

#### **Regression With Correlated Errors**

Feb 9 2016

Charlotte Wickham



Projects

Groups of 3-5, I will assign and announce next week.

In weeks 8, 9 & 10 we will have some in class time for working/getting advice on projects.

On form, indicate:

- preference for type of project
- comfort/skill level
- anything else I should know when assigning groups

Proposal due Thursday Feb 25th Final project due Thursday March 10th

#### SARIMA models

Good if you are just interested in a short term forecast.

Doesn't result directly in estimates of parameters of the past, i.e. seasonal means, trends, regression parameters...

**Today:** regression models with correlated errors



Can mortality be explained by temperature and particulate matter?

## Regression with correlated errors

Use a regression model to explain the non-stationarity in mean.

Extend the usual regression model to allow the errors to be an ARMA process.

## Linear Regression Review

A linear regression model, models the response, y<sub>t</sub>, as a linear combination of p covariates,  $x_{t1}$ ,  $x_{t2}$ , ...,  $x_{tp}$ , and noise,  $\varepsilon_t$ .  $y_t = \beta_0 + \beta_1 x_{t1} + \beta_2 x_{t2} + ... + \beta_p x_{tp} + \varepsilon_t$ t = 1, ..., n you might be used to seeing i

#### $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

(in matrix notation, y is a nx1 vector of the responses, X is a nxp matrix of covariates,  $\beta$  a px1 vector of parameters,  $\varepsilon$  is a nx1 vector of errors)

# Your turn

# What are the assumptions of ordinary least squares regression?



Suggests the regression model: mortality<sub>t</sub> =  $\beta_0 + \beta_1 t + \beta_2 temp_t + \beta_3 temp_t^2 + \beta_4 part_t + \epsilon_t$ 

```
> fit_lm <- lm(mortality ~ time0 + temp_sc + temp_2 + part, data = mort)</pre>
                                     1
                       I started time at 0, centered temp about it's mean
                       and found the square of temp
             > summary(fit_lm)
             Call:
             lm(formula = mortality ~ time0 + temp_sc + temp_2 + part, data = mort)
             Residuals:
                          10 Median
                  Min
                                           30
                                                  Max
             -19.0760 -4.2153 -0.4878 3.7435 29.2448
             Coefficients:
                         Estimate Std. Error t value Pr(>|t|)
             (Intercept) 81.565394 1.101148 74.07 < 2e-16 ***
             time0
                        -1.395901
                                  0.101009 -13.82 < 2e-16 ***
                       -0.472469 0.031622 -14.94 < 2e-16 ***
             temp_sc
                     0.022588 0.002827 7.99 9.26e-15 ***
             temp_2
             part
                      0.255350
                                  0.018857 13.54 < 2e-16 ***
              ___
             Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
             Residual standard error: 6.385 on 503 degrees of freedom
             Multiple R-squared: 0.5954, Adjusted R-squared: 0.5922
             F-statistic: 185 on 4 and 503 DF, p-value: < 2.2e-16
```

### Assumptions

Everything looks good, except....

#### examine\_corr(residuals(fit\_lm))



Assumptions say this should be white noise, looks like ?

## Generalized Least Squares

- Instead of variance  $\sigma^2 I$ , allow errors to have covariance matrix  $\Sigma$ .
- If  $\Sigma$  is known multiplying by  $\Sigma^{-1/2}$ , reduces the problem to the usual regression problem.

If  $\Sigma$  is unknown, start with estimating the  $\beta$  using OLS, use the residuals to estimate  $\Sigma$  and iterate. Should converge to the MLEs (under Normal errors).

### Correlated errors model

 $y_t = \beta_0 + \beta_1 x_{t1} + \beta_2 x_{t2} + \dots + \beta_p x_{tp} + z_t$ z<sub>t</sub> is a stationary ARMA(p,q) process

Or equivalently  $z_t \sim N(0, \Sigma)$  where,  $\Sigma_{ij} = Cov(z_i, z_j) = \Upsilon(|i - j|) = \text{some function of } \sigma^2, \beta$ and  $\alpha$ 

Plus usual assumptions:

linearity, zt independent of X.

#### ARMA errors

If the errors are an ARMA process,  $\Sigma$ , can be written in terms of our parameters  $\beta$  and  $\alpha$ . (Since  $\Sigma_{ij} = Cov(z_i, z_j) = \gamma(|i - j|)$ ) If we specify p and q,  $\Sigma$  is known up to  $\beta$  and

α.

and can be estimated by GLS.

The standard errors on the regression coefficients,  $\beta$ , depend on X,  $\sigma^2$ ,  $\beta$  and  $\alpha$ .



- 1. Fit the model for the mean using OLS.
- 2. Examine the residuals to identify an appropriate ARMA(p, q) process.
- 3. Fit the GLS model.
- 4. Model diagnostics.
- 5. Interpret (forecast?).

### $\ln R$

Either (a little more general, i.e. not just for time series):
library(nlme)
gls(y ~ x1 + x2,
 correlation = corARMA(p = p, q = q),

```
method = "ML")
```

Or (faster and can handle seasonal arima models):
arima(y, order = c(p, d, q), xreg = X)

We already identified the errors in the mortality series as AR(2).

#### Or

AR(2): 
$$z_t = \alpha_1 z_{t-1} + \alpha_2 z_{t-2} + w_t$$

Diagnostics

residuals(arima\_fit) gives estimates of  $w_t$ , the white noise.

residuals(gls\_fit) gives estimates of  $z_t$ , the ARMA process.

use residuals(gls\_fit, type = "normalized") to
get wt



#### wrong

#### Ordinary linear regression Assume white noise errors

Linear regression with correlated errors Assume AR(2) errors

right

<pre>&gt; round(confint(fit_lm),</pre>			2)
	2.5 %	97.5 %	
(Intercept)	79.40	83.73	
time0	-1.59	-1.20	
temp_sc	-0.53	-0.41	
temp_2	0.02	0.03	
part	0.22	0.29	

How would you interpret the coefficient on time0?

# Your turn

Where does most of the variation in temperature come from?

What about particulates?

What about mortality?

