

Multiple Imputation

Multiple Imputation

Multiple imputation generates many complete data sets (e.g., $M > 20$), each with different imputations

X	Y	Z									
4	4	3	4	4	3	4	4	3	4	4	3
3	NA	5	3	3.3	5	3	4.7	5	3	2.6	5
7	1	6	7	1	6	7	1	6	7	1	6
NA	1	6	2.4	1	6	1.3	1	6	2.1	1	6
5	9	3	5	9	3	3	2.1	1.9	3	6.5	3.5
3	NA	NA	1	6	7	1	6	7	1	6	7
1	6	7	9	4	9	9	4	9	9	4	9
9	4	9	2	5.3	6	2	4.2	6	2	4.6	6
2	NA	6									

Analyzing Multiply Imputed Data

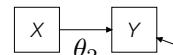
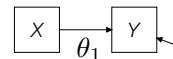
The model of interest is fit to each data set

X	Y	Z									
4	4	3	4	4	3	4	4	3	4	4	3
3	3.3	5	3	4.7	5	3	2.6	5	3	2.1	1.9
7	1	6	7	1	6	7	1	6	7	1	6
2.4	1	6	1.3	1	6	2.1	1	6	2.1	1	6
5	9	3	5	9	3	5	9	3	5	9	3
3	2.1	1.9	3	6.5	3.5	3	3.9	3.0	3	3.9	3.0
1	6	7	1	6	7	1	6	7	1	6	7
9	4	9	9	4	9	9	4	9	9	4	9
2	5.3	6	2	4.2	6	2	4.6	6	2	4.6	6

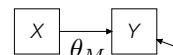
$X \xrightarrow{\theta_1} Y$ $X \xrightarrow{\theta_2} Y$... $X \xrightarrow{\theta_M} Y$

Combining Analysis Results

The pooling phase averages estimates and standard errors into a single set of analysis results



$$\hat{\theta} = (\theta_1 + \theta_2 + \dots + \theta_M) / M$$



Benefits of Multiple Imputation

Accurate with MCAR and MAR mechanisms

Maximizes power by using all available data

Analyzing several data sets gives a method for incorporating missing data error into standard errors and significance tests

Readily accommodates different levels of measurement

Imputation Phase of Multiple Imputation

Bivariate Example

The substantive analysis is a simple regression model, where the covariate is incomplete

$$Y = \beta_0 + \beta_1 (X) + e$$

For example, efficacy to quit predicted by the number of cigarettes smoked, which is incomplete

$$\text{Efficacy} = \beta_0 + \beta_1 (\text{Cigs}) + e$$

Bivariate Example

The imputation model is a reverse regression

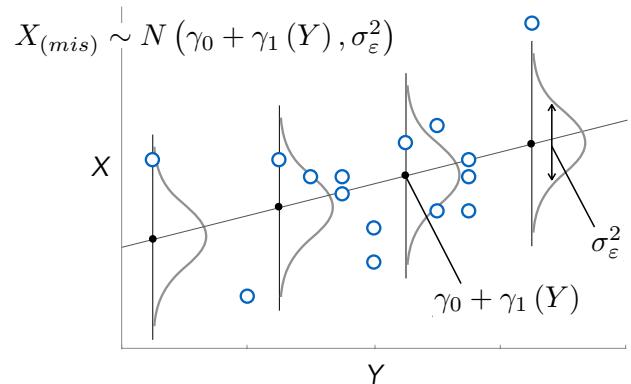
$$X_{(mis)} = \gamma_0 + \gamma_1 (Y) + \varepsilon$$

$$\phi = (\gamma_0, \gamma_1, \sigma_\varepsilon^2)$$

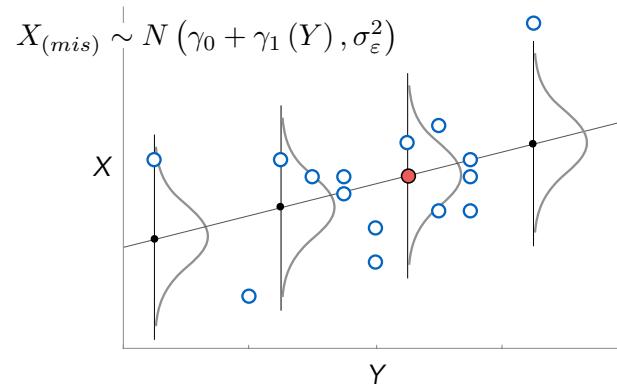
The imputation model parameters define a distribution of plausible replacement values

