# Marginal Effects in Multiply Imputed Datasets

Daniel Klein

daniel.klein@uni-kassel.de

University of Kassel

14th German Stata Users Group meeting

GESIS Cologne

June 10, 2016

Outline       Motivation       Marginal effects in multiply imputed datasets       Summary
000000         000
                            00000
                            0000000

**1** Motivation

**2** Marginal effects in multiply imputed datasets
  - Some considerations
  - A general approach
  - The mimrgns command

**3** Summary

# Multiple Imputation
**The mi commands**

mi

- introduced in Stata11
- imputes missing values
    - mi impute
- fits models to imputed data
    - mi estimate
- some postestimation commands
    - mi test
- performs data management tasks
    - mi xeq, mi passive, . . .

# Marginal effects and adjusted predictions

**The** margins **command**

margins

- introduced in Stata11
- estimates marginal effects
    - as derivatives or (semi-)elasticities
- estimates adjusted predictions
- replaces old mfx and adjust

**Factor variable notation**

- introduced with and required by margins
- creates indicator variables, higher order terms, interactions
- replaces old xi

# Subsequent enhancements

`mi`

- chained equations
- user imputation methods
- predictions

`margins`

- `marginsplot` visualizes results
- pairwise comparisons
    - `pwcompare`
- `contrasts`

# So, what is missing?
**Setting up example data**

```
version 12.1

webuse nhanes2d , clear

svyset , clear
drop if (hlthstat > 5)

set seed 42

mi set mlong
mi register imputed vitaminc zinc
mi impute chained ///
        (pmm , knn(5)) vitaminc zinc ///
        = hlthstat i.sex i.ageg i.psu finalwgt i.strata ///
        , add(5) noisily

mi svyset psu [pweight = finalwgt] , strata(strata)

save nhanes2d_imputed5.dta , replace
```

# So, what is missing?

### Combining `mi` ...

```
. mi estimate : svy : mlogit hlthstat i.sex vitaminc i.ageg zinc
```

| | | |
|---|---|---|
| Multiple-imputation estimates | Imputations = | 5 |
| Survey: Multinomial logistic regression | Number of obs = | 10335 |
| | | |
| Number of strata = 31 | Population size = | 116997257 |
| Number of PSUs = 62 | | |
| | Average RVI = | 0.2784 |
| | Largest FMI = | 0.1845 |
| | Complete DF = | 31 |
| DF adjustment: Small sample | DF: min = | 21.03 |
| | avg = | 27.89 |
| | max = | 29.17 |
| Model F test: Equal FMI | F( 32, 28.1) = | 197.60 |
| Within VCE type: Linearized | Prob > F = | 0.0000 |

| hlthstat | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| **1** | | | | | | |
| 2.sex | -.2535739 | .0757789 | -3.35 | 0.002 | -.4085658 | -.098582 |
| vitaminc | .3783082 | .0836045 | 4.52 | 0.000 | .207102 | .5495144 |
| | | | | | | |
| **agegrp** | | | | | | |
| 2 | -.2417474 | .1129682 | -2.14 | 0.041 | -.4727367 | -.0107581 |
| 3 | -.5148524 | .116651 | -4.41 | 0.000 | -.75337 | -.2763347 |
| 4 | -1.158049 | .0944494 | -12.26 | 0.000 | -1.351177 | -.9649211 |
| 5 | -1.464487 | .1040847 | -14.07 | 0.000 | -1.677314 | -1.251659 |
| 6 | -1.236679 | .1044329 | -11.84 | 0.000 | -1.45025 | -1.023107 |
| | | | | | | |
| zinc | .0069707 | .0029981 | 2.33 | 0.028 | .0008261 | .0131153 |
| _cons | -.3710475 | .2991727 | -1.24 | 0.225 | -.9837576 | .2416626 |

(*output omitted*)

