# Chap 15: Missing Data and Other Opportunities

```
1  md"# Chap 15: Missing Data and Other Opportunities"
```

```
1  versioninfo()
```

```
Julia Version 1.11.0                                    ?
Commit 501a4f25c2b (2024-10-07 11:40 UTC)
Build Info:
  Official https://julialang.org/ release
Platform Info:
  OS: Linux (x86_64-linux-gnu)
  CPU: 32 × Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz
  WORD_SIZE: 64
  LLVM: libLLVM-16.0.6 (ORCJIT, haswell)
Threads: 16 default, 0 interactive, 8 GC (on 32 virtual cor
es)
Environment:
  JULIA_PKG_SERVER = https://mirrors.tuna.tsinghua.edu.cn/j
ulia
  JULIA_REVISE_WORKER_ONLY = 1
```

```
 1  html"""
 2  <style>
 3      main {
 4          margin: 0 auto;
 5          max-width: max(1800px, 75%);
 6          padding-left: max(5px, 1%);
 7          padding-right: max(350px, 10%);
 8      }
 9  </style>
10  """
```

## Table of Contents

```
1  begin
2     using Pkg, DrWatson
3     using PlutoUI
4     TableOfContents()
5  end
```

```
1  begin
2      using Turing
3      using Turing
4      using DataFrames
5      using CSV
6      using Random
7      using Dagitty
8      using Distributions
9      #using StatisticalRethinking
10     using StatisticalRethinking: link
11     using StatisticalRethinkingPlots
12     using StatsPlots
13     using StatsBase
14     using Logging
15     using LinearAlgebra
16 end
```

# Code 15.1

```
1  md"## Code 15.1"
```

0.6617857711284418

```
1  begin
2      Random.seed!(2)
3
4      function sim_pancake()
5          pancake = [[1, 1], [1, 0], [0, 0]]
6          sides = sample(pancake)
7          sample([sides, reverse(sides)])
8      end
9
10     @time pancakes = vcat([sim_pancake() for _ in
       1:100_000]'...)
11     up = pancakes[:,1]
12     down = pancakes[:,2]
13
14     num_11_10 = sum(up .== 1)
15     num_11 = sum((up .== 1) .& (down .== 1))
16     num_11 / num_11_10
17 end
```

> 0.103752 seconds (1.65 M allocations: 64.906 MiB, 58.
> 63% compilation time)   ⑦

pancake =   [[1, 1], [1, 0], [0, 0]]
```
1  pancake = [[1, 1], [1, 0], [0, 0]]
```

sides =   [0, 0]
```
1  sides = sample(pancake)
```

[0, 0]
```
1  sample([sides, reverse(sides)])
```
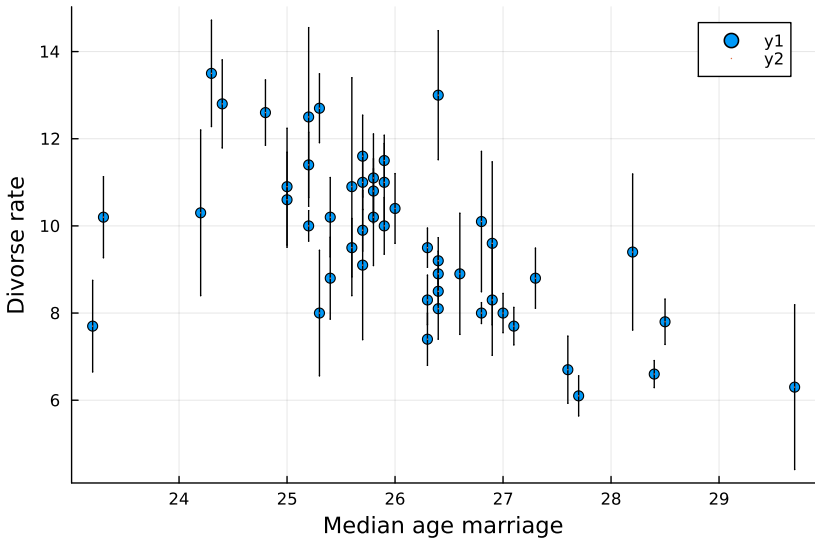
[[0, 0], [0, 0]]
```
1  [sides, reverse(sides)]
```