$$X_{(mis)} \sim N(\gamma_0 + \gamma_1 (Y), \sigma_\varepsilon^2)$$

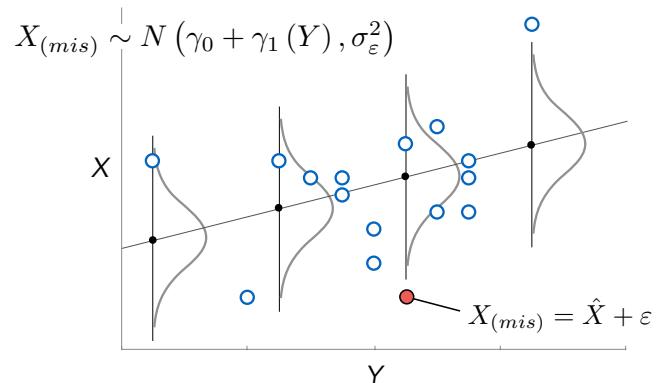
Distribution Of Missing Values



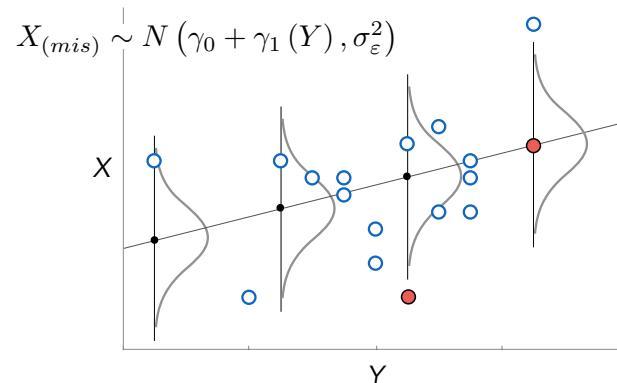
Sampling Imputations



Imputation = Predicted Score + Noise



Sampling Imputations



Imputation = Predicted Score + Noise

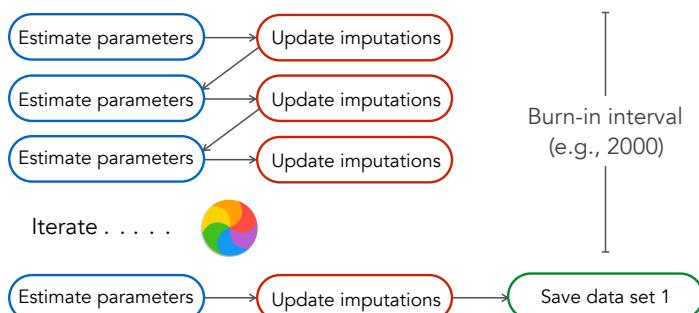
Markov Chain Monte Carlo

Each iteration t of the MCMC algorithm generates estimates of the regression model parameters, after which it samples imputations from a distribution based on those estimates

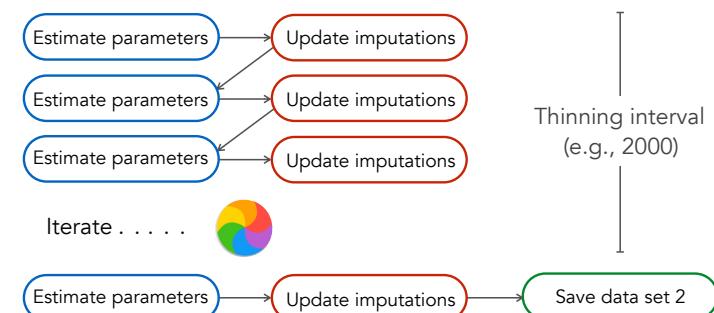
$\phi^{(t)} \sim f(\gamma_0, \gamma_1, \sigma_\varepsilon^2 | Y, X^{(t-1)})$ ————— Estimation

$X_{(mis)}^{(t)} \sim f(X_{(mis)}|Y, \phi^{(t)})$ ————— Imputation

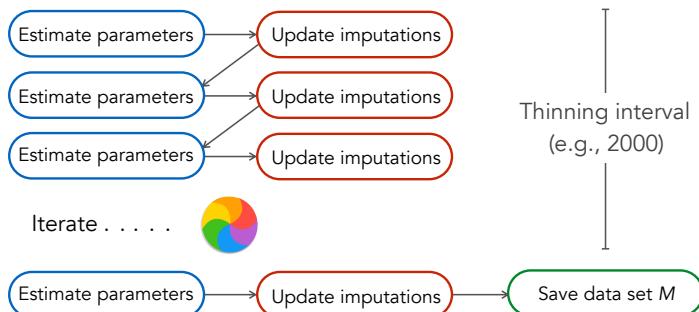
MCMC Burn-In Period



MCMC Thinning Period



Keep Thinning Until Finished

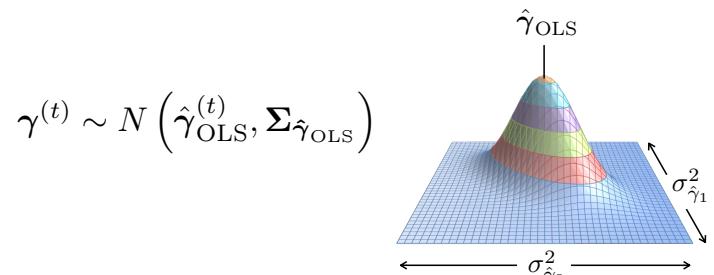


Recipe For One MCMC Iteration

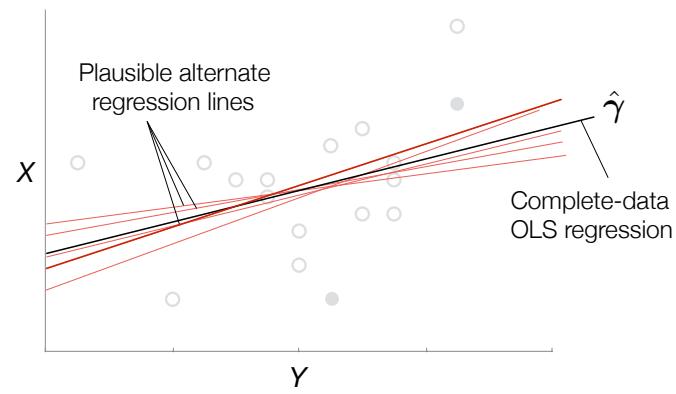
- Step 1: Sample regression coefficients from a distribution that conditions on the imputed data and the current residual variance estimate
- Step 2: Sample a residual variance from a distribution that conditions on the imputed data and the current regression coefficients
- Step 3: Sample imputations from a distribution that conditions on the current regression coefficients and residual variance

Estimating Regression Coefficients

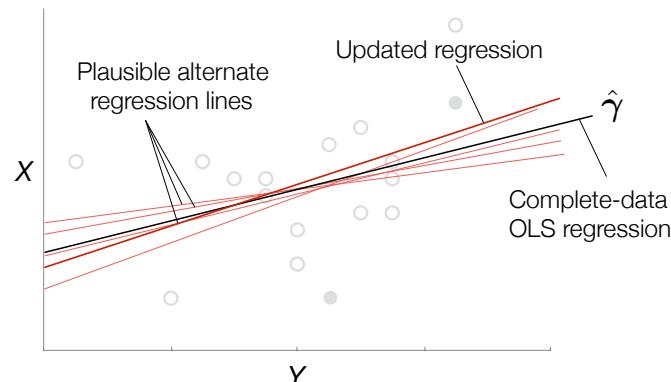
Monte Carlo simulation “samples” coefficients from a normal distribution with center and spread equal to OLS estimates and covariance matrix