# So, what is missing?
### . . . and margins

(*output omitted*)

| 3 | (base outcome) |

(*output omitted*)

| 5 | | | | | | |
|---|---|---|---|---|---|---|
| 2.sex | -.0199766 | .1179985 | -0.17 | 0.867 | -.2613745 | .2214214 |
| vitaminc | -.9080216 | .1067354 | -8.51 | 0.000 | -1.129973 | -.6860703 |
| | | | | | | |
| agegrp | | | | | | |
| 2 | .2911128 | .3417395 | 0.85 | 0.401 | -.4076501 | .9898758 |
| 3 | 1.109931 | .2787963 | 3.98 | 0.000 | .5398748 | 1.679988 |
| 4 | 1.577982 | .2523318 | 6.25 | 0.000 | 1.062038 | 2.093925 |
| 5 | 2.092809 | .2107573 | 9.93 | 0.000 | 1.66184 | 2.523779 |
| 6 | 2.483614 | .2154178 | 11.53 | 0.000 | 2.043079 | 2.924149 |
| | | | | | | |
| zinc | -.0107241 | .0038114 | -2.81 | 0.010 | -.0186396 | -.0028087 |
| _cons | -1.316412 | .3994441 | -3.30 | 0.003 | -2.139385 | -.493439 |

```
.
. margins , vce(unconditional) dydx(*) predict(outcome(1))
last estimates not found
r(301);
```

# General considerations
### How to?

Technically,

- run both, the estimation command and margins on each imputed dataset
- post estimates to e(b) and e(V)
- combine results according to Rubin's rules
    - better yet: let mi estimate do the work

# Rubin's rules

### A brief review

**Point estimate**

$$\bar{\mathbf{q}}_{mi} = \frac{1}{M} \sum_{i=1}^{M} \hat{\mathbf{q}}_i$$

**Variance-covariance matrix**

$$\mathbf{VCE}_{mi} = \bar{\mathbf{W}} + \left(1 + \frac{1}{M}\right) \mathbf{B}$$

with

$$\bar{\mathbf{W}} = \frac{1}{M} \sum_{i=1}^{M} \hat{\mathbf{W}}_i$$

$$\mathbf{B} = \left(\frac{1}{M-1}\right) \sum_{i=1}^{M} \left(\hat{\mathbf{q}}_i - \bar{\mathbf{q}}_{mi}\right) \left(\hat{\mathbf{q}}_i - \bar{\mathbf{q}}_{mi}\right)'$$

# Requirements
**Applicability of Rubin's rules**

Statistically, q ...

- estimates population parameter
- does not depend on sample size
- is (asymptotic) normal

Satisfied by, e.g.,

- regression coefficients
- linear predictor
- average marginal effects (in large samples)

# The proposed solution

**Isabel Cañette and Yulia Marchenko**

```
program emargins , eclass properties(mi)
    version 11
    args outcome
    svy : mlogit hlthstat i.sex vitaminc i.ageg zinc
    margins , vce(unconditional) dydx(*) ///
        predict(outcome(`outcome´)) post
end
```

# Let's try

```
. forvalues j = 1/5 {
  2.         mi estimate : emargins `j´
  3. }
```

Multiple-imputation estimates                      Imputations      =         5
Average marginal effects                           Number of obs    =     10335

Number of strata  =          31
Number of PSUs    =          62

                                                   Average RVI      =    0.0209
                                                   Largest FMI      =    0.0906
                                                   Complete DF      =        31
DF adjustment:   Small sample                      DF:    min       =     25.74
                                                          avg       =     28.64
Within VCE type:   Linearized                             max       =     29.17

