

Methods of Applied Statistics I

STA2101H F LEC9101

Week 7

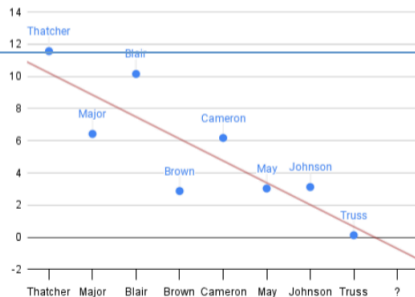
October 26 2022



Rob Sansom
@Sansom_Rob

Following current trends, the next PM will be in office for approximately minus 200 days

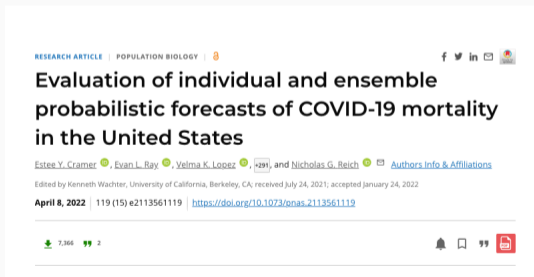
PM Tenure (years)



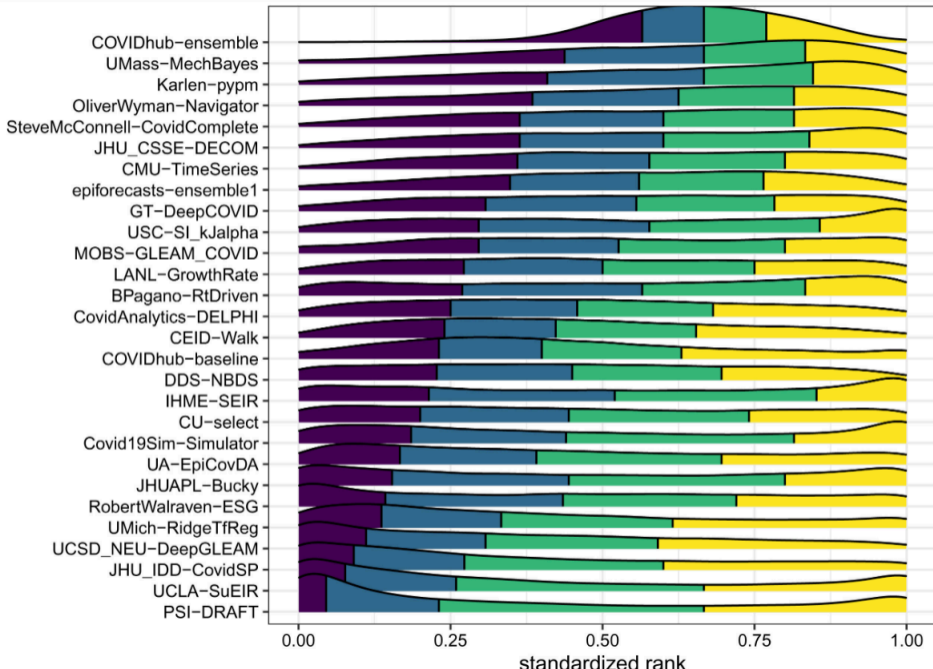
1:59 PM · Oct 20, 2022 · Twitter Web App

1. Upcoming events
2. Recap
3. Likelihood theory and logistic regression
4. Observational studies and causality
5. In the News
6. Hour 3: Comments on HW 1-6 estimates of effect size, missing data
7. Office Hour Wednesday October 26 4-5 pm on Zoom

- Monday October 24 3.30-4.30 : DoSS Seminar Room 9014 (Hydro Building)
- Data Science Seminar Series
- **Daniel McDonald, U Chicago**
- Markov-Switching State Space Models for Uncovering Musical Interpretation



model



Upcoming

- October 27 3.30-4.30 Statistical Sciences Seminar; Room 9014, Hydro Building and [online](#)
- **Mireille Schnitzer, U Montreal**
“Outcome-Adaptive LASSO for Confounder Selection With Time-Varying Treatments”



