

PUBH 7405: Final Presentation

Annika Fredrickson, Mykhaylo M. Malakhov, and Daniel Whitford



UNIVERSITY OF MINNESOTA
*Driven to Discover*SM

Predicting Length of Stay From an Electronic Patient Record System: a Primary Total Knee Replacement Example

Evelene M Carter and Henry WW Potts

BMC Medical Informatics and Decision Making, Volume 14, Issue 26 (2014)



UNIVERSITY OF MINNESOTA
Driven to DiscoverSM

Background: Length of Stay

Hospital length of stay is an important metric for quality of care.

Previous studies:

- Gender, age, social deprivation, comorbidity, discharge destination, and consultant are known to impact length of stay

Use in the United Kingdom:

- Reported to Department of Health
- Published on NHS Choices website
- Efforts toward reducing length of stay

Background: Electronic Patient Record Systems

Technological advances allow for unprecedented access to healthcare data, yet these data are not used to their full potential.

Advantages of electronic systems:

- In England, most patient records are submitted to the national Hospital Episode Statistics database
- Standardized reporting of measures and outcomes
- Easy for researchers to access
- Ethics approval is not needed (in the UK)

No previous studies had modeled length of stay using all of the factors available in electronic patient records.

STUDY AIMS

Can factors be identified that significantly affect hospital length of stay from those available in an electronic patient record system?

Can a model be produced to predict the length of stay based on these factors?

Design and Data Collection

The study adheres to a retrospective design.

Electronic patient records obtained ($n = 2,130$) that satisfy:

- Data reported to NHS by Nuffield Orthopaedic Centre hospital
- Patients discharged between January 2007 and December 2011
- Primary total knee operations (OPCS codes W401, W411, and W421)
- Not cancer patients
- Only adults (age >15)

Factors Considered

- season of admission
- season of discharge
- day of the week admitted
- country of residence
- distance between residence and the hospital
- deprivation of postal code
- commissioner reimbursing the cost of treatment
- lead consultant
- discharge destination
- discharge method
- gender
- age
- ethnicity
- comorbidities:
 - diabetes
 - renal failure
 - heart failure
 - retention of urine
 - difficulty swallowing
 - pulmonary embolism
 - respiratory failure

Methods

- Descriptive analyses
- Non-parametric tests for effect on length of stay
 - Mann-Whitney test - two groups
 - Kruskal-Wallis test - more than two groups
 - Spearman's rank correlation coefficient - continuous
 - P-value < 0.05 was considered significant

Modelling Techniques

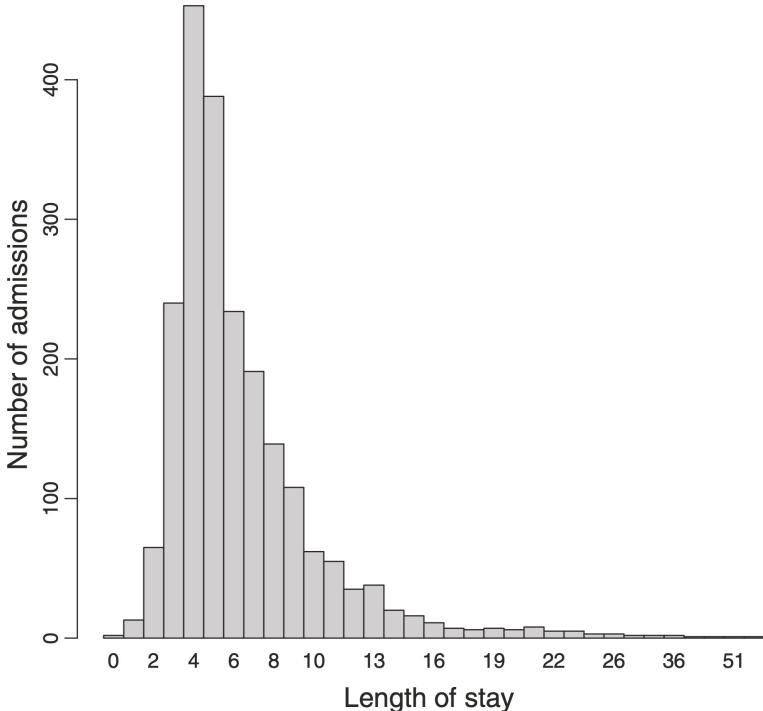


Figure 1 LoS distribution - primary total knees.

- Poisson regression (PRM)
 - Models count data
- Negative binomial regression (NBM)
 - Variance can be different from mean - better for over-dispersed data
- Authors found no other papers using these methods to analyse length of stay for primary total knee replacements (or any specialty for NBM)
- Forward selection procedure at 5% level of significance for PRM
- Evidence of over-dispersion for PRM
- Same variables included for NBM

Model Performance

- Data randomly split to avoid over-fitting
 - 90% of admissions for derivation
 - 10% of admissions for validation
 - Predicted total days compared to actual total days
- Residual plots to check models' fit
 - Heteroscedasticity
 - Data points with high leverage
- Akaike Information Criterion (AIC)
- Analysis carried out using R

Table 1 LoS analysis significant factor results

Group	N	Proportion	Median	Mean	Median average deviation	Proportion %	Median deviation range
<2011	1748	82.1%	5	6.6	1.48	82.1	3.5–6.5
2011	382	17.9%	4	5.7	1.48	17.9	2.5–5.5
Monday	449	21.1%	5	6.0	1.5	21.1	3.5–6.5
Tuesday	186	8.7%	6	7.4	3.0	8.7	3–9
Wednesday	356	16.7%	6	6.7	3.0	16.7	3–9
Thursday	394	18.5%	5.5	6.7	2.2	18.5	3.3–7.7
Friday	359	16.9%	5	6.2	1.5	16.9	3.5–6.5
Saturday	322	15.1%	4	5.2	1.5	15.1	2.5–5.5
Sunday	64	3.0%	9	10.7	4.4	3.0	4.6–13.4
Female	1274	59.8%	6	6.8	2.97	59.8	3–9
Male	856	40.2%	5	5.9	2.97	40.2	2 - 8
C58	274	12.9%	5	6.3	1.5	12.9	3.5–6.5
C38	256	12.0%	5	6.2	1.5	12.0	3.5–6.5
C49	246	11.5%	5	6.3	3.0	11.5	2–8
C46	211	9.9%	5	6.6	3.0	9.9	2–8
C42	200	9.4%	6	7.0	3.0	9.4	3–9
C59	179	8.4%	4	5.1	1.5	8.4	2.5–5.5
C26	171	8.0%	5	5.9	1.5	8.0	3.5–6.5
C62	156	7.3%	6	7.7	3.0	7.3	3–9
C64	148	6.9%	6	6.3	3.0	6.9	3–9
Other consultant	289	13.6%	5	6.9	2.97	13.6	2–8
NHS hospital provider	88	4.1%	8	9.8	7.4	4.1	0.6– 15.4
Other discharge destination	26	1.2%	10	11.7	5.2	1.2	4.8–15.2
Usual place of residence	2016	94.6%	5	6.2	1.5	94.6	3.5–6.5
White, declined and unknown	2074	97.6%	5	6.4	1.5	97.6	3.5–6.5
Other ethnicity	50	2.4%	6.5	7.6	3.7	2.4	2.8–10.2
Factor	Mean	First Quartile	Median	Third Quartile	Median LoS at 40	Median LoS at 60	Median LoS at 80
Age	70	64	71	77	7.5	5	6
Factor	Mean	First Quartile	Median	Third Quartile	IMD at 3 Days LoS	IMD at 5 Days LoS	IMD at 7 Days LoS
Indices of deprivation	14720	7526	15353	22400	24368	24189	24011