|            | Coef.     | Std. Err. | t      | P>\|t\| | [95% Conf. Interval] |           |
|-----------:|-----------|-----------|--------|---------|----------------------|-----------|
|      2.sex | -.0506149 | .011062   | -4.58  | 0.000   | -.073242             | -.0279878 |
|   vitaminc | .0746023  | .0103436  | 7.21   | 0.000   | .0534397             | .0957648  |
|            |           |           |        |         |                      |           |
|     agegrp |           |           |        |         |                      |           |
|          2 | -.0283339 | .0185412  | -1.53  | 0.137   | -.0662458            | .0095781  |
|          3 | -.0703962 | .0191397  | -3.68  | 0.001   | -.1095317            | -.0312607 |
|          4 | -.1800159 | .0154733  | -11.63 | 0.000   | -.2116561            | -.1483758 |
|          5 | -.2412945 | .0155433  | -15.52 | 0.000   | -.2730762            | -.2095128 |
|          6 | -.2335327 | .0144687  | -16.14 | 0.000   | -.2631215            | -.203944  |
|            |           |           |        |         |                      |           |
|       zinc | .0010539  | .0004302  | 2.45   | 0.021   | .0001692             | .0019386  |

(*output omitted*)

# Let's try

(*output omitted*)

```
Multiple-imputation estimates                    Imputations      =          5
Average marginal effects                         Number of obs    =      10335

Number of strata  =        31
Number of PSUs    =        62
                                                 Average RVI      =     0.0683
                                                 Largest FMI      =     0.2394
                                                 Complete DF      =         31
DF adjustment:   Small sample                    DF:       min    =      18.37
                                                           avg    =      26.98
Within VCE type:  Linearized                               max    =      29.17
```

|          | Coef.     | Std. Err. | t     | P>\|t\| | [95% Conf. Interval]    |
|---------:|-----------|-----------|-------|-------|-------------------------|
| 2.sex    | -.0006929 | .0044593  | -0.16 | 0.878 | -.0098181    .0084322   |
| vitaminc | -.0396292 | .0052561  | -7.54 | 0.000 | -.0504941   -.0287643   |
| agegrp   |           |           |       |       |                         |
| 2        | .005949   | .0042071  | 1.41  | 0.168 | -.0026535    .0145515   |
| 3        | .0314414  | .0075778  | 4.15  | 0.000 | .015947    .0469358     |
| 4        | .0666996  | .0083927  | 7.95  | 0.000 | .0495386    .0838605    |
| 5        | .103145   | .0093551  | 11.03 | 0.000 | .0840132    .1222768    |
| 6        | .1286984  | .0152025  | 8.47  | 0.000 | .0976045    .1597924    |
| zinc     | -.0005411 | .0001548  | -3.50 | 0.003 | -.0008658   -.0002164   |

# Adapting the general approach
### A simple logit model

```
mi xeq : generate byte goodhlth = (hlthstat < 3)

program emargins2 , eclass properties(mi)
    version 11
    svy : logit goodhlth i.sex vitaminc i.ageg zinc
    margins , vce(unconditional) dydx(*) post
end

. mi estimate : emargins2
varlist specification required
r(198);
```

# Adapting the general approach

**Pairwise comparisons**

```
capture program drop emargins2
program emargins2 , eclass properties(mi)
    version 11
    svy : logit goodhlth i.sex vitaminc i.ageg zinc
    margins i.ageg , vce(unconditional) ///
        post pwcompare
end

. mi estimate : emargins2 whatever
```

```
Multiple-imputation estimates                    Imputations      =         5
Pairwise comparisons of predictive margins

Number of strata  =         31
Number of PSUs    =         62
                                                 Average RVI      =    0.0017
                                                 Largest FMI      =    0.0078
                                                 Complete DF      =        31
DF adjustment:    Small sample                   DF:     min      =     29.05
                                                         avg      =     29.14
Within VCE type:  Linearized                             max      =     29.16
```

| | Coef. | Std. Err. | t | P>|t| | [95% Conf. Interval] |
|---|---|---|---|---|---|
| agegrp | | | | | | |
| 1 | .7059828 | .0106223 | 66.46 | 0.000 | .684263 | .7277026 |
| 2 | .6568665 | .0195888 | 33.53 | 0.000 | .6168121 | .6969209 |
| 3 | .5459818 | .0152304 | 35.85 | 0.000 | .5148394 | .5771241 |
| 4 | .4056276 | .014571 | 27.84 | 0.000 | .375833 | .4354221 |
| 5 | .3318254 | .01227 | 27.04 | 0.000 | .3067364 | .3569145 |
| 6 | .3026305 | .0155862 | 19.42 | 0.000 | .2707557 | .3345053 |

