

Assignment One: The MAUP and Multilevel Modelling

Contents

1	Demonstrating the MAUP	1
1.1	Background	1
1.2	Demonstrating the MAUP	2
2	Multilevel Modelling	7
2.1	Model specification	9
2.2	Interpretation	13
3	Bibliography	13
4	Appendix	13

1 Demonstrating the MAUP

1.1 Background

Areal units in zoning systems amalgamate into objects that constitute the basic units for the observation and analysis of spatial phenomena (Openshaw 2015). Yet, no gold standard for guiding the spatial aggregation process exists, with the validity of zonal objects subject to the arbitrary and modifiable decision-making of quantitative geographers. Problematically, the analysis of socioeconomic data involving areal units is encumbered by the modifiable areal unit problem (MAUP): “the sensitivity of analytical results to the definition of units for which data are collected.” According to the literature, the MAUP constrains the reliability of analyses for aggregated spatial data, as findings have shown varying results with the scale of aggregation and configuration of zoning systems (Avery and Clark 2015).

In practice, the MAUP is condensed into two issues of scale and zoning sensitivity which this paper will attempt to demonstrate in Section 1.2. The first issue, described as the *scale problem*, is the variation in findings when data for zonal units are progressively aggregated. This has been demonstrated empirically by Avery and Clark (2015) who found that whilst correlation coefficients did not increase monotonically with aggregation¹, a general increase in data aggregation corresponds to an increase in correlation coefficients.

The second issue, the *zoning problem*, pertains to the variation in findings when alternative combinations of zonal units are analysed with the scale or number of units held constant (Openshaw 2015). Zoning sensitivity in multivariate analysis has been demonstrated empirically in Fotheringham and Wong (2015) who simulated the aggregation of 871 block groups into 218 zones in 150 different iterations. They highlight the severity of the zoning problem by demonstrating the possibility of concluding with one iteration of zones no association between the percentage of blue-collar workers and mean family income, with another iteration finding a unit increase in blue-collar worker percentages as reducing mean family income by \$20,000. In all, ignoring scale and zoning sensitivity in model calibration can lead to inferential conclusions that a researcher’s areal data is applicable to the constituents who form the zones under study - the ecological fallacy problem (Openshaw 2015).

¹At higher levels of aggregation, there is smaller adjacency of zonal units meaning groupings are more heterogenous leading to lower correlation coefficients.

1.2 Demonstrating the MAUP

1.2.1 Data

To demonstrate the scaling and zoning sensitivities of the MAUP, we calculate the bivariate strength of association between two open data variables. For our first variable, *crime_count*, we submit a HTTP request using the POST verb to send a custom polygon for retrieving all street-level crimes occurring in 2012.

```
# download geojson
u <- "http://statistics.data.gov.uk/boundaries/E08000012.json"
# store in temporary directory
downloader::download(url = u, destfile = "/tmp/lpool.geojson")
lpool <- readOGR(dsn = "/tmp/lpool.geojson", layer = "OGRGeoJSON")
# access coords slot
lpool <- lpool@polygons[[1]]@Polygons[[1]]@coords
# build lat/lon + date string to send with postrequest
curl.string <- paste0('poly=', paste0(sprintf('%s,%s', lpool[,2], lpool[,1]), collapse = ':'))

# build dates list for loop
dates = c("2012-01", "2012-02", "2012-03", "2012-04", "2012-05", "2012-06",
          "2012-07", "2012-08", "2012-09", "2012-10", "2012-11", "2012-12")

document <- lapply(dates, function(month) {
  # format acceptable packet for http request
  curl.string <- list(poly=c(curl.string), date=c(month))
  # post custom polygon to police api
  r <- httr::POST("https://data.police.uk/api/crimes-street/all-crime",
                  body = curl.string, encode="multipart", verbose())
  json <- content(r, "text", encoding = "ISO-8859-1")
  # return as data.frame
  jsonlite::fromJSON(txt=json)
})

# cast lat/lon columns to numeric data type
document$lat <- as.numeric(document$lat)
document$lon <- as.numeric(document$lon)

# convert data.frame to shapefile in rgdal
coordinates(document) <- ~lon+lat
proj4string(document) <- CRS("+init=epsg:4326") # WGS 84

# subsetting operations
document <- document[document$category == "violent-crime",]

# write once
writeOGR(document, "/Users/samcomber/Documents/spatial_analysis/shp/crimes",
          "crimes", driver = "ESRI Shapefile")

crimes <- readOGR(dsn = "/Users/samcomber/Documents/spatial_analysis/shp/crimes",
                  layer = "crimes")
crimes <- crimes[sample(nrow(crimes), 500), ]
```