Recap

- regression models for binomial and binary data
- examples: O-ring failure; heart disease; nodal involvement
- inference, residuals, diagnostics, analysis of deviance, nested models oct19.pdf: 20–24
- covariate classes; binary data
- model selection with

$$AIC = -2\ell(\hat{\beta}; y) + 2p$$

$$BIC = 2\ell(\hat{\beta}; y) + \log(n)p$$

Likelihood theory

- model: $y_i \sim f(y_i; \theta), i = 1, \dots, n; \quad \theta \in \Theta \subset \mathbb{R}^p$
- joint density: $f(\underline{y}; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- likelihood function $L(\theta; \underline{y}) = f(\underline{y}; \theta)$

independent

Likelihood theory

- model: $y_i \sim f(y_i; \theta), i = 1, \dots, n; \quad \theta \in \Theta \subset \mathbb{R}^p$ independent
- joint density: $f(\underline{y}; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- likelihood function $L(\theta; \underline{y}) = f(\underline{y}; \theta)$
- log-likelihood function $\ell(\theta; \underline{y}) = \log L(\theta; \underline{y}) = \sum_{i=1}^n \log f(y_i; \theta)$
- maximum likelihood estimate $\hat{\theta} = \arg \sup \ell(\theta; \underline{y});$ $\ell'(\hat{\theta}) = \mathbf{0}$
- Fisher information $j(\hat{\theta}) = -\ell''(\hat{\theta})$

Likelihood theory

- model: $y_i \sim f(y_i; \theta), i = 1, \dots, n; \quad \theta \in \Theta \subset \mathbb{R}^p$ independent
- joint density: $f(\underline{y}; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- likelihood function $L(\theta; \underline{y}) = f(\underline{y}; \theta)$

- log-likelihood function $\ell(\theta; \underline{y}) = \log L(\theta; \underline{y}) = \sum_{i=1}^n \log f(y_i; \theta)$
- maximum likelihood estimate $\hat{\theta} = \arg \sup \ell(\theta; \underline{y});$ $\ell'(\hat{\theta}) = \mathbf{0}$
- Fisher information $j(\hat{\theta}) = -\ell''(\hat{\theta})$

- two theorems:

$$j^{1/2}(\hat{\theta})(\hat{\theta} - \theta) \xrightarrow{d} N_p(\mathbf{0}, I)$$

asymptotically normal

- likelihood ratio statistic

$$w(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{d} \chi_p^2$$

p is dimension of θ

\xrightarrow{d} is convergence in distribution

- two theorems:

$$\begin{aligned} j^{1/2}(\hat{\theta})(\hat{\theta} - \theta) &\xrightarrow{d} N(0, I) \\ w(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} &\xrightarrow{d} \chi_p^2 \end{aligned}$$

- two theorems:

$$\begin{aligned} j^{1/2}(\hat{\theta})(\hat{\theta} - \theta) &\xrightarrow{d} N(0, I) \\ w(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} &\xrightarrow{d} \chi_p^2 \end{aligned}$$

- two approximations

$$\begin{aligned} \hat{\theta} &\sim N_d\{\theta, j^{-1}(\hat{\theta})\} \\ \hat{\theta}_k &\sim N(\{\theta_k, j^{-1}(\hat{\theta})_{kk}\}) \end{aligned}$$

$$w(\theta) \sim \chi_p^2$$

- two theorems:

$$j^{1/2}(\hat{\theta})(\hat{\theta} - \theta) \xrightarrow{d} N(0, I)$$
$$w(\theta) = 2\{\ell(\hat{\theta}) - \ell(\theta)\} \xrightarrow{d} \chi_p^2$$