# Generalizing the proposed solution

mimrgns

**Goals**

- syntax similar to margins
- output similar to margins
- no hard coded estimation command

# Replicating the original example

```
. mi estimate : svy : mlogit hlthstat i.sex vitaminc i.ageg zinc
  (output omitted )

. mimrgns , vce(unconditional) dydx(*) predict(outcome(1))
```

Multiple-imputation estimates          Imputations      =          5
Average marginal effects               Number of obs    =      10335

Number of strata  =         31
Number of PSUs    =         62

                                       Average RVI      =     0.0209
                                       Largest FMI      =     0.0906
                                       Complete DF      =         31
DF adjustment:    Small sample         DF:    min       =      25.74
                                              avg       =      28.64
Within VCE type:    Linearized                max       =      29.17

Expression   : Pr(hlthstat==1), predict(outcome(1))
dy/dx w.r.t. : 2.sex vitaminc 2.agegrp 3.agegrp 4.agegrp 5.agegrp 6.agegrp zinc

|          | dy/dx     | Std. Err. | t      | P>\|t\| | [95% Conf. Interval]  |
|----------|-----------|-----------|--------|-------|-----------------------|
| 2.sex    | -.0506149 | .011062   | -4.58  | 0.000 | -.073242    -.0279878 |
| vitaminc | .0746023  | .0103436  | 7.21   | 0.000 | .0534397    .0957648  |
|          |           |           |        |       |                       |
| agegrp   |           |           |        |       |                       |
| 2        | -.0283339 | .0185412  | -1.53  | 0.137 | -.0662458   .0095781  |
| 3        | -.0703962 | .0191397  | -3.68  | 0.001 | -.1095317   -.0312607 |
| 4        | -.1800159 | .0154733  | -11.63 | 0.000 | -.2116561   -.1483758 |
| 5        | -.2412945 | .0155433  | -15.52 | 0.000 | -.2730762   -.2095128 |
| 6        | -.2335327 | .0144687  | -16.14 | 0.000 | -.2631215   -.203944  |
|          |           |           |        |       |                       |
| zinc     | .0010539  | .0004302  | 2.45   | 0.021 | .0001692    .0019386  |

Note: dy/dx for factor levels is the discrete change from the base level.

# The logit model

```
. mi estimate : svy : logit goodhlth i.sex vitaminc i.ageg zinc

Multiple-imputation estimates              Imputations       =          5
Survey: Logistic regression                Number of obs     =      10335

Number of strata =         31              Population size    =  116997257
Number of PSUs   =         62
                                           Average RVI       =     0.0340
                                           Largest FMI       =     0.1144
                                           Complete DF        =         31
DF adjustment:   Small sample              DF:       min      =      24.57
                                                     avg       =      27.90
                                                     max       =      29.16
Model F test:       Equal FMI              F(   8,   29.0)    =      93.99
Within VCE type:    Linearized             Prob > F          =     0.0000
```

| goodhlth | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 2.sex | -.21535 | .061994 | -3.47 | 0.002 | -.3421322 | -.0885678 |
| vitaminc | .5040865 | .0639858 | 7.88 | 0.000 | .3727249 | .6354481 |
| | | | | | | |
| agegrp | | | | | | |
| 2 | -.2307692 | .0919848 | -2.51 | 0.018 | -.4188591 | -.0426792 |
| 3 | -.7052038 | .0684112 | -10.31 | 0.000 | -.8450875 | -.5653201 |
| 4 | -1.284035 | .0787597 | -16.30 | 0.000 | -1.445081 | -1.122989 |
| 5 | -1.608576 | .073101 | -22.00 | 0.000 | -1.758052 | -1.459101 |
| 6 | -1.746134 | .0925375 | -18.87 | 0.000 | -1.935379 | -1.556888 |
| | | | | | | |
| zinc | .0082014 | .0020081 | 4.08 | 0.000 | .0040671 | .0123356 |
| _cons | -.2285908 | .1915549 | -1.19 | 0.244 | -.6234565 | .166275 |