# 15.1 Measurement error

```
1  md" # 15.1 Measurement error"
```
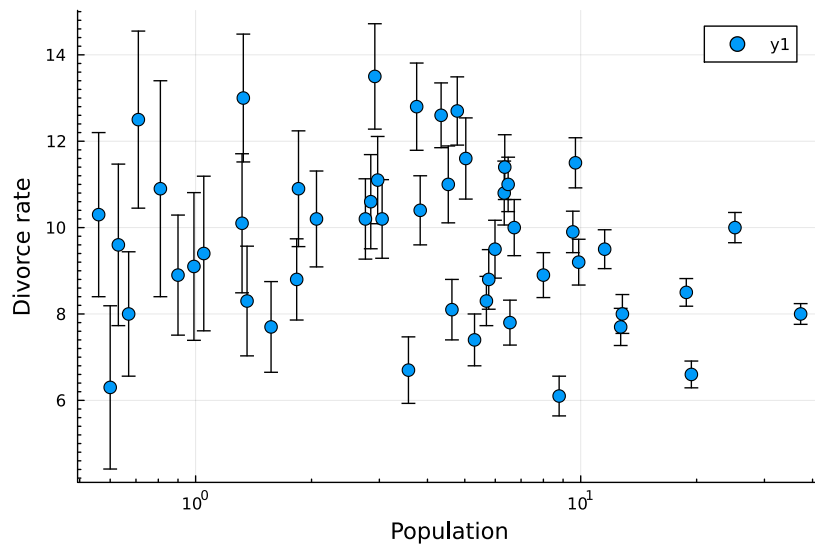
## Code 15.2

```
1  md"## Code 15.2"
```



```
1  begin
2      d_divorce = DataFrame(CSV.File("data/WaffleDivorce.csv"))
3
4      scatter(d_divorce.MedianAgeMarriage, d_divorce.Divorce,
5          xlab="Median age marriage", ylab="Diverse rate")
6      scatter!(d_divorce.MedianAgeMarriage, d_divorce.Divorce,
        yerror=d_divorce."Divorce SE", ms=0)
7  end
```

| | Location | Loc | Population | MedianAgeMarriage | Marriage | Ma |
|---|---|---|---|---|---|---|
| 1 | "Alabama" | "AL" | 4.78 | 25.3 | 20.2 | 1. |
| 2 | "Alaska" | "AK" | 0.71 | 25.2 | 26.0 | 2. |
| 3 | "Arizona" | "AZ" | 6.33 | 25.8 | 20.3 | 0. |

```
1  first(d_divorce,3)
```

```
1  begin
2      scatter(d_divorce.Population, d_divorce.Divorce,
       xaxis=:log10,
3          xlab="Population", ylab="Divorce rate",
4          xminorticks=9, yminorticks=10,
5          yerror=d_divorce."Divorce SE", ms=5)
6      #scatter!(d_divorce.Population, d_divorce.Divorce,
       yerror=d_divorce."Divorce SE", ms=0)
7  end
```

## Code 15.3 model m15_1

```
1  md"## Code 15.3 model `m15_1`"
```

|       | D_true[10] | D_true[11] | D_true[12] | D_true[13] | D_true[14] | D_t |
|-------|------------|------------|------------|------------|------------|-----|
| 1     | −0.784182  | 0.591328   | −0.412152  | 0.12795    | −0.820124  | 0.9 |
| 2     | −0.418201  | 1.43575    | −0.850778  | 0.108591   | −0.808901  | 0.5 |
| 3     | −0.638492  | 0.392377   | −0.361222  | 1.00863    | −0.986006  | 0.6 |
| 4     | −0.509696  | 1.10383    | −0.518881  | 0.27218    | −0.814202  | 0.5 |
| 5     | −0.691987  | 0.769628   | −0.666241  | 0.826797   | −0.71571   | 0.4 |
| 6     | −0.521323  | 0.589482   | −0.991947  | −0.00510027 | −1.00075  | 0.6 |
| 7     | −0.586819  | 0.689822   | −0.252866  | 0.751031   | −0.657953  | 0.4 |
| 8     | −0.582998  | 0.455282   | −0.310074  | 0.56882    | −0.711395  | 0.5 |
| 9     | −0.630576  | 0.790008   | −0.888668  | 1.30604    | −0.891928  | 0.4 |
| 10    | −0.630576  | 0.790008   | −0.888668  | 1.30604    | −0.891928  | 0.4 |
| more  |            |            |            |            |            |     |
| 1000  | −1.02493   | 0.883494   | −0.144567  | −0.00200293 | −0.840648 | 0.6 |

```julia
1  begin
2      d_divorce_ls = (
3          D_obs = standardize(ZScoreTransform, d_divorce.Divorce),
4          D_sd = d_divorce."Divorce SE" ./ std(d_divorce.Divorce),
5          M = standardize(ZScoreTransform, d_divorce.Marriage),
6          A = standardize(ZScoreTransform,
       d_divorce.MedianAgeMarriage),
7          N = nrow(d_divorce),
8      )
9
10     @model function m15_1(D_obs, D_sd, M, A, N)
11         a ~ Normal(0, 0.2)
12         bA ~ Normal(0, 0.5)
13         bM ~ Normal(0, 0.5)
14         μ = @. a + bA * A + bM * M
15         σ ~ Exponential()
16         D_true ~ MvNormal(μ, σ)
17         @. D_obs ~ Normal(D_true, D_sd)
18     end
19
20     Random.seed!(1)
21     @time m15_1_ch = sample(m15_1(d_divorce_ls...), NUTS(),
       1000)
22     m15_1_df = DataFrame(m15_1_ch);
23 end
```