Age: Non-Linear Relationship with Length of Stay

Table 7 LoS for the average patient by age bands

Age	LoS for the average patient
20	8.4
30	6.5
40	5.5
50	5.1
60	5.1
70	5.5
80	6.6
90	8.5

Table 2 Summary of univariately significant independent variables on LoS

Variable	p-value	Test type	Test value	Order of modelling
Admission year	<0.0001	Mann-Whitney	U = 104842283	1
Age at admission	<0.0001	Spearman's	r = 0.26	2
Age ²	<0.0001	Spearman's	r = 0.26	3
Gender	<0.0001	Mann-Whitney	U = 652862	4
Consultant	<0.0001	Kruskal-Wallis	$\chi^2(9) = 75.76$	5
Admission day of week	<0.0001	Kruskal-Wallis	$\chi^2(6) = 146.15$	6
Discharge destination	<0.0001	Kruskal-Wallis	$\chi^2(2) = 35.37$	7
IMD	0.01	Spearman's	r = 0.06	8
Ethnicity	0.03	Mann-Whitney	U = 68095	9

Factors not included in modelling

- There was no seasonal effect on length of stay
- Commissioning area of patient's residence ($p=0.20$)
- Distance patient lived from hospital ($p=0.52$)
- Patients who died in hospital had much longer length of stay
 - whether or not a patient will die is not the focus of this paper's prediction
- Comorbidities
 - Retention of urine ($p<0.001$, 1.2% prevalence)
 - Difficulty swallowing ($p=0.081$, <0.1% prevalence)
 - Pulmonary embolism ($p=0.038$, 0.9% prevalence)
 - Criterion for reliable test was a minimum of 5% of the cohort per group
 - Other comorbidities were insignificant

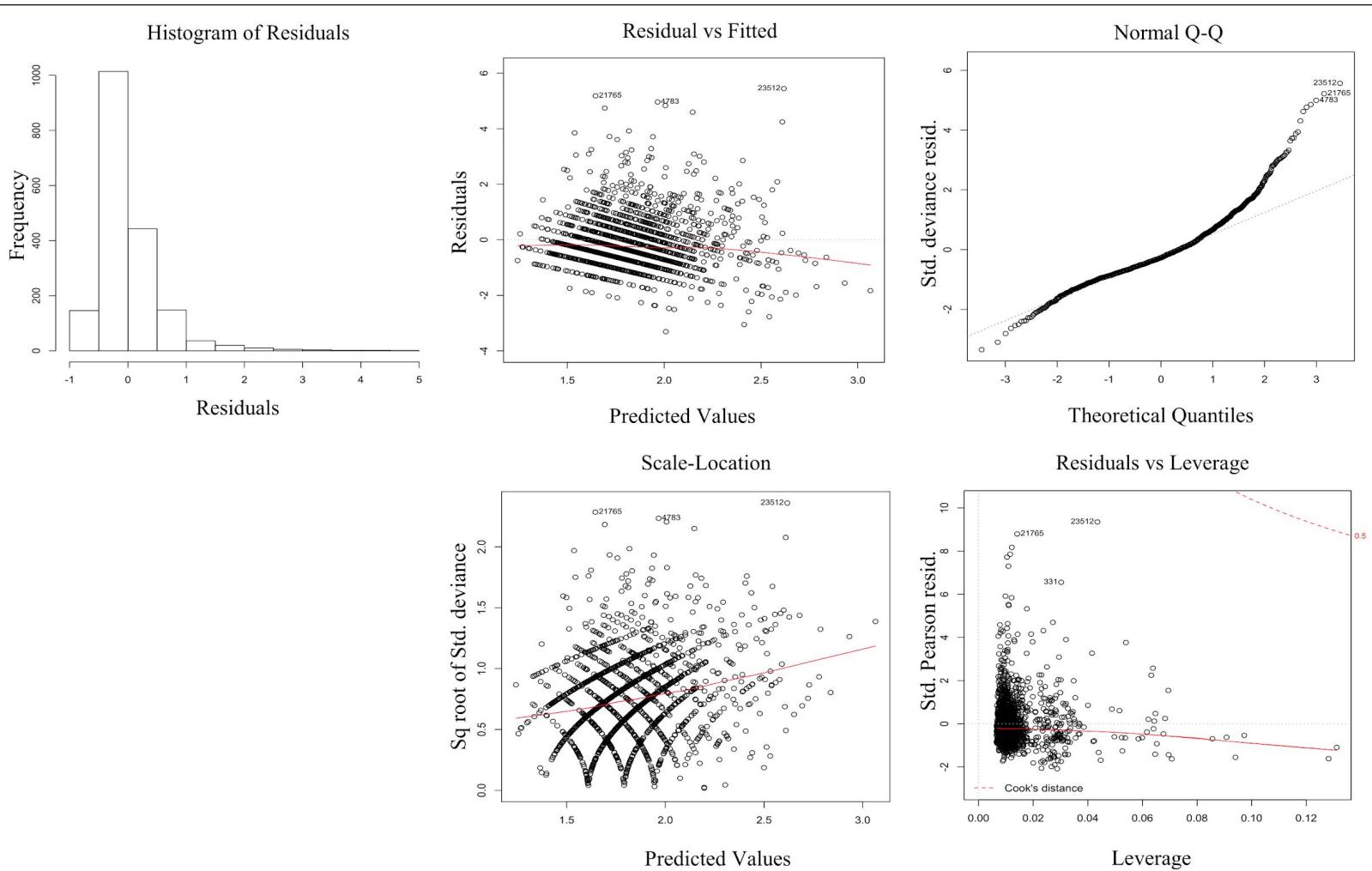


Figure 2 Negative binomial model residual plots.

