# Thinking about Cross-Sectional and Longitudinal Data

## Dr Vernon Gayle

Longitudinal data analysis is important because it permits insights into the processes of change. Davies (1994) states this glib claim is inadequate and certainly fails to convince many social science researchers who are concerned with substantive rather than methodological challenges. What is required is an understanding of the limitations of cross-sectional analysis.

This handout outlines four central issues regarding longitudinal data analysis. These are **age and cohort effects**, **direction of causality**, **state dependence** and **residual heterogeneity**. These issues are complicated and only a cursory explanation is given in this handout!

## Age & Cohort Effects

Table 1 shows some data on the average number of days spent off of the road for a particular make and model of car. It appears that older cars, on average, spend more days per year off the road. Buying one of these cars new looks like a bad idea because they appear to 'age' badly. There appears to be an 'ageing effect'.

Table 1 Cross-sectional Data on Car Reliability and Age

| Age of car (years) | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Average number of days off the road | 3 | 4 | 15 | 16 | 18 |

However, the manufacturer tells me that cars made more recently are better. For example the ones made in the last two years are more reliable and only spend 3 or 4 days of off the road on average. The manufacturer suggests that there is a 'cohort effect'. The cohort of cars made 3–5 years ago was less reliable. Could this be true?

Without longitudinal data we cannot unravel whether age or cohort (or a combination of each) provides the correct explanation. Below in Table 2 and Table 3 are some longitudinal data on reliability for two cars, the Ford Viva and the Vauxhall Capri. If you were buying a new car, which would you choose on reliability?

Table 2 Longitudinal Data on Car Reliability and Age (Ford Viva)

| | | Year of Manufacture | | | |
|---|---|---|---|---|---|
| Age of Car (years) | 1997 | 1998 | 1999 | 2000 | 2001 |
| 1 | 3 | 4 | 3 | 3 | 4 |
| 2 | 4 | 4 | 3 | 3 | |
| 3 | 10 | 10 | 10 | | |
| 4 | 16 | 16 | | | |
| 5 | 18 | | | | |

Table3LongitudinalDataonCarReliabilityandAge(VauxhallCapri)

| | YearofManufacture | | | | |
|---|---|---|---|---|---|
| | 1997 | 1998 | 1999 | 2000 | 2001 |
| AgeofCar (years) | | | | | |
| 1 | 18 | 18 | 15 | 3 | 4 |
| 2 | 18 | 16 | 15 | 3 | |
| 3 | 18 | 16 | 15 | | |
| 4 | 18 | 16 | | | |
| 5 | 18 | | | | |

InTable2wecanseeaclear'ageingeffect'.AstheFordVivacarsgetol derthey becomelessreliable.
InTable3wecanseeaclear'cohorteffect'.SofartheVauxhallsmadein2000and 2001appeartobemuchmorereliablethantheonesmanufacturedbetween1997and 1999.

TheVauxhallCapriappearstobeabetterpurchaseintermsofreliabilitytha nthe FordViva.

Cross-sectionaldataareuninformativeaboutageand/orcohorteffects–tountangl e theseeffectsweneedlongitudinaldata.


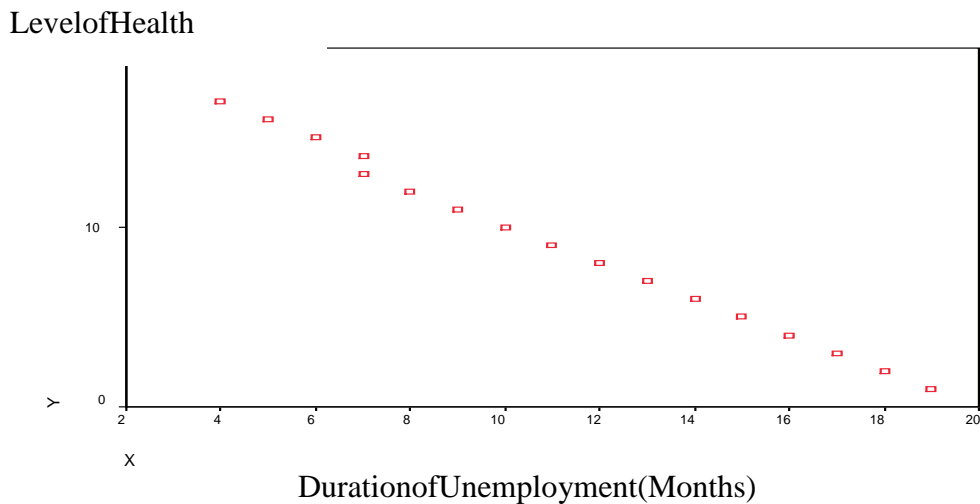**DirectionofCausality**

Thereisunequivocalevidencefromcross-sectionaldatathat,overall,theunemploy ed havepoorerhealth.Thisisconsistentwithtwohypotheses;1)unemploymentcauses illhealth2)illhealthcausesunemployment.

ImaginethatIhadasurveythatcollecteddataonpeople'semploymentstat us(e.g. whetherornottheywereunemployed)andalsoontheirlevelofhealth.Itwouldalso beplausiblethatwemightaskpeoplehowlongtheyhadbeenunemployed. Generally,wemightexpecttofindanegativerelationship.Thosethathadbeen unemployedforlongerwouldhavelowerlevelsofhealth.