Sampling  [100%]

Found initial step size
ε: 0.2

11.256042 seconds (16.66 M allocations: 6.157 GiB, 10.
68% gc time, 55.28% compilation time)   ⓘ

# Code 15.4

```julia
1  md"## Code 15.4"
```

| | variable | mean | min | median | max |
|---|---|---|---|---|---|
| 1 | Symbol("D_true[10]") | -0.622426 | -1.17513 | -0.621466 | -0.059846 |
| 2 | Symbol("D_true[11]") | 0.752743 | -0.167793 | 0.764524 | 1.76655 |
| 3 | Symbol("D_true[12]") | -0.54162 | -2.09472 | -0.538969 | 1.42518 |
| 4 | Symbol("D_true[13]") | 0.191023 | -1.80048 | 0.197183 | 1.54803 |
| 5 | Symbol("D_true[14]") | -0.86873 | -1.59464 | -0.878422 | -0.135698 |
| 6 | Symbol("D_true[15]") | 0.563774 | -0.450136 | 0.559766 | 1.55619 |
| 7 | Symbol("D_true[16]") | 0.269308 | -0.855484 | 0.282876 | 1.57184 |
| 8 | Symbol("D_true[17]") | 0.505615 | -0.78145 | 0.504514 | 1.83022 |
| 9 | Symbol("D_true[18]") | 1.25328 | 0.14058 | 1.25724 | 2.48261 |
| 10 | Symbol("D_true[19]") | 0.428978 | -0.812482 | 0.441281 | 1.63373 |
| | more | | | | |
| 54 | :σ | 0.579131 | 0.30084 | 0.575787 | 1.00322 |

```
1  describe(m15_1_df)
```

## Code 15.5 model `m15_2`

```
1  md"## Code 15.5 model `m15_2`"
```

```
(D_obs = [1.65421, 1.54436, 0.610716, 2.09357, -0.927058, 1.05008, -1.
```

```
1  begin
2      dlist2 = (
3          D_obs = standardize(ZScoreTransform, d_divorce.Divorce),
4          D_sd = d_divorce."Divorce SE" ./ std(d_divorce.Divorce),
5          M_obs = standardize(ZScoreTransform,
       d_divorce.Marriage),
6          M_sd = d_divorce."Marriage SE" ./
       std(d_divorce.Marriage),
7          A = standardize(ZScoreTransform,
       d_divorce.MedianAgeMarriage),
8          N = nrow(d_divorce),
9      )
10 end
```

