

PhD course in Basic Biostatistics - Day 7

Erik Parner, Department of Biostatistics, Aarhus University®

Post term delivery and parity

A two by two table

Odds and odds ratios

Odds ratio and the relative risk

Estimating the odds ratio

Relative measures - why inference on the log-scale?

Odds ratios, Confidence Intervals and testing

Odds ratio and logistic regression

Post term delivery and age of the woman

A simple logistic regression model

Understanding the model and the parameters

The estimates

Formulations

The risk versus age curves

Erik Parner

Basic Biostatistics - Day 7

1

Post term delivery and parity - adjusting for age

Two simple logistic regression models

Comparing the slopes

A logistic regression model with the same slope

Understanding the parameters

The estimates

Formulation

A risk versus age and parity plot

Why use logistic regression

Linear and logistic regression models - a comparison

Why do we need regression models

Adjustment

Effect modification

Prediction

Erik Parner

Basic Biostatistics - Day 7

2

Example: Post term delivery and parity**Question:** How does the risk of post term delivery depend on parity?**Data:** Parity and gestational age for 12,311 women in the age of 20 to 39. Post term delivery defined as a gestational age larger than 40 weeks.

Parity	N	Postterm	Risk
First child	5,938	1,722	29.0 (28.8; 30.2)%
Not first child	6,373	1,677	26.3 (25.2; 27.4)%
Total	12,311	3,399	27.6 (26.8; 28.4)%

Model: Independent samples from two binomial distributions.Let π_0 and π_1 be the probability (risk) of post term delivery among women giving birth to their first child or not, respectively.

Erik Parner

Basic Biostatistics - Day 7

3

Example: Post term delivery and parity

The assumptions behind the model was discussed on day 4.

On that day we also looked at three different measures of associations: Risk Difference, Relative Risk and Odds Ratio. And the chi-squared test for no association.

Today we will look closer at the Odds Ratio.

Risk difference	-2.7 (-4.3; -1.1)%
Relative risk	0.91 (0.86; 0.96)
Odds ratio	0.87 (0.81; 0.95)

$$X^2=11.09 \quad p=0.001$$

In the table above we compare π_1 to π_0 , i.e. women giving birth to their first child is the **reference group**.

We see that the risk is (statistically significant) smaller if the woman already had a child.

Erik Parner

Basic Biostatistics - Day 7

4

Odds and risk

The **odds** is defined as $\pi/(1-\pi)$, i.e. the probability of post term delivery divided by the probability of not having a post term delivery.

$$\text{odds} = \frac{\pi}{1-\pi}$$

If the odds is equal to 0.5=1/2, then the risk of post term delivery is only **half** of the risk of not having a post term delivery.

We can also go from odds to risk:

$$\pi = \frac{\text{odds}}{1 + \text{odds}}$$

We see that

$$\text{odds} = 0.5 \text{ gives } \pi = 0.5/(1+0.5)=0.3333.$$

$$\text{odds} = 1 \text{ gives } \pi = 1/(1+1)=0.5.$$

Erik Parner

Basic Biostatistics - Day 7

5

Odds and odds ratios

The odds ratio comparing parity>0 to the reference is given by

$$OR_{10} = \frac{\text{odds}_1}{\text{odds}_0} = \frac{\pi_1/(1-\pi_1)}{\pi_0/(1-\pi_0)} = \frac{\pi_1 \cdot (1-\pi_0)}{\pi_0 \cdot (1-\pi_1)}$$

It is easily seen that $\pi_1 = \pi_0 \Leftrightarrow \text{odds}_1 = \text{odds}_0 \Leftrightarrow OR = 1$

OR has **nice properties**:

Switching reference group or event will just lead to 1/OR, e.g.

$$OR_{01} = \frac{\text{odds}_0}{\text{odds}_1} = \frac{1}{OR_{10}}$$

And of course the estimates and confidence intervals will transform similarly.

$$OR_{01} : \frac{1}{0.87} \left(\frac{1}{0.95}; \frac{1}{0.81} \right) = 1.14(1.06; 1.24)$$

Erik Parner

Basic Biostatistics - Day 7

6

Odds ratios and relative risks

The odds ratio is related to the relative risk:

$$OR_{10} = \frac{\pi_1 \cdot (1-\pi_0)}{\pi_0 \cdot (1-\pi_1)} = RR_{10} \cdot \frac{(1-\pi_0)}{(1-\pi_1)}$$

We can see that if **the event is rare**, i.e. both π_1 and π_0 are small, then the last ratio is close to 1/1=1.