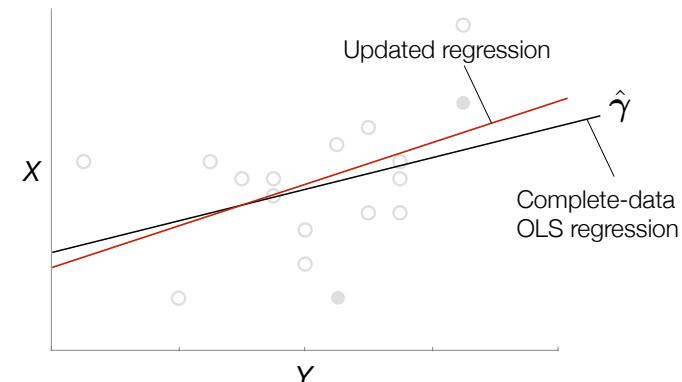
Distribution Of Plausible Regressions



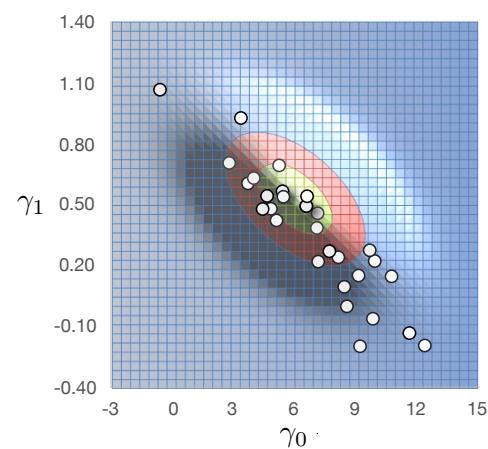
Updated Regression Line



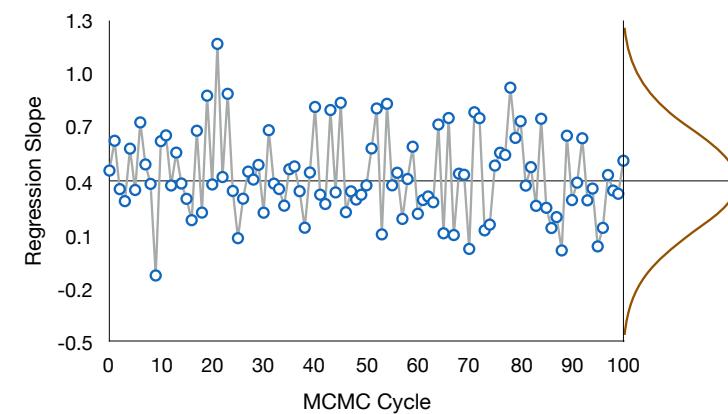
Updated Regression Line



Regression Coefficients from 30 MCMC Iterations



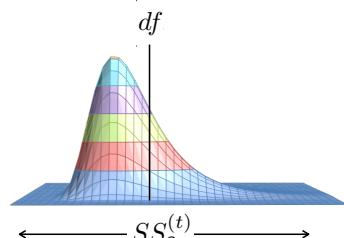
Trace Plot of Slopes from 100 MCMC Iterations



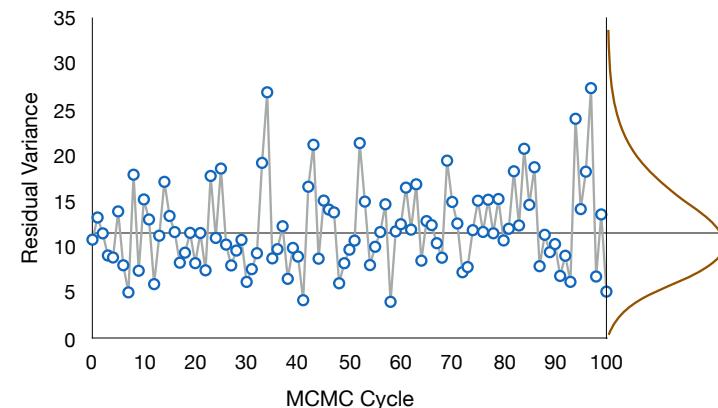
Estimating The Residual Variance

Monte Carlo simulation “samples” coefficients from an inverse gamma distribution with center and spread determined by the degrees of freedom and residual sum of squares

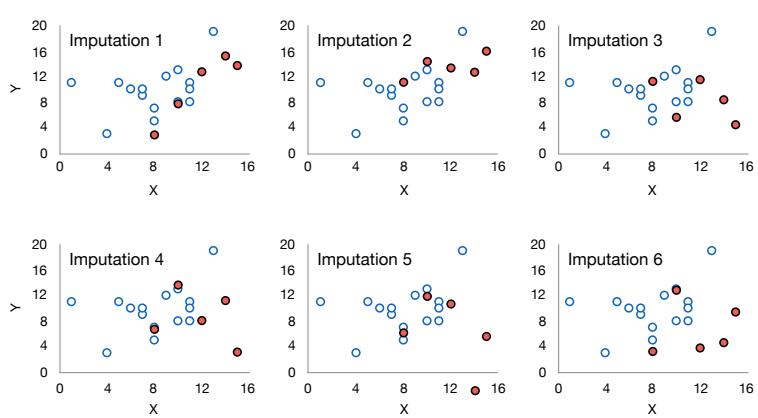
$$\sigma_{\varepsilon}^{2(t)} \sim IG(df, SS_{\varepsilon}^{(t)})$$



Trace Plot of Residual Variances from 100 Iterations



Imputed Data Sets



Multivariate Missing Data

The substantive analysis is a multiple regression model, where any variable could be incomplete

$$Y = \beta_0 + \beta_1 (X_1) + \beta_2 (X_2) + e$$

For example, efficacy to quit predicted by years smoking and number of cigarettes smoked

$$Efficacy = \beta_0 + \beta_1 (Years) + \beta_2 (Cigs) + e$$

Fully Conditional Specification (FCS)

Fully conditional specification (FCS) imputes variables one at a time in a sequence

FCS uses a series of univariate imputation models, each of which consists of an estimation step that updates parameters and an imputation step

An imputed variable from one model functions as a predictor in all other models

Y Imputation Model

The Y imputation model conditions on the current values of X_1 and X_2

$$Y_{(mis)}^{(t)} = \gamma_0(Y) + \gamma_1(Y) (X_1^{(t-1)}) + \gamma_2(Y) (X_2^{(t-1)}) + \varepsilon_{(Y)}$$

Imputations are drawn from a normal distribution

$$Y_{(mis)}^{(t)} \sim N(\gamma_0(Y) + \gamma_1(Y) (X_1^{(t-1)}) + \gamma_2(Y) (X_2^{(t-1)}), \sigma_{\varepsilon(Y)}^2)$$