[0.083057, 1.01903, 0.0594721, 1.41732, -0.266635, 0.830463, -0.76543

```
1  begin
2
3      @model function m15_2(D_obs, D_sd, M_obs, M_sd, A, N)
4          a ~ Normal(0, 0.2)
5          bA ~ Normal(0, 0.5)
6          bM ~ Normal(0, 0.5)
7          M_true ~ filldist(Normal(), N)
8
9          μ = @. a + bA * A + bM * M_true
10         σ ~ Exponential()
11         D_true ~ MvNormal(μ, σ)
12         @. D_obs ~ Normal(D_true, D_sd)
13         @. M_obs ~ Normal(M_true, M_sd)
14     end
15
16     Random.seed!(1)
17     @time m15_2_ch = sample(m15_2(dlist2...), NUTS(), 1000)
18     m15_2_df = DataFrame(m15_2_ch);
19     D_true = [mean(m15_2_df[!, "D_true[$i]"]) for i ∈
           1:dlist2.N]
20     M_true = [mean(m15_2_df[!, "M_true[$i]"]) for i ∈
           1:dlist2.N]
21 end
```

```
 46.280141 seconds (67.04 M allocations: 42.705 GiB, 1    ⑦
7.65% gc time, 36.69% compilation time)
 35.882352 seconds (53.92 M allocations: 42.077 GiB, 19.90%
gc time, 21.32% compilation time)
```
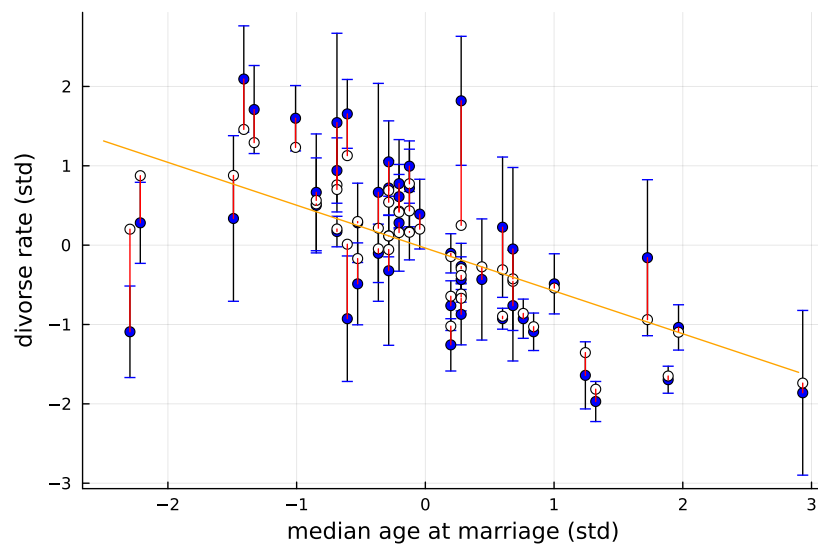
Sampling  `100%`

```
Found initial step size
ϵ: 0.4
```

| | variable | mean | min | median | max |
|---|---|---|---|---|---|
| 1 | Symbol("D_true[10]") | -0.616598 | -1.09836 | -0.616169 | -0.10918 |
| 2 | Symbol("D_true[11]") | 0.773391 | -0.153289 | 0.772106 | 1.59042 |
| 3 | Symbol("D_true[12]") | -0.455932 | -1.96422 | -0.469349 | 1.27627 |
| 4 | Symbol("D_true[13]") | 0.201203 | -1.44406 | 0.204312 | 1.67876 |
| 5 | Symbol("D_true[14]") | -0.860255 | -1.57298 | -0.85922 | -0.15458 |
| 6 | Symbol("D_true[15]") | 0.540992 | -0.540644 | 0.543722 | 1.62189 |
| 7 | Symbol("D_true[16]") | 0.297736 | -0.943139 | 0.293591 | 1.44996 |
| 8 | Symbol("D_true[17]") | 0.519618 | -1.31079 | 0.522772 | 2.32168 |
| 9 | Symbol("D_true[18]") | 1.23177 | 0.22005 | 1.22341 | 2.29087 |
| 10 | Symbol("D_true[19]") | 0.431547 | -0.906202 | 0.416142 | 1.98877 |
| | more | | | | |
| 104 | :σ | 0.563163 | 0.242072 | 0.558362 | 0.974335 |

```
1  describe(m15_2_df)
```

# Figure 15.2

```
1  md"## Figure 15.2"
```

```
1  begin
2
3      p1 = scatter(dlist2.A, dlist2.D_obs, mc=:blue,
       yerror=dlist2.D_sd,
4          label="observed", xlab="median age at marriage (std)",
       ylab="divorse rate (std)")
5      scatter!(dlist2.A, D_true, mc=:white, label="true")
6
7      for i ∈ 1:dlist2.N
8          plot!([dlist2.A[i], dlist2.A[i]], [dlist2.D_obs[i],
       D_true[i]], c=:red, legend=false)
9      end
10     x = -2.5:0.2:3
11     y = -0.0368595 .+  -0.540089 .* x
12     plot!(x,y, c=:orange, label="m15_2 estimate")
13     p1
14 end
```
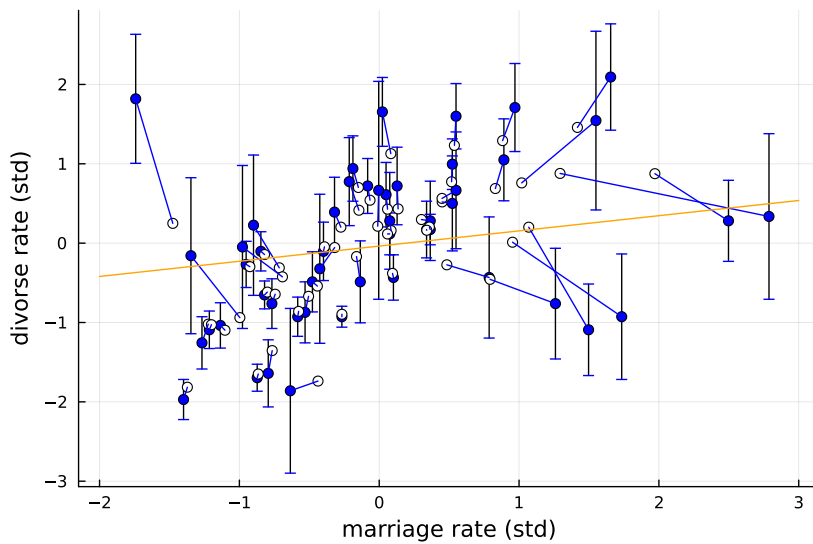
```
1  Enter cell code...
```