Regarding our second variable, *tweet_count*, we aggregate geo-referenced tweets containing timestamps

relating to Twitter postings within the municipality of Liverpool for 2012².

```
# read twitter shapefile
tweets <- readOGR(dsn = "/Users/samcomber/Documents/spatial_analysis/shp/tweets",
                  layer = "tweets_mersey")

# sample 500 tweets
tweets <- tweets[sample(nrow(tweets), 500), ]
```

To demonstrate the MAUP, we aggregate a count of each crime and tweet into separate variables for each region in the shapefile, before computing correlations between the two vectors of values in the dataframe. We investigate the scaling problem using an iterative process that increments the number of regions in a hex-binned lattice from 10 to 100 (stepping each iteration by 10 regions). To demonstrate the zoning problem we constrain cardinality - i.e. the number of regions - to 30, and calculate correlation coefficients using 10 different permutations of 30 regions. To build random regions, we use the multiple aggregatons feature in IMAGE studio (???) which uses a depth-first search algorithm for a given contiguity matrix to create aggregated spatial regions. We then dissolve regions with the `unionSpatialPolygons` using the generated ZoneID for aggregations.

1.2.2 Scale problem

```
# directory loop
folders <- dir("/Users/samcomber/Documents/spatial_analysis/shp/img_stud/run1/")

# we have folders + file names and are in run1 directory
setwd("/Users/samcomber/Documents/spatial_analysis/shp/img_stud/run1/")

# create master dataframe to hold correlation coefficients
cors.scale.master <- data.frame()

# loop over aggregation folder - i.e. 10 regions, 20, 30
for (folder in folders) {
  setwd(paste0("/Users/samcomber/Documents/spatial_analysis/shp/img_stud/run1/",
              ,folder))

  # reread hex shapefile everytime we move folder
  hex <- readOGR(dsn = "/Users/samcomber/Documents/spatial_analysis/shp/hex",
                layer = "hex_1_clip")

  out <- data.frame()
  files <- list.files(pattern = "*.zon")

  # pick first iteration from each scale of aggregation
  img <- read.csv(files[1])

  # join regions to hex shapefile on area id
  hex@data$id <- as.numeric(hex@data$id)
  hex@data <- left_join(hex@data, img, by = c("id"= "AreaID"))
  hex@data <- hex@data[order(hex@data$id),]

  # dissolve by area id
```

²This dataset was data mined by Guy Lansley from UCL, and processed by Dani Arribas-Bel.

```

tes <- unionSpatialPolygons(hex, hex@data$ZoneID)

# convert spatial polygons to spatialpolygonsdataframe
oldw <- getOption("warn")

# silence irrelevant warnings
options(warn = -1)
df <- data.frame(id = getSpPPolygonsIDSlots(tes))
row.names(df) <- getSpPPolygonsIDSlots(tes)
options(warn = oldw)

spdf <- SpatialPolygonsDataFrame(tes, data = df)

# points in polygon spatial join
pts.crimes <- point.in.poly(crimes, spdf)
pts.tweets <- point.in.poly(tweets, spdf)

# aggregate crimes/tweets data
pts.agg.crimes <- pts.crimes@data %>% group_by(id.1) %>% summarise(num = n())
%>% arrange(id.1)
pts.agg.tweets <- pts.tweets@data %>% group_by(id) %>% summarise(num = n())
%>% arrange(id)

# cast ids as factors for left_join
pts.agg.crimes$id.1 <- as.integer(pts.agg.crimes$id.1)
pts.agg.tweets$id <- as.integer(pts.agg.tweets$id)

# join summed crimes/tweets back to data.frame
hex@data <- left_join(hex@data, pts.agg.tweets, by = c("ZoneID" = "id"))
hex@data <- left_join(hex@data, pts.agg.crimes, by = c("ZoneID" = "id.1"))

# correlate unique zone id sums
uniques <- hex@data[!duplicated(hex@data$num.x), ]
cor.tc <- cor(uniques$num.x, uniques$num.y, use="complete.obs")

# rbind to master frame
cors.scale.master <- rbind(cors.scale.master, cor.tc)
}

```

From Table 1, we demonstrate the scaling problem. Generally, we observe rising correlations with increasing scale. Had the MAUP not existed, we would observe little systematic variation in the correlation between violent crime and tweets - we observe a difference of 0.36 between the upper and lower coefficients. Interestingly, we observe a correlation coefficient of 0.61 at 80 regions which lends evidence to Avery and Clark (2015) findings that the low adjacency of zonal units at higher aggregations leads to lower correlation.