So for a **rare event** we have:

$$OR \approx RR$$

Erik Parner

Basic Biostatistics - Day 7

7

Estimating the odds ratios

The odds ratio is of course estimated by:

$$\widehat{OR}_{10} = \frac{\hat{\pi}_1 \cdot (1-\hat{\pi}_0)}{\hat{\pi}_0 \cdot (1-\hat{\pi}_1)}$$

Another way to find the estimate is to make the 'classical' 2x2 table:

Exposed	Event	
	Yes	No
Yes	a	b
No	c	d

$$\widehat{OR}_{10} = \frac{a \cdot d}{b \cdot c}$$

Parity>0	Post term	
	Yes	No
Yes	1,677	4,696
No	1,722	4,216

$$\widehat{OR}_{10} = \frac{1,677 \cdot 4,216}{1,722 \cdot 4,696} = 0.8743$$

Erik Parner

Basic Biostatistics - Day 7

8

Odds ratios - why inference on the log-scale

The odds ratio is limited to be positive.

A value in the interval 0 to 1 corresponds to lower risk among the Parity>0.

A value from 1 to infinity corresponds to higher risk among the Parity>0

If we switch "exposed" and "unexposed" we get

Erik Parner Basic Biostatistics - Day 7 9

Odds ratios - why inference on the log-scale

The log-odds ratio is not limited.

A value in the interval -infinity to 0 corresponds to lower risk among the Parity>0.

A value from 0 to infinity to higher risk among the Parity>0

If we switch "exposed" and "unexposed" we get

Symmetry on the log scale!!!

Erik Parner Basic Biostatistics - Day 7 10

Odds ratios - Approx. CI (Woolf/Wald)

So on the log scale we have a symmetric measure of association.

On the log scale it makes sense to find the CI as 'usual', i.e. as estimate ± 1.96 * se .

Using the notation from page 8 we have:

$$se(\ln(\widehat{OR})) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

$$se(\ln(\widehat{OR})) = \sqrt{\frac{1}{1,677} + \frac{1}{4,696} + \frac{1}{1,711} + \frac{1}{4,216}} = 0.0403$$

$$CI \ln(OR) = \ln(0.8743) \pm 1.96 \cdot 0.0403 = (-0.2134; -0.0552)$$

$$CI OR = (\exp(-0.2134); \exp(-0.0552)) = (0.81; 0.95)$$

Erik Parner Basic Biostatistics - Day 7 11

Odds ratios - Testing

If one wants to test a hypothesis that the odds ratio has a specific value: $OR = OR_0$,

then this is also done on the log-scale:

$$z_{obs} = \frac{\ln(\widehat{OR}) - \ln(OR_0)}{se(\ln(\widehat{OR}))}$$

Could the odds be reduced by 10%, i.e. H: $OR = 0.9$?

$$z_{obs} = \frac{\ln(0.8743) - \ln(0.9)}{0.0403} = -0.719$$

$$p = 2 \cdot \Pr(z > |z_{obs}|) = 2 \cdot \Pr(z < -|z_{obs}|) = 0.47$$

The hypothesis cannot be rejected.

```
disp ( ln(0.8743)-ln(0.9) )/ 0.0403
disp 2*normal( abs(-0.719) )
```

Erik Parner Basic Biostatistics - Day 7 12

Odds ratios - logistic regression

Here we will see how one can find the odds ratio by **logistic regression**.

Let **Par1** be an **indicator variable** for Parity>0 ,
i.e. **Par1**=1 if parity>0 and **Par1**=0 if Parity=0.

