

Week 4: Misclassification and Measurement Error

Video 5: Misclassification and Measurement Error Terminology

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Week 4: Misclassification and Measurement Error

Video 5: Misclassification and Measurement Error Terminology

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health

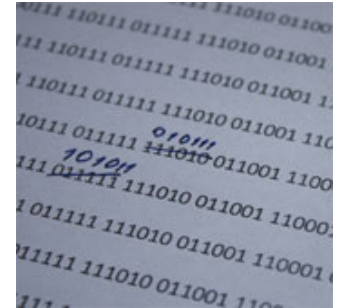


HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Key Concepts

- Terminology
 - Misclassification vs. measurement error
 - Nondifferential vs. differential misclassification
 - Independent vs. dependent misclassification
- Indices of measurement accuracy
- Formulae for misclassification of a binary exposure
 - Truth \rightarrow Observed
 - Observed \rightarrow Truth
- Misclassification of a polytomous exposure
- Impact of misclassification
 - Estimation and hypothesis tests
 - Confounding
 - Effect Measure Modification



Perfect Measurements

- So far, we have assumed that the exposure, outcome and any confounders have been measured with no error.
- However, measurement is not always perfect.
- These errors may result in incorrect estimates of the underlying association in studies.



Misclassification vs. Measurement Error

- Misclassification

- ☐ error in a categorical disease determinant or outcome
- ☐ e.g. ever exposed to DES (yes/no)
- ☐ e.g. cancer incidence by the end of the follow-up period (yes/no)

- Measurement Error

- ☐ error in a continuous exposure, confounder or outcome
- ☐ e.g. blood pressure, cholesterol or sex hormone levels



Nondifferential Versus Differential Misclassification

- **Nondifferential** misclassification occurs in the same proportion in each group
 - Equal amount of misclassification of exposure in cases and non-cases
 - Equal amount of misclassification of outcome in exposed and unexposed
- **Differential** misclassification occurs in different proportions in each group
 - Different amount of misclassification of exposure in cases compared to non-cases
 - Different amount of misclassification of outcome in exposed compared to unexposed



Independent vs. Dependent Misclassification

- **Independent** misclassification is not dependent on the probability of errors in classifying other variables
 - Probability of misclassifying exposure does not depend on probability of misclassifying outcome
 - Probability of misclassifying outcome does not depend on probability of misclassifying exposure
- **Dependent** misclassification depends on the errors in measuring or classifying other variables
 - Probability of misclassifying exposure depends on probability of misclassifying outcome
 - Probability of misclassifying outcome depends on probability of misclassifying exposure
- If only one of exposure or disease is misclassified, but not the other, the independence assumption is not needed since, by definition, it is always present.



Misclassification of a Dichotomous Exposure or Outcome

	Non-Differential	Differential
Exposure	<ul style="list-style-type: none"> ■ Misclassification of exposure is similar for cases and non-cases ■ Generally leads to a bias towards the null* 	<ul style="list-style-type: none"> ■ Misclassification of exposure is different for cases and non-cases ■ Lead to a bias in either direction, depending on situation
Outcome	<ul style="list-style-type: none"> ■ Misclassification of outcome is similar for exposed and unexposed ■ Generally leads to a bias towards the null* 	<ul style="list-style-type: none"> ■ Misclassification of outcome is different for exposed and unexposed ■ Lead to a bias in either direction, depending on situation

*This holds when exposure and outcome are binary



25% of Exposed Misclassified as Unexposed

Truth (No Misclassification)

	E+	E-
D+	80	40
D-	20	50

$$OR_{\text{true}} = 5.0$$

Observed (25% of E+ Misclassified as E-)

	E+	E-
D+	80-20=60	40+20=60
D-	20-5=15	50+5=55

$$OR_{\text{observed}} = 3.6$$



BREAK

Week 4: Misclassification and Measurement Error

Video 6: Indices of Measurement Accuracy

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Validation Study

- In order to determine the amount of misclassification, conduct validation study to compare a
 - gold standard measure (“truth”)
 - e.g. electronic records of medication use
 - alternative feasible measure (“observed ” or “test”)
 - e.g. self-reported medication use
- 2 x 2 table to compare the performance of the gold standard to the feasible measure