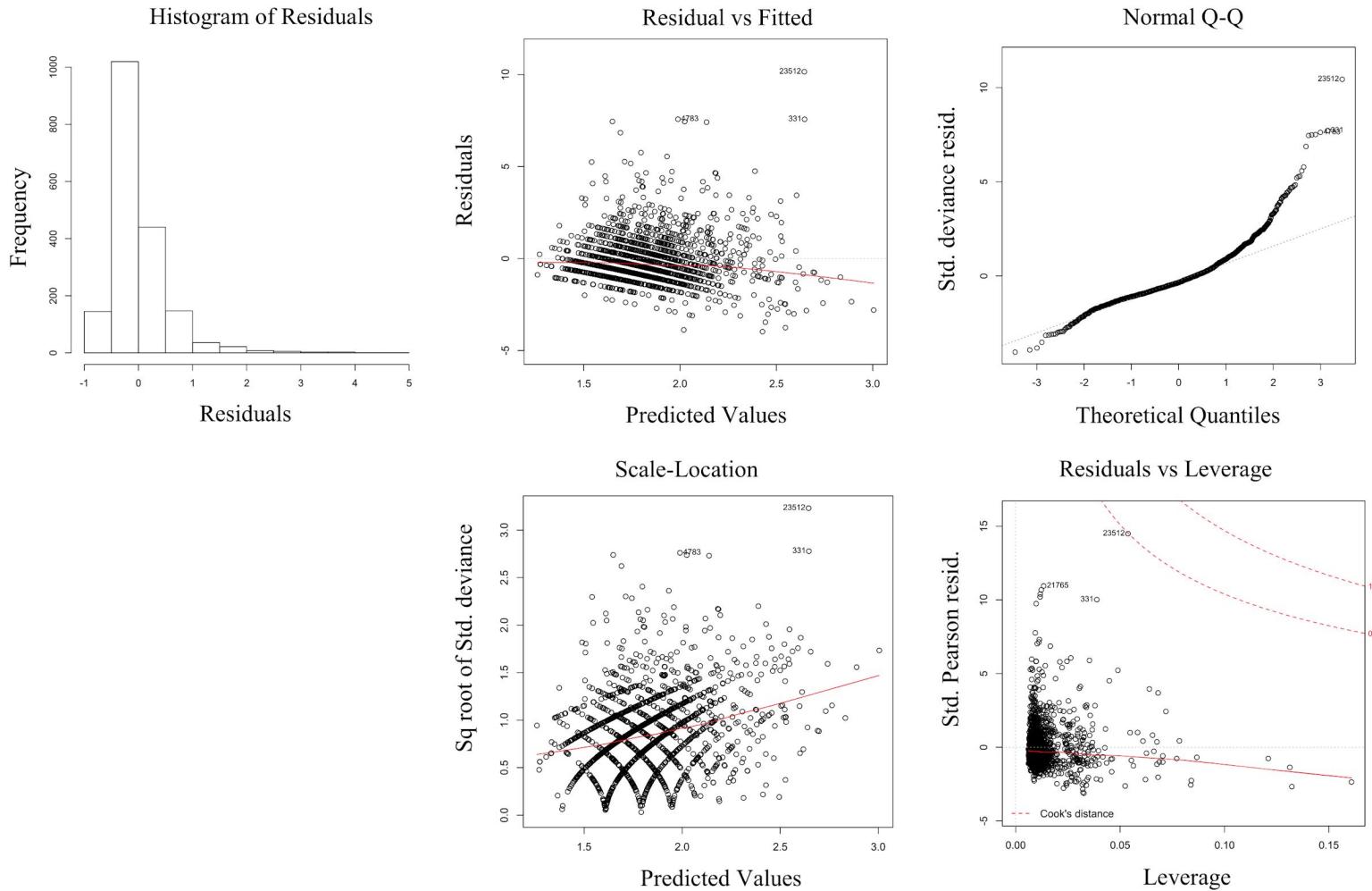


Figure 3 Poisson model residual plots.

Table 3 Akaike information criterion for the models

Model	AIC	Diff AIC	Model likelihood	AIC weight
Negative binomial	9170	0	1.0	1.0
Poisson	9668	498	0.0	0.0
Total model Likelihood			1.0	

Final Model - Negative Binomial Regression

Model Equation

$$LoS = \exp^{(\text{Intercept} + b(x_1) * \text{cf}(Year=2011) + \text{Age} * \text{cf}(Age) + \text{Age}^2 * \text{cf}(Age^2) + b(x_2) * \text{cf}(Gender=Male) + b(x_3) * \text{cf}(Cons=C38) + \text{etc.})}$$

$\text{cf}(variable)$: regression coefficient for that variable

$b(x_n)$: binomial variable with value of 1 or 0

Incidence Rate Ratios (IRR) for interpretation of variables

Table 4 Summary of coefficients and incident rate ratios (IRR)

	Coefficient	Std. error	z value	Pr (> z)	Significance	IRR
(Intercept)	3.51500	0.32630	10.772	<2e-16	***	33.62
Admission.Year.Group2011	-0.0999	0.03409	-2.930	0.00339	**	0.90
Age.at.Admission	-0.04659	0.00933	-4.996	0.00000	***	0.95
Age.Squared	0.00043	0.00007	6.160	0.00000	***	1.00
GenderMale	-0.13760	0.02525	-5.450	0.00000	***	0.87
Consultant.Pseudo.Code.PK.GroupC38	0.07489	0.05752	1.302	0.19292		1.08
Consultant.Pseudo.Code.PK.GroupC42	0.11630	0.05977	1.946	0.05164	.	1.12
Consultant.Pseudo.Code.PK.GroupC46	0.13800	0.06180	2.234	0.02551	*	1.15
Consultant.Pseudo.Code.PK.GroupC49	0.01046	0.05734	0.182	0.85523		1.01
Consultant.Pseudo.Code.PK.GroupC58	0.06061	0.05859	1.034	0.30093		1.06
Consultant.Pseudo.Code.PK.GroupC59	-0.07585	0.06385	-1.188	0.23481		0.93
Consultant.Pseudo.Code.PK.GroupC62	0.17290	0.06781	2.550	0.01077	*	1.19
Consultant.Pseudo.Code.PK.GroupC64	0.03859	0.06718	0.574	0.56574		1.04
Consultant.Pseudo.Code.PK.GroupOther consultant	0.12170	0.05685	2.142	0.03223	*	1.13
Admission.DayMonday	-0.01045	0.04263	-0.245	0.80631		0.99
Admission.DaySaturday	-0.16110	0.04504	-3.578	0.00035	***	0.85
Admission.DaySunday	0.45220	0.07280	6.211	0.00000	***	1.57
Admission.DayThursday	0.00404	0.04522	0.089	0.92877		1.00
Admission.DayTuesday	0.10370	0.05123	2.024	0.04300	*	1.11
Admission.DayWednesday	0.01796	0.04357	0.412	0.68018		1.02
Discharge.Destination.PK.GroupOther discharge dest	0.10560	0.11270	0.937	0.34855		1.11
Discharge.Destination.PK.GroupUsual place of resid	-0.32660	0.05522	-5.915	0.00000	***	0.72
Rank.of.IMD.Score	-0.000005	0.00000	-2.901	0.00372	**	1.00
Ethnicity.Common.GroupWhite, declined and unknown	-0.13000	0.07777	-1.671	0.09467	.	0.88

0 = **** : 0.001 = *** : 0.01 = ** 0.05 : ' = 0.1 : ' = 1.