X_1 Imputation Model

The X_1 imputation model conditions on the current values of Y and X_2

$$X_{1(mis)}^{(t)} = \gamma_0(X_1) + \gamma_1(X_1) (Y^{(t)}) + \gamma_2(X_1) (X_2^{(t-1)}) + \varepsilon_{(X_1)}$$

Imputations are drawn from a normal distribution

$$X_{1(mis)}^{(t)} \sim N(\gamma_0(X_1) + \gamma_1(X_1) (Y^{(t)}) + \gamma_2(X_1) (X_2^{(t-1)}), \sigma_{\varepsilon(X_1)}^2)$$

X_2 Imputation Model

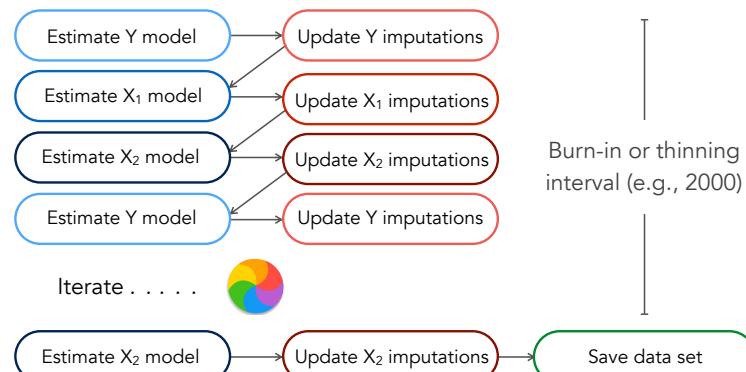
The X_2 imputation model conditions on the current values of Y and X_1

$$X_{2(mis)}^{(t)} = \gamma_0(X_2) + \gamma_1(X_2) (Y^{(t)}) + \gamma_2(X_2) (X_1^{(t)}) + \varepsilon_{(X_2)}$$

Imputations are drawn from a normal distribution

$$X_{2(mis)}^{(t)} \sim N(\gamma_0(X_2) + \gamma_1(X_2) (Y^{(t)}) + \gamma_2(X_2) (X_1^{(t)}), \sigma_{\varepsilon(X_2)}^2)$$

MCMC Steps



MCMC Convergence

MCMC gives parameters that constantly change across iterations, estimates never converge on a single value

MCMC converges when the parameter values achieve a stable distribution (the mean and variance of the distribution no longer changes across iterations)

Determining the number of cycles needed to achieve convergence is an important diagnostic step

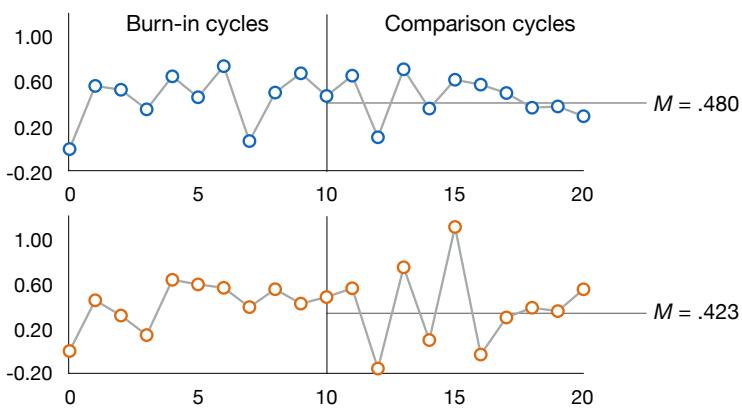
Evaluating Convergence

Run MCMC twice and determine how many iterations are needed for the two runs to give similar distributions

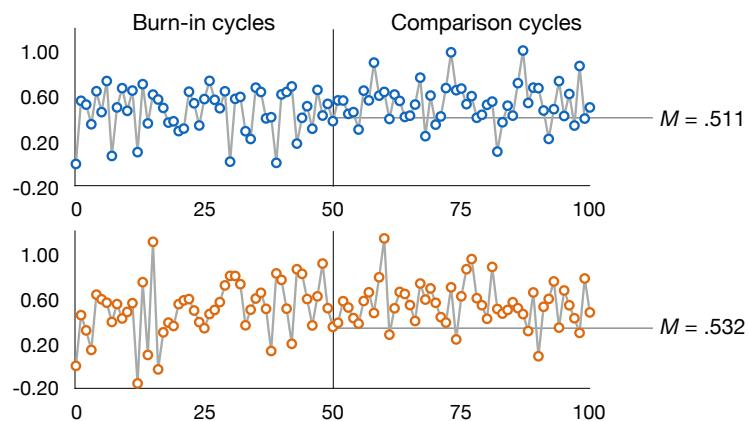
Similarity is assessed by looking at mean differences between the two runs

The potential scale reduction (PSR) factor is a common diagnostic based on the mean parameter values from two or more separate MCMC chains