1.2.3 Zoning problem

```

# we choose 20 regions to illustrate zoning problem
setwd("/Users/samcomber/Documents/spatial_analysis/shp/img_stud/run1/z070")

# get all 10 iterations of 20 regions
files.20 <- list.files(pattern = "*.zon")

```

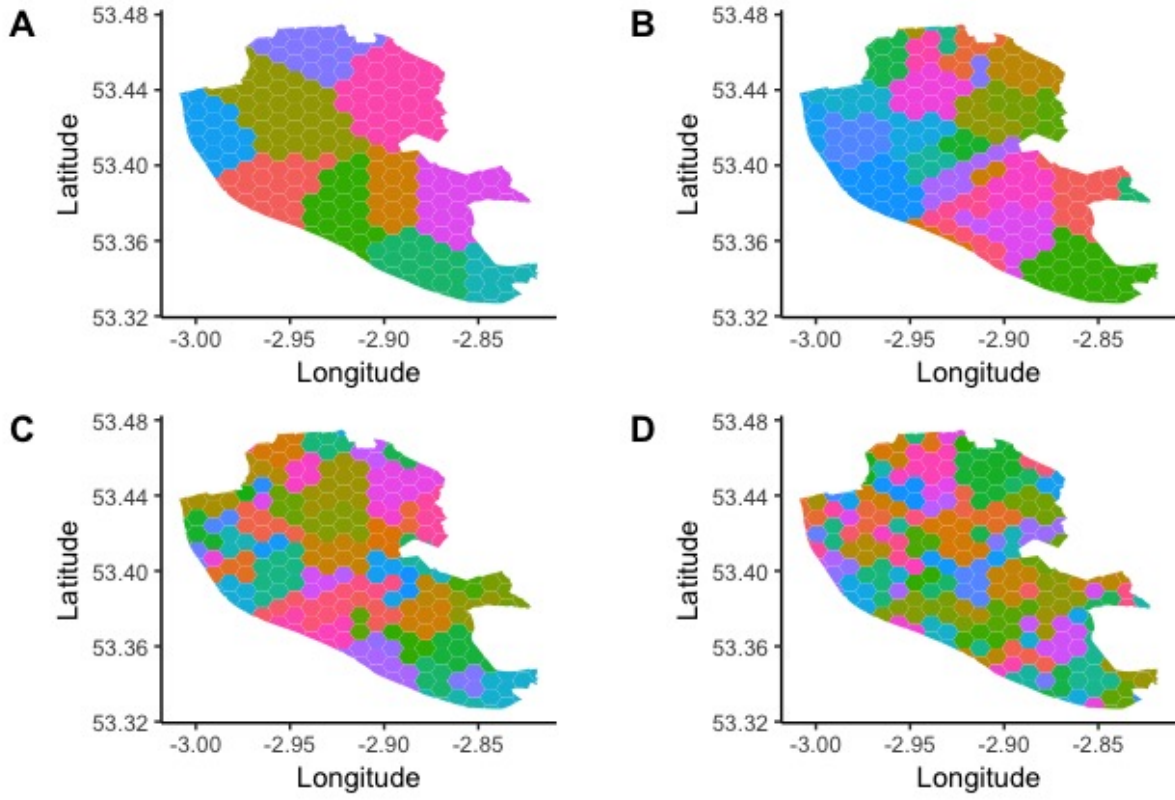


Figure 1: Scale aggregations. A shows 10 regions, B shows 30 regions, C shows 60 regions and D shows 100 regions.

Table 1: Correlations between crime and tweet counts.

<i>Number of regions</i>	<i>Correlation coeff.</i>
10	0.84
20	0.88
30	0.88
40	0.82
50	0.73
60	0.89
70	0.89
80	0.61
90	0.97
100	0.63

```

cors.zone.master
# build master data.frame
cors.zone.master <- data.frame()
for (file in files.20) {
  hex <- readOGR(dsn = "/Users/samcomber/Documents/spatial_analysis/shp/hex",
                layer = "hex_1_clip")

  # read file to df
  zon <- read.csv(file)

  # join regions to hex shapefile on area id
  hex@data$id <- as.numeric(hex@data$id)
  hex@data <- left_join(hex@data, zon, by = c("id" = "AreaID"))
  hex@data <- hex@data[order(hex@data$id),]

  # dissolve by area id
  tes <- unionSpatialPolygons(hex, hex@data$ZoneID)

  # convert spatial polygons to spatialpolygonsdataframe
  oldw <- getOption("warn")

  # silence irrelevant warnings
  options(warn = -1)
  df <- data.frame(id = getSpPPolygonsIDSlots(tes))
  row.names(df) <- getSpPPolygonsIDSlots(tes)
  options(warn = oldw)

  spdf <- SpatialPolygonsDataFrame(tes, data = df)

  # points in polygon spatial join
  pts.crimes <- point.in.poly(crimes, spdf)
  pts.tweets <- point.in.poly(tweets, spdf)

  # aggregate crimes/tweets data
  pts.agg.crimes <- pts.crimes@data %>% group_by(id.1) %>% summarise(num = n())
  %>% arrange(id.1)
  pts.agg.tweets <- pts.tweets@data %>% group_by(id) %>% summarise(num = n())
  %>% arrange(id)

  # cast ids as factors for left_join
  pts.agg.crimes$id.1 <- as.integer(pts.agg.crimes$id.1)
  pts.agg.tweets$id <- as.integer(pts.agg.tweets$id)

  # join summed crimes/tweets back to data.frame
  hex@data <- left_join(hex@data, pts.agg.tweets, by = c("ZoneID" = "id"))
  hex@data <- left_join(hex@data, pts.agg.crimes, by = c("ZoneID" = "id.1"))

  # correlate unique zone id sums
  uniques <- hex@data[!duplicated(hex@data$num.x), ]
  cor.tc <- cor(uniques$num.x, uniques$num.y, use="complete.obs")

  # rbind to master frame
  cors.zone.master <- rbind(cors.zone.master, cor.tc)