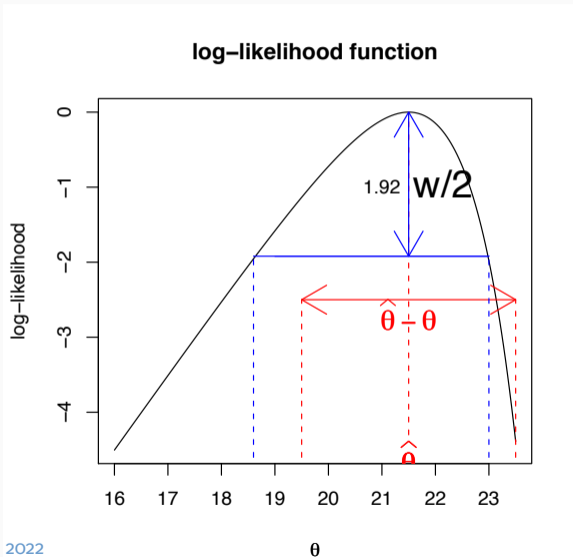
- two approximations

$$\hat{\theta} \sim N_d\{\theta, j^{-1}(\hat{\theta})\}$$
$$\hat{\theta}_k \sim N(\{\theta_k, j^{-1}(\hat{\theta})_{kk}\})$$

$$w(\theta) \sim \chi_p^2$$

- compare two models using **change in likelihood ratio statistic**

nested models



STA2101: Likelihood Cheatsheet

$Y = Y_1, \dots, Y_n$ independently distributed with densities $f(y_i | x_i; \theta), \theta \in \Theta \subset \mathbb{R}^p; y_i \in \mathbb{R}$. The observations are independent, but not identically distributed, due to the dependence on the $p \times 1$ vector x_i . Independence is critical, but i.d. can usually be handled, so the dependence on x_i below is often suppressed.

Likelihood function is the joint probability of the observations, considered as a function of the parameter

$$L(\theta; y) \propto \prod_{i=1}^n f(y_i | x_i; \theta)$$

Nested models

- Comparing two models:
- likelihood ratio test

$$2\{\ell_A(\hat{\beta}_A) - \ell_B(\hat{\beta}_B)\}$$

compares the maximized log-likelihood function under model A and model B

- example

model A: $\text{logit}(p_i) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$, $\beta_A = (\beta_0, \beta_1, \beta_2)$

model B: $\text{logit}(p_i) = \beta_0 + \beta_1 x_{1i}$, $\beta_B = (\beta_0, \beta_1)$

- when model B is **nested** in model A, LRT is approximately χ^2_ν distributed under model B
- $\nu = \dim(A) - \dim(B)$ theory of profile likelihood

... nested models

```
> logitmodcorrect2 <- glm(cbind(r,m-r) ~ temperature + pressure, family = binomial, data = shuttle2)
```

```
> summary(logitmodcorrect2)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	2.520195	3.486784	0.723	0.4698
temperature	-0.098297	0.044890	-2.190	0.0285 *
pressure	0.008484	0.007677	1.105	0.2691

Null deviance: 24.230 on 22 degrees of freedom

Residual deviance: 16.546 on 20 degrees of freedom

AIC: 36.106

Number of Fisher Scoring iterations: 5

... nested models

```
> logitmodcorrect2 <- glm(cbind(r,m-r) ~ temperature + pressure, family = binomial, data = shuttle2)
```

