# **scientific** reports

OPEN

# Improving the design stage of air pollution studies based on wind patterns

Léo Zabrocki[1]✉, Anna Alari[2] & Tarik Benmarhnia[3]

A growing literature in economics and epidemiology has exploited changes in wind patterns as a source of exogenous variation to better measure the acute health effects of air pollution. Since the distribution of wind components is not randomly distributed over time and related to other weather parameters, multivariate regression models are used to adjust for these confounding factors. However, this type of analysis relies on its ability to correctly adjust for all confounding factors and extrapolate to units without empirical counterfactuals. As an alternative to current practices and to gauge the extent of these issues, we propose to implement a causal inference pipeline to embed this type of observational study within an hypothetical randomized experiment. We illustrate this approach using daily data from Paris, France, over the 2008–2018 period. Using the Neyman–Rubin potential outcomes framework, we first define the treatment of interest as the effect of North-East winds on particulate matter concentrations compared to the effects of other wind directions. We then implement a matching algorithm to approximate a pairwise randomized experiment. It adjusts nonparametrically for observed confounders while avoiding model extrapolation by discarding treated days without similar control days. We find that the effective sample size for which treated and control units are comparable is surprisingly small. It is however reassuring that results on the matched sample are consistent with a standard regression analysis of the initial data. We finally carry out a quantitative bias analysis to check whether our results could be altered by an unmeasured confounder: estimated effects seem robust to a relatively large hidden bias. Our causal inference pipeline is a principled approach to improve the design of air pollution studies based on wind patterns.

A growing literature in economics and epidemiology has recently re-examined the short-term effects of air pollution on mortality and emergency admissions using causal inference methods. Among these techniques, instrumental variable strategies have been very popular since they can overcome the biases caused by unmeasured confounders and measurement errors in air pollution exposure[1–6]. Daily changes in wind directions are such instrumental variables since they arguably meet two of the three main requirements for the method to be valid: they can strongly affect air pollutant concentrations while having no direct effects on health outcomes[7–9]. This strategy however rests on the remaining assumption that changes in wind directions occur randomly, which is often not credible without further statistical adjustments. One could unfortunately fear that the resulting analysis would depend on the quality of the model[10,11]. Does the model take into account all relevant confounding factors, and if so, are they adjusted for with the correct functional forms? Is the model also able to extrapolate when there is little overlap in covariate distributions?

To illustrate these issues, imagine that we are interested in estimating the influence of particulate matters on daily mortality in Paris, France, over the 2008–2018 period. Research in atmospheric science has shown that winds blowing from the North-East could transport particulate matters due to wood burning in the region but also from other sources located in North-Eastern Europe[12–14]. We could therefore use the comparison of winds blowing from the North-East to those from other directions as an instrumental variable for particulate matters.

In Panel A of Fig. 1, we display polar plots of air pollutant concentrations that were predicted using a Generalized Additive Model (GAM) and wind components as inputs[15]. We clearly see that winds blowing from the North-East are associated with higher $PM_{10}$ and $PM_{2.5}$ concentrations. These patterns could however be confounded by other variables such as the weather parameters or a shared seasonality in air pollution and

[1]Paris School of Economics and École des Hautes Etudes en Sciences Sociales, 48 Boulevard Jourdan, 75014 Paris, France. [2]Barcelona Institute for Global Health (ISGlobal), Carrer del Rosselló, 132, 08036 Barcelona, Spain. [3]Department of Family Medicine and Public Health, Scripps Institution of Oceanography, University of California, 8622 Kennel Way, La Jolla, San Diego, CA, USA. ✉email: leo.zabrocki@psemail.eu
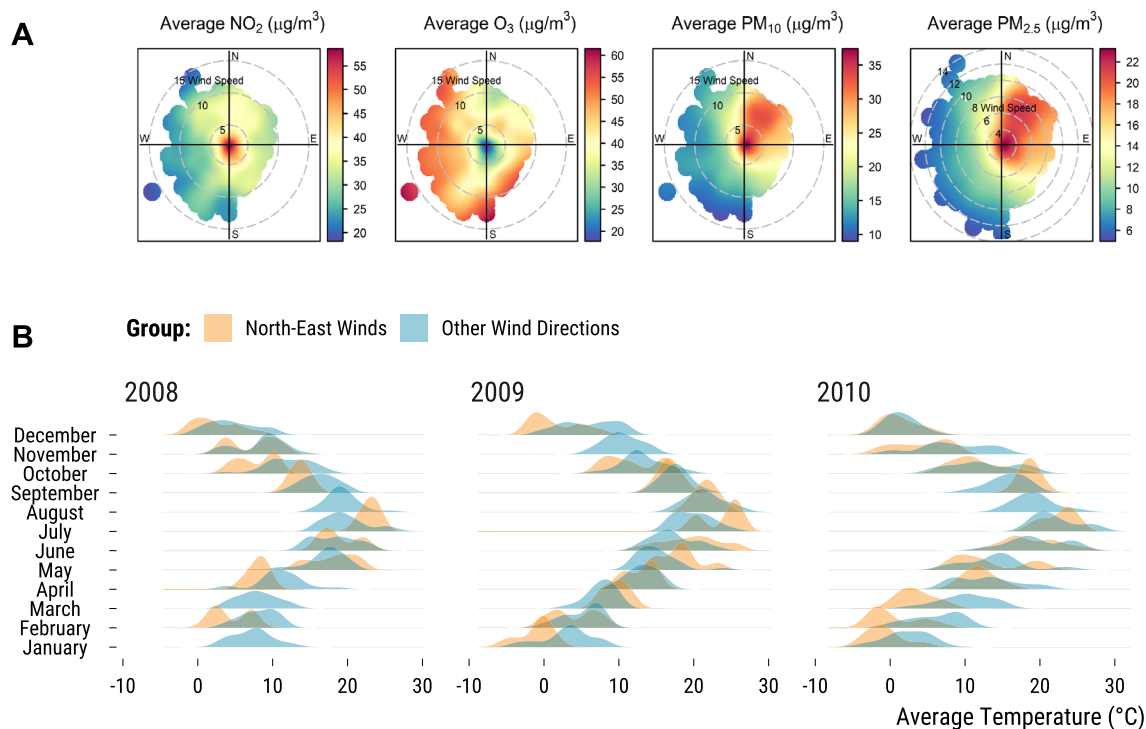
**Figure 1.** Polar plots of air pollutant concentrations predicted by wind components and average temperature imbalance of wind directions by year and month. In panel (**A**), each plot represents the concentrations (in $\mu g/m^3$) of an air pollutant that were predicted using a generalized additive model based on a smooth isotropic function of the two wind components $u$ and $v$[15]. The direction from which the wind blows is described on a 360° compass rose and wind speed (in m/s) is represented by a series of increasing circles starting from the intersection of the two cardinal directions axes where wind speed is null: the farther the circle is away from the intersection, the faster the wind speed is. In panel (**B**), the density distribution of the average temperature (in °C) is drawn for North-East winds (orange colour) and other wind directions (blue colour). The figure is divided into subplots by month and year (2008–2010).