```

```
}
```

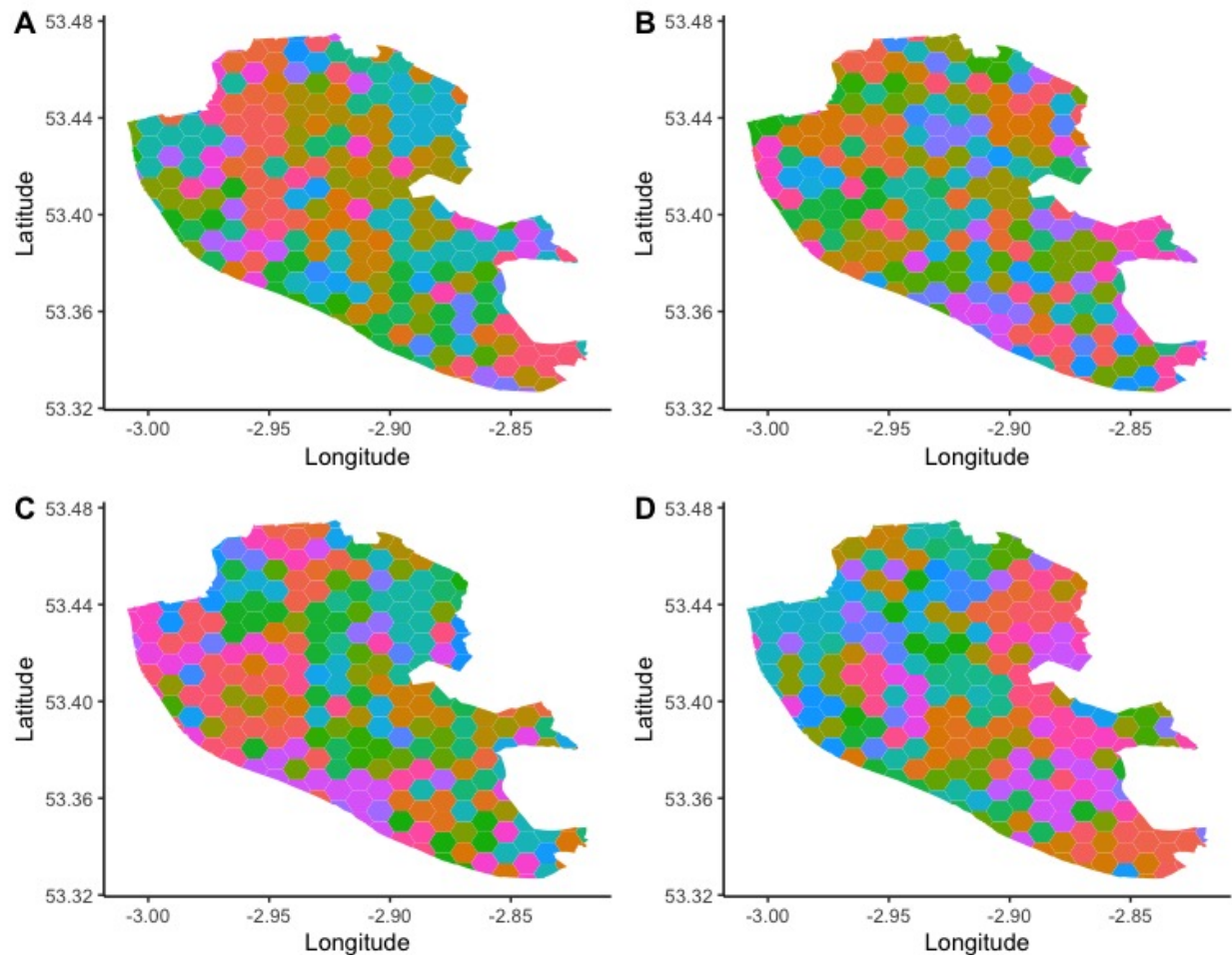


Figure 2: Zone aggregations. Four iterations displayed A-D for 70 regions.