		Truth	
		E+	E-
Observed	T+		
	T-		



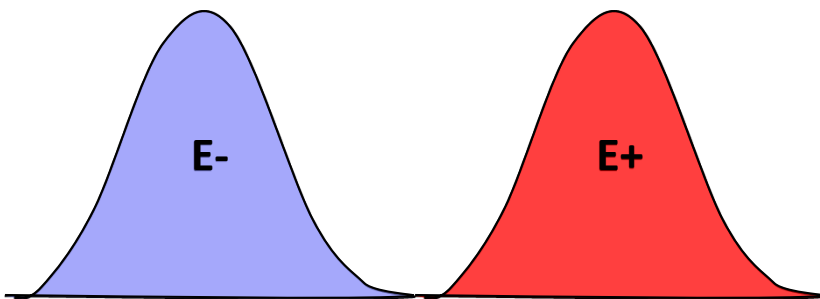
Ideal vs. Observed Exposure Data

Ideal

- Everyone who is truly E+ is classified as E+
- Everyone who is truly E- is classified as E-

Truth

		E+	E-
Observed	T+	80	0
	T-	0	50

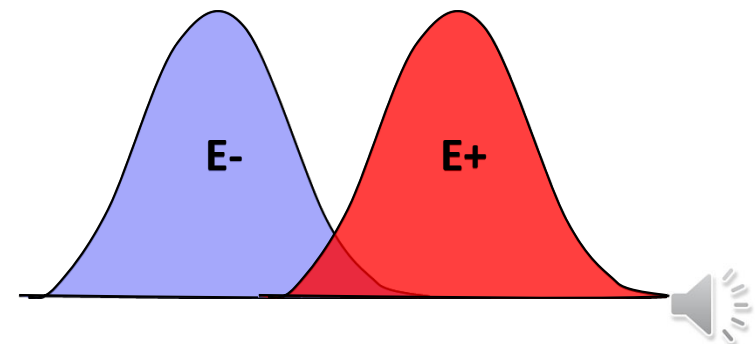


Actual

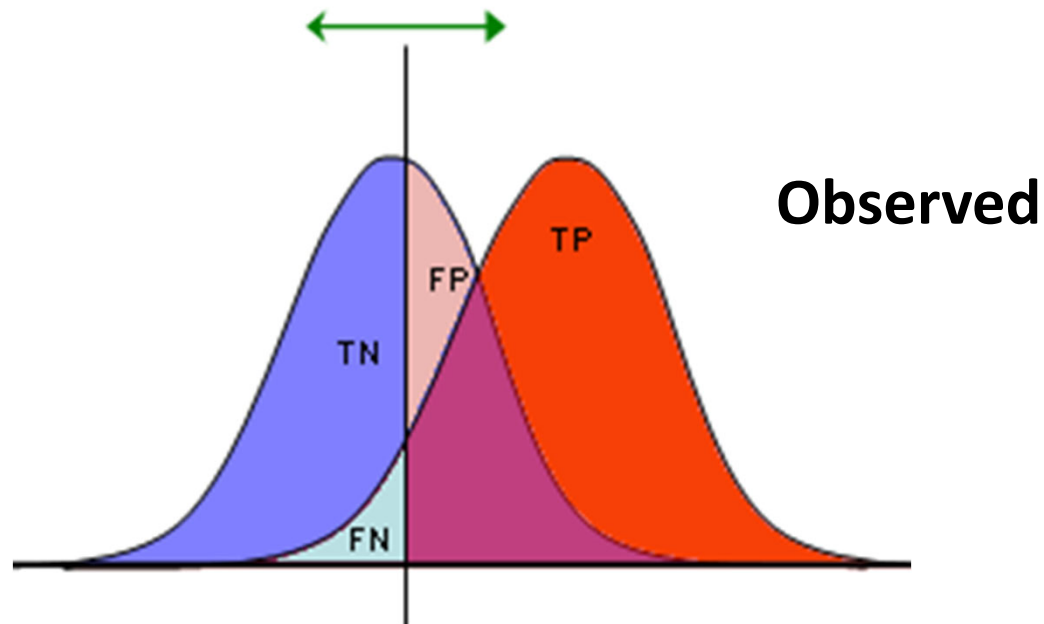
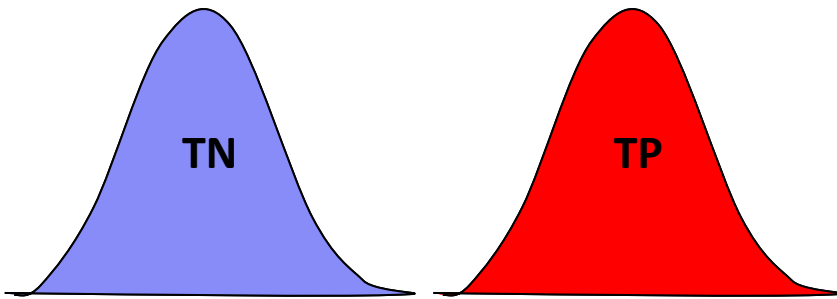
- Some of the truly E+ are incorrectly classified as E-
- Some of the truly E- are incorrectly classified as E+