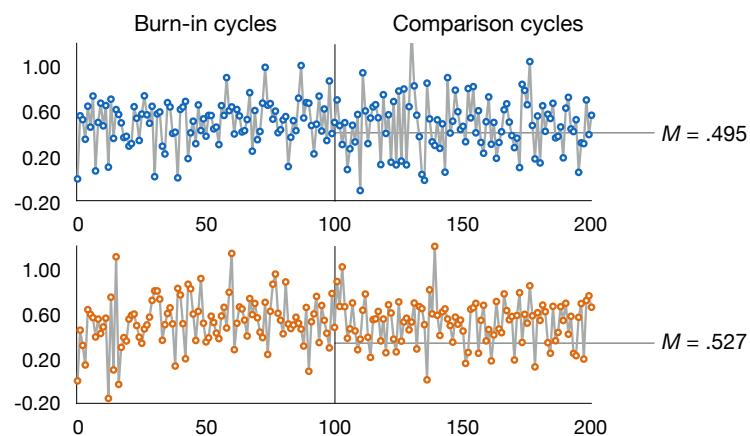
Two Chains after 20 Iterations



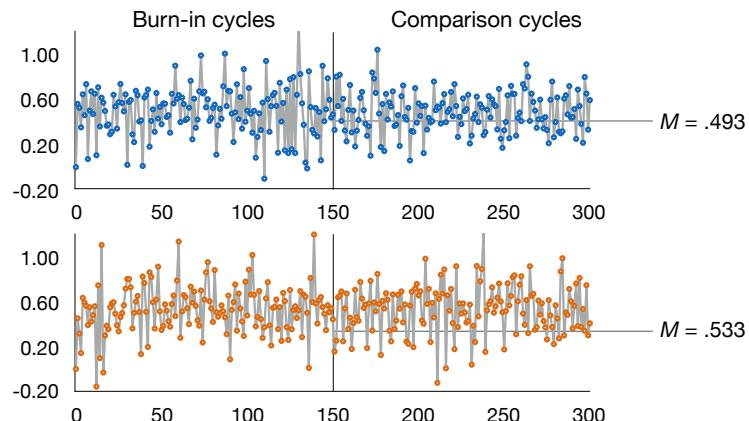
Two Chains after 100 Iterations



Two Chains after 200 Iterations



Two Chains after 300 Iterations



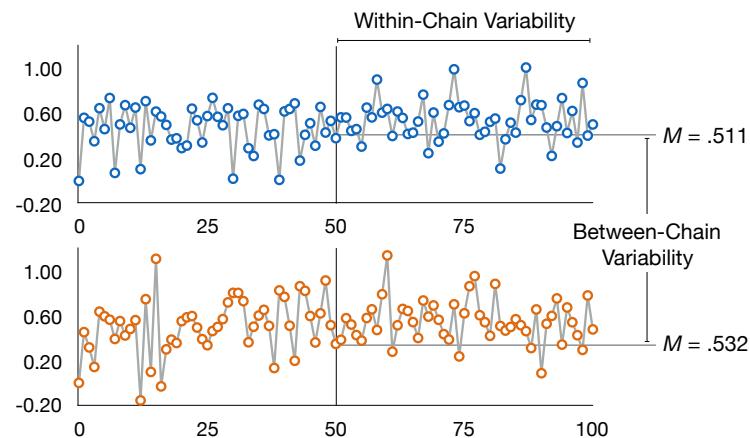
Reading the Trace Plots

The chains have noticeably different means at iteration 20, but the mean difference dissipates by cycle 100

The distribution of synthetic slope values stabilizes (converges) by 100 iterations and additional iterations do not affect the mean difference

The burn-in and/or thinning interval can safely be set to a value of 100 or greater

Within- and Between-Chain Variation



Potential Scale Reduction (PSR) Factor

The potential scale reduction (PSR) factor (Gelman & Rubin, 1992) uses ANOVA mean squares formulas to compute between- and within-chain variance

$$PSR = \sqrt{\frac{\sigma_{\text{Within}}^2 + \sigma_{\text{Between}}^2}{\sigma_{\text{Within}}^2}} = \sqrt{\frac{1}{1 - R^2}}$$

The between-chain variance decreases as mean differences diminish, giving PSR values closer to 1

Recommendations

Parameters converge at different rates, so PSR values vary across parameters

Determine the number of MCMC iterations required for the worst (highest) PSR to fall below 1.05 to 1.10 (common rules of thumb in the literature)

Set the burn-in and thinning intervals to a value that exceeds the number of iterations required to converge

Analysis Example

Motivating Example

Analysis involving number of years smoking (Years), number of cigarettes smoked per week (Cigs), and self-efficacy to quit smoking (Efficacy)

Imputation assumes that values are missing at random

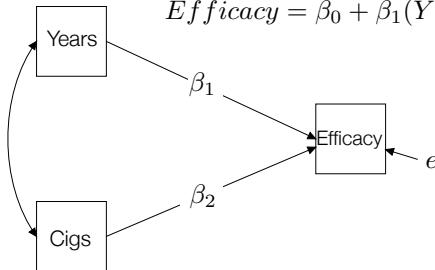
What does MAR require in this context?

Years	Cigs	Efficacy
7	9	NA
8	NA	NA
1	11	16
4	3	21
6	10	17
8	5	10
8	7	13
10	NA	10
15	NA	11
5	11	13
9	12	11
11	11	16
14	NA	10
13	19	9
12	NA	5
11	8	7
10	13	10
10	8	NA
7	10	7
11	10	6

Analysis Model

The analysis model is a multiple regression predicting self-efficacy to quit based on years smoking and number of cigarettes smoked

$$Efficacy = \beta_0 + \beta_1(Years) + \beta_2(Cigs) + e$$



Imputation Highlights

The analysis model is linear and additive (e.g., no interactions), so FCS imputation based on reverse regression models is appropriate

e.g., Number of cigarettes is imputed from a linear regression with efficacy and years smoking as predictors

The imputation process should include all variables in the analysis, although extra variables can be used

Ex3.1.imp Blimp Diagnostic Script

```
DATA: ~/desktop/examples/smoking.csv;
VARNAMES: id quitmeth male age years cigs heavycig
efficacy stress;
MISSING: -99;
MODEL: ~ years cigs efficacy;
SEED: 90291;
BURN: 3000;
THIN: 1;
NIMPS: 2;
OUTFILE: ~/desktop/examples/imp*.csv;
OPTIONS: separate psr;
CHAINS: 2 processors 2;
```

Diagnostic Output

POTENTIAL SCALE REDUCTION (PSR) OUTPUT:

Comparing iterations 51 to 100 for 2 chains.

	Fix Eff	Ran Eff Var	Err Var	Threshold
Max PSR	1.027	nan	1.008	nan
Missing Variable	cigs		cigs	

Comparing iterations 101 to 200 for 2 chains.

	Fix Eff	Ran Eff Var	Err Var	Threshold
Max PSR	1.043	nan	1.002	nan
Missing Variable	efficacy		cigs	

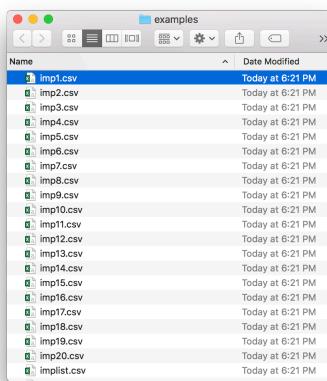
Ex3.2.imp Blimp Imputation Script (Mplus Format)

```
DATA: ~/desktop/examples/smoking.csv;
VARNAMES: id quitmeth male age years cigs heavycig
efficacy stress;
MISSING: -99;
MODEL: ~ years cigs efficacy;
SEED: 90291;
BURN: 100;
THIN: 100;
NIMPS: 20;
OUTFILE: ~/desktop/examples/imp*.csv;
OPTIONS: separate;
CHAINS: 2 processors 2;
```

Imputed Data Sets

The imputation phase generates a set of imputed data sets

The next step is to analyze the data ...