From Table 2, we demonstrate the zoning problem: with 10 aggregation iterations for 70 regions we observe positive correlation between violent crime and tweets ranging from 0.42 (moderate association) to 0.96 (very strong association).

2 Multilevel Modelling

Having demonstrated the MAUP, we implement a two-level multilevel model that accounts for spatial heterogeneity effects between the scale geographies. Firstly, we aggregate the median price paid per property for each LSOA:

```
# save location to tmp folder
temp <- "/tmp/lpool_pp.csv"

# download from source
download.file(paste0("http://landregistry.data.gov.uk/app/ppd/ppd_data.csv
                    ?et%5B%5D=lrcommon%3Afreehold&et%5B%5D=lrcommon%3Aleasehold
```


Table 2: Correlation coefficients for 10 iterations of 70 regions.

<i>Iteration</i>	<i>Correlation coeff.</i>
1	0.89
2	0.66
3	0.61
4	0.89
5	0.42
6	0.96
7	0.75
8	0.92
9	0.82
10	0.84

```

&header=true&limit=all&max_date=31+December+2015
&min_date=1+January+2015&nb%5B%5D=true&nb%5B%5D=false
&ptype%5B%5D=lrcommon%3Adetached&ptype%5B%5D=lrcommon%3Asemi-detached
&ptype%5B%5D=lrcommon%3Aterraced&ptype%5B%5D=lrcommon%
3Aflat-maisonette&ptype%5B%5D=lrcommon%3AotherPropertyType
&tc%5B%5D=ppd%3AstandardPricePaidTransaction
&tc%5B%5D=ppd%3AadditionalPricePaidTransaction&town=",
"Liverpool"), temp)

# load LR data
lr <- read.csv("/tmp/lpool_pp.csv",
               stringsAsFactors = FALSE)

# select only relevant columns
lr <- dplyr::select(lr, price_paid, property_type, new_build, saon, paon, street,
                   locality, town, postcode)

# add concatenated column for geocoding
lr.geocode <- mutate(lr, address = paste(saon, paon, street, locality, town, postcode))

# drop previous columns
lr.geocode <- dplyr::select(lr.geocode, -saon, -paon, -street,
                           -locality, -town, -postcode)

# take sample to prevent API query limit from Google (2,500)
lr.geocode.sample <- lr.geocode[sample(nrow(lr.geocode), 2450),]

# loop over address, get lat/lon
geo.done <- data.frame()
for (i in 1:nrow(lr.geocode.sample)) {
  g <- geocode(lr.geocode.sample$address[i])
  g <- mutate(g, price=lr.geocode.sample$price_paid[i])
  geo.done <- rbind(geo.done, g)
}

# export to csv for backup
write.csv(geo.done, file = "/Users/samcomber/Documents/spatial_analysis/data/

```



```

    geocoded_lr.csv", row.names=FALSE)

# convert to spatial points df
lr <- read.csv("lpool_pp.csv", header = TRUE, sep=",")
coordinates(lr) <- ~lon + lat
proj4string(lr) <- CRS("+proj=longlat +datum=WGS84 +no_defs +ellps=WGS84 +towgs84=0,0,0")

# write once
writeOGR(lr, "/Users/samcomber/Documents/spatial_analysis/shp/spatial_join",
        "lr", driver = "ESRI Shapefile")

```

We then tidy the data, before log transforming median housing price per LSOA to relax the assumption of normality:

```

setwd("/Users/samcomber/documents/spatial_analysis/shp/spatial_join")
# LSOA and MSOA geographies joined by location in QGIS as
# R lacks library for polygon.in.polygon joins
lsoa.join <- readOGR(dsn = ".", layer = "lsoa_join", stringsAsFactors = FALSE)
# take only dataframe
lsoa.join <- lsoa.join@data
# cast as integer
lsoa.join$md_2015 <- as.numeric(lsoa.join$md_2015)
# log housing prices
lsoa.join$md_2015 <- log(lsoa.join$md_2015)
# take only complete cases
lsoa.join <- lsoa.join[complete.cases(lsoa.join),]

# proportion over 65