wind patterns. For instance, in Panel B of Fig. 1, the density distribution of the average temperature (°C) is not similar for the groups of wind directions. We must take into account this confounding variable if we want to make the as-if random distribution of North-East wind more credible. Multivariate linear regression have been the standard approach to help achieve this goal but more flexible methods such as generalized additive models and machine learning algorithms could also be used[16,17]. Yet, even a very flexible model will not overcome the second issue visible in Panel B of Fig. 1: as for January 2008, the model will sometimes depend on extrapolation since there are no empirical counterfactuals to estimate what would have happened had the wind blown from the North-East. Finally, it could be argued that we fail to adjust for a confounding variable which we have not measured. In addition to explaining with qualitative arguments why it is not likely the case, we should also try to quantify the bias induced by an unmeasured confounder.

In this paper, we show how we can evaluate the extent to which studies exploiting wind directions as instrumental variables could be prone to the issues raised above. To achieve this goal, we follow the four consecutive stages of the causal inference pipeline proposed by[18,19] that explicitly embed the design of this type of observational study within an hypothetical randomized experiment[20–23].

First, in a *conceptual stage*, we clearly state the causal question of interest using the Neyman–Rubin potential outcomes framework[24,25]. Our treatment of interest is the effect of North-East winds on air pollution compared to other wind directions. To estimate this effect, for treated days with winds blowing from the North-East, we need to impute the concentrations that would have been observed had winds blown from other directions. The issue is that wind patterns are not randomly assigned: control days with wind blowing from other directions are not similar to treated days.

We therefore implement a *design stage* where we approximate a pairwise randomized experiment using a matching algorithm recently designed for air pollution studies[26]. Matching is a transparent method to adjust for confounders without making parametric assumption and directly looking at observed outcomes[27,28]. Given a set of chosen covariate distances, each treated day is matched to its closet control day. This method also avoids model extrapolation since treated days for which no control days exist in the data are discarded from the analysis.

The third step is an *analysis stage* where we estimate the influence of North-East winds on air pollutant concentrations. We simply compute the average difference in concentrations between matched treated and control days and rely on Neymanian inference to compute an estimate of the sampling variability[22]. The last and fourth step is to carry out a *sensitivity analysis*. Throughout the previous steps, we must make the strong assumption that no unmeasured variables could be related both to wind patterns and air pollutant concentrations. Quantitative
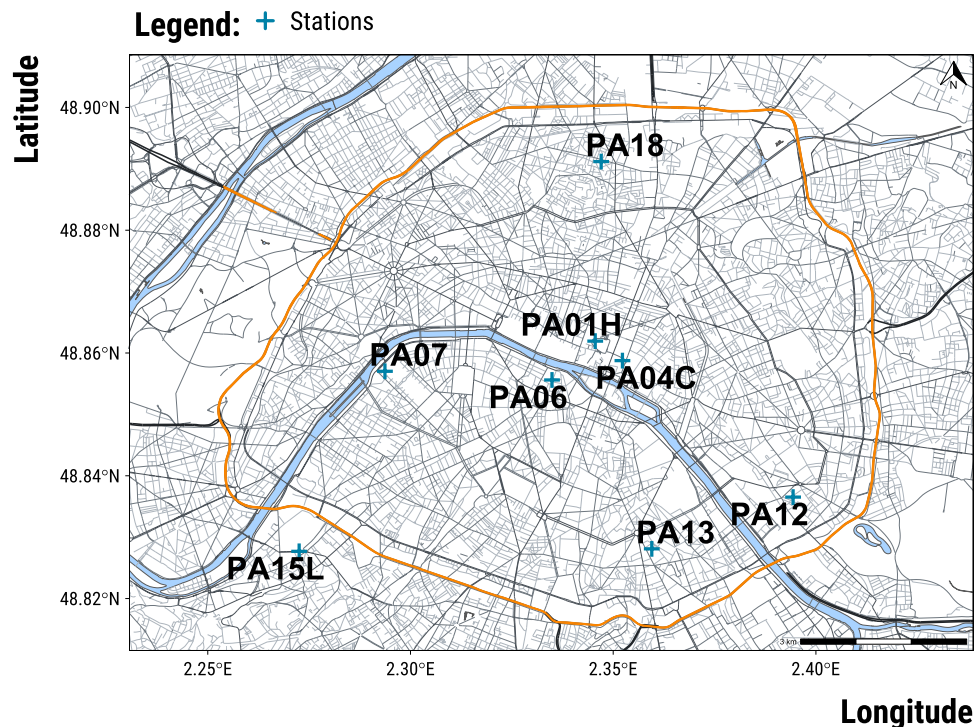
**Figure 2.** Map of road network and location of air pollution measuring stations in Paris, France . Grey lines represent the road network. The orange line is the orbital ring surrounding Paris. Blue crosses are the locations of air pollution measuring stations. $NO_2$ concentrations are measured at stations PA07, PA12, PA13, PA18; $O_3$ concentrations at PA13, PA18; $PM_{10}$ at PA18; $PM_{2.5}$ at PA01H and PA04C. The map was created with the R programming language (version 4.1.0)[36], data were provided by OpenStreetMap[37] and retrieved with the `osmdata` package[38].

bias analysis was initially proposed by[29] to assess which magnitude of hidden bias would be required to alter observed results. We follow here the method developed by[21,30].

With this study, we aim to bring two contributions to the causal inference literature on the acute health effects of air pollution. First, we show that using wind directions as instrumental variables requires more caution to make the assumption that they are "as-if" randomly distributed according to observed covariates convincing. The effective sample size where treated and control units are similar on a set of observed covariates is actually small. The standard approach used in the literature based on multivariate regression models will therefore rely on its ability to adjust correctly for the functional forms of covariates and extrapolate to units without empirical counterfactuals. Second, our quantitative bias analysis reveals that the estimated increase in particulate matter concentrations due to North-East winds is relatively robust to the presence of hidden bias. Even if an unobserved confounding factor is twice more common among days with winds blowing from the North-East than among days with winds from other directions, the large range of estimates consistent with the data remains positive.