## Code 15.6 Figure 15.3

```
1  md"## Code 15.6 Figure 15.3"
```
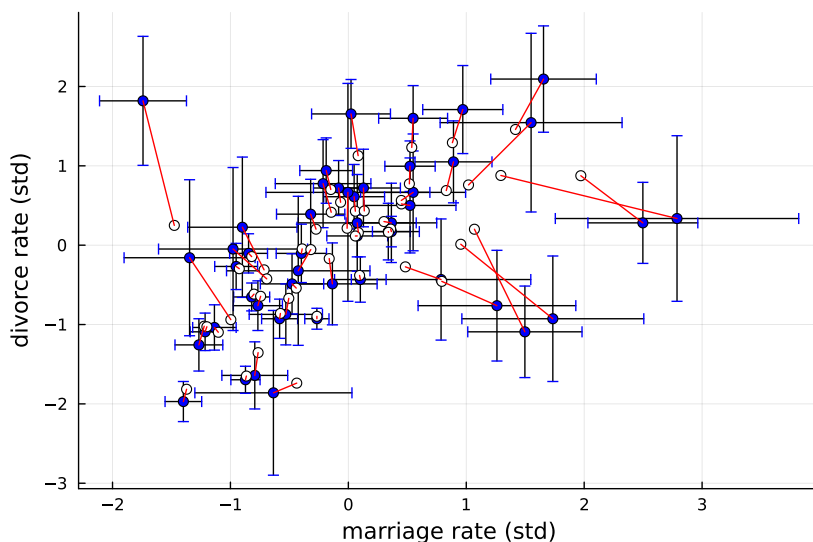
```
1  begin
2      p2 = scatter(dlist2.M_obs, dlist2.D_obs, mc=:blue,
       yerror=dlist2.D_sd,
3          label="observed", xlab="marriage rate (std)",
       ylab="divorse rate (std)",
4          legend=true)
5      scatter!(M_true, D_true, mc=:white, label="true",
       legend=true)
6
7      for i ∈ 1:dlist2.N
8          plot!([dlist2.M_obs[i], M_true[i]], [dlist2.D_obs[i],
       D_true[i]], c=:blue, legend=false)
9      end
10     x2 = -2:0.2:3
11     y2 = -0.0368595 .+  0.1915 .* x2
12     plot!(x2,y2, c=:orange, label="m15_2 estimate")
13     p2
14  end
```



```
1  begin
2      p3 = scatter(dlist2.M_obs, dlist2.D_obs, mc=:blue,
       xerror=dlist2.M_sd, yerror=dlist2.D_sd,
3          label="observed", xlab="marriage rate (std)",
       ylab="divorce rate (std)")
4      scatter!(M_true, D_true, mc=:white, label="true")
5
6      for i ∈ 1:dlist2.N
7          plot!([dlist2.M_obs[i], M_true[i]], [dlist2.D_obs[i],
       D_true[i]], c=:red, legend=false)
8      end
9      p3
10  end
```

```
1 Enter cell code...
```

## Code 15.7

```
1 md"## Code 15.7"
```

[-0.366839, -2.48606, 0.579584, -0.588886, -1.54843, -2.13782, -1.197

```
1 let
2     N = 500
3     A = rand(Normal(), N)
4     M = rand.(Normal.(-A))
5     D = rand.(Normal.(A))
6     A_obs = rand.(Normal.(A));
7 end
```

## Code 15.7

```
1 md"## Code 15.7"
```

# 15.2 Missing data

```
1 md"# 15.2 Missing data"
```

```
1 Enter cell code...
```

## Code 15.8

```
1 md"## Code 15.8"
```

```
[4, 8, 0, 7, 6, 5, 4, 5, 5, 6, 8, 7, 6, 9, 4, 4, 8, 8, 8, 1,    more ,4, 6
```

```
1 let
2    N = 100
3    S = rand(Normal(), N)
4    H = rand.([BinomialLogit(10, l) for l in S]);
5 end
```