Truth

		E+	E-
Observed	T+	60	5
	T-	20	45



Indices of Measurement Accuracy (1)



Truth			
		E+	E-
T+		True Positive TP	False Positive FP
T-		False Negative FN	True Negative TN



Indices of Measurement Accuracy (2)

		Truth		
		E+	E-	
Observed	T+	True Positive TP	False Positive FP	TP+FP
	T-	False Negative FN	True Negative TN	FN+TN
		TP+FN	FP+TN	

Sensitivity	$\Pr[T+ E+]$	$TP / (TP + FN)$
Specificity	$\Pr[T- E-]$	$TN / (FP + TN)$
False Negative Rate (1-Sensitivity)	$\Pr[T- E+]$	$FN / (TP + FN)$
False Positive Rate (1-Specificity)	$\Pr[T+ E-]$	$FP / (FP + TN)$
Positive Predictive Value	$\Pr[E+ T+]$	$TP / (TP + FP)$
Negative Predictive Value	$\Pr[E- T-]$	$TN / (TN + FN)$



Sensitivity and Specificity

Ideal

		Truth	
		E+	E-
Observed	T+	80	0
	T-	0	50

Sensitivity

$$TP / (TP + FN) = 80 / (80 + 0) = 100\%$$

Specificity

$$TN / (FP + TN) = 50 / (0 + 50) = 100\%$$

Observed

		Truth	
		E+	E-
Observed	T+	60	5
	T-	20	45

Sensitivity

$$TP / (TP + FN) = 60 / (60 + 20) = 75\%$$

Specificity

$$TN / (FP + TN) = 45 / (5 + 45) = 90\%$$



BREAK

Week 4: Misclassification and Measurement Error

Video 7: Misclassification Correction

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Exposure Misclassification Impact:

Truth \rightarrow Observed
Notation

Truth			
	E+	E-	
D+	a	b	m ₁
D-	c	d	m ₀



Observed			
	E+	E-	
D+	A	B	m ₁
D-	C	D	m ₀

Cases	A:	= (sensitivity)*a + (1-specificity)*b = (TP*a) + (FP*b)
	B:	= (specificity)*b + (1-sensitivity)*a = (TN*b) + (FN*a)
Non-cases	C:	= (sensitivity)*c + (1-specificity)*d = (TP*c) + (FP*d)
	D:	= (specificity)*d + (1-sensitivity)*c = (TN*d) + (FN*c)

Exposure Misclassification Impact:

Truth → Observed
Calculation

Truth		
	E+	E-
D+	a=50	b=500
D-	c=20	d=800

Sensitivity=0.9

Specificity=0.8



Observed		
	E+	E-
D+		
D-		

$$OR_{\text{true}} = \frac{50 \times 800}{500 \times 20} = 4.0$$



Exposure Misclassification Impact:

Observed → Truth

Notation

Based on the formulae above:

Observed			
	E+	E-	
D+	A	B	m_1
D-	C	D	m_0



Truth			
	E+	E-	
D+	a	b	m_1
D-	c	d	m_0

Cases	a:	$= \frac{(\text{specificity}) * m_1 - \mathbf{B}}{\text{sensitivity} + \text{specificity} - 1}$
	b:	$= \frac{(\text{sensitivity}) * m_1 - \mathbf{A}}{\text{sensitivity} + \text{specificity} - 1}$
Non-cases	c:	$= \frac{(\text{specificity}) * m_0 - \mathbf{D}}{\text{sensitivity} + \text{specificity} - 1}$
	d:	$= \frac{(\text{sensitivity}) * m_0 - \mathbf{C}}{\text{sensitivity} + \text{specificity} - 1}$

Note: Since we are assuming no misclassification of disease,

$$m_1 = \mathbf{A} + \mathbf{B} = \mathbf{a} + \mathbf{b}$$

and

$$m_0 = \mathbf{C} + \mathbf{D} = \mathbf{c} + \mathbf{d}$$



Exposure Misclassification Impact:

Observed → Truth

Notation

Observed			
	E+	E-	
D+	90	60	150
D-	300	500	800

Sensitivity=0.96

Specificity=0.70



Truth			
	E+	E-	
D+	68	82	150
D-	91	709	800

$$a = \frac{(0.70)150 - 60}{0.96 + 0.70 - 1} = \frac{45}{0.66} = 68$$

$$c = \frac{(0.70)800 - 500}{0.96 + 0.70 - 1} = \frac{60}{0.66} = 91$$

$$b = \frac{(0.96)150 - 90}{0.96 + 0.70 - 1} = \frac{54}{0.66} = 82$$

$$d = \frac{(0.96)800 - 300}{0.96 + 0.70 - 1} = \frac{468}{0.66} = 709$$

$$OR_{\text{observed}} = \frac{90 \times 500}{60 \times 300} = 2.5$$

$$OR_{\text{true}} = \frac{68 \times 709}{82 \times 91} = 6.5$$



BREAK

Week 4: Misclassification and Measurement Error

Video 8: A Worked Example

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Amoxicillin and Oral Clefts

- In a study to evaluate the association between maternal exposure to amoxicillin and the risk of oral clefts (*Epidemiology*. 2012;23(5):699-705), amoxicillin use was ascertained via interview within 6 months of delivery.
- Based on how exposure information was collected, do you think misclassification of the exposure is likely?



Amoxicillin and Oral Clefts

- In a study to evaluate the association between maternal exposure to amoxicillin and the risk of oral clefts in the baby (*Epidemiology*. 2012;23(5):699-705), amoxicillin use was ascertained via interview within 6 months of delivery.
- Based on how exposure information was collected, do you think misclassification of the exposure is likely?

Yes, because the information on prior exposure had to be recalled and it was collected after the occurrence of the outcome.

In this situation, misclassification of the exposure is a particularly serious threat to the study's validity because mothers with malformed babies may report their exposure to amoxicillin more completely than mothers of children without cleft lip/palate (recall bias), which would tend to make amoxicillin appear associated with the outcome, even in the absence of a true effect of amoxicillin use on the risk of oral clefts.



Amoxicillin and Oral Clefts

- Calculate the odds ratio based on the observed data below:

Amoxicillin Use in First Trimester			
	Yes	No	
Cases	28	810	838
Controls	144	6379	6523
	172	7189	7631

- $$\widehat{OR}_{observed} = \frac{(28 \times 6379)}{(810 \times 144)} = 1.53$$



Amoxicillin and Oral Clefts

- Suppose the authors were concerned about exposure misclassification, so they conducted a validation study among a subgroup of their participants, comparing self-reported amoxicillin use during pregnancy (as done in this paper) to true amoxicillin use during pregnancy based on a review of medical and pharmacy records. The table below shows the results.