Now we will look at the (logistic regression) model:

$$\log(\text{odds}) = \beta_0 + \beta_1 \cdot \text{Par1}$$

This is equivalent to:

$$\text{odds} = \exp(\beta_0 + \beta_1 \cdot \text{Par1}) = \exp(\beta_0) \cdot \exp(\beta_1)^{\text{Par1}}$$

and

$$\pi = \frac{\exp(\beta_0 + \beta_1 \cdot \text{Par1})}{1 + \exp(\beta_0 + \beta_1 \cdot \text{Par1})}$$

Odds ratios - logistic regression

$$\log(\text{odds}) = \beta_0 + \beta_1 \cdot \text{Par1} \quad \text{odds} = \exp(\beta_0) \cdot \exp(\beta_1)^{\text{Par1}}$$

We see that if **Parity=0** then we have:

$$\log(\text{odds}) = \beta_0 \quad \boxed{\text{odds} = \exp(\beta_0)}$$

and if **Parity>0** then we have

$$\log(\text{odds}) = \beta_0 + \beta_1 \quad \text{odds} = \exp(\beta_0) \cdot \exp(\beta_1)$$

Combining we have

$$\boxed{OR_{10} = \frac{\text{odds}(\text{if parity} > 0)}{\text{odds}(\text{if parity} = 0)} = \frac{\exp(\beta_0) \cdot \exp(\beta_1)}{\exp(\beta_0)} = \exp(\beta_1)}$$

Odds ratios - logistic regression

$$\log(\text{odds}) = \beta_0 + \beta_1 \cdot \text{Par1} \quad \text{odds} = \exp(\beta_0) \cdot \exp(\beta_1)^{\text{Par1}}$$

In summary we have that in the model:

The "intercept" β_0 is the **log odds in the "reference group"**.

The "slope" β_1 is the **log OR**.

That is, we can find the odds ratio from before by what is called a **logistic regression model**.

So the computer will give us **estimates** and **confidence intervals** for the **odds in the reference group** and the **odds ratio** comparing the 'exposed' to the reference.

Post term delivery and age

We know that the age distribution among the two groups of women is different - the women giving birth for the first time will on average be younger!

It might be relevant to compare the two groups after "adjustment for age".

We will start by modeling the association between post term delivery and age among the **women with Parity=0**.

The simplest **logistic regression model** is:

$$\log(\text{odds}) = \alpha_0 + \alpha_1 \cdot \text{Age} \quad \text{odds} = \exp(\alpha_0) \cdot \exp(\alpha_1)^{\text{Age}}$$

To get a sensible reference age:

$$\log(\text{odds}) = \alpha_0 + \alpha_1 \cdot (\text{Age} - 30) \quad \text{odds} = \exp(\alpha_0) \cdot \exp(\alpha_1)^{\text{Age}-30}$$

Post term delivery and age (Parity==0)

$\log(odds) = \alpha_0 + \alpha_1 \cdot (Age - 30) \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)^{Age-30}$

Age = 30: $\log(odds) = \alpha_0 \quad odds = \exp(\alpha_0)$

Age = 31: $\log(odds) = \alpha_0 + \alpha_1 \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)$

Age = 18: $\log(odds) = \alpha_0 - 12 \cdot \alpha_1 \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)^{-12}$

Age = 19: $\log(odds) = \alpha_0 - 11 \cdot \alpha_1 \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)^{-11}$

Age = 25: $\log(odds) = \alpha_0 - 5 \cdot \alpha_1 \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)^{-5}$

$OR_{31 \text{ vs } 30} = \exp(\alpha_0) \cdot \exp(\alpha_1) / \exp(\alpha_0) = \exp(\alpha_1)$

$OR_{19 \text{ vs } 18} = \exp(\alpha_0) \cdot \exp(\alpha_1)^{-11} / (\exp(\alpha_0) \cdot \exp(\alpha_1)^{-12}) = \exp(\alpha_1)$

$OR_{25 \text{ vs } 19} = \exp(\alpha_0) \cdot \exp(\alpha_1)^{-5} / (\exp(\alpha_0) \cdot \exp(\alpha_1)^{-11}) = \exp(\alpha_1)^6$