Thecross-sectionaldatareportedinFigure1supportthehypothesisthat unemploymentcausesillhealth(illhealthismeasuredona20pointscale;lower scoreindicatingpoorerhealth).Thepeoplethathavebeenunemployedforlonger periodshavelowerlevelsofhealth.

LevelofHealth

DurationofUnemployment(Months)

However,thesecross-sectionaldataaarealsoconsistentwiththesecondhypothesi      s; thatillhealthcausesunemployment.Ifillhealthpreventsapersonfromworkin      g thosewithlesssevereillhealthwillrecoverandreturntowork.Withtheincrea      sing durationofunemploymentthosewithlesssevereillhealthwillbeprogressivel      y under-represented,whilethosewithmoresevereillhealthwilltheoverrepresent      ed. Thisisknownas'sampleselectionbias'.Thissampleselectionbiascouldtherefo      re explainthecross-sectionalpictureofddeclininghealthwithdurationofunemployme      nt.

Withoutlongitudinaldatawecannotsolvethispuzzle.Togettogripswiththiswe wouldneedlongitudinaldata.Tables4and5report(hypothetical)dataonthelevel ofhealth(measuredona20pointscale)andemploymentstatusfortwoindividuals. **Thisisasimplifieddepictionoftheissue.**

Table4LevelofHealthandEmploymentStatus(IndividualA)

| Month | LevelofHealth | EmploymentStatus |
|---|---|---|
| 1 | 17 | Employed |
| 2 | 17 | Employed |
| 3 | 17 | Employed |
| 4 | 17 | Unemployed |
| 5 | 17 | Unemployed |
| 6 | 10 | Unemployed |
| 7 | 6 | Unemployed |
| 8 | 5 | Unemployed |
| 9 | 4 | Unemployed |
| 10 | 3 | Unemployed |
| 11 | 2 | Unemployed |
| 12 | 1 | Unemployed |

ThesedatainTable4(IndividualA)suggestthatunemploymenthasprecededill healthandastheduorationoftheirunemploymentincreasedtheirlevelofhealth declined.Thisisconsistentwiththeideathatunemploymentleadstoillhealth.

Table 5 Level of Health and Employment Status (Individual B)

| Month | Level of Health | Employment Status |
|---|---|---|
| 1 | 17 | Employed |
| 2 | 1 | Employed |
| 3 | 1 | Employed |
| 4 | 1 | Unemployed |
| 5 | 1 | Unemployed |
| 6 | 1 | Unemployed |
| 7 | 1 | Unemployed |
| 8 | 1 | Unemployed |
| 9 | 1 | Unemployed |
| 10 | 1 | Unemployed |
| 11 | 1 | Unemployed |
| 12 | 1 | Unemployed |

These data in Table 5 (Individual B) suggest that ill health preceded unemployment. This is consistent with the idea that ill health causes unemployment. *Note*: if we had undertaken a cross-section survey in month 12 for these two individuals

* A would have been unemployed for 9 months and have a health score of 1
* B would have been unemployed for 9 months and have a health score of 1

However, we can see that with longitudinal data we can start to untangle this.

* If more individuals experienced the same relationship between health and employment as Individual A, this would lend support to the hypothesis that unemployment causes to ill health.

* If more individuals experienced the same relationship between health and employment as Individual B, this would lend support to the hypothesis that ill health causes unemployment.

**State Dependence**

This is the idea that current behaviour is influenced by past or previous behaviour. Consider the following two examples.

* The chances that you are married in any given month are highly contingent on whether or not you are married in the previous month. Most people don't change their marital status with great frequency.

* Your chances of being employed in any month are contingent on whether or not you were employed in the previous month. The labour market tends not to be sufficiently volatile to cause the switching of employment states.

Therefore we can see why we might want to take account of prior information when examining current situations. This is only possible with longitudinal information.

## ResidualHeterogeneity

Thisisacomplicatedissue.Putsimply,longitudinaldataallowsustoincrease control inouranalysisforresidualheterogeneity.Residualheterogeneityisaterm thatrefers totheomissionofexplanatoryvariablesinouranalysis.Theseomittedexplanatory variablesareeitherunmeasuredorun-measurable.

Ifyouwanttounderstandthismorefully…
*Cross-sectionaldataallowsustoundertakeanalysisbetweencases(usually individuals).Theadvantageoflongitudinaldataisthat,inadditiontobetweencases analysis,itallowsustoundertakeanalysiswithincasesbecausewehavemeasures forthesameindividualatdifferenttimepoints.*
*Inourdatasettheremightbeindividualswithsimilarcharacteristicsbuttheybehave differentlyatdifferenttimepoints.Thiswouldsuggestthatsomeoftheirbehaviour mightbeexplainedbyunmeasured(andpossiblyun-measurable)variables(e.g. motivation).Thepossibilityofsubstantialvariationduetounmeasuredandpossibly un-measurablevariablesisknownas'residualheterogeneity'.*

ReferencesandFurtherReading

Davies,R.B.(1994)'FromCross-SectionaltoLongitudinalAnalysis',inDale,A.and Davies,R.B. *AnalyzingSocial&PoliticalChange:ACasebookofMethods* ,Sage.

Plewis,I.(1985) *AnalysingChange:MeasurementandExplanationUsing LongitudinalData,* Wiley.

Ruspini,E.(2000)'LongitudinalResearchintheSocialSciences', *SocialResearch Update,* 28.<http://www.soc.surrey.ac.uk/search/search.htm>

Taris,T.(2000) *LongitudinalDataAnalysis* ,Sage.