		Truth (medical/pharmacy records)	
		Yes	No
Observed (self-report)	Yes	63	4
	No	21	196

- Calculate the following indices of measurement accuracy. Round all results to two decimal points.
 - ☐ Sensitivity
 - ☐ Specificity
 - ☐ False Positive Rate
 - ☐ False Negative Rate
 - ☐ Positive Predictive Value
 - ☐ Negative Predictive Value



Amoxicillin and Oral Clefts

- Suppose the authors were concerned about exposure misclassification, so they conducted a validation study among a subgroup of their participants, comparing self-reported amoxicillin use during pregnancy (as done in this paper) to true amoxicillin use during pregnancy based on a review of medical and pharmacy records. The table below shows the results.

		Truth (medical/pharmacy records)	
		Yes	No
Observed (self-report)	Yes	63	4
	No	21	196

- Calculate the following indices of measurement accuracy. Round all results to two decimal points.
 - ☐ Sensitivity $= \Pr(T+ | E+) = 63/(63+21) = 0.75$
 - ☐ Specificity $= \Pr(T- | E-) = 196/(4+196) = 0.98$
 - ☐ False Positive Rate $= 1 - \text{specificity} = 0.02$
 - ☐ False Negative Rate $= 1 - \text{sensitivity} = 0.25$
 - ☐ Positive Predictive Value $= \Pr[E+ | T+] = 63/(63+4) = 0.94$
 - ☐ Negative Predictive Value $= \Pr[E- | T-] = 196/(196+21) = 0.90$



Amoxicillin and Oral Clefts

Using the observed data and the indices of measurement accuracy from the validation study, construct a 2x2 table for the expected true frequencies and calculate the corrected (true) odds ratio for the association between maternal exposure to amoxicillin and oral clefts, assuming that there is nondifferential independent exposure misclassification and no other sources of bias

Observed			
Amoxicillin Use in First Trimester			
	Yes	No	
Cases	28	810	838
Controls	144	6379	6523
	172	7189	7631

$OR_{\text{observed}} =$

Truth		
Amoxicillin Use in First Trimester		
	Yes	No
Cases		
Controls		

$OR_{\text{true}} =$

Cases	$a = \frac{(\text{specificity}) * m_1 - \text{B}}{\text{sensitivity} + \text{specificity} - 1} =$
	$b = \frac{(\text{sensitivity}) * m_1 - \text{A}}{\text{sensitivity} + \text{specificity} - 1} =$
Non-cases	$c = \frac{(\text{specificity}) * m_0 - \text{D}}{\text{sensitivity} + \text{specificity} - 1} =$
	$d = \frac{(\text{sensitivity}) * m_0 - \text{C}}{\text{sensitivity} + \text{specificity} - 1} =$



Amoxicillin and Oral Clefts

Using the observed data and the indices of measurement accuracy from the validation study, construct a 2x2 table for the expected true frequencies and calculate the corrected (true) odds ratio for the association between maternal exposure to amoxicillin and oral clefts, assuming that there is nondifferential independent exposure misclassification and no other sources of bias

Observed			
Amoxicillin Use in First Trimester			
	Yes	No	
Cases	28	810	838
Controls	144	6379	6523
	172	7189	7631

$$OR_{\text{observed}} = (AD)/(BC) = (28*6379)/(810*144) = 1.53$$

Truth			
Amoxicillin Use in First Trimester			
	Yes	No	
Cases	15.4	822.6	838
Controls	18.5	6504.5	6523
	33.9	7327.1	7361