```
> anova(logitmodcorrect,logitmodcorrect2)
```

Analysis of Deviance Table

Model 1: cbind(r, m - r) ~ temperature

Model 2: cbind(r, m - r) ~ temperature + pressure

	Resid. Df	Resid. Dev	Df	Deviance
1	21	18.086		
2	20	16.546	1	1.5407

...nested models

- Model A: $\text{logit}(p_i) = \beta_0 + \beta_1 \text{temp}_i + \beta_2 \text{pressure}_i$
- Model B: $\text{logit}(p_i) = \beta_0 + \beta_1 \text{temp}_i$
- **nested**: Model B is obtained by setting $\beta_2 = 0$
- Under Model B, the **change in deviance** is (approximately) an observation from a χ^2_1
- $\Pr(\chi^2_1 \geq 1.5407) = 0.22$; this is a p -value for testing $H_0 : \beta_2 = 0$
- so is $1 - \Phi\left\{\frac{\hat{\beta}_2}{\widehat{\text{s.e.}}(\hat{\beta}_2)}\right\} = 1 - \Phi(1.105) = 0.27$

ELM-1 p.30

Binomial likelihood

- model

Binomial likelihood

- model
- likelihood

Binomial likelihood

- model
- likelihood
- log-likelihood

Binomial likelihood

- model
- likelihood
- log-likelihood
- score function

Binomial likelihood

- model
- likelihood
- log-likelihood
- score function
- maximum likelihood estimate

Binomial likelihood

- model
- likelihood
- log-likelihood
- score function
- maximum likelihood estimate
- Fisher information

... binomial likelihood

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$

no regression

... binomial likelihood

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$
- likelihood

no regression

... binomial likelihood

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$
- likelihood
- log-likelihood

no regression

... binomial likelihood

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$
- likelihood
- log-likelihood
- score function

no regression

... binomial likelihood

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$
- likelihood
- log-likelihood
- score function
- maximum likelihood estimate

no regression

- model $y_i \sim \text{Bin}(n_i, p_i), i = 1, \dots, m$
- likelihood
- log-likelihood
- score function
- maximum likelihood estimate
- maximized log-likelihood function

no regression

- regression model is nested in saturated model

-

$$w = 2[\ell(\hat{p}) - \ell\{p(\hat{\beta})\}] \sim \chi_{m-p}^2$$

- logistic regression model $p_i = p_i(\beta) = \text{expit}(x_i^T \beta)$, $\hat{p}_i = p_i(\hat{\beta})$
is **nested** in the **saturated** model $\tilde{p}_i = y_i/n_i$

- logistic regression model $p_i = p_i(\beta) = \text{expit}(\mathbf{x}_i^T \beta)$, $\hat{p}_i = p_i(\hat{\beta})$ is **nested** in the **saturated** model $\tilde{p}_i = y_i/n_i$
- **residual deviance** compares fitted model to saturated model

- logistic regression model $p_i = p_i(\beta) = \text{expit}(x_i^T \beta)$, $\hat{p}_i = p_i(\hat{\beta})$ is **nested** in the **saturated** model $\tilde{p}_i = y_i/n_i$
- **residual deviance** compares fitted model to saturated model
- under the fitted model, approximately distributed as χ^2_{n-p}

- logistic regression model $p_i = p_i(\beta) = \text{expit}(\mathbf{x}_i^T \beta)$, $\hat{p}_i = p_i(\hat{\beta})$ is **nested** in the **saturated** model $\tilde{p}_i = y_i/n_i$
- **residual deviance** compares fitted model to saturated model
- under the fitted model, approximately distributed as χ^2_{n-p} if each n_i “large”

ELM-2 §3.2

```
> summary(Ex1018.glm)
```

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 40.710  on 22  degrees of freedom  
Residual deviance: 18.069  on 17  degrees of freedom  
AIC: 41.69
```

- if $n_i \equiv 1$, then

```
> summary(Ex1018binom.glm)
```

Call:

```
glm(formula = cbind(r, m - r) ~ ., family = binomial, data = nodal2)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.4989	-0.7726	-0.1265	0.7997	1.4351

```
> summary(Ex1018binom.glm)
```

Call:

```
glm(formula = cbind(r, m - r) ~ ., family = binomial, data = nodal2)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.4989	-0.7726	-0.1265	0.7997	1.4351

Deviance: $2 \sum_{i=1}^n [y_i \log\{y_i/n_i p_i(\hat{\beta})\} + (n_i - y_i) \log\{(n_i - y_i)/(n_i - n_i p_i(\hat{\beta}))\}]$

```
> summary(Ex1018binom.glm)
```