We also hope that the approach we propose in this paper could be of interest to atmospheric scientists. The fact that wind patterns play a key role in the variation of air pollution concentrations is obviously not new[31–34]. Yet, causal inference methods have rarely been implemented in atmospheric science to estimate the influence of weather parameters on air pollution. We believe that mimicking a randomized experiment corresponds to an intuitive approach and could complement source apportionment and emission inventory approaches. While wind is non manipulable, emission sources are and our framework could also serve as a stepping-stone to evaluate potential interventions to control emissions—if a source is shut-down in the North-East of Paris, would wind blowing from this direction influence less specific air pollutant concentrations?

We took great care to make our work fully reproducible to help researchers implement but also improve and criticize our approach. Data and detailed **R** codes are available at https://lzabrocki.github.io/design_stage_wind_air_pollution/ and backed-up in an Open Science Framework repository[35].

## Methods

**Data.** We built a dataset combining daily time series of air pollutant concentrations and weather parameters in Paris over the 2008–2018 period. We chose to carry out an analysis at the daily level as done in studies on the acute health effects of air pollution[3,4,6].

First, we obtained hourly air quality data from AirParif, the local air quality monitoring agency. Figure 2 displays the location of the selected measuring stations. Using a 2.5% trimmed mean, we first averaged at the daily level the concentrations ($\mu g/m^3$) of background measuring stations for $NO_2$, $O_3$ and $PM_{10}$. For a given

day, if more than 3 hourly readings were missing, the average daily concentration was set to missing. The proportion of missing values for stations ranged from 2.8% up to 9.1%. We also computed the average daily concentrations of $PM_{2.5}$ but 25% of the recordings were missing: the air pollutant was not measured by Airparif between 2009/09/22 and 2010/06/23. It is important to note that we did not retrieve data from traffic monitors but only from background monitors as they are used to assess the residential exposure of a city population in epidemiological studies.

We then retrieved meteorological data from the single monitoring station located in the South of the city and ran by the French national meteorological service Météo-France. We extracted daily observations on wind speed (m/s), wind direction (measured on a 360° wind rose where 0° is the true North), the average temperature (°C), and the rainfall duration (min). Weather parameters had very few missing values (e.g., at most 2.5% of observations were missing for the rainfall duration).

Finally, to avoid working with a reduced sample size, we imputed missing values for all variables but $PM_{2.5}$. There were no clear patterns in the missingness of $NO_2$, $O_3$ and $PM_{10}$ concentrations. We used the chained random forest algorithm implemented by the **R** package missRanger[39]. A small simulation exercise showed that it had good performance for imputing $NO_2$ concentrations (the absolute difference between observed and imputed values was equal to $3.2\,\mu g/m^3$ for an average concentration of $37.6\,\mu g/m^3$) but was much less effective for imputing $PM_{10}$ concentrations (the absolute difference between observed and imputed values was equal to 6.1 $\mu g/m^3$ for an average concentration of $23.4\,\mu g/m^3$). Once the data were imputed, we averaged the air pollutant concentrations at the city level as it is the spatial level of analysis used in[3,4].

Further details on data wrangling and an exploratory analysis of the data can be found in the supplementary materials (https://lzabrocki.github.io/design_stage_wind_air_pollution, tab Data). We were not allowed to share weather data from Météo-France so we added some noise to the weather parameters.

**A causal inference pipeline.** We present below the four stages of the causal inference pipeline we advocate to use for improving the design of air pollution studies based on wind patterns. Its implementation was done with the R programming language (version 4.1.0)[36].

*Stage 1: Defining the treatment of interest.* The first step of our causal inference approach is to clearly state the question we are trying to answer: *What is the effect of North-East winds on particulate matter in Paris over the 2008–2018 period?* This question is motivated by the exploratory analysis of Fig. 1 and research in atmospheric science on the sources of particulate matter located in the North-East of the city. Our treatment of interest is therefore defined as the comparison of air pollutant concentrations when winds are blowing from the North-East (10°–90°) with concentrations when wind come from other directions. We frame this question in the Rubin–Neyman causal framework[24,25]. Our units are 4018 days indexed by $i$ ($i=1,..., I$). For each day, we define our treatment indicator $W_i$ which takes two values. It is equal to 1 if the unit is treated (the wind blows from the North-East), and 0 if the unit belongs to the control group (the wind is blowing from another direction). Under the Stable Unit Treatment Value Assumption (STUVA), we assume that each day can have two potential concentrations in $\mu g/m^3$ for an air pollutant: $Y_i(1)$ if the wind blows from the North-East and $Y_i(0)$ if the wind blows from another direction.

The fundamental problem of causal inference states that we can only observe for each day one of these two potential outcomes: it is a missing data problem[40,41]. The observed concentration of an air pollutant $Y^{obs}$ is defined as $Y^{obs} = (1-W_i) \times Y_i(0) + W_i \times Y_i(1)$. If the unit is treated, we observe $Y_i(1)$. If it is a control, we observe $Y_i(0)$. To estimate the effect of North-East winds on air pollutant concentrations, we therefore need to impute the missing potential outcomes of treated units—what would have been the air pollutant concentrations if the wind had blown from another direction?

*Stage 2: Designing the hypothetical randomized experiment.* The second stage of our causal inference pipeline is to embed our non-randomized study within an hypothetical randomized experiment. We are dealing with an observational study where North-East winds are not randomly distributed through a year and are correlated with other weather parameters influencing air pollutant concentrations. In Fig. 3, we plot, for each month, the absolute standardized mean differences between treated and control units for the average temperature, relative humidity and wind speed: most differences are superior to 0.1, which is often considered as a threshold to assess the imbalance of covariates.

To better approximate a randomized experiment, we must therefore find the subset of treated units which are similar to control units. Formally, we want to make plausible for this subset of units the assumption that the treatment assignment is independent from the potential outcomes of units given their covariates **X**: Pr(**W** | **X**, **Y**(0), **Y**(1)) = Pr(**W** | **X**). The issue is that some units' covariates are observed while other are not. Unlike a randomized experiment where both observed and unobserved covariates will be, on average, balanced across treatment and control groups, we must assume that no unobserved covariates affect the treatment assignment.