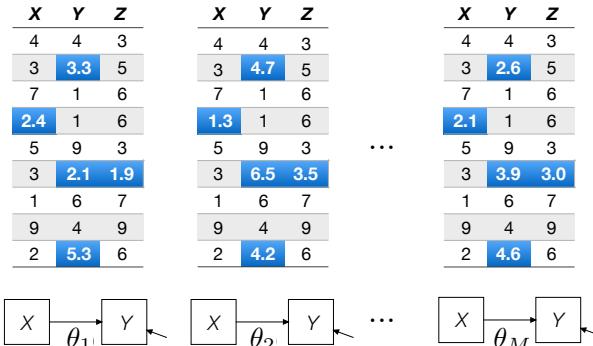
Ex3.3.imp Blimp Imputation Script (R, SAS, SPSS, and Stata Format)

```
DATA: ~/desktop/examples/smoking.csv;
VARNAMES: id quitmeth male age years cigs heavycig
efficacy stress;
MISSING: -99;
MODEL: ~ years cigs efficacy;
SEED: 90291;
BURN: 100;
THIN: 100;
NIMPS: 20;
OUTFILE: ~/desktop/examples/imps.csv;
OPTIONS: stacked;
CHAINS: 2 processors 2;
```

Multiple Imputation Analysis and Pooling Phase

Analyzing Multiply Imputed Data

The model of interest is fit to each data set



Analysis Phase

The multiple regression analysis is fit to each imputed data set (e.g., 20 multiple regression analyses)

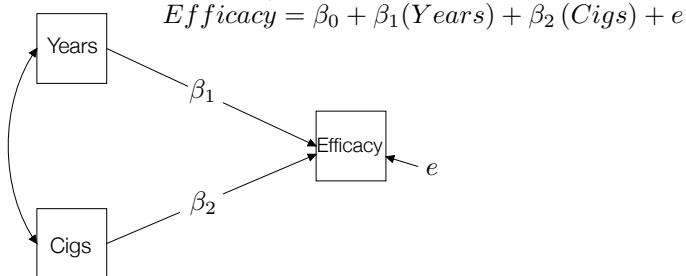


Illustration With Three Imputed Data Sets

Years	Cigs	Efficacy
7	9	15.50
8	8.96	15.78
1	11	16
4	3	21
6	10	17
8	5	10
8	7	13
10	9.92	10
15	13.62	11
5	11	13
9	12	11
11	11	16
14	14.42	10
13	19	9
12	18.04	5
11	8	7
10	13	10
10	8	9.18
7	10	7
11	10	6

Years	Cigs	Efficacy
7	9	15.38
8	8.28	9.25
1	11	16
4	3	21
6	10	17
8	5	10
8	7	13
10	13.41	10
15	6.99	11
5	11	13
9	12	11
11	11	16
14	15.31	10
13	19	9
12	12.75	5
11	8	7
10	13	10
10	8	14.85
7	10	7
11	10	6

Years	Cigs	Efficacy
7	9	5.86
8	10.04	13.88
1	11	16
4	3	21
6	10	17
8	5	10
8	7	13
10	12.95	10
15	14.40	11
5	11	13
9	12	11
11	11	16
14	14.47	10
13	19	9
12	11.46	5
11	8	7
10	13	10
10	8	6.79
7	10	7
11	10	6

Pooling Estimates

The multiple imputation point estimate is the arithmetic average of the M complete-data estimates

$$\hat{\theta} = \frac{\sum_{m=1}^M \hat{\theta}_m}{M}$$

Estimate from one data set
 Pooled estimate
 Number of data sets

Pooling Descriptive Statistics

Data Set 1			Data Set 2			Data Set 3					
	M	SD	N		M	SD	N		Years	Cigs	SE
Years	9.00	3.45	20	Years	9.00	3.45	20	Years	9.00	3.45	20
Cigs	10.59	3.84	20	Cigs	10.19	3.61	20	Cigs	10.52	3.51	20
SE	11.62	4.19	20	SE	11.57	4.14	20	SE	10.93	4.28	20

$$\hat{\theta} = \frac{\sum_{m=1}^M \hat{\theta}_m}{M} = \frac{10.59 + 10.19 + 10.52}{3} = 10.43$$

Pooled Estimates			
	M	SD	N
Years	9.00	3.45	20
Cigs	10.43	3.65	20
SE	11.37	4.20	20

Pooling Correlations

Data Set 1			Data Set 2			Data Set 3		
Years	Cigs	SE	Years	Cigs	SE	Years	Cigs	SE
Years	1.00		Years	1.00		Years	1.00	
Cigs	0.54	1.00	Cigs	0.38	1.00	Cigs	0.54	1.00
SE	-0.60	-0.45	SE	-0.57	-0.38	SE	-0.52	-0.26

Pooled Estimates		
Years	Cigs	SE
Years	1.00	
Cigs	0.49	1.00
SE	-0.57	-0.37

$$\hat{\theta} = \frac{\sum_{m=1}^M \hat{\theta}_m}{M} = \frac{-0.45 - 0.57 - 0.26}{3} = -0.37$$

Pooling Regression Parameters

	Data Set 1	Data Set 2	Data Set 3	Pooled
B ₀ (Intercept)	19.20	19.16	18.30	18.88
B ₁ (Years)	-0.62	-0.59	-0.63	-0.61
B ₂ (Cigarettes)	-0.19	-0.22	0.04	-0.12
Residual Variance	12.07	12.41	14.81	13.10
R ²	0.39	0.35	0.28	0.34

$$\hat{\theta} = \frac{\sum_{m=1}^M \hat{\theta}_m}{M} = \frac{-0.19 - 0.22 + 0.04}{3} = -0.12$$

Pooling Standard Errors

Averaging standard errors underestimates sampling variability because the component standard errors are based on complete data sets

The pooling procedure incorporates an adjustment that reflects the additional error from missing data

The within-imputation component estimates complete-data sampling error, and the between-imputation component is additional missing data error

Within-Imputation Variance

The within-imputation variance is the average squared standard error (sampling variance)

$$Var_W = \frac{\sum_{m=1}^M SE_m^2}{M}$$

Within-imputation variance estimates sampling error in the hypothetically complete data

