

# **Die Imputation von Löhnen oberhalb der Beitragsbemessungsgrenze in der IAB-Beschäftigten-Stichprobe**

Hermann Gartner

Methodenreport

Institut für Arbeitsmarkt- und Berufsforschung, Nürnberg

5. November 2004

# 1 Einführung

Ein häufiges Problem, mit dem man bei Untersuchungen mit der Beschäftigten-Stichprobe konfrontiert ist, besteht darin, dass die Einkommen der Beschäftigten nur bis zur Beitragsbemessungsgrenze gemeldet werden. Die Daten sind daher ab dieser Grenze zensiert (eine Dokumentation des Datensatzes findet sich in Bender et al. (1996)). Um dadurch bedingte Verzerrungen der zu schätzenden Koeffizienten zu vermeiden, können die zensierten Werte durch Schätzungen imputiert werden. Hierzu wird der Lohn mit einem Tobit-Modell geschätzt, wobei die in der Beschäftigten-Stichprobe zur Verfügung stehenden Variablen (Qualifikation, Region, Branche usw.) als erklärende Faktoren verwendet werden können. Als Schätzung für den Lohn sollte jedoch nicht der Erwartungswert

$$E(\ln w) = x' \hat{\beta} \quad (1)$$

verwendet werden ( $\ln w$  ist der logarithmierte Lohn). Erstens würde durch so imputierte Werte eine zu große Korrelation mit den erklärenden Variablen vorgetäuscht (zur Diskussion von Problemen verschiedener Verfahren zur Ersetzung von fehlenden Werten sei auf Schafer (1997) verwiesen). Zweitens kann es Fälle geben, in denen der Erwartungswert unter der Beitragsbemessungsgrenze liegt, obwohl bekannt ist, dass er darüber liegen muss.

Um die korrekte Korrelation des Lohnes mit den Kovariablen zu erhalten, wird einen zufälligen Wert  $\ln w^{imp}$  aus einer Normalverteilung  $N(x' \hat{\beta}, \sigma^2)$  gezogen<sup>1</sup>. Dem Erwartungswert ist also ein Störterm  $\eta$  mit der Standardabweichung  $\sigma$  hinzuzufügen. Für  $\sigma$  wird der geschätzte Wert aus der Tobit-Regression zugrundegelegt.

$$\ln w^{imp} = x' \hat{\beta} + \eta \quad (2)$$

Da der wahre Wert, der geschätzt wird, über der Beitragsbemessungsgrenze liegt, wird der imputierte Lohn aus einer gestutzten Verteilung gezogen, so dass der Lohn über dieser Schwelle liegt. Es ist also ein Verfahren zur Ziehung aus einer gestutzten Verteilung erforderlich.

## 2 Ziehungen einer gestutzten Zufallsvariable

Gegeben sei eine normalverteilte Zufallsvariable  $e \sim N(\mu, \sigma^2)$ . Die Untergrenze sei  $a$ . Zur Vereinfachung der Schreibweise wird normiert: Es gelte  $\alpha = \frac{a-\mu}{\sigma}$  und  $\epsilon = \frac{e-\mu}{\sigma}$ .  $\epsilon$  ist somit standardnormalverteilt:

$$g(\epsilon) = \phi(\epsilon) \quad (3)$$

$\phi(x)$  ist die Dichtefunktion der Standardnormalverteilung. Aus dieser Verteilung sollen Werte gezogen werden, die über der Grenze  $\alpha$  liegen. Es ist also eine gestutzte Verteilung

---

<sup>1</sup>Vorausgesetzt wird, dass die Verteilung der Löhne einer Lognormal-Verteilung folgt. Empirisch ist dies eine gute Approximation.

zugrunde zu legen. Deren Dichtefunktion lautet

$$g(\epsilon|\epsilon > \alpha) = \frac{f(\epsilon)}{1 - \Phi(\alpha)}, \quad \epsilon > \alpha. \quad (4)$$

$\Phi(x)$  ist die Verteilungsfunktion einer standardnormalverteilten Zufallsvariable.