Matching methods are particularly convenient to design hypothetical randomized experiments. Contrary to standard regression approaches, matching is a non-parametric way to adjust for observed covariates while avoiding model extrapolation since units without counterfactuals in the data are discarded from the analysis. Specifically, we use a constrained matching algorithm to design a pairwise randomized experiment where, for each pair, the probability of receiving the treatment is equal to 0.5 (see[26] for further details on the algorithm). Each treated unit is matched to its closest unit given a set of covariate constraints which represent the maximum distance, for each covariate, allowed between treated and control units. We match on the two sets of covariates influencing both wind directions and air pollutant concentrations.
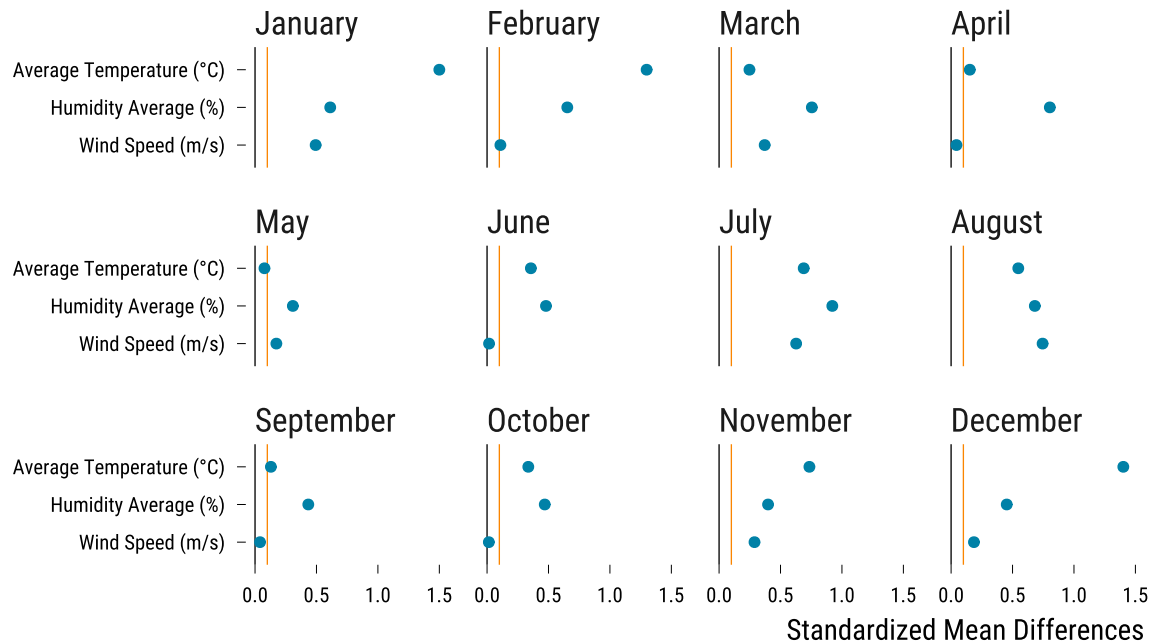
**Figure 3.** Evidence of imbalance for weather covariates . For each month, we compute the absolute standardized differences for continuous weather covariates between treated and control groups. These differences are represented as blue points. The vertical orange line is the 0.1 threshold which is used in the matching literature to spot covariates imbalance. The vertical black line is at 0.

First, we match on calendar variables such as the Julian date, weekend, holidays and bank days indicators. A treated unit could be matched up to a control unit with a maximum distance of 60 days. If we extend this distance, it would be easier to match treated units to control units but the treatment effect could be biased by seasonal variation in air pollutant concentrations. We match exactly treated and control units for the other calendar indicators.

Second, we match on weather variables. The average temperature between treated and control units could not differ by more than 5°. The difference in wind speed must be less than 0.5 m/s. The rainfall duration (divided in four ordinal categories) needs to be the same and the absolute difference in average humidity could be up to 12 percentage points. We also force the absolute difference in $PM_{10}$ concentrations in the previous day to be less or equal to $8 \, \mu g/m^3$. The thresholds we set up were chosen through an iterative process were we checked (1) that they led to balanced sample of treated and control units and (2) that there were enough matched pairs to draw our inference upon.

Finally, the Stable Unit Treatment Value Assumption (SUTVA) requires that there is no interference between units and no hidden variation of the treatment. To make this assumption more plausible, we discard from the analysis the matched pairs for which the distance in days is inferior to 4 days and make sure that the first lag of the treatment indicator for treated and control units.

*Stage 3: Analyzing the experiment using Neymanian inference.* In the third stage, we proceed to the analysis of our hypothetical pairwise randomized experiment. Several modes of statistical inference such as Fisherian, Neymanian or Bayesian could be implemented[42]. Here, we take a Neymanian perspective where the potential outcomes are assumed to be fixed and the treatment assignment is the basis of inference. Our goal is to measure the average causal effect for the sample of matched units. We assume that each of the two units of a matched pair $j$ has two potential concentrations for an air pollutant. If we were able to observe these potential outcomes, we could simply measure the effect of North-East winds on air pollutant concentrations by computing the finite-sample average treatment effect for matched treated units $\tau_{fs}$. We would first compute for each pair the mean difference in concentrations and then average the differences over the $J$ pairs. While we only observe one potential outcome for each unit, we can nonetheless estimate $\tau_{fs}$ with the average of observed pair differences $\hat{\tau}$:

$$\hat{\tau} = \frac{1}{J} \sum_{j=1}^{J} (Y_{t,j}^{obs} - Y_{c,j}^{obs}) = \overline{Y}_t^{obs} - \overline{Y}_c^{obs}$$

Here, the subscripts $t$ and $c$ respectively indicate if the unit in a given pair is treated or not. Since there are only one treated and one control unit within each pair, the standard estimate for the sampling variance of the average of pair differences is not defined. We can however compute a conservative estimate of the variance[22]:

$$\hat{\mathbb{V}}(\hat{\tau}) = \frac{1}{J(J-1)} \sum_{j=1}^{J} (Y_{t,j}^{obs} - Y_{c,j}^{obs} - \hat{\tau})^2$$

We finally compute an asymptotic 95% confidence interval using a Gaussian distribution approximation:

$$\text{CI}_{0.95}(\tau_{fs}) = \left( \hat{\tau} - 1.96 \times \sqrt{\hat{\mathbb{V}}(\hat{\tau})}, \ \hat{\tau} + 1.96 \times \sqrt{\hat{\mathbb{V}}(\hat{\tau})} \right)$$

The obtained 95% confidence interval gives the set of effect sizes compatible with our data[43].