Within-Imputation Variance Example

Data Set 1		
Est.	SE	SE ²
B ₀	19.20	2.559
B ₁	-0.62	0.275
B ₂	-0.19	0.247

Data Set 2		
Est.	SE	SE ²
B ₀	19.16	2.762
B ₁	-0.59	0.253
B ₂	-0.22	0.242

Data Set 3		
Est.	SE	SE ²
B ₀	16.54	2.970
B ₁	-0.67	0.304
B ₂	0.04	0.299

Pooled Estimates		
Est.	Var _w	
B ₀	18.30	10.745
B ₁	-0.63	0.080
B ₂	-0.12	0.097

$$Var_W = \frac{\sum_{m=1}^M SE_m^2}{M} = \frac{.061 + .059 + .089}{3} = .097$$

Between-Imputation Variance

Variability in the estimates across the M data sets is a result of using different imputations (missing data error)

$$Var_B = \frac{\sum_{m=1}^M (\hat{\theta}_m - \hat{\theta})^2}{M - 1}$$

Between-imputation variance applies the sample variance formula to the M estimates

Between-Imputation Variance Example

Data Set 1		
Est.	SE	SE ²
B ₀	19.20	2.559
B ₁	-0.62	0.275
B ₂	-0.19	0.247

Data Set 2		
Est.	SE	SE ²
B ₀	19.16	2.762
B ₁	-0.59	0.253
B ₂	-0.22	0.242

Data Set 3		
Est.	SE	SE ²
B ₀	16.54	2.970
B ₁	-0.67	0.304
B ₂	0.04	0.299

$$\begin{aligned} Var_B &= \frac{\sum_{m=1}^M (\hat{\theta}_m - \bar{\theta})^2}{M-1} \\ &= \frac{(-.19 - .12)^2 + (-.22 + .12)^2 + (.04 + .12)^2}{3-1} = .021 \end{aligned}$$

Pooled Estimates		
Est.	Var _W	Var _B
B ₀	18.30	10.745
B ₁	-0.63	0.080
B ₂	-0.12	0.097

Pooled Standard Error

The pooled standard error combines complete-data sampling error and additional missing data error

$$SE = \sqrt{Var_W + Var_B + \frac{Var_B}{M}}$$

Complete data sampling error
Missing data error
Sampling error of the pooled estimate

Standard Error Example

	Est.	Var _W	Var _B	SE
B ₀	18.30	10.745	2.311	3.72
B ₁ (Years)	-0.63	0.080	0.002	0.29
B ₂ (Cigs)	-0.12	0.097	0.021	0.35

$$\begin{aligned} SE &= \sqrt{Var_W + Var_B + \frac{Var_B}{M}} \\ &= \sqrt{.097 + .021 + \frac{.012}{3}} = .35 \end{aligned}$$

Significance Test

The single-parameter test statistic is based on the pooled estimate and standard error

$$t \text{ (or } z) = \frac{\hat{\theta} - \theta_0}{SE}$$

Pooled estimate
Hypothesized value (e.g., 0)

Software packages use different reference distributions for p-values (e.g., t with various df adjustments, z)

Degrees of Freedom Adjustments

Rubin's (1987) original degrees of freedom expression for the t statistic often exceed the sample size

Use Barnard and Rubin (1999) or Reiter (2007) expressions for degrees of freedom when possible

Degrees of freedom expressions are complex functions of the complete-data degrees of freedom, sample size, missing data rates, and correlations among the variables

p -Value Comparison

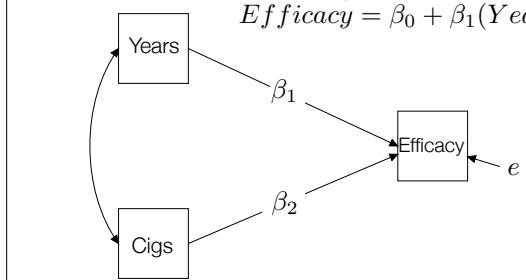
	Est.	SE	Est / SE	df	p
t Distribution with Rubin (1987) degrees of freedom					
B_0	18.16	2.95	6.17	1498.59	< .01
B_1 (Years)	-0.64	0.28	-2.31	8416.41	0.01
B_2 (Cigs)	-0.10	0.26	-0.38	1414.63	0.35
t Distribution with Barnard and Rubin (1999) degrees of freedom					
B_0	18.16	2.95	6.17	13.46	< .01
B_1 (Years)	-0.64	0.28	-2.31	14.55	0.02
B_2 (Cigs)	-0.10	0.26	-0.38	13.40	0.36
Normal z Distribution					
B_0	18.16	2.95	6.17	NA	< .01
B_1 (Years)	-0.64	0.28	-2.31	NA	0.01
B_2 (Cigs)	-0.10	0.26	-0.38	NA	0.35

Analysis Example

Analysis Model

The analysis model is a multiple regression predicting self-efficacy to quit based on years smoking and number of cigarettes smoked

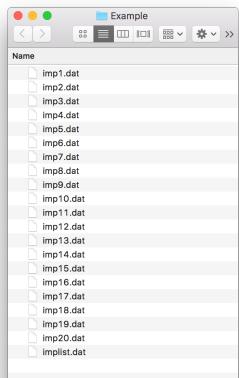
$$Efficacy = \beta_0 + \beta_1(Years) + \beta_2(Cigs) + e$$



Mplus Imputation Format

Mplus requires each imputed data set as a separate file

Blimp creates a list file (e.g., implist.dat) containing the names of the data sets, and this file serves as the input data for subsequent analyses



Ex3.2.imp Blimp Imputation Script (Mplus Format)

```
DATA: ~/desktop/examples/smoking.csv;
VARNAMES: id quitmeth male age years cigs heavycig
efficacy stress;
MISSING: -99;
MODEL: ~ years cigs efficacy;
SEED: 90291;
BURN: 100;
THIN: 100;
NIMPS: 20;
OUTFILE: ~/desktop/examples/imp*.csv;
OPTIONS: separate;
CHAINS: 2 processors 2;
```

Blimp Output

```
-----  
VARIABLE ORDER IN SAVED DATA:  
id quitmeth male age years cigs heavycig efficacy stress  
-----
```