Die zugehörige Verteilungsfunktion ist eine Abbildung der Zufallsvariable  $\epsilon$  auf ein Intervall zwischen null und eins.  $G: \epsilon \rightarrow Y$  mit  $Y \in [0, 1]$  ist

$$G(\epsilon) = \int_{\alpha}^{\epsilon} \frac{\phi(z)}{1 - \Phi(\alpha)} dz. \quad (5)$$

Umformen führt zu

$$G(\epsilon) = \frac{1}{1 - \Phi(\alpha)} \left( \int_{-\infty}^{\epsilon} \phi(t) dt - \int_{-\infty}^{\alpha} \phi(t) dt \right) \quad (6)$$

oder

$$G(\epsilon) = \frac{1}{1 - \Phi(\alpha)} (\Phi(\epsilon) - \Phi(\alpha)) \quad (7)$$

Um die Zufallswerte in STATA generieren zu können, ist die Inverse dieser Funktion  $G^{-1}(Y) = \epsilon$  erforderlich. Auflösen der Gleichung  $Y = \frac{1}{1 - \Phi(\alpha)} (\Phi(\epsilon) - \Phi(\alpha))$  nach  $\Phi(\epsilon)$  führt zu:

$$Y(1 - \Phi(\alpha)) + \Phi(\alpha) = \Phi(\epsilon) \quad (8)$$

Daraus folgt

$$\Phi^{-1}(Y(1 - \Phi(\alpha)) + \Phi(\alpha)) = \epsilon. \quad (9)$$

In STATA lässt sich die Variable  $\epsilon$  somit generieren mit:

$$\epsilon = \text{invnorm}(\text{uniform}() * (1 - \text{norm}(\alpha)) + \text{norm}(\alpha)) \quad (10)$$

Wobei  $Y \in [0, 1]$  durch die Funktion `uniform()` ersetzt wurde, welche eine gleichverteilte Zufallsvariable im Intervall  $[0, 1]$  erzeugt.

### 3 Anwendung auf die Imputation des Lohnes

$e$  ist bei uns der Logarithmus des zu imputierenden Lohnes  $\ln w_i^{imp}$ ,  $\mu$  ist dessen Erwartungswert  $x_i' \hat{\beta}$ .  $a$  ist die logarithmierte Betragsbemessungsgrenze  $\ln(Beitr)$ . Statt  $\sigma$  verwenden wir den Schätzer  $\hat{\sigma}$ . Es gilt damit:

$$\epsilon_i = \frac{\ln w_i^{imp} - x_i' \hat{\beta}}{\hat{\sigma}} \quad (11)$$

bzw.

$$\ln w_i^{imp} = \epsilon_i \hat{\sigma} + x_i' \hat{\beta} \quad (12)$$

Die zensierten Werte lassen sich dann imputieren durch den Befehl:

$$\ln w_i^{imp} = \text{invnorm}(\text{uniform}() * (1 - \text{norm}(\alpha_i)) + \text{norm}(\alpha_i)) \hat{\sigma} + x_i' \hat{\beta}, \quad (13)$$

wobei gilt

$$\alpha_i = \frac{\ln(\text{Beitr}) - x_i' \hat{\beta}}{\hat{\sigma}}. \quad (14)$$

## 4 Das Programm zur Imputation

Die einzelnen Schritte des Imputationsverfahrens lassen sich mit einem STATA-Programm folgender Gestalt durchführen:

```
cens = 1 wenn Lohn zensiert ist, sonst 0
lnw = beobachteter logarithmierter Lohn
grenze = Betragsbemessungsgrenze
```

```
/* Schätzung des Tobit-Modell*/
cnreg lnw variablenliste, censored(cens)
/*Generieren der Prognose  $x_i' \hat{\beta}$  */
predict xb, xb
```

```
/*standardisierte Grenze für jede Beobachtung wird generiert*/
gen alpha=(ln(grenze)-xb)/_b[_se] if cens==1
gen lnw_i=lnw
```

```
/*Imputation, wenn Wert zensiert ist*/
replace lnw_i = invnorm(uniform()*(1-norm(alpha))+norm(alpha)) * _b[_se] +
xb if cens==1
```

## Literatur

**Bender, S., Hilzendegen, J., Rohwer, G. und Rudolph, H.** (1996). Die IAB-Beschäftigtenstichprobe. *Beiträge zur Arbeitsmarkt- und Berufsforschung*, (197).

**Schafer, J. L.** (1997). *Analysis of Incomplete Multivariate Data*. Chapman & Hall, London.