*Stage 4: Sensitivity analysis.* The fourth step of our causal inference pipeline is to explore how sensitive our analysis is to violation of the assumptions it relies upon. We carry out three types of robustness checks.

First, we make the strong assumption that the treatment assignment is as-if random: winds blowing from the North-East occur randomly conditional on a set of measured covariates. Other researchers could however argue that we fail to adjust for unmeasured variables influencing both the occurrence of North-East winds and air pollutant concentrations. Within matched pairs, these unobserved counfounders could make the treated day more likely to have wind blowing from the North-East than the control day. We therefore implement the quantitative bias analysis, also called sensitivity analysis, that was developed by[21,30]. It allows us to explore how our results would be altered by the effect of an unobserved confounder on the treatment odds, denoted by $\Gamma$. In our matched pairwise experiment, we assume that within each pair, control and treated days have the odds to see the wind blowing from the North-East: the odds of treatment is such that $\Gamma = 1$. The quantitative bias analysis allows to compute the 95% confidence intervals obtained for different values of bias the unmeasured confounder has on the treatment assignment. For instance, if we assume that an unmeasured confounder has a small effect on the odds of treatment (i.e., for a $\Gamma > 1$ and close to 1) but the resulting 95% confidence interval becomes completely uninformative, it would imply that our results are highly sensitive to hidden bias. Conversely, if we assume that an unmeasured confounder has a strong effect on the odds of treatment (i.e., for a large $\Gamma$) and we find that the resulting 95% confidence interval remains similar, it would imply that our results are very robust to hidden bias. In a complementary manner, we also check whether unmeasured biases could be present by using the first daily lags of air pollutant concentrations as control outcomes[44]. If our matched pairs are indeed similar in terms of unobserved covariates, the treatment occurring in $t$ should not influence concentration of air pollutants in $t-1$.

Second, for many matched pairs, air pollutant concentrations were imputed using the chained random forest algorithm[39]. We check whether the results are sensitive to the imputation by re-running the analysis for the non-missing concentrations.

Third, we make sure that the treatment assignment within pairs was effective to increase the precision of estimates. We compare the estimate of the sampling variance of a pairwise randomized experiment to the one of a completely randomized experiment. If the estimate of sampling variability for the pairwise experiment is smaller than the estimate of sampling variability for a complete experiment, it means that our matching procedure was successful to match similar units within pairs compared to randomly selected units[22].

## Results

### Performance of the matching procedure.
Our initial dataset consists in 4018 daily observations, divided into 912 treated units and 3106 control units. The matching procedure results in 121 pairs of matched treated-control units—only 13% of treated units could be matched to similar control units given the constraints we set. In the supplementary materials (https://lzabrocki.github.io/design_stage_wind_air_pollution/4_comparing_initial_to_matched_data.html), we show that the matched sample has different characteristics from the initial sample: observations belong more to the period ranging from May to October, their average temperature is higher and their relative humidity is lower.

In Fig. 4, we display how the balance of continuous and categorical covariates improves after the matching procedure. Blue dots represent either the absolute mean differences between treated and control units for continuous variables or the absolute differences in percentage points for categorical variables. For continuous covariates, the average standardized mean differences between treated and control days is 0.26 before matching and reduces to 0.07 after the procedure. For categorical covariates, the average difference in percentage points diminishes from 6.2 to 1.8 after matching. Our matching procedure therefore leads to a consequent reduction of our sample size but allows us to compare treated units that are more similar to control units. A complete analysis of the balance improvement for each covariate is available in the supplementary materials (https://lzabrocki.github.io/design_stage_wind_air_pollution/6_checking_balance_improvement.html).

### North-east wind effects on air pollutant concentrations.
For each air pollutant, we plot in Fig. 5 the estimated average difference in concentration ($\mu g/m^3$) between North-East winds and other wind directions. We also display the estimated differences for the previous day and the following day. Thick lines represent the 95% confidence intervals while thin lines are the 99% confidence intervals. The third panel of Fig. 5 confirms the exploratory analysis of the polar plot. When wind blows from the North-East, $PM_{10}$ concentrations increase by $4.4\,\mu g/m^3$, with the lower and upper bounds of the 95% confidence being respectively equal to an increase by $1.7\,\mu g/m^3$ and $7.2\,\mu g/m^3$. The estimated difference represents an 18% increase in the average concentration of $PM_{10}$. We also observe a positive difference of 25% in $PM_{10}$ concentrations the following day (point estimate of 4.9; 95% CI 1.8, 8.1).
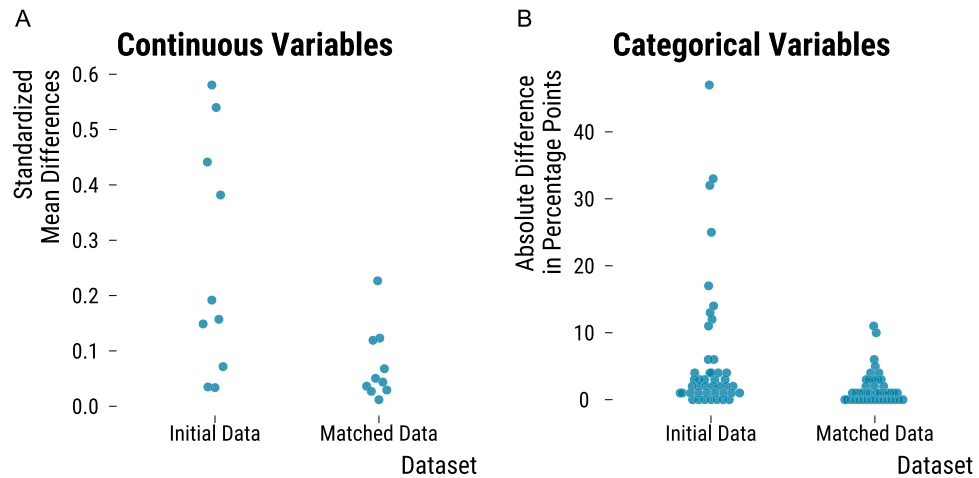
**Figure 4.** Overall balance improvement in continuous and categorical covariates . In Panel (**A**), we plot, before and after matching, the absolute standardized differences in continuous covariates between treated and control groups. Each blue dot represents an absolute mean difference for a given covariate. In panel (**B**), we plot, before and after matching, the absolute difference in percentage points for categorical covariates.
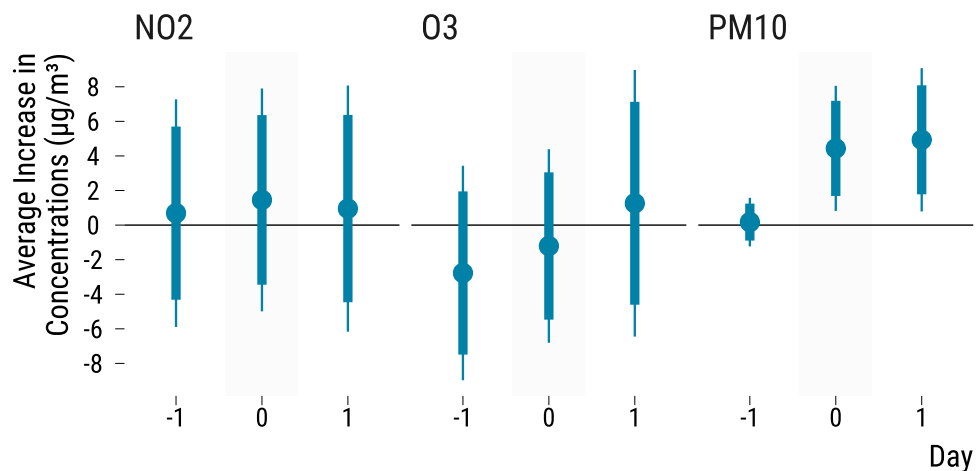