Erik Parner Basic Biostatistics - Day 7 17

Post term delivery and age (Parity==0)

$\log(odds) = \alpha_0 + \alpha_1 \cdot (Age - 30) \quad odds = \exp(\alpha_0) \cdot \exp(\alpha_1)^{Age-30}$

We saw: **Age = 30:** $odds = \exp(\alpha_0)$

$OR_{31 \text{ vs } 30} = \exp(\alpha_1)$

$OR_{19 \text{ vs } 18} = \exp(\alpha_1)$

$OR_{25 \text{ vs } 19} = \exp(\alpha_1)^6$

That is,

$\exp(\alpha_0)$ is the odds in the reference (**Age==30**)

$\exp(\alpha_1)$ is the **OR** for 1 year difference.

and

$OR_{6 \text{ years}} = OR_{1 \text{ year}}^6$

Erik Parner Basic Biostatistics - Day 7 18

Post term delivery and age (Parity==0)

Using a computer we get:

log odds scale	CI				H = 0	
	est	se	lower	upper	z	p
Const	-0.8446	0.0332	-0.9096	-0.7795	-25.44	<0.0001
Age-30	0.0207	0.0071	0.0069	0.0345	2.93	0.003

Exp odds scale	CI				H = 1	
	est	se	lower	upper	z	p
Const	0.4297		0.4027	0.4586	-25.44	<0.0001
Age-30	1.0209		1.0069	1.0351	2.93	0.003

Note, only the estimates and the confidence intervals should be transformed!

Erik Parner Basic Biostatistics - Day 7 19

Post term delivery and age (Parity==0)

	est	CI		H = 1	
		lower	upper	z	p
Odds if Age=30	0.4297	0.4027	0.4586	-25.44	<0.0001
One years age dif.	1.0209	1.0069	1.0351	2.93	0.003

From odds to probability:

$$\Pr(\text{post term if Age==30}) = \frac{0.4297}{1 + 0.4297} \left(\frac{0.4027}{1 + 0.4027}; \frac{0.4586}{1 + 0.4586} \right)$$

$= 30.1(28.7; 31.4)\%$

Five years age difference:

$$OR_{5 \text{ years}} = 1.0209^5 (1.0069; 1.0351)^5$$

$= 1.11(1.03; 1.19)$

Erik Parner Basic Biostatistics - Day 7 20

Post term delivery and age (Parity==0) - Formulations

Methods

The risk of post term delivery among women giving birth for the first time was described by a logistic regression model with age as a continuous variable. ...

Results

We found that five year age difference corresponds to an odds ratio of 1.11(1.03; 1.19). A 30 year old woman giving birth for the first time has 30(29;31)% risk of post term delivery.

Conclusion

The risk of post delivery among women giving birth for the first time increases with the age of the woman....

Logistic regression checking the model

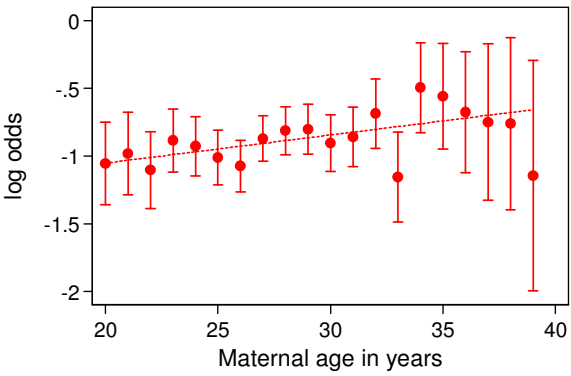
It is outside the scope of this course to go into details on how to check the model, so we will just state the assumptions behind the model:

1. All the observations should be **independent**.
2. There is exactly the same **two possible outcomes** for each observation.
3. The log odds is a **linear function** of age.

The last assumption can to some extent be checked by plotting the fitted regression line and the observed odds (with 95% CI) for each distinct age.