Component + Residual Plots

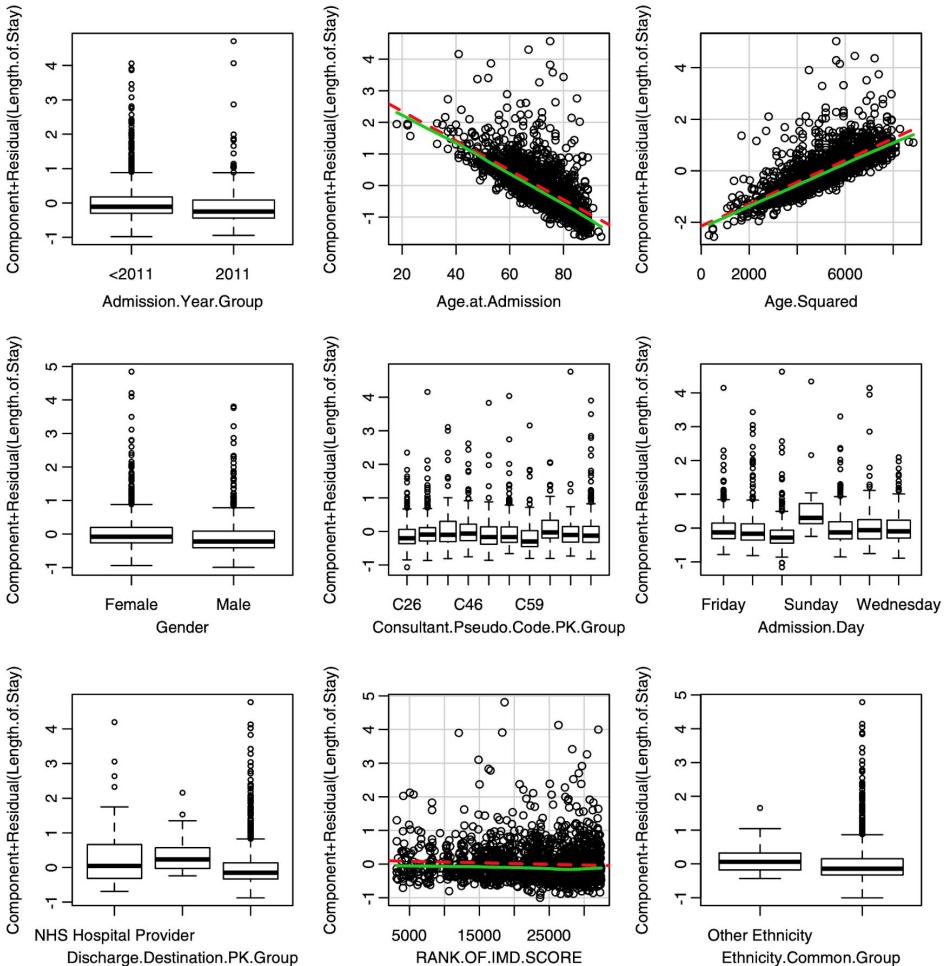


Figure 4 Negative binomial model residual plots by variable.

Table 5 Model results – model data

LoS grouping	Proportion of admissions	Total number of actual days stayed	Total number of days difference (Model vs Actual)	% difference in total number of days	Predicting within 1 days accuracy	Predicting within 2 days accuracy	Predicting within 3 days accuracy
4 to 6 days	50.8%	4469	1193	26.7%	41.1%	74.7%	91.4%
Shorter LoS	14.8%	729	780	107.0%	2.2%	30.4%	61.1%
Longer LoS	34.4%	6607	-1968	-29.8%	24.3%	42.1%	59.2%
Total	100.0%	11805	5	0.0%	29.6%	56.9%	75.8%

Table 6 Model results – test data

LoS grouping	Proportion of admissions	Total number of actual days stayed	Total number of days difference (Model vs Actual)	% difference in total number of days	Predicting within 1 days accuracy	Predicting within 2 days accuracy	Predicting within 3 days accuracy
4 to 6 days	50.9%	500	158	31.6%	32.4%	67.6%	88.6%
Shorter LoS	17.3%	95	110	115.6%	0.0%	30.6%	58.3%
Longer LoS	31.8%	667	-180	-27.0%	25.4%	52.2%	68.7%
Total	100.%	1262	88	6.9%	24.5%	56.3%	76.9%

Example - Mrs. Everett

- 70 years old
- female
- admitted under consultant C58
- admitted on a Monday
- discharged to her normal place of residence
- IMD Rank of 24871.5
- white ethnic origin

$$\text{LoS} = \exp(3.515 + (-0.0999) + (-0.04659 * 70) + (0.0004 * 70^2) + (0) + (0.06061) + (-0.01045) + (-0.3266) + (2.48715 * -0.000005) + (-0.13))$$

$$= 5.5 \text{ days}$$

Example - Mrs. Everett Continued

Showing application of incident rate ratio:

Coefficient for gender male = -0.13760

IRR for gender male = $\exp(-0.13760):1 = 0.87:1$

A man will only stay 87% of the time of a woman. If a man was admitted with the same admission and discharge attributes of Mrs Everett his LoS reduces by 0.7 days to 4.8 days stay.

Predictive Tool

Estimate Your Patients Length of Stay for a Total Primary Knee Admission

Please pick from the drop down lists for all entries except age.

Age (Years):

Gender: F for Female, M for Male

Consultant Admitted To: (This is Pseudo Consultant Code for this paper. It could easily be consultant name or code)

Admission Day of Week: If no date date for admission yet please enter "No date"

Discharge Destination:

Postcode: To look up the IMD Rank
(This postcode is for illustration only and pulls back the average IMD rank)

Ethnicity:

Estimated Length of Stay:
(If LoS is N/A please ensure all values have been entered above)

NB1 Patients from Malta tend to stay longer as they are more complex cases due to our contract with them.

NB2 Patients admitted over the Christmas and New Year period are likely to have a shorter stay due to attempts to get patients home for Christmas and New Year.

Strengths

- Large dataset, $N > 2,000$
- Very thorough/robust
 - Univariate testing using non-parametric tests
 - Poisson and Negative Binomial Regression to account for outcome distribution
 - Diagnostic plots and AIC for model fit
 - Model performance on training and test data
- Extensive explanation of interpretation of model results

Limitations

- Only one medical center - cannot extrapolate results.
- Variable selection methods for multiple regression analysis used stepwise selection, but only for variables significant in univariate analysis.
- Did not check for multicollinearity, assumed all variables were independent.
- Residuals of age and age squared follow a pattern, but nothing done to remedy this.