**Figure 5.** Effects of North-East winds on air pollutant concentrations . In each panel, we plot the estimated effects of North-East winds on air pollutant concentrations for the previous, current and following days. Point estimates are depicted by blue points; blue thick lines are 95% confidence intervals and thin lines are 99% confidence intervals. The 95% and 99% confidence intervals associated with the estimated average difference in $PM_{10}$ in the first lag are smaller than other intervals for the following days since we added a constraint in the matching procedure for this lag of the air pollutant.

North-East winds do not seem to influence $NO_2$ (point estimate of 1.5; 95% CI − 3.4, 6.4), and $O_3$ (point estimate of − 1.2; 95% CI − 5.5, 3.1) concentrations on the current day. This is also the case for the concentrations of these two air pollutants on the following day.

Regarding the effects of North-East winds on $PM_{2.5}$, we restrain our analysis to pairs without missing concentrations. For the current and following days, we respectively find an average increase of 1.4 µg/m$^3$ (95% CI − 0.6, 3.4) and 2.7 µg/m$^3$ (95% CI 0.8, 4.5). These point estimates respectively represent a 8.8% and a 17% relative increases in $PM_{2.5}$ concentrations.

**Sensitivity analysis.** Our quantitative bias analysis reveals that if we have failed to adjust for an unobserved confounder twice more common among treated days, the resulting 95% confidence intervals for the estimated effects of North-East winds on $PM_{10}$ would be equal to (0.5, 9) for the current day and to (− 0.2, 10) for the the following day. Confidence intervals are still consistent with mostly positive effects but are relatively wide. As a complementary test for unobserved confounders, we also check that the occurrence of North-East winds on the current day does not have any effect on concentrations measured in the previous day. Reassuringly, for $NO_2$ and $O_3$, 95% confidence intervals do not suggest clear negative or positive average differences in concentrations as shown in Fig. 5 (for $PM_{2.5}$, the estimated average difference is − 0.1 µg/m$^3$ (95% CI − 1.2, 1)).

In the supplementary materials (https://lzabrocki.github.io/design_stage_wind_air_pollution/7_analyzing_results.html), we check whether the imputation of missing air pollutant concentrations did not drive our results. For $NO_2$, $O_3$ and $PM_{10}$, 13%, 8% and 7% of concentrations were respectively imputed. We replicate our analysis on the subset of pairs without missing observations: point estimates remain very similar but confidence intervals are a bit larger due to the sample size loss. This robustness check implies that our imputation did not bias our estimates.

Finally, the pairwise design of our hypothetical experiment does not help increase the precision of the estimated differences in $PM_{10}$ concentrations. The standard error under a completely randomized assignment is equal to 1.35 while the one of a pairwise randomized assignment is 1.4. The pairwise design however increases the precision estimates for $O_3$ by 23% for $O_3$ but decreases the precision by 42% for $NO_2$.

## Discussion

In our study, we follow a causal inference pipeline to craft a hypothetical experiment for measuring the effects of North-East winds on daily particulate matter concentrations in Paris. Our constrained pair matching algorithm enables us to find the subset of treated days that were similar to control days for a set of calendar and weather confounding factors. Compared to a statistical adjustment based on a multivariate regression model, matching is non-parametric and avoids to extrapolate to units without empirical counterfactuals. At the very heart of this method, graphical displays of covariates balance allow to check in a transparent manner whether the as-if random distribution of the treatment was achieved conditional on observed confounders. We were surprised that covariates balance could only be achieved for 13% of treated units. It would be an interesting question for future research to see if alternative methods such as cardinality matching or bayesian additive regression trees lead to similar results[45–47]. The relevant structure of the hypothetical experiment to target should also be of interest since our pair matching algorithm failed to increase the precision of estimates compared to a completely randomized assignment of the treatment.

The difficulty to find similar treated and control units could lead researchers interested in the acute health effects of air pollution to worry that instrumental variable strategies exploiting wind patterns and based on multivariate regression models might suffer from extrapolation bias[10,27]. In the supplementary materials (https://lzabrocki.github.io/design_stage_wind_air_pollution/7_analyzing_results.html), we show that results based on an outcome regression approach, even if they are based on the entire sample, are consistent with those found with the matched data. This may increase the confidence in the capability of a multivariate regression model to correctly extrapolate. Matching estimates are however much less precise. Further research is therefore needed to better understand if improving the design stage of instrument variable studies with matching methods is feasible given the small sample size it entails[48–51]. If it is the case, could matching methods actually lead to different results[52–54]?

In addition to providing evidence on the effective sample size for which covariates balance was achievable, our study was the occasion to assess whether the estimated effects of North-East wind on particulate matters were robust hidden bias. It would require an unmeasured confounder twice more common among treated days to raise doubt on the direction of the estimated effects. This raises our confidence in the assumption that North-East wind are also randomly distributed according to unobserved variables. To the best of our knowledge, this assumption was waiting to be quantitatively evaluated. This could be explained by the fact that the sensitivity analysis we rely on was developed for pairwise matched data[30]. As an alternative, researchers wishing to keep working with a regression approach could implement the new method developed by[55,56].