$$OR_{\text{true}} = (ad)/(bc) = (15.4*6504.5)/(822.6*18.5) = 6.56$$

Cases	a	$= \frac{(\text{specificity}) * m_1 - \mathbf{B}}{\text{sensitivity} + \text{specificity} - 1} = \frac{0.98 * 838 - 810}{0.75 + 0.98 - 1} = 15.4$
	b	$= \frac{(\text{sensitivity}) * m_1 - \mathbf{A}}{\text{sensitivity} + \text{specificity} - 1} = \frac{0.75 * 838 - 28}{0.75 + 0.98 - 1} = 822.6$
Non-cases	c	$= \frac{(\text{specificity}) * m_0 - \mathbf{D}}{\text{sensitivity} + \text{specificity} - 1} = \frac{0.98 * 6523 - 6379}{0.75 + 0.98 - 1} = 18.5$
	d	$= \frac{(\text{sensitivity}) * m_0 - \mathbf{C}}{\text{sensitivity} + \text{specificity} - 1} = \frac{0.75 * 6523 - 144}{0.75 + 0.98 - 1} = 6504.5$



BREAK

Week 4: Misclassification and Measurement Error

Video 9: Impact of a Misclassified Exposure

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Direction of Bias from Exposure Misclassification

- Differential misclassification can either overestimate or underestimate the true association, depending on the situation.
- On the other hand, the direction of bias produced by independent non-differential misclassification of a dichotomous exposure is toward the null value of no association.
- If the independent non-differential misclassification is worse than randomly classifying people as exposed or unexposed, the bias can go beyond the null value and reverse the direction of the association.

Non-differential Misclassification of a Dichotomous Exposure

- If the sum of the sensitivity and specificity is greater than 1, the expected estimate will be biased towards the null.
- If the sum of the sensitivity and specificity is less than 1, the expected estimate will be in the opposite direction of the actual effect.
- If sensitivity and specificity sum to 1, the resulting expected estimate will be null, regardless of the magnitude of the effect.
- If both sensitivity and specificity are zero, the resulting expected odds ratio is the inverse of the correct value.

Table 9-2 Nondifferential Misclassification with Two Exposure Categories			
	E+	E-	
Correct data			OR = 3.0
Cases	240	200	
Controls	240	600	
Sensitivity = 0.8			OR = 2.6
Specificity = 1.0			
Cases	192	248	
Controls	192	648	
Sensitivity = 0.8			OR = 1.9
Specificity = 0.8			
Cases	232	208	
Controls	312	528	
Sensitivity = 0.4			OR = 1.0
Specificity = 0.6			
Cases	176	264	
Controls	336	504	
Sensitivity = 0.0			OR = 0.33
Specificity = 0.0			
Cases	200	240	
Controls	600	240	

Non-differential Misclassification May Lead to Estimates Farther from the Null

- Non-differentiality alone does not guarantee bias toward the null.
- Non-differential exposure or disease misclassification can sometimes produce bias away from the null
 - if the exposure or disease variable has more than two levels or
 - if the classification errors depend on errors made in other variables
- In a given study, random fluctuations in the errors produced by a method may lead to estimates that are farther from the null than what they would be if no error were present, even if the method satisfies all the conditions that guarantee that the expected value of the bias is toward the null.

Non-differential Misclassification

Exposure with >2 Categories

- When the exposure is polytomous (i.e., >2 categories), the bias from independent non-differential misclassification of the exposure for a given comparison may be away from the null value.
- When the exposure is polytomous and there is non-differential misclassification between two of the categories and no others, the estimates for those two categories will be biased toward one another.

Table 9-3: Misclassification with Three Exposure Categories

	No Exposure	Low Exposure	High Exposure
Correct data			
Cases	100	200	600
Controls	100	100	100
		OR = 2	OR = 6
40% of high exposure → low exposure			
Cases	100	440	360
Controls	100	140	60
		OR = 3.1	OR = 6

The bias in the estimate for the low-exposure category is toward that of the high-exposure category and away from the null value.