Limitations continued

- Interpretation of final model coefficients and incident rate ratios is not easily accessible. - Needed to add long example to make this concept clear.
- Results presented in a way that make the model look like it performs well, but it is not clear if it would be clinically useful.
- Not clear in research focus/purpose - waffled between inference and prediction.

Association of Short-term Exposure to Air Pollution With Mortality in Older Adults

Qian Di, Lingzhen Dai, Yun Wang, Antonella Zanobetti, Christine Choirat, Joel D. Schwartz, and
Francesca Dominici
JAMA. Volume 318, Issue 24 (2017)



UNIVERSITY OF MINNESOTA
Driven to DiscoverSM

Research Question

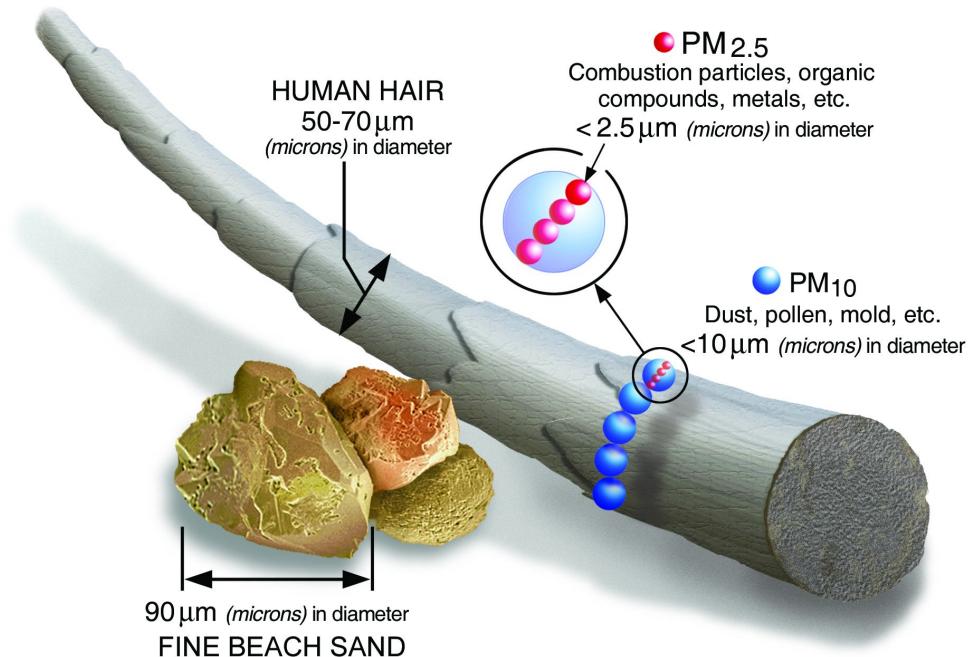
What is the association between short-term exposure to air pollution below current air quality standards and all-cause mortality?

Background - Air Pollution

Ground-Level Ozone and PM_{2.5}

Largely from:

- Power plants
- Factories
- Vehicles
- Burning fuels



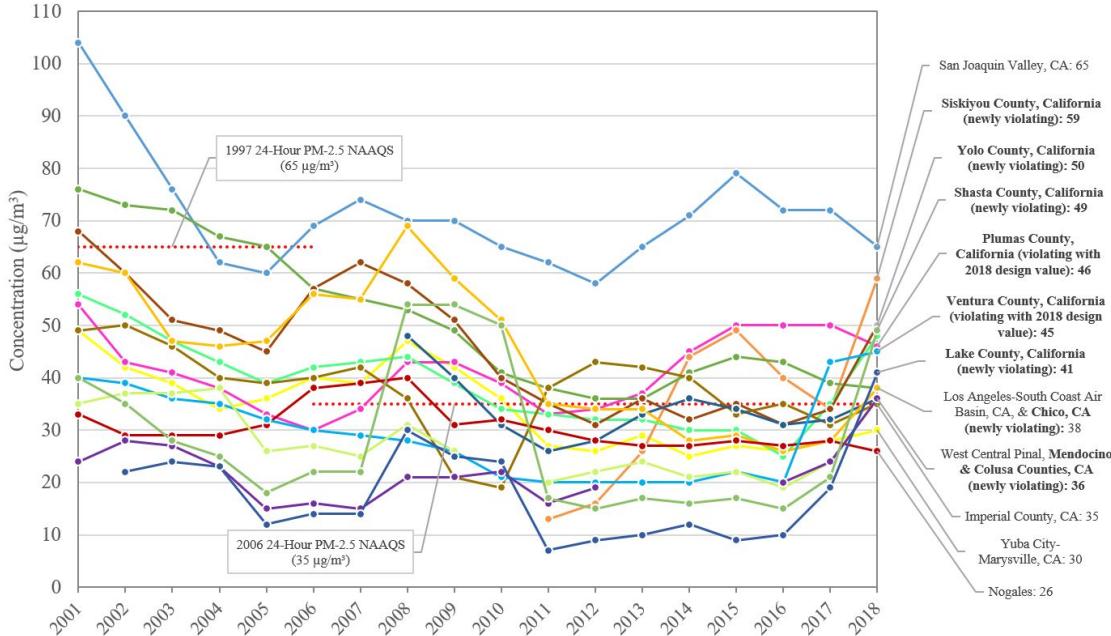
PM_{2.5} Levels

The plot shows an example of PM_{2.5} trends in several locations.

Max: San Jose, CA in 2011 with ~105 µg/m³

Min: Lake County, CA in 2011 with ~5 µg/m³

US EPA REGION 9 AIR QUALITY TRENDS, 2001-2018
24-HOUR FINE PARTICULATE MATTER (PM_{2.5}) DESIGN VALUES
IN DESIGNATED AND OTHER VIOLATING AREAS



Source: US EPA's Air Quality Systems (AQS) database (July 18, 2019).

The 2006 national ambient air quality standard (NAAQS) for 24-hour fine particulate matter (24-hour PM_{2.5}) is 35 micrograms per cubic meter (µg/m³). The design value for 24-hour PM_{2.5} is the three-year average of third-highest daily values. X-axis values represent the last year of a monitoring site's three year time period. All incomplete, invalid data and exceptional event data (e.g., high winds and wildfires) that EPA has concurred on have been excluded from design value calculations.

AIR19100 - 2018 annual air quality update - PM25.xlsx (September 26, 2019)

National Ambient Air Quality Standards

- Clean Air Act required NAAQS be set by the US EPA
- Reviewed every 5 years
- State, local, and tribal agencies use these standards for designing emission reduction plans

Standards for 2012:

PM_{2.5}: Annual - 12 µg/m³; 24-hour - 35 µg/m³

Ozone: 8-hour - 70 parts per billion (ppb).