Call:

```
glm(formula = cbind(r, m - r) ~ ., family = binomial, data = nodal2)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.4989	-0.7726	-0.1265	0.7997	1.4351

Deviance: $2 \sum_{i=1}^n [y_i \log\{y_i/n_i \hat{p}_i\} + (n_i - y_i) \log\{(n_i - y_i)/(n_i - n_i \hat{p}_i)\}]$

approximately χ^2_{n-q}

$$r_{Di} = \pm \sqrt{(2[y_i \log\{y_i/n_i \hat{p}_i\} + (n_i - y_i) \log\{(n_i - y_i)/(n_i - n_i \hat{p}_i)\}])}$$

- $Y_i \sim \text{Bin}(n_i, p_i) \Rightarrow E(Y_i) = n_i p_i, \quad \text{Var}(Y_i) = n_i p_i (1 - p_i)$
- variance is determined by the mean

- $Y_i \sim \text{Bin}(n_i, p_i) \Rightarrow E(Y_i) = n_i p_i, \quad \text{Var}(Y_i) = n_i p_i (1 - p_i)$
- variance is determined by the mean
- ```
bmod <- glm(cbind(survive,total-survive) ~ location + period, family = binomial,
 data = troutegg)
```

```
summary(bmod)
```

```
Null deviance: 1021.469 on 19 degrees of freedom
```

```
Residual deviance: 64.495 on 12 degrees of freedom
```

```
AIC: 157.03
```

- $Y_i \sim \text{Bin}(n_i, p_i) \Rightarrow E(Y_i) = n_i p_i, \quad \text{Var}(Y_i) = n_i p_i (1 - p_i)$
- variance is determined by the mean
- ```
bmod <- glm(cbind(survive,total-survive) ~ location + period, family = binomial,  
             data = troutegg)
```

```
summary(bmod)
```

```
Null deviance: 1021.469  on 19  degrees of freedom
```

```
## Residual deviance:   64.495  on 12  degrees of freedom
```

```
## AIC: 157.03
```

- quasi-binomial: $E(Y_i) = n_i p_i, \quad \text{Var}(Y_i) = \phi n_i p_i (1 - p_i)$
- estimate ϕ ?
- usually use $X^2/(n - p)$, where

over-dispersion parameter

$$\chi^2 = \sum \frac{(y_i - n_i \hat{p}_i)^2}{n \hat{p}_i (1 - \hat{p}_i)}$$

`overdisp.Rmd`; `overdisp.html`

```
> step(EX1018binom.glm)
```

Coefficients:

(Intercept)	stage1	xray1	acid1
-3.05	1.65	1.91	1.64

Degrees of Freedom: 22 Total (i.e. Null); 19 Residual

Null Deviance: $\hat{\chi}^2$ 40.7

Residual Deviance: 19.6 $\hat{\chi}^2$ AIC: 39.3

- we can drop age and grade without affecting quality of the fit
- in other words the model can be simplified by setting two regression coefficients to zero
- **several mistakes** in text on pp. 491,2;
- deviances in Table 10.9 are incorrect as well <http://statwww.epfl.ch/davison/SM/> has corrected version

- step implements stepwise regression
- evaluates each fit using $\text{AIC} = -2\ell(\hat{\beta}; y) + 2p$
- penalizes models with larger number of parameters
- we can also compare fits by comparing deviances \longrightarrow [binaryELM2.html](#)

- step implements stepwise regression
- evaluates each fit using $AIC = -2\ell(\hat{\beta}; y) + 2p$
- penalizes models with larger number of parameters

- we can also compare fits by comparing deviances

→ [binaryELM2.html](#)

- ```
> update(Ex1018binom.glm, .~. - aged - grade)
```

```
Call: glm(formula = cbind(rtot, total - rtot) ~ stage + xray + acid,
 family = binomial, data = nodal2)
```

Coefficients:

|             |        |       |       |
|-------------|--------|-------|-------|
| (Intercept) | stage1 | xray1 | acid1 |
| -3.05       | 1.65   | 1.91  | 1.64  |

Degrees of Freedom: 22 Total (i.e. Null); 19 Residual

Null Deviance: 40.7

Residual Deviance: 19.6 AIC: 39.3