```

2.1 Model specification

2.1.1 Null Model

Generally, our goal of estimation is the partial-pooling estimates of the median logged housing prices among all LSOAs in MSOA j . In our null model of no predictors, we approximate the multilevel estimate for MSOA j as a weighted average of the mean of properties sold in the MSOA - an unpooled estimate \hat{y}_j - and the mean across all MSOAs - the completely pooled estimate, \hat{y}_{all} . As the weighted average uses information available about individual LSOAs, averages from MSOAs with a small number of constituent LSOAs carry less weighting, and so the multilevel estimate is pulled closer to Liverpool's overall average³ (Gelman 2017). To allow this 'soft constraint' on α_j , partial-pooling estimates are derived by assigning a probability distribution to the α_j 's which pulls estimates of α_j some way towards their mean level μ_α as below:

$$\alpha_j \sim \mathcal{N}(\mu_\alpha, \sigma_\alpha^2) \text{ for } j = 1, \dots, J,$$

where the mean μ_α and standard deviation σ_α derived from the data.

To begin decomposing how median log housing prices vary by LSOA and MSOA geographies, we specify the following null multilevel model:

$$y_{i,j} = \alpha_{0,j} + \epsilon_{i,j}$$

³In this way, the relative weights are determined by the sample size in the group and the variation within and between groups.

$$\alpha_{0,j} = \alpha_0 + \mu_{0,j}$$

where in the first level, $y_{i,j}$ is the median logged housing price for the i^{th} LSOA in MSOA j , $\alpha_{0,j}$ is the MSOA-level mean of $y_{i,j}$ (varying intercept) for the j^{th} MSOA, and $\epsilon_{i,j}$ is the residual error term assumed i.i.d; where in the second level, α_0 is the overall intercept, and finally $\mu_{0,j}$ is the random error component for the deviation of the intercept of the j^{th} group from the overall intercept.

```
# ----- NULL MODEL -----

# (1|MSOA11CD) allows intercept (coefficient of 1 is the constant
# term in the regression) to vary by MSOA
model.null <- lmer(md_2015 ~ 1 + (1|MSOA11CD), data = lsoa.join)

# get variance
vc <- VarCorr(model.null)
print(vc, comp=c("Variance"))

# calculate variance partition coefficient
print(paste0(round(0.18341 / (0.18341 + 0.12261), 3), "%"))
```

Overall, we estimate mean attainment (across MSOAs) as *. The mean for the j^{th} MSOA is estimated as $30.60 + \mu_{0,j}$, where $\mu_{0,j}$ is the MSOA-level residual⁴.

Having observed a Variance Partition Coefficient (VPC) of 0.599, we approximate ~60% of the variation in median log housing prices as explained by inequalities between MSOAs. This infers the presence of spatial heterogeneity between MSOAs, which justifies the use of a multilevel model to simultaneously capture outcomes at LSOA and MSOA geographies.

```
# caterpillar graph to visualise mean of each MSOA random effect
msoa.rand.eff <- RESim(model.null)
p <- plotRESim(msoa.rand.eff)
```

2.1.2 Random-slopes model

To rationalise a varying intercept and slope model, we assess the variability between median housing prices and the percentage of neighbourhood income deprivation.

As Figure 3 demonstrates variance between median logged housing prices and percentage of neighbourhood income deprivation, we are able to motivate the use of a random slope multilevel model⁵. As observed, given random intercept models assume parallel grouping lines, they perform poorly at fitting the data. For this reason, we allow coefficients of explanatory variables to vary for each MSOA, which can be understood as interactions between group-level indicators and an individual-level predictor (Gelman 2017). This is achieved by adding a random term to the coefficient of x_k which relaxes the constraint of β_k as fixed across each MSOA, allowing the coefficient to vary randomly across groups. We specify our random-slope model as follows:

$$y_{i,j} = \alpha_{ij} + \beta_1 x_{1,ij} + \beta_2 x_{2,ij} + \epsilon_{0,ij}$$

$$\alpha_{0j} = \alpha_0 + \mu_{0j}$$

⁴For MSOA ** DO EXAMPLE HERE

⁵-2 log likelihood test also concluded **.