[Environmental Topics](#)

[Laws & Regulations](#)

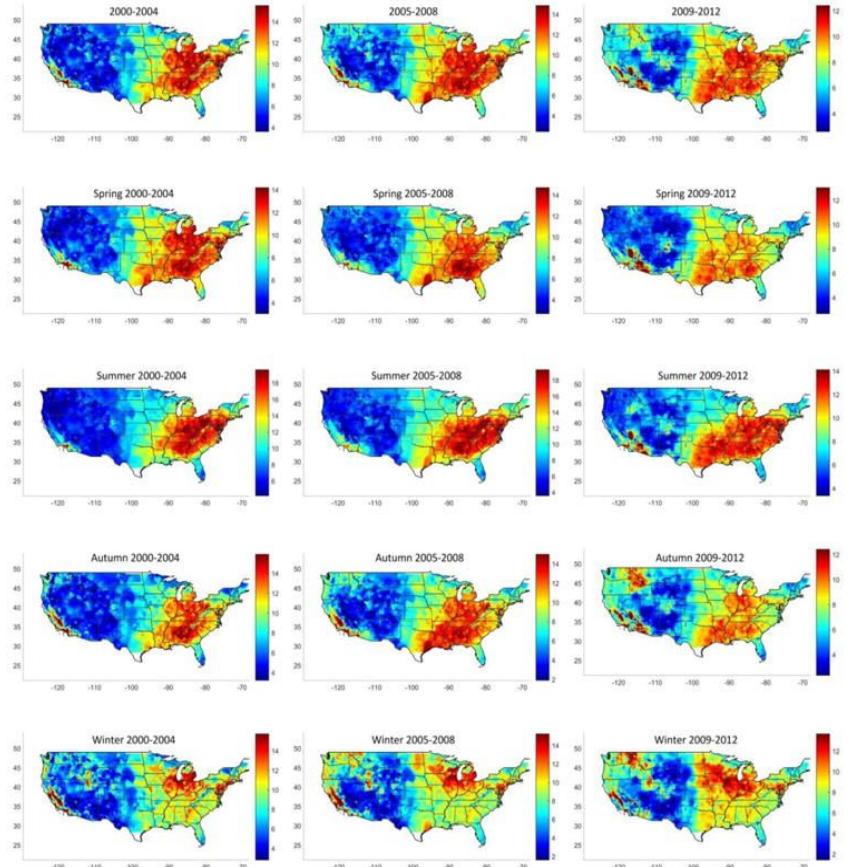
[About EPA](#)

Related Topics: [Criteria Air Pollutants](#)

NAAQS Table

Environmental Data

- Used published and validated air pollution neural network prediction models
- Air pollution was estimated for each 1 km square nationwide, then assigned to ZIP codes
 - Meteorological satellite data
 - EPA monitoring data
- $R^2 = 0.84$ for $PM_{2.5}$; $R^2 = 0.76$ for ozone



Cohort and Variables

Medicare members with a date of death between Jan. 1, 2000 and Dec. 31, 2012

N = >22 million members

Predictors matched by zip code for day of death and control days:

- Ozone and PM_{2.5} exposure levels
- Air and dew point temperatures

Subgroup analysis:

- age, sex, race, ethnicity, and eligibility for Medicaid

Outcome:

- All-cause mortality

Review: Study Designs

Case-control design:

- Observational study
- Case group containing subjects with specified outcome (e.g. disease)
- Control group containing subjects without specified outcome

Crossover design:

- Controlled experiment
- Each subject serves as his/her own control “group”
- Used in clinical trials studying treatments with temporary effects

Case-Crossover Design

“A hybrid between a matched case-control design and a traditional crossover design”

- Each subject serves as his/her own control
- The control data is from a different time period when the event that defines case status was not experienced

In this study:

- Case day defined as day of death
- Each case matched to 3-4 control days

Table 1. Baseline Characteristics of Study Population (2000-2012)

Baseline Characteristic	Value
Case days, No.	22 433 862
Control days, No.	76 143 209
Among All Cases (n = 22 433 862), %	
Age at death, y	
≤69	10.38
70-74	13.37
75-84	38.48
≥85	37.78
Sex	
Male	44.73
Female	55.27
Race/ethnicity	
White	87.34
Black	8.87
Asian	1.03
Hispanic	1.51
Native American	0.31
Medicaid Eligibility (n = 22 433 862), %	
Ineligible	77.36
Eligible	22.64

Conditional Logistic Regression

An extension of logistic regression that takes matching into account.

Here, each subject has 1 case day and 3-4 control days, denoted t_1, t_2, \dots, t_M .
Then the probability that subject i dies at time t_k is

$$P_{ik} = \frac{\exp(\boldsymbol{\beta}^T \mathbf{X}_{ik})}{\sum_{j=1}^M \exp(\boldsymbol{\beta}^T \mathbf{X}_{ij})}$$

where predictors are PM_{2.5}, ozone, and natural splines of air temperature and dew point temperature with 3 d.f..

Estimating Relative Risk

The relative risk increase for all-cause mortality associated with a $10\text{-}\mu\text{g}/\text{m}^3$ daily increase in $\text{PM}_{2.5}$ was calculated as $RR_{PM2.5} = \exp(10 * \beta_{PM2.5})$

Similarly, the relative risk increase for all-cause mortality associated with a 10-ppb daily increase in ozone was calculated as $RR_{ozone} = \exp(10 * \beta_{ozone})$

Subgroup analyses were conducted by sex, race/ethnicity, age group, eligibility for Medicaid, and population density quartile. For example, to test the null hypothesis $H_0: RR_{male} = RR_{female}$ the z statistic is calculated:

$$Z = \frac{RR_{male} - RR_{female}}{\sqrt{se(RR_{male})^2 + se(RR_{female})^2}}$$

