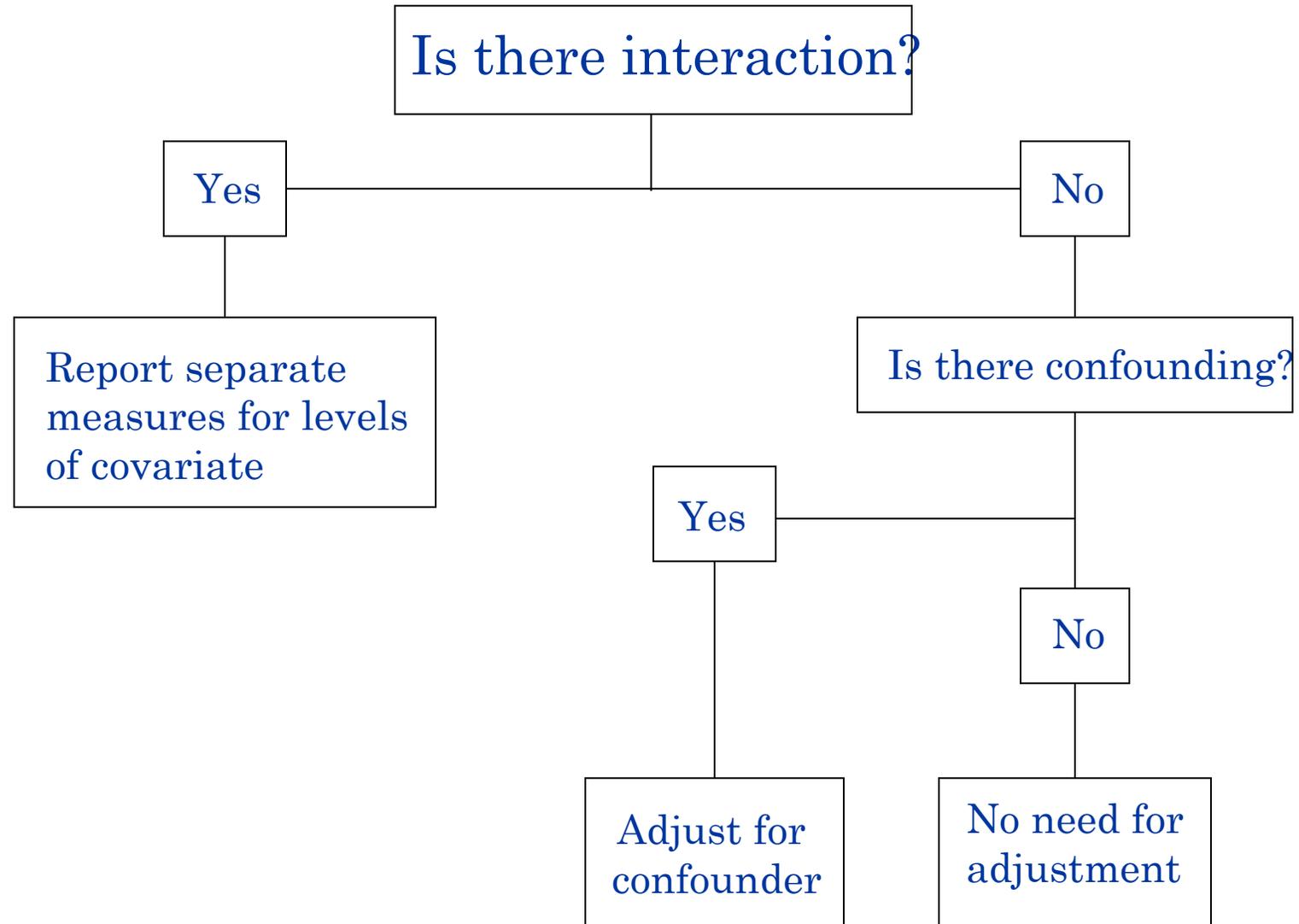
$$\begin{aligned} & \log[\text{Pr}(\text{EverSmoker})] \\ &= \beta_0 + \beta_1 * \text{EduSecHigh} + \beta_2 * \text{Age3549} + \beta_3 * \text{Age5065} + \beta_4 * \text{Age65pl} + \beta_5 * \text{SexFemale} + \beta_6 \\ & * \text{Alcohol17} + \beta_7 * \text{Alcohol7pl} + \beta_8 * \text{UnemployedYes} + \beta_9 * \text{EduSecHigh:Age3549} + \beta_{10} \\ & * \text{EduSecHigh:Age5065} + \beta_{11} * \text{EduSecHigh:Age65pl} + \beta_{12} * \text{EduSecHigh:SexFemale} + \beta_{13} \\ & * \text{Age3549:SexFemale} + \beta_{14} * \text{Age5065:SexFemale} + \beta_{15} * \text{Age65pl:SexFemale} \end{aligned}$$

<i>Characteristics</i>	<i>OR (95% CI)</i>	<i>p</i>
<b>Alcohol consumption (glasses/week)</b>		
1–7	1.41 (1.22–1.64)	<0.001
>7	2.52 (1.97–3.23)	0.001
<b>Unemployed/employed</b>	1.42 (1.16–1.73)	0.001
<b>Age (18–34.9 years)</b>		
<b>Women</b>		
Secondary/Higher vs Primary	0.84 (0.43–1.66)	0.614
<b>Men</b>		
Secondary/Higher vs Primary	0.44 (0.23–0.84)	0.013
<b>Primary education</b>		
Women/Men	0.37 (0.23–0.59)	<0.001
<b>Secondary/Higher education</b>		
Women/Men	0.71 (0.54–0.94)	0.015

<b>Age (35–49.9 years)</b>		
<b>Women</b>		
Secondary/Higher vs Primary	1.28 (0.85–1.92)	0.24
<b>Men</b>		
Secondary/Higher vs Primary	0.67 (0.44–1.03)	0.06
<b>Primary education</b>		
Women/Men	0.35 (0.23–0.53)	<0.001
<b>Secondary/Higher education</b>		
Women/Men	0.67 (0.52–0.87)	0.001
<b>Age (50–64.9 years)</b>		
<b>Women</b>		
Secondary/Higher vs Primary	1.00 (0.75–1.34)	0.99
<b>Men</b>		
Secondary/Higher vs Primary	0.52 (0.35–0.79)	0.001
<b>Primary education</b>		
Women/Men	0.17 (0.12–0.26)	<0.001
<b>Secondary/Higher education</b>		
Women/Men	0.33 (0.25–0.43)	<0.001



# APPROACH TO INTERACTION AND CONFOUNDING





Thank you for your attention