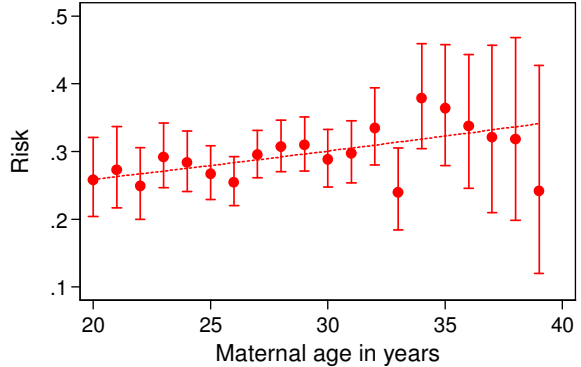
Post term delivery and age (Parity==0)

log(odds) = alpha_0 + alpha_1 * (Age - 30)



Post term delivery and age (Parity==0)

risk = exp(alpha_0 + alpha_1 * (Age - 30)) / (1 + exp(alpha_0 + alpha_1 * (Age - 30)))



Not linear! But almost

Post term delivery and age (Parity==0)

Extrapolating we see the non-linearity.

$$risk = \frac{\exp(\alpha_0 + \alpha_1 \cdot (Age - 30))}{1 + \exp(\alpha_0 + \alpha_1 \cdot (Age - 30))}$$

Erik Parner Basic Biostatistics - Day 7 25

Post term delivery parity adjusting for age

We now know that

- The risk of post term delivery **increases with age** (among women with **parity==0**).
- The risk of post term delivery is **smaller for Parity>0**.
- Women with **Parity>0** are older.

From this we can deduce that **adjusting for age** (if reasonable) will **increase** the difference between the two parity groups.

We now show how to find an age adjusted estimate, when we assume a linear "effect" of age on log odds.

Erik Parner Basic Biostatistics - Day 7 26

Post term delivery and age

We fit the same model to the Parity>0 group and then look at the difference:

log odds	Slope				log odds Age==30			
	est	se	lower	upper	est	se	lower	upper
Parity>0	0.025	0.007	0.012	0.039	-1.037	0.029	-1.093	-0.981
Parity==0	0.021	0.007	0.007	0.035	-0.845	0.033	-0.910	-0.780
Difference	0.005	<u>0.010</u>	-0.015	0.024	-0.193	<u>0.044</u>	-0.279	-0.107

The **standard errors** of the differences are found as usual:

$$se(est_{parity>0} - est_{parity==0}) = \sqrt{se^2_{parity>0} + se^2_{parity==0}}$$

We see that we can assume the slopes to be identical, we could also test the hypothesis:

$$z = \frac{0.005 - 0}{0.010} = 0.5 \quad p = 64\%$$

Erik Parner Basic Biostatistics - Day 7 27

Post term delivery and age - assuming identical slopes

If we assume identical the slopes, then we can write the model:

$$\log(odds) = \gamma_0 + \gamma_1 \cdot (Age - 30) + \gamma_2 \cdot Par1$$

$$odds = \exp(\gamma_0) \cdot \exp(\gamma_1)^{(Age-30)} \cdot \exp(\gamma_2)^{Par1}$$

We see that

- $\exp(\gamma_0)$ is the odds among 30-year old with **Parity==0**.
- $\exp(\gamma_0) \cdot \exp(\gamma_2)$ is the odds among 30-year old with **Parity>0**.
- $\exp(\gamma_1)^A$ is the odds ratio for the age difference **A** years among women in the same **Parity** group.
- $\exp(\gamma_2)$ is the odds ratio comparing **Parity>0** to **Parity==0**, at the same age.