1.05%

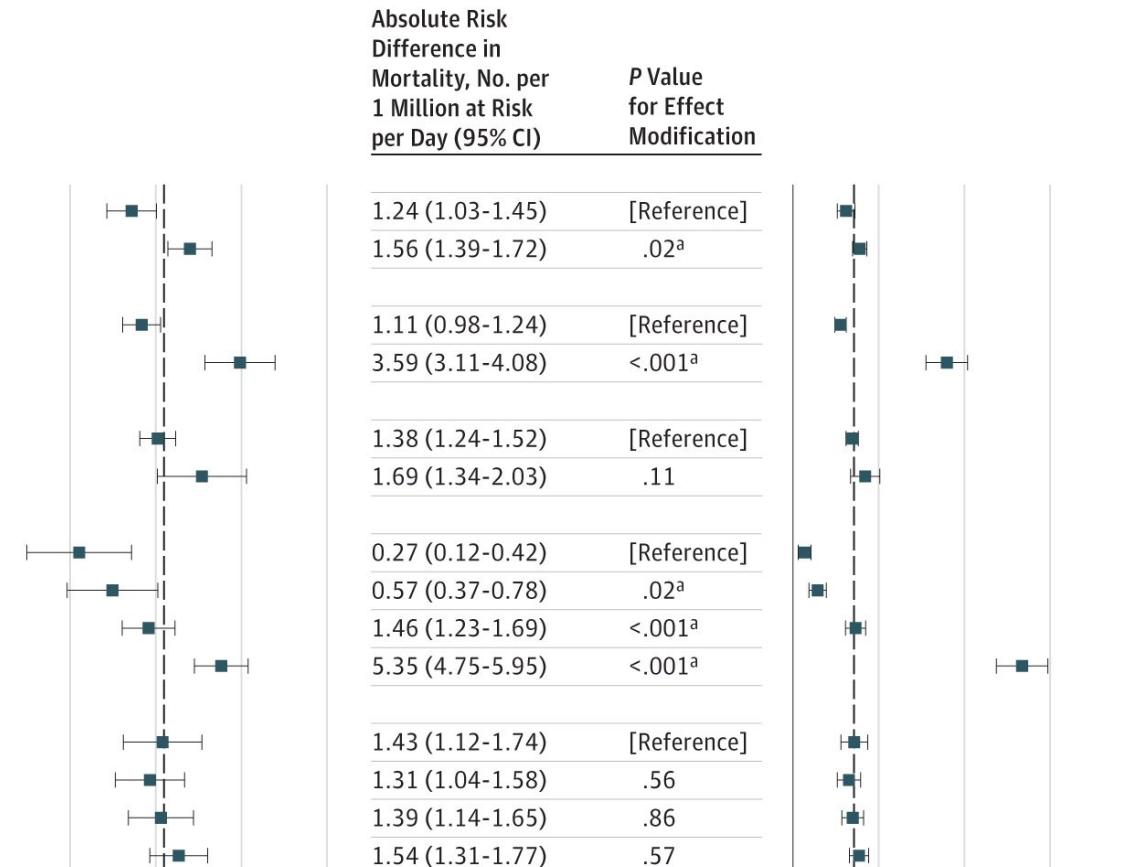
Relative risk increase
associated with a $10\text{-}\mu\text{g}/\text{m}^3$
daily increase in $\text{PM}_{2.5}$

0.51%

Relative risk increase
associated with a 10-ppb daily
increase in ozone

Results of Subgroup Analyses

Model	Relative Risk Increase in Mortality per 10- $\mu\text{g}/\text{m}^3$ Increase in PM _{2.5} , % (95% CI)	P Value for Effect Modification
Sex		
Male	0.86 (0.72-1.00)	[Reference]
Female	1.20 (1.07-1.33)	<.001 ^a
Medicaid eligibility		
Noneligible	0.92 (0.81-1.03)	[Reference]
Eligible	1.49 (1.29-1.70)	<.001 ^a
Race/ethnicity		
White	1.01 (0.91-1.12)	[Reference]
Nonwhite	1.27 (1.01-1.53)	.07
Age, y		
≤69	0.55 (0.25-0.86)	[Reference]
70-74	0.75 (0.48-1.01)	.35
75-84	0.96 (0.80-1.11)	.02 ^a
≥85	1.38 (1.23-1.54)	<.001 ^a
Population density		
Low	1.04 (0.81-1.27)	[Reference]
Medium low	0.97 (0.76-1.17)	.64
Medium high	1.03 (0.84-1.22)	.95
High	1.13 (0.97-1.30)	.52



Results of Subgroup Analyses (Continued)

Whites

Sex

Male	0.83 (0.67-0.99)	[Reference]
Female	1.16 (1.02-1.30)	.002 ^a

Medicaid eligibility

Noneligible	0.88 (0.77-1.00)	[Reference]
Eligible	1.58 (1.34-1.83)	<.001 ^a

Nonwhites

Sex

Male	1.03 (0.65-1.42)	[Reference]
Female	1.47 (1.12-1.82)	.01

Medicaid eligibility

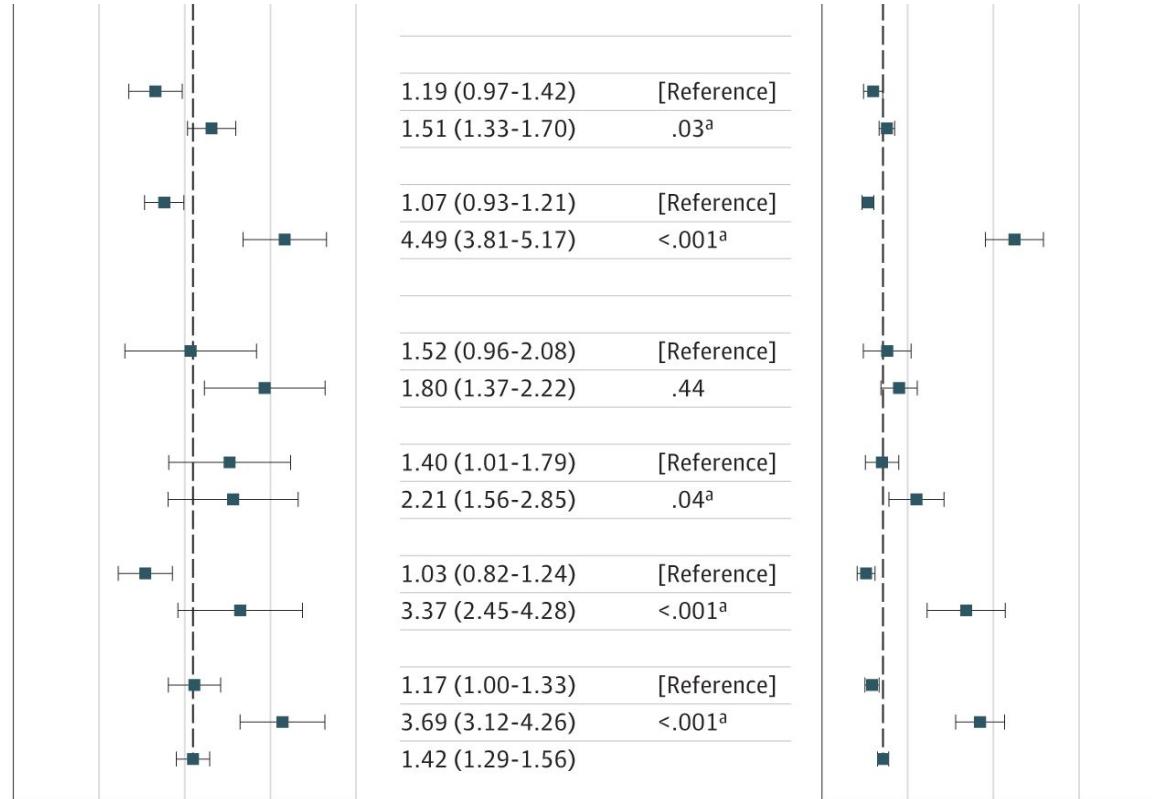
Noneligible	1.26 (0.91-1.62)	[Reference]
Eligible	1.28 (0.90-1.66)	.94