```
> deviance(ex1018binom)
```

```
[1] 18.06869
```

```
> pchisq(19.6-18.07, df = 2, lower = F)
```

```
[1] 0.4653
```

- as terms are added to the model, deviance always decreases
- because log-likelihood function always increases
- similar to residual sum of squares

- as terms are added to the model, deviance always decreases
- because log-likelihood function always increases
- similar to residual sum of squares
- Akaike Information Criterion penalizes models with more parameters
- 

$$AIC = 2\{-\ell(\hat{\beta}; y) + p\}$$

SM (4.57)

- comparison of two model fits by difference in *AIC*

- see **posted handout** on case-control studies
- consider for simplicity binomial responses with a single binary covariate:

$$\text{logit}(p_i) \sim \beta_0 + \beta_1 z_i, \quad i = 1, \dots, n$$

- see **posted handout** on case-control studies
- consider for simplicity binomial responses with a single binary covariate:

$$\text{logit}(p_i) \sim \beta_0 + \beta_1 z_i, \quad i = 1, \dots, n$$

- no difference between groups  $\iff$  odds-ratio  $\equiv 1$

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1 - p_0)}{p_0(1 - p_1)}$
- also called **relative risk**

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1-p_0)}{p_0(1-p_1)}$
- also called **relative risk**
- if  $p_1$  and  $p_0$  are both small, ( $y = 1$  is rare), then

$$\frac{p_1}{p_0} \approx \frac{p_1(1-p_0)}{p_0(1-p_1)}$$

- sometimes  $p_1/p_0$  can be large but if  $p_1$  and  $p_0$  are both small the difference  $p_1 - p_0$  might also be very small

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1 - p_0)}{p_0(1 - p_1)}$
- also called **relative risk**
- if  $p_1$  and  $p_0$  are both small, ( $y = 1$  is rare), then

$$\frac{p_1}{p_0} \approx \frac{p_1(1 - p_0)}{p_0(1 - p_1)}$$

- sometimes  $p_1/p_0$  can be large but if  $p_1$  and  $p_0$  are both small the difference  $p_1 - p_0$  might also be very small
- in order to estimate the **risk difference** we need to know the baseline risk  $p_0$

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1-p_0)}{p_0(1-p_1)}$
- also called **relative risk**
- if  $p_1$  and  $p_0$  are both small, ( $y = 1$  is rare), then

$$\frac{p_1}{p_0} \approx \frac{p_1(1-p_0)}{p_0(1-p_1)}$$

- sometimes  $p_1/p_0$  can be large but if  $p_1$  and  $p_0$  are both small the difference  $p_1 - p_0$  might also be very small
- in order to estimate the **risk difference** we need to know the baseline risk  $p_0$
- bacon sandwiches [www.youtube.com/watch?v=4szyEbU94ig](https://www.youtube.com/watch?v=4szyEbU94ig)
- risk calculator <https://realrisk.wintoncentre.uk/p1>

## Results



### Risk for Usual care

Out of 100 UK patients receiving mechanical ventilation for COVID-19, we would expect around 41 to die after 28 days

Edit Text



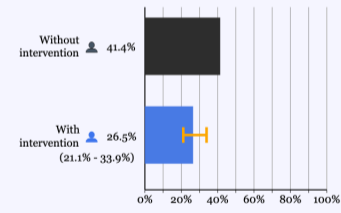
### Risk for Usual care plus dexamethasone

Out of 100 UK patients receiving mechanical ventilation for COVID-19, we would expect around 26 to die after 28 days

Edit Text

Barchart

Icon Array



<< Reset

< Back

FAQs

Download

Share

## Results summary

### PAPER TITLE

[Dexamethasone and 28 day mortality for COVID-19 patients on ventilation](#)