Erik Parner Basic Biostatistics - Day 7 28

Post term delivery and age - assuming identical slopes

The model is easily fitted by a computer:

log odds	Slope				log odds Age==30			
	est	se	lower	upper	est	se	lower	upper
Parity>0	0.023	0.005	0.013	0.033	-1.036	0.029	-1.092	-0.981
Parity==0					-0.839	0.031	-0.900	-0.778
Difference	0				-0.197	0.043	-0.281	-0.114

The **age adjusted OR** comparing Parity>0 to Parity==0:

$OR_{10} : \exp(-0.197) (\exp(-0.281); \exp(-0.114)) = 0.821 (0.755; 0.892)$

and if we compare Parity==0 to Parity>0:

$OR_{01} : \frac{1}{0.821} \left(\frac{1}{0.892}; \frac{1}{0.755} \right) = 1.22 (1.12; 1.32)$

Erik Parner Basic Biostatistics - Day 7 29

Post term delivery and Parity - Formulations

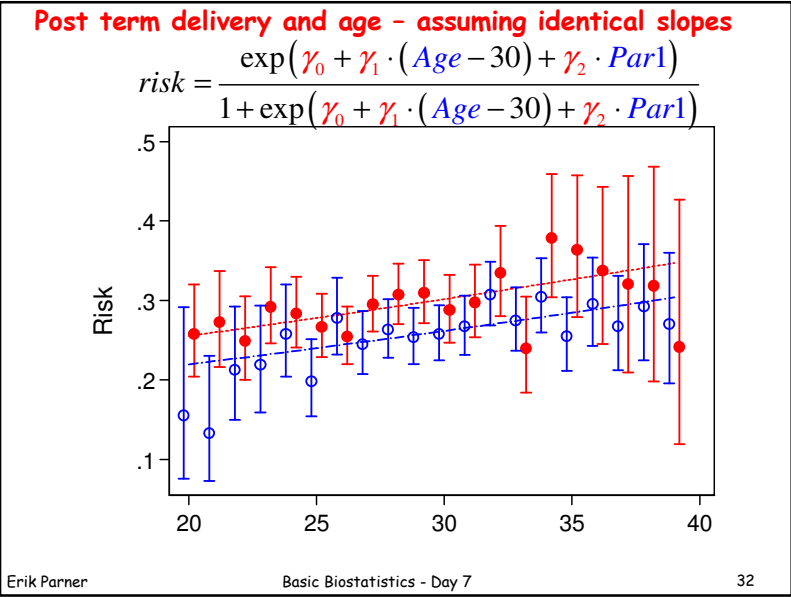
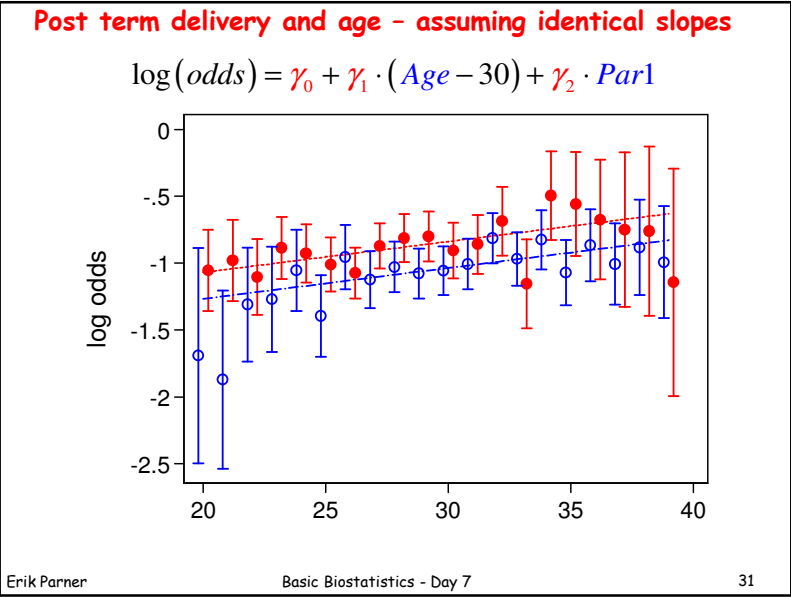
Methods
The risk of post term delivery among women was modeled by a logistic regression model with age as a continuous variable....

Results version1
Comparing Parity >0 to Parity==0 the crude odds ratio was 0.87(0.81;0.95). The age adjusted odds ratio was 0.82(0.76;0.89).

Results version2
Comparing Parity==0 to Parity>0 the crude odds ratio was 1.14(1.06;1.24). The age adjusted odds ratio was 1.22(1.12;1.32).