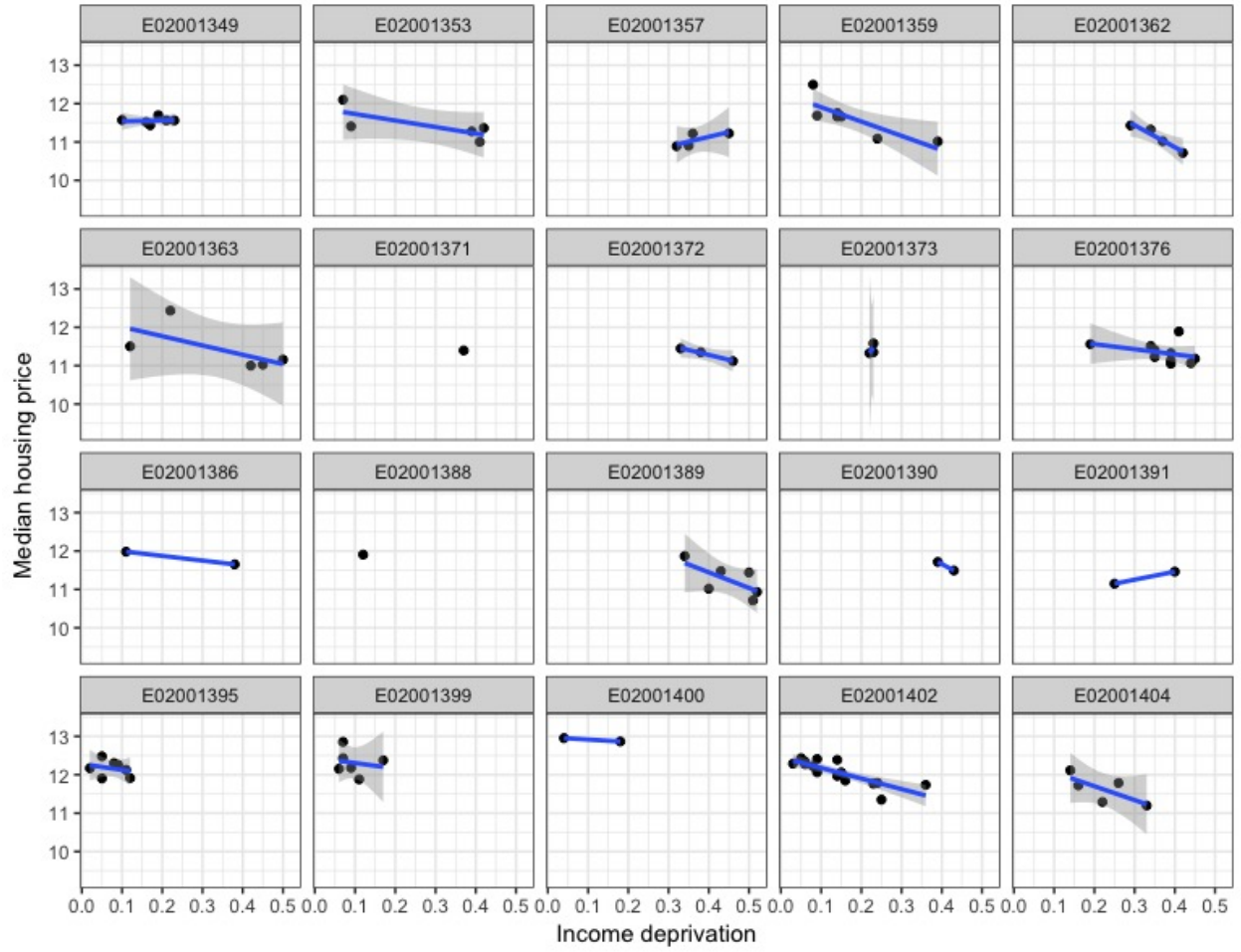


Figure 3: Variability in income deprivation by median log housing price for 20 sampled MSOAs.

$$\beta_{2j} = \beta + 1 + \mu_{1j}$$

which can be re-expressed into a single level by substituting formulae for $\alpha_{0,j}$ and $\beta_{1,j}$ into $y_{i,j}$ as:

$$y_{i,j} = \alpha_0 + \beta_1 x_{ij} + \mu_{0j} + \mu_{1j} x_{ij} + \epsilon_{ij}$$

Conceptually, this equation is similar to the null model, but the $\mu_{1j}x_{ij}$ term has been added to reflect the interaction between the j^{th} MSOA and predictor x . The random effects μ_{0j} and μ_{1j} are assumed:

$$\begin{bmatrix} \mu_{0j} \\ \mu_{1j} \end{bmatrix} \sim \mathcal{N}(0, \begin{bmatrix} \sigma_{\mu 0}^2 & \\ \sigma_{\mu 0 \mu 1} & \sigma_{\mu 1}^2 \end{bmatrix})$$