Finally, our study presents two main limits regarding the improvement of the design stage of air pollution studies based on wind directions. The first limit concerns the definition of the contrast of interest, that is to say the difference of air pollutant concentrations between North-East winds and other wind directions. If this comparison is easy to understand, the treatment we defined is not manipulable contrary to those found in randomized controlled trials. It might lack a certain appeal to policy-makers as our estimates only indicate whether North-East winds lead to higher particulate matter concentrations than other wind directions[57,58], without determining the origin of the sources emitting the air pollutant. To overcome this limit, a study exploiting variations in wind directions should be combined with a clear shock on one of the sources emitting an air pollutant. For instance, in a recent paper in Southern California[34], it was shown that Santa Ana winds have a predominant ventilation effect on $PM_{2.5}$ but when inland wildfires occur, Santa Ana winds are instead increasing $PM_{2.5}$ levels on the coast.

The second limit revolves around the assumption that, for wind direction to be a valid instrument, its effects on a health outcome must be fully mediated by a single air pollutant[7–9]. As recognized by researchers, studies exploiting wind patterns could violate this assumption if changes in wind direction affect simultaneously several air pollutants. In our study, once the data are matched, it seems that North-East winds only influence particulate matter, which could reinforce the credibility of the assumption. Yet, this should not be always the case as it would be highly dependent on the city and air pollutant investigated. Methodological work is much needed to understand in which cases the air pollutants co-variance structure could lead to biased dose-response. In a recent work[59], propose to run a multi-pollutant model where each air pollutant concentration is predicted by selecting the optimal set of instrumental variables using least absolute shrinkage and selection operator (lasso). The authors show that results of an instrumented multi-pollutant model can be very different from those found by single-pollutant models. It remains to be studied if matching could also help limit this well-known issue.

## References

1. Schlenker, W. & Walker, W. R. Airports, air pollution, and contemporaneous health. *Rev. Econ. Stud.* **83**(2), 768–809 (2016).

2. Arceo, E., Hanna, R. & Oliva, P. Does the effect of pollution on infant mortality differ between developing and developed countries? Evidence from Mexico City. *Econ. J.* **126**(591), 257–280 (2016).
3. Schwartz, J., Bind, M.-A. & Koutrakis, P. Estimating causal effects of local air pollution on daily deaths: Effect of low levels. *Environ. Health Perspect.* **125**(1), 23–29 (2017).
4. Schwartz, J., Fong, K. & Zanobetti, A. A national multicity analysis of the causal effect of local pollution, no 2, and pm 2.5 on mortality. *Environ. Health Perspect.* **126**(8), 087004 (2018).
5. Halliday, T. J., Lynham, J. & de Paula, A. Vog: Using volcanic eruptions to estimate the health costs of particulates. *Econ. J.* **129**(620), 1782–1816 (2019).
6. Deryugina, T., Heutel, G., Miller, N. H., Molitor, D. & Reif, J. The mortality and medical costs of air pollution: Evidence from changes in wind direction. *Am. Econ. Rev.* **109**(12), 4178–4219 (2019).
7. Angrist, J. D., Imbens, G. W. & Rubin, D. B. Identification of causal effects using instrumental variables. *J. Am. Stat. Assoc.* **91**(434), 444–455 (1996).
8. Angrist, J. D. & Pischke, J.-S. *Mostly Harmless Econometrics* (Princeton University Press, 2008).
9. Baiocchi, M., Cheng, J. & Small, D. S. Instrumental variable methods for causal inference. *Stat. Med.* **33**(13), 2297–2340 (2014).
10. King, G. & Zeng, L. The dangers of extreme counterfactuals. *Polit. Anal.* **14**(2), 131–159 (2006).
11. Stuart, E. A. & Rubin, D. B. Best practices in quasi-experimental designs. *Best Pract. Quantitative Methods* **20**, 155–176 (2008).
12. Bressi, M. *et al.* Sources and geographical origins of fine aerosols in Paris (France). *Atmos. Chem. Phys.* **14**(16), 8813–8839 (2014).
13. Petetin, H. *et al.* A novel model evaluation approach focusing on local and advected contributions to urban pm 2.5 levels-application to Paris, France. *Geosci. Model Dev.* **7**(4), 1483–1505 (2014).
14. Stirnberg, R. *et al.* Meteorology-driven variability of air pollution (pm 1) revealed with explainable machine learning. *Atmos. Chem. Phys.* **21**(5), 3919–3948 (2021).
15. Carslaw, D. C. & Ropkins, K. Openair-An r package for air quality data analysis. *Environ. Modell. Softw.* **27**, 52–61 (2012).
16. Grange, S. K., Carslaw, D. C., Lewis, A. C., Boleti, E. & Hueglin, C. Random forest meteorological normalisation models for swiss pm 10 trend analysis. *Atmos. Chem. Phys.* **18**(9), 6223–6239 (2018).
17. Grange, S. K. & Carslaw, D. C. Using meteorological normalisation to detect interventions in air quality time series. *Sci. Total Environ.* **653**, 578–588 (2019).
18. Bind, M.-A.C & Rubin, D. B. Bridging observational studies and randomized experiments by embedding the former in the latter. *Stat. Methods Med. Res.* **28**(7), 1958–1978 (2019).
19. Bind, M.-A.C & Rubin, D. B. The importance of having a conceptual stage when reporting non-randomized studies. *Biostatist. Epidemiol.* **20**, 1–10 (2021).
20. Rubin, D. B. For objective causal inference, design trumps analysis. *Ann. Appl. Stat.* **2**(3), 808–840 (2008).
21. Rosenbaum, P. R., Rosenbaum, P. R. & Briskman. *Design of Observational Studies* Vol. 10 (Springer, 2010).
22. Imbens, G. W. & Rubin, D. B. *Causal Inference in Statistics, Social, and Biomedical Sciences* (Cambridge University Press, 2015).
23. Hernán, M. A. & Robins, J. M. Using big data to emulate a target trial when a randomized trial is not available. *Am. J. Epidemiol.* **183**(8), 758–764 (2016).
24. Neyman, J. Sur les applications de la théorie des probabilités aux experiences agricoles: Essai des principes. *Roczniki Nauk Rolniczych* **10**, 1–51 (1923).
25. Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.* **66**(5), 688 (1974).
26. Sommer, A. J., Leray, E., Lee, Y. & Bind, M.-A.C. Assessing environmental epidemiology questions in practice with a causal inference pipeline: An investigation of the air pollution-multiple sclerosis relapses relationship. *Stat. Med.* **40**(6), 1321–1335 (2021).
27. Ho, D. E., Imai, K., King, G. & Stuart, E. A. Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Polit. Anal.* **15**(3), 199–236 (2007).
28. Stuart, E. A. Matching methods for causal inference: A review and a look forward. *Stat. Sci. Rev. J. Inst. Math. Stat.* **25**(1), 1 (2010).
29. Cornfield, J. *et al.* Smoking and lung cancer: Recent evidence and a discussion of some questions. *J. Natl. Cancer Inst.* **22**(1), 173–203 (1959).
30. Fogarty, C. B. Studentized sensitivity analysis for the sample average treatment effect in paired observational studies. *J. Am. Stat. Assoc.* **115**(531), 1518–1530 (2020).
31. Wilson, W. E. & Suh, H. H. Fine particles and coarse particles: Concentration relationships relevant to epidemiologic studies. *J. Air Waste Manage. Assoc.* **47**(12), 1238–1249 (1997).
32. Hoek, G. *et al.* A review of land-use regression models to assess spatial variation of outdoor air pollution. *Atmos. Environ.* **42**(33), 7561–7578 (2008).
33. Tai, A. P. K., Mickley, L. J. & Jacob, D. J. Correlations between fine particulate matter (pm 2.5) and meteorological variables in the united states: Implications for the sensitivity of pm 2.5 to climate change. *Atmos. Environ.* **44**(32), 3976–3984 (2010).
34. Aguilera, R., Gershunov, A., Ilango, S. D., Guzman-Morales, J. & Benmarhnia, T. Santa ana winds of Southern California impact pm 2.5 with and without smoke from wildfires. *GeoHealth* **4**(1), e2019GH000225 (2020).
35. Zabrocki, L. Improving the design stage of air pollution studies based on wind patterns. https://osf.io/7x23u/, 2022.
36. R Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, 2021).
37. OpenStreetMap contributors. Planet dump retrieved from https://planet.osm.org. https://www.openstreetmap.org, 2017.
38. Padgham, M., Rudis, B., Lovelace, R. & Salmon, M. Osmdata. *J. Open Source Softw.* **2**, 14 (2017).
39. Mayer, M. missranger: Fast imputation of missing values (2.1. 0). https://CRAN.R-project.org/package=missRanger. R package version, 2(0), 2019.
40. Holland, P. W. Statistics and causal inference. *J. Am. Stat. Assoc.* **81**(396), 945–960 (1986).
41. Ding, P. & Li, F. Causal inference: A missing data perspective. *Stat. Sci.* **33**(2), 214–237 (2018).
42. Rubin, D. B. Practical implications of modes of statistical inference for causal effects and the critical role of the assignment mechanism. *Biometrics* **20**, 1213–1234 (1991).
43. Amrhein, V., Trafimow, D. & Greenland, S. Inferential statistics as descriptive statistics: There is no replication crisis if we don't expect replication. *Am. Stat.* **73**(sup1), 262–270 (2019).
44. Rosenbaum, P. *Observation and Experiment* (Harvard University Press, 2018).
45. Hill, J. L. Bayesian nonparametric modeling for causal inference. *J. Comput. Graph. Stat.* **20**(1), 217–240 (2011).
46. Hill, J. & Su, Y.-S. Assessing lack of common support in causal inference using bayesian nonparametrics: Implications for evaluating the effect of breastfeeding on children's cognitive outcomes. *Ann. Appl. Stat.* **20**, 1386–1420 (2013).
47. Visconti, G. & Zubizarreta, J. R. Handling limited overlap in observational studies with cardinality matching. *Observ. Stud.* **4**(1), 217–249 (2018).
48. Small, D. S. & Rosenbaum, P. R. War and wages: The strength of instrumental variables and their sensitivity to unobserved biases. *J. Am. Stat. Assoc.* **103**(483), 924–933 (2008).
49. Baiocchi, M., Small, D. S., Yang, L., Polsky, D. & Groeneveld, P. W. Near/far matching: A study design approach to instrumental variables. *Health Serv. Outcomes Res. Methodol.* **12**(4), 237–253 (2012).
50. Kang, H., Kreuels, B., May, J. & Small, D. S. Full matching approach to instrumental variables estimation with application to the effect of malaria on stunting. *Ann. Appl. Stat.* **10**(1), 335–364 (2016).
51. Keele, L. & Morgan, J. W. How strong is strong enough? Strengthening instruments through matching and weak instrument tests. *Ann. Appl. Stat.* **10**(2), 1086–1106 (2016).