Medicaid eligibility, males

Noneligible	0.77 (0.61-0.93)	[Reference]
Eligible	1.32 (0.96-1.69)	.006

Medicaid eligibility, females

Noneligible	1.06 (0.90-1.21)	[Reference]
Eligible	1.57 (1.32-1.82)	<.001 ^a
Overall	1.05 (0.95-1.15)	

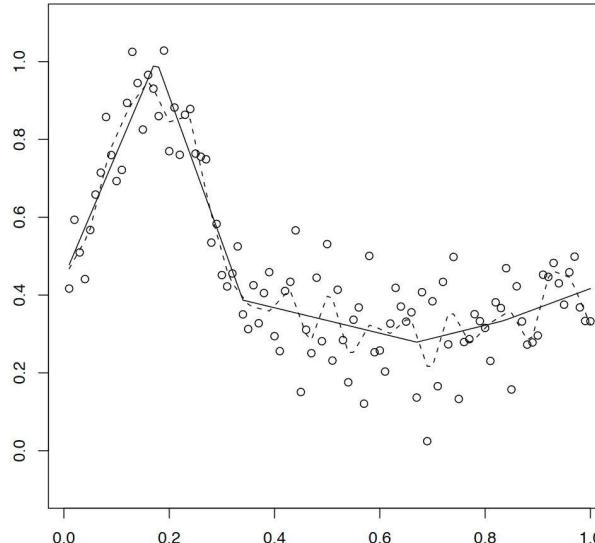


Meta-analysis: Inverse Variance Weighting

To estimate exposure-response curves, linear terms for the two pollutants were replaced with penalized splines in the regression model.

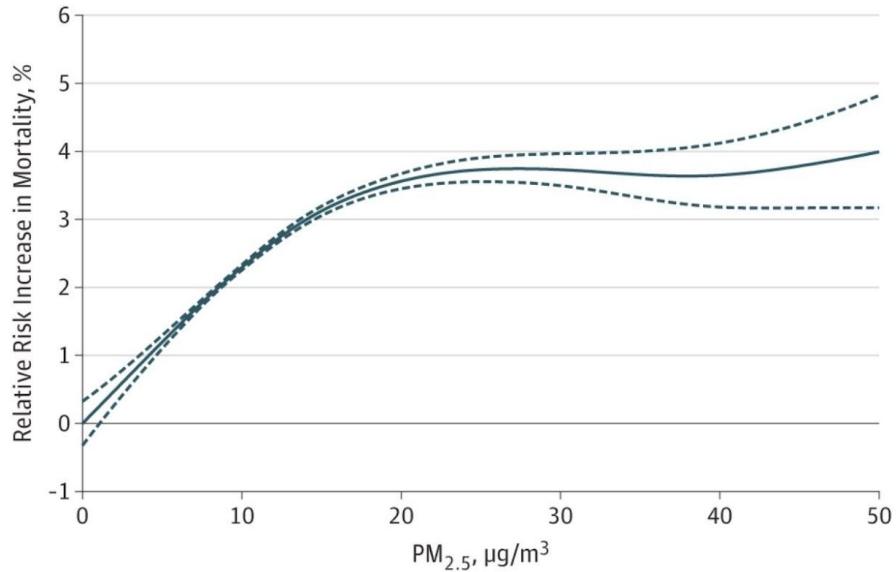
Computational limitations render full analysis infeasible.

1. Data was split randomly into 50 groups
2. Meta-analysis used to combine group-level effect estimates
3. Regressing them against indicator variables for each exposure level with inverse variance weights

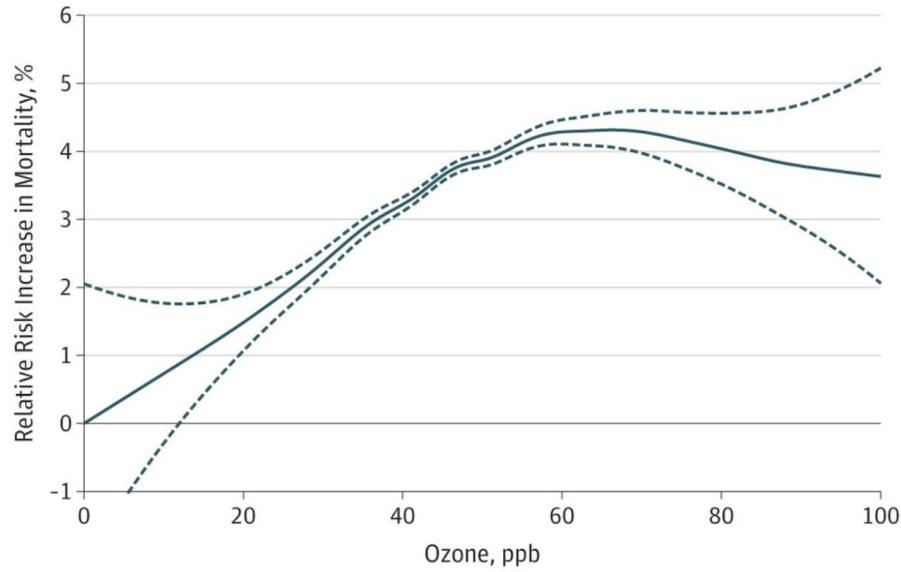


Exposure-Response Curves

A Exposure-response curve for PM_{2.5}



B Exposure-response curve for ozone



Findings

- Both PM_{2.5} and ozone levels are associated with an increased risk of mortality
 - 10- $\mu\text{g}/\text{m}^3$ daily increase in PM_{2.5} associated with an increase of 1.42 deaths per 1 million per day
 - 10-ppb daily increase in warm-season ozone exposures associated with an increase of 0.66 deaths per 1 million per day
 - Remains significant at pollution levels much lower than current daily NAAQS
 - Effect even more pronounced for subpopulations such as medicaid-eligible individuals, black individuals, females, and the elderly
 - No apparent threshold for safe PM_{2.5} and ozone levels

Strengths

- Very large study population - over 22 million deaths
- Air pollution prediction models allowed for accurate estimates across the United States, including unmonitored areas
- Differences in mortality rates compared across several vulnerable populations
- Case-crossover design allowed for matching of potential confounders

Weaknesses

- The variables for PM_{2.5} and ozone levels are estimates from another model, rather than directly measured data
- Case-crossover design does not allow estimation of the effect of long-term exposure to pollution on mortality (this was done in another study)
- Limited generalizability to younger populations
- Cause of mortality was not included
- More detailed and reliable subpopulation data needed (i.e. Medicaid eligibility varies by state)

Any Questions?