Ex3.4.inp Mplus Analysis Script

```
DATA:  
file = implist.csv;  
type = imputation;  
VARIABLE:  
names = id quitmeth male age years  
cigs heavycig efficacy stress;  
usevariables = years cigs efficacy;  
MODEL:  
efficacy on years (b1)  
cigs (b2);  
MODEL TEST:  
b1 = 0; b2 = 0;  
OUTPUT:  
standardized(stdyx);
```

Mplus Analysis Output

MODEL RESULTS

	Estimate	S.E.	Est./S.E.	Two-Tailed P-Value	Rate of Missing
EFFICACY ON YEARS	-0.652	0.269	-2.423	0.015	0.129
CIGS	-0.069	0.278	-0.249	0.803	0.322
Intercepts					
EFFICACY	17.815	2.818	6.323	0.000	0.178
Residual Variances					
EFFICACY	11.240	3.972	2.830	0.005	0.185

Mplus Analysis Output

STANDARDIZED MODEL RESULTS

STDYX Standardization

	Estimate	S.E.	Est./S.E.	Two-Tailed P-Value	Rate of Missing
EFFICACY ON YEARS	-0.540	0.201	-2.684	0.007	0.185
CIGS	-0.055	0.254	-0.216	0.829	0.351
Intercepts					
EFFICACY	4.393	0.735	5.975	0.000	0.263
Residual Variances					
EFFICACY	0.677	0.178	3.791	0.000	0.087

Ex3.3.imp Blimp Imputation Script (R, SAS, SPSS, and Stata Format)

```
DATA: ~/desktop/examples/smoking.csv;
VARNAMES: id quitmeth male age years cigs heavycig
efficacy stress;
MISSING: -99;
MODEL: ~ years cigs efficacy;
SEED: 90291;
BURN: 100;
THIN: 100;
NIMPS: 20;
OUTFILE: ~/desktop/examples/imps.csv;
OPTIONS: stacked;
CHAINS: 2 processors 2;
```

Blimp Output

VARIABLE ORDER IN SAVED DATA:

```
imp# id quitmeth male age years cigs heavycig efficacy stress
```

Ex3.5.r R Analysis Script

```
# Required packages
library(mitml)

# Read data
filepath <- "~/desktop/examples/imps.csv"
impdata <- read.csv(filepath, header = F)
names(impdata) <- c("imputation", "id", "quitmeth", "male", "age",
"years", "cigs", "heavycig", "efficacy", "stress")

# Analyze data and pool estimates
implist <- as.mitml.List(split(impdata, impdata$imputation))
analysis <- with(implist, lm(efficacy ~ years + cigs))
estimates <- testEstimates(analysis, var.comp = T, df.com = 17)
estimates

# Test full model with Wald test
emptymodel <- with(implist, lm(efficacy ~ 1))
testModels(analysis, emptymodel, method = "D1")
```

R Analysis Output

Final parameter estimates and inferences obtained from 20 imputed data sets.

	Estimate	Std. Error	t.value	df	P(> t)	RIV	FMI
(Intercept)	17.815	3.016	5.907	12.760	0.000	0.180	0.155
years	-0.652	0.289	-2.255	13.497	0.041	0.124	0.111
cigs	-0.069	0.294	-0.235	10.514	0.818	0.391	0.287

Estimate
Residual--Residual 13.224

Hypothesis test adjusted for small samples with df=[17]
complete-data degrees of freedom.

R Analysis Output

Model comparison calculated from 20 imputed data sets.
Combination method: D1

F.value	df1	df2	P(>F)	RIV
3.215	2	948.316	0.041	0.222

Unadjusted hypothesis test as appropriate in larger samples.

Ex3.6.sps SPSS Analysis Script

```
data list free file = '/users/craig/desktop/examples/imps.csv'
/ imputation_ id quitmeth male age years cigs heavycig
efficacy stress.
exe.

* Initiate pooling routines.
sort cases by imputation_.
split file layered by imputation_.

* Analysis and pooling.
regression
 /descriptives mean stddev corr sig n
 /dependent efficacy
 /method enter years cigs.
```

SPSS Analysis Output

Coefficients ^a						
imputation_	Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1.00	1	(Constant)	18.639	2.931	6.359	.000
		years	-.645	.278	-.519	.2318
		cigs	-.106	.279	-.085	.381
...						
20.00	1	(Constant)	17.606	3.246	5.424	.000
		years	-.578	.282	-.454	-.2053
		cigs	-.135	.286	-.105	.473
Pooled	1	(Constant)	17.815	3.016	5.907	.000
		years	-.652	.289	-.255	.024
		cigs	-.069	.294	-.235	.814

a. Dependent Variable: efficacy

Ex3.7.do Stata Analysis Script

```
// Import and save original data
import delimited "~/desktop/examples/smoking.csv"
rename (v1 - v9)(id quitmeth male age years cigs heavycig
efficacy stress)
generate imp=0

// Recode missing values
foreach var of varlist id - stress {
    replace `var' = . if `var'== -99
}
save original, replace

// Import and save imputed data
clear
import delimited "~/desktop/examples/imps.csv"
rename (v1 - v10)(imp id quitmeth male age years cigs heavycig
efficacy stress)
save imputed, replace
```

Ex3.7.do Stata Analysis Script

```
// Append original and imputed data
use original, clear
append using imputed

// Convert to mi data
mi import flong, m(imp) id(id) imputed(quitmeth - stress) clear

// Analyze data and pool results
mi estimate, cmdok: regress efficacy years cigs
```

Stata Analysis Output

```
Multiple-imputation estimates
Linear regression
Imputations      =      20
Number of obs   =      20
Average RVI    =  0.1942
Largest FMI    =  0.3196
Complete DF     =      17
DF adjustment: Small sample
DF:      min      =  10.51
                    avg      =  12.26
                    max      =  13.50
Model F test: Equal FMI
F(  2,    14.0) =    3.22
Within VCE type: OLS
Prob > F        =  0.0709

efficacy |   Coef.  Std. Err.      t  P>|t| [95% Conf. Interval]
-----+-----+-----+-----+-----+-----+-----+
years |  -.651715  .2889468  -2.26  0.041  -1.273616  -.0298138
cigs |  -.069252  .2941369  -0.24  0.818  -.7203108  .5818068
_cons |  17.8151  3.015846  5.91  0.000  11.28726  24.34295
```