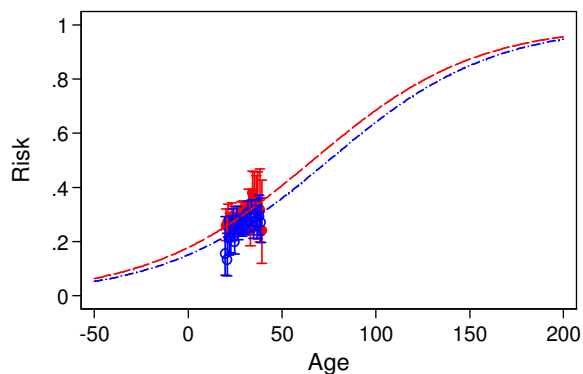
Conclusion ??
Women giving birth for the first time have up to 32% higher odds (risk) of post term delivery compared to other women of the same age.

Erik Parner Basic Biostatistics - Day 7 30



Post term delivery and age - assuming identical slopes

$$\text{risk} = \frac{\exp(\gamma_0 + \gamma_1 \cdot (\text{Age} - 30) + \gamma_2 \cdot \text{Par1})}{1 + \exp(\gamma_0 + \gamma_1 \cdot (\text{Age} - 30) + \gamma_2 \cdot \text{Par1})}$$



Erik Parner

Basic Biostatistics - Day 7

33

Why use logistic regression

There is a **long tradition** for using logistic regression when considering binary outcome. Some of the reasons are:

It is the **mathematical nicest** model for binary outcome, and hence the **first** type of models that was included in the statistical software.

If you have a **case-control** design, then you want to work with odds ratios.

If the **event is rare**, then it will give you **relative risk** estimates.

It is one of the few models for binary data that ensures that the estimated probability is **between zero and one**.

Erik Parner

Basic Biostatistics - Day 7

34

Linear and logistic regression - a comparison

In a **linear regression** the **outcome** is **continuous**:
Lung function, Blood pressure, BMI, concentrations...

In a **logistic regression** the **outcome** is **binary**:
Post term delivery, gender, dead/alive, sick/ well, BMI>30..

Neither of the models make **any assumptions** about the **explanatory variable**!!

In both models they can be continuous, binary or categorical.

In **both models** we have to assume **independence** between observations.

In **both models** we assume linearity -
of expected value or the log odds.

Both models are readily fitted by standard statistical packages.

Erik Parner

Basic Biostatistics - Day 7

35

Regression models in general - why? Adjustment

You have now looked at two of the most commonly used regression models in their most simple forms, involving one continuous and one binary explanatory variable.

You have seen how one can use such models for **adjustment**:
What is the 'effect' of the binary 'exposure' when **adjusting** for the continuous variable?

Exactly the same models could answer the question:
What is the 'effect' (slope) of the **continuous** 'exposure' when **adjusting** for the binary variable?

E.g. what is the increase in risk of post term delivery associated with age when we **adjust** for parity?

Often one has several explanatory variables, a mixture of continuous, binary and categorical and the purpose is to **adjust for more than one**.

In such case one might apply a **multiple regression** model.

Erik Parner

Basic Biostatistics - Day 7

36

Regression models in general - why? Effect modification

We have also seen how we could compare the 'effect' of one explanatory variable for subgroups described by another explanatory variable (**effect modification**):

What is the difference in the PEFR-height relationship for men and women?

What is the difference in the Risk-age relationship for the two parity groups?

Typically by **comparing the slopes**.

Often one has several explanatory variables, a mixture of continuous, binary and categorical and the purpose is to **model effect modification** between explanatory variables.

In such case one might apply a **multiple regression** model.

Regression models in general - why? Prediction

We could also have used exactly the same models for **prediction/prognosis**:

What is the expected PEFR for a person with a given sex and a given height?

What is the risk of post term delivery for women of a given age having her first child?

Often one has several explanatory variables, a mixture of continuous, binary and categorical and the purpose is to make **prediction** for a person with a given set of characteristics.

In such a case one might apply a **multiple regression** model.