meaning the slope of regression line for MSOA j is $\beta_1 + \mu_{1j}$, which allows the coefficient to vary by group.

```
model.var.slope <- lmer(md_2015 ~ lr_crime
                        + (1 + lr_income | MSOA11CD), data=lsoa.join)

model.var.slope.fit <- lmer(md_2015 ~ lr_crime
                           + (1 | MSOA11CD), data=lsoa.join)

anova(model.var.slope, model.var.slope.fit)

## refitting model(s) with ML (instead of REML)

## Data: lsoa.join
## Models:
## model.var.slope.fit: md_2015 ~ lr_crime + (1 | MSOA11CD)
## model.var.slope: md_2015 ~ lr_crime + (1 + lr_income | MSOA11CD)
##
##      Df      AIC      BIC    logLik deviance Chisq Chi Df
## model.var.slope.fit  4 267.23 281.68 -129.616   259.23
## model.var.slope      6 203.76 225.44  -95.881   191.76 67.47    2
##
##      Pr(>Chisq)
## model.var.slope.fit
## model.var.slope      2.233e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Firstly, we use a likelihood ratio test derived from an ANOVA table to estimate whether the income deprivation effect varies across MSOAs. As the deduction of the log-likelihood values computes a likelihood ratio value of 67.47 on 2 degrees of free, we can infer the income deprivation effect varies across MSOAs, justifying the use of a random-slopes model.

We graphically project the plot of income deprivation slopes by intercept ($\hat{\mu}_{1j}$ by $\hat{\mu}_{0j}$) in Figure 4. From this, we can derive *.

```
# data object containing random slope and intercepts
inc.random <- ranef(model.var.slope, condVar = TRUE)

plot(inc.random[[1]], xlab = "Intercept (u0j)", ylab = "Slope of lr_income (u1j)")
abline(h = 0, col = "red")
abline(v = 0, col = "red")
```

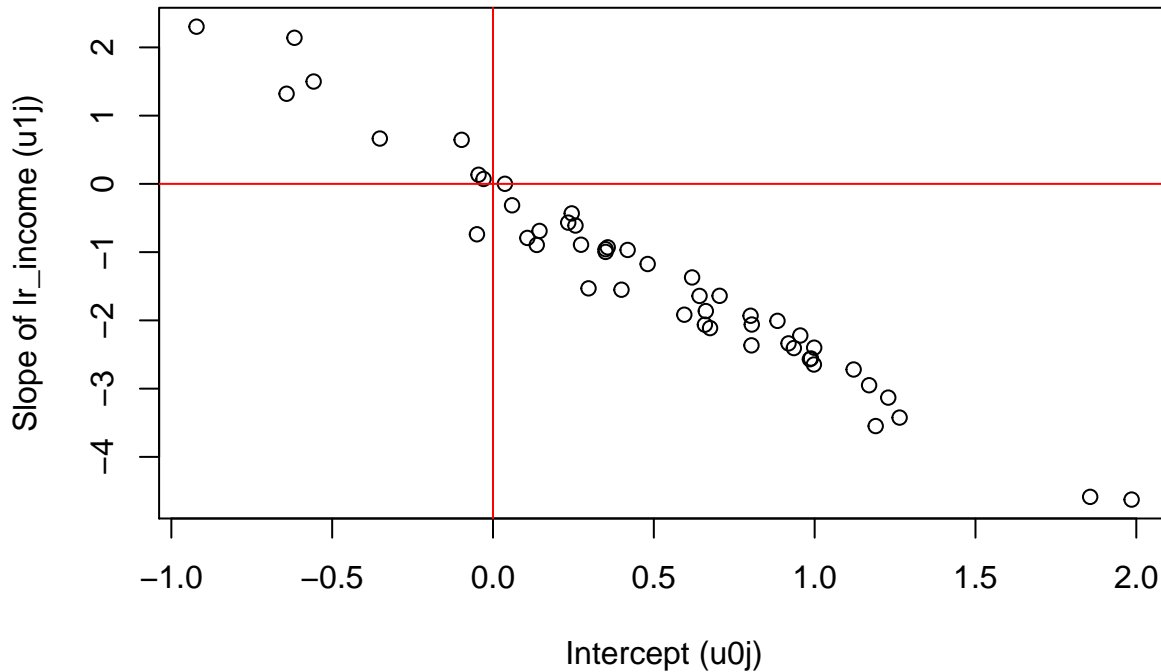


Figure 4: Income deprivation slopes versus intercept.

To produce

```
predicted <- fitted(model.var.slope)
```

```
# create caterpillar graph
msoa.rand <- REsim(model.var.slope)
plotREsim(msoa.rand)
```

2.2 Interpretation

3 Bibliography

4 Appendix

Avery and Clark. 2015. “R: A Language and Environment for Statistical Computing.” Journal Article. <http://www.R-project.org>.

Fotheringham and Wong. 2015. “R: A Language and Environment for Statistical Computing.” Journal Article. <http://www.R-project.org>.

Gelman. 2017. “R: A Language and Environment for Statistical Computing.” Journal Article. <http://www.R-project.org>.

Openshaw. 2015. “R: A Language and Environment for Statistical Computing.” Journal Article. <http://www.R-project.org>.