# Pairwise comparisons

```
. mimrgns i.ageg , predict(pr) vce(unconditional) pwcompare

Multiple-imputation estimates                     Imputations      =        5
Pairwise comparisons of predictive margins        Number of obs    =    10335

Number of strata =        31
Number of PSUs   =        62
                                                  Average RVI      =   0.0017
                                                  Largest FMI      =   0.0078
                                                  Complete DF      =       31
DF adjustment:   Small sample                     DF:    min       =    29.05
                                                         avg       =    29.14
Within VCE type:  Linearized                             max       =    29.16
Expression  : Pr(goodhlth), predict(pr)
```
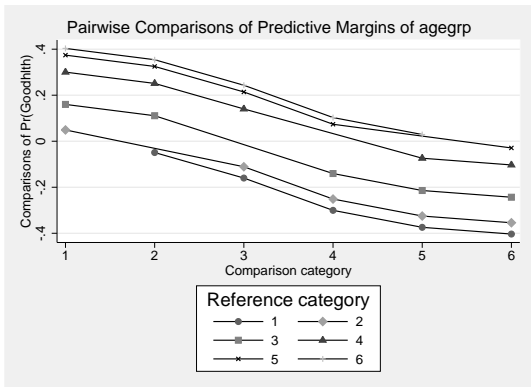
|              | Contrast  | Std. Err. | [95% Conf. Interval] |          |
|-------------:|----------:|----------:|---------------------:|---------:|
| agegrp       |           |           |                      |          |
| 2 vs 1       | -.0491163 | .0201039  | -.0902249            | -.0080076 |
| 3 vs 1       | -.1600011 | .0156144  | -.191928             | -.1280741 |
| 4 vs 1       | -.3003552 | .0176029  | -.3363501            | -.2643604 |
| 5 vs 1       | -.3741574 | .0152539  | -.4053472            | -.3429676 |
| 6 vs 1       | -.4033523 | .0189403  | -.4420859            | -.3646187 |
| 3 vs 2       | -.1108848 | .0276638  | -.1674514            | -.0543182 |
| (*output omitted*) |     |           |                      |          |
| 6 vs 4       | -.1029971 | .019452   | -.1427737            | -.0632204 |
| 6 vs 5       | -.0291949 | .0156898  | -.0612829            | .002893  |

# Plotting results

. mimrgns i.ageg , predict(pr) vce(unconditional) pwcompare cmdmargins
  (*output omitted*)

. marginsplot , noci
  Variables that uniquely identify margins: _pw1 _pw0

       _pw enumerates all pairwise comparisons; _pw0 enumerates the reference
           categories; _pw1 enumerates the comparison categories.



Pairwise Comparisons of Predictive Margins of agegrp

# What is the catch?

**Unresolved issues**

Statistical

- non-linear (adjusted) predictions
    - Rubin's rules applicable?
- higher order terms and interactions
    - passive imputation vs. JAV
- . . .

# What is the catch?

**Unresolved issues**

**Technical**

- non-linear (adjusted) predictions
  - transformations not easily implemented
- higher order terms and interactions
  - margins relies on factor variables
  - JAV approach not feasible
- execution time
  - need to run estimation command and margins $M$ times
  - cf. Miles (2015)
- . . .

# Summary

- some of `margins`' results can be combined using Rubin's rules
  - e.g., linear predictions, average marginal effects
- `mimrgns` makes this easy
- there are still unresolved statistical and technical issues
- produced results are <u>not</u> guaranteed to be valid

## References

- Cañette, I and Marchenko, Y. (2010). Re: st: Average marginal effects for a multiply imputed complex survey. Statalist.
- Miles, A. (2015). Obtaining Predictions from Models Fit to Multiply-Imputed Data. Sociological Methods and Research 45(1):175-185.