### DOI

<https://www.nejm.org/doi/10.1056/NEJMoa2021436>

### STUDY GROUP

UK patients receiving mechanical ventilation for COVID-19

### STUDY TYPE

experimental

### RISK FACTOR

taking dexamethasone

### OUTCOME

die after 28 days

### MEASURE OF CHANGE

Relative risk 0.64 (0.51 - 0.82)

### BASELINE CONDITION

Usual care

### EXPERIMENTAL CONDITION

Usual care plus dexamethasone

### BASELINE RISK

41.4%

Odds ratio 0.64; baseline risk 41.4%

## Results



### Risk for Usual care

Out of 100 UK patients receiving mechanical ventilation for COVID-19, we would expect around 41 to die after 28 days

Edit Text



### Risk for Usual care plus dexamethasone

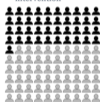
Out of 100 UK patients receiving mechanical ventilation for COVID-19, we would expect around 26 to die after 28 days

Edit Text

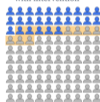
Barchart

Icon Array

41 out of 100 without  
intervention



26 (22 - 33) out of 100  
with intervention



<< Reset

< Back

FAQs

Download

Share

## Results summary

### PAPER TITLE

[Dexamethasone and 28 day mortality for COVID-19 patients on ventilation](#)

### DOI

<https://www.nejm.org/doi/10.1056/NEJMoa2021436>

### STUDY GROUP

UK patients receiving mechanical ventilation for COVID-19

### STUDY TYPE

experimental

### RISK FACTOR

taking dexamethasone

### OUTCOME

die after 28 days

### MEASURE OF CHANGE

Relative risk 0.64 (0.51 – 0.82)

### BASELINE CONDITION

Usual care

### EXPERIMENTAL CONDITION

Usual care plus dexamethasone

### BASELINE RISK

41.4%

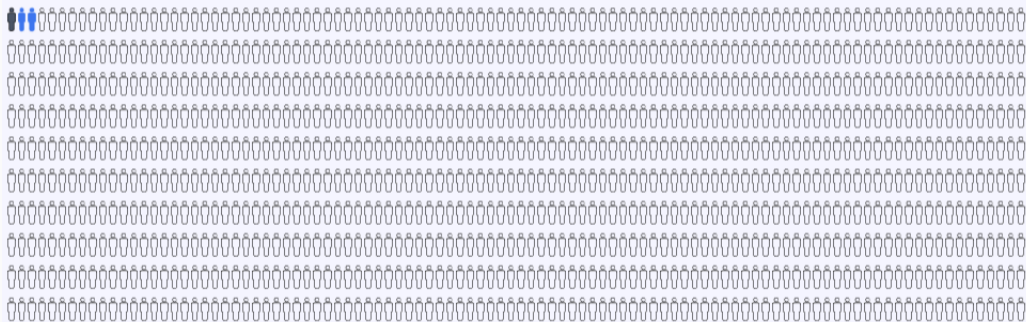
Odds ratio 0.64; baseline risk 41.4%



1 / 1000



3 / 1000 (2 extra cases)



Odds ratio 2.91; baseline risk 1/1000

Whether we sample **prospectively** or **retrospectively**, the odds ratio is the same

|                 | Lung cancer |          |
|-----------------|-------------|----------|
|                 | 1           | 0        |
|                 | cases       | controls |
| smoke = 1 (yes) | 688         | 650      |
| smoke = 0 (no)  | 21          | 59       |
|                 | 709         | 709      |

$$\text{retro: OR} = \frac{(688/709)/(21/709)}{(650/709)/(59/709)} = \frac{688 \times 59}{650 \times 21} = 2.97$$

$$\text{prosp: OR} = \frac{\{688/(688 + 650)\}/\{650/(688 + 650)\}}{21/(21 + 59)/\{59/(21 + 59)\}} = \frac{688 \times 59}{650 \times 21} = 2.97$$

# Types of observational studies

- secondary analysis of data collected for another purpose
- estimation of some feature of a defined population
  - could in principle be found exactly
- tracking across time of such features
- study of a relationship between features, where individuals may be examined
  - at a single time point
  - at several time points for different individuals
  - at different time points for the same individual
- census
- meta-analysis: statistical assessment of a collection of studies on the same topic



## E-commerce domain survival rates, by platform, 2019–2021

Percentage of domains that survive by number of days after sign-up

Shopify Wix Squarespace WooCommerce PrestaShop