52. Schwartz, J., Austin, E., Bind, M.-A., Zanobetti, A. & Koutrakis, P. Estimating causal associations of fine particles with daily deaths in boston. *Am. J. Epidemiol.* **182**(7), 644–650 (2015).
53. Baccini, M., Mattei, A., Mealli, F., Bertazzi, P. A. & Carugno, M. Assessing the short term impact of air pollution on mortality: A matching approach. *Environ. Health* **16**(1), 1–12 (2017).
54. Forastiere, L., Carugno, M. & Baccini, M. Assessing short-term impact of pm 10 on mortality using a semiparametric generalized propensity score approach. *Environ. Health* **19**(1), 1–13 (2020).
55. Cinelli, C. & Hazlett, C. Making sense of sensitivity: Extending omitted variable bias. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **82**(1), 39–67 (2020).
56. Cinelli, C. & Hazlett, C. An omitted variable bias framework for sensitivity analysis of instrumental variables. Working Paper (2020).
57. Zigler, C. M. & Dominici, F. Point: Clarifying policy evidence with potential-outcomes thinking-beyond exposure-response estimation in air pollution epidemiology. *Am. J. Epidemiol.* **180**(12), 1133–1140 (2014).
58. Dominici, F. & Zigler, C. Best practices for gauging evidence of causality in air pollution epidemiology. *Am. J. Epidemiol.* **186**(12), 1303–1309 (2017).
59. Godzinski, A. & Castillo, M. S. Disentangling the effects of air pollutants with many instruments. *J. Environ. Econ. Manage.* **20**, 102489 (2021).

## Acknowledgements

## Author contributions

L.Z., A.A. and T.B. conceptualize the study and wrote the main manuscript text. L.Z. and A.A. carried out the statistical analysis. L.Z. prepared the replication materials. T.B. supervised L. Z. and A.A.

## Competing interests

The authors declare no competing interests.

## Additional information

**Correspondence** and requests for materials should be addressed to L.Z.

**Reprints and permissions information** is available at www.nature.com/reprints.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.