Impact of Non-Differential Exposure Misclassification

When **sensitivity + specificity** > 1, binary non-differential misclassification leads, on average, to:

- bias in the odds ratio
- confidence intervals which are incorrectly centered and too narrow
- reduction in the **power** of hypothesis tests
- sample size calculations which have less than the specified power (Type II error is inflated)
- valid Neyman-Pearson tests of the null hypothesis (i.e. when the null is true, the hypothesis test will have the correct Type I error rate and alpha will retain its precise meaning)
 - Note that despite the valid N-P hypothesis tests, the p-value **does not** retain its precise meaning



BREAK

Week 4: Misclassification and Measurement Error

Video 10: Impact of a Misclassified Confounder

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Non-differential Misclassification

Residual Confounding

- Errors in confounder measurement compromise our ability to control for confounding, leaving **residual confounding** (i.e., confounding left after control of the available confounder measurements).
- Independent non-differential misclassification of a dichotomous confounding variable will reduce the degree to which the confounder can be controlled, and thus causes a bias in the direction of the confounding by the variable.
- This will result in non-exchangeability, even when the DAG is correctly specified and the analysis is designed correctly.
- The expected result will lie between the unadjusted association and the correctly adjusted association that would have obtained if the confounder had not been misclassified.
- Hypothesis tests of the association between exposure and outcome are not valid.
- Confidence intervals are incorrectly centered and of incorrect width.



Non-differential Misclassification

Effect Measure Modification

- The degree of residual confounding left within strata of the misclassified confounder often will differ across strata, which will distort the apparent degree of heterogeneity (effect modification) across strata.
- Independent non-differential misclassification of either the confounder or exposure can therefore give rise to the appearance of effect-measure modification when in fact there is none or mask the appearance of such modification when in fact it is present.



A WORKED EXAMPLE

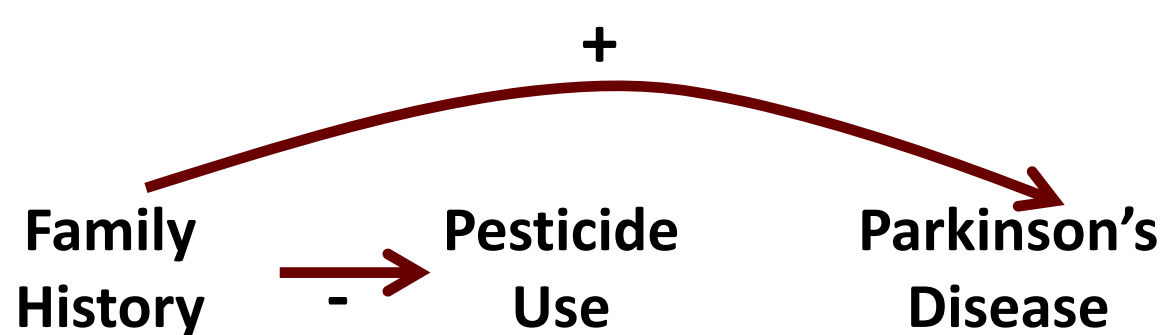


Pesticides and Parkinson's Disease

Confounding by Family History

- People with a family history of Parkinson's disease may be more careful to avoid pesticide use and they are more likely to develop Parkinson's even in the absence of pesticide exposure.

- DAG



Pesticides and Parkinson's Disease

Correctly Classified Pesticide Use

Sensitivity=1.0 and Specificity=1.0

Crude

	Pesticide	No Pesticide
Cases	900	1,890
Controls	3,150	13,500

OR = 2.04

Family History

	Pesticide	No Pesticide
Cases	450	1,350
Controls	900	6,750

OR = 2.50

No Family History

	Pesticide	No Pesticide
Cases	450	540
Controls	2,250	6,750

OR = 2.50

P-heterogeneity=1.0



Pesticides and Parkinson's Disease

Nondifferentially Misclassified Pesticide Use

Sensitivity=1.0 and Specificity=0.8

Crude

	Pesticide	No Pesticide
Cases	1,278	1,512
Controls	5,850	10,800

OR = 1.56

Family History

	Pesticide	No Pesticide
Cases	720	1,080
Controls	2,250	5,400

OR = 1.60

No Family History

	Pesticide	No Pesticide
Cases	558	432
Controls	3,600	5,400

OR = 1.94

P-heterogeneity